

# Inverse Covariance Estimation with Structured Groups

Shaozhe Tao<sup>\*</sup> and Yifan Sun<sup>†</sup> and Daniel Boley<sup>‡</sup>

## Abstract

Estimating the inverse covariance matrix of  $p$  variables from  $n$  observations is challenging when  $n \ll p$ , since the sample covariance matrix is singular and cannot be inverted. A popular solution is to optimize for the  $\ell_1$  penalized estimator; however, this does not incorporate structure domain knowledge and can be expensive to optimize. We consider finding inverse covariance matrices with group structure, defined as potentially overlapping principal submatrices, determined from domain knowledge (e.g. categories or graph cliques). We propose a new estimator for this problem setting that can be derived efficiently via the conditional gradient method, leveraging chordal decomposition theory for scalability. Simulation results show significant improvement in sample complexity when the correct group structure is known. We also apply these estimators to 14,910 stock closing prices, with noticeable improvement when group sparsity is exploited.

## 1 Introduction

The inverse covariance matrix is of interest to statisticians in biology, finance, machine learning, *etc.* In finance, it is a key ingredient for computing value-at-risk, a factor in portfolio optimization. In graphical models, for  $p$  random variables with true covariance matrix  $C$ , the sparsity pattern of  $C^{-1}$  gives the conditional independence between each pair of variables. However, if  $n \ll p$ , then the sample covariance matrix  $\hat{C}$  is invertible, and the pseudoinverse  $\hat{C}^\dagger$  is inaccurately dense. The most popular alternative is the graphical LASSO (G-LASSO) estimator [Yuan and Lin, 2007; Banerjee *et al.*, 2008], the solution to

$$\begin{aligned} & \underset{X}{\text{minimize}} && -\log\det(X) + \text{tr}(\hat{C}X) + \rho\|X\|_1 \\ & \text{subject to} && X \succeq 0 \end{aligned} \quad (1)$$

<sup>\*</sup>University of Minnesota - Twin Cities, MN  
taoxx120@umn.edu

<sup>†</sup>Technicolor Research - Los Altos, CA,  
yifan.sun@technicolor.com

<sup>‡</sup>University of Minnesota - Twin Cities, MN  
boley@cs.umn.edu

for some regularization parameter  $\rho > 0$ . By adding a sparsity-inducing regularizer, the effective degrees of freedom are reduced, and as these works show, the resulting estimator has a much lower sample complexity than inverting  $\hat{C}$ . However, this estimator does not incorporate any *prior structural knowledge* from the problem domain. Additionally, in general solving (1) is computationally challenging if  $p$  is large.

Most existing methods for solving (1) require a sequence of eigenvalue decompositions (EDs) [Banerjee *et al.*, 2006; Friedman *et al.*, 2008; d’Aspremont *et al.*, 2008; Yuan, 2009; Rolfs *et al.*, 2012]. This is expensive if  $p$  is large; a dense ED requires  $O(p^3)$  computations, and sparse EDs (like Lanczos type methods) can be even slower when the full eigenvalue spectrum is needed. There are some exceptions, for example [Scheinberg and Rish, 2009] at each step updates a row in a block coordinate descent fashion, and maintains inverses using only rank-2 updates; [Dahl *et al.*, 2008] uses chordal decomposition to compute Newton steps efficiently in an interior point solver; and [Meinshausen and Bühlmann, 2006] uses neighborhood selection, which enforces the conditional independence condition one variable at a time. These methods are more or less intuitive, relying on general convex optimization principles; however, their scalability is limited. On the other end of the spectrum is BIG-QUIC [Hsieh *et al.*, 2013] which can solve up to 1 million variables. This breakthrough method simultaneously makes estimates of the matrix sparsity while also optimizing for it, and updating via block coordinate descent with carefully chosen (non-principle) submatrices. However, it demonstrates the tradeoff between simplicity and scalability; there are many intricate details for a successful implementation.

At the same time, there has been growing interest in the statistics community to exploit *group structure* in the estimators [Bach *et al.*, 2011; Negahban *et al.*, 2009; Chandrasekaran *et al.*, 2012; Obozinski *et al.*, 2011]. For example, [Danaher *et al.*, 2014] proposes a group graphical LASSO, but where groups are defined as membership in  $K$  classes. And, [Mazumder and Hastie, 2012] proposes thresholding the sample covariance matrix in order to identify fully-connected components of the graphical model, effectively decomposing (1). More recently, [Hosseini and Lee, 2016] learns overlapping submatrix groups probabilistically and penalizes accordingly. To our knowledge, this is the only work that addresses overlapping group sparsity in matrices; however, iterative full

eigenvalue decompositions are still needed to find the inverse covariance estimate.

We propose an estimator that exploits group structure, where a matrix group is described as either a principal submatrix or the matrix diagonal in Section 2. The solution  $X$  is then described as a sum of these possibly overlapping components. We then apply the Frank-Wolfe method to derive the estimator in Section 3. The algorithm at each iteration decomposes into parallelizable eigenvalue computations on the submatrices. In this way, unlike [Mazumder and Hastie, 2012; Hosseini and Lee, 2016], this estimator *explicitly* uses the predetermined groups as components for decomposition, thus using group structure to improve both performance and computation time. In Section 4 we give simulation results, which demonstrate that knowing and exploiting group structure significantly improves sample complexity. Finally, Section 5, we show the performance of our model on the stock datasets.

## 2 Group Norm Constrained Estimator

For an index set  $\gamma \subset \{1, \dots, p\}$  and a vector  $u \in \mathbb{R}^p$ , define  $\gamma$  as the subvector of  $u$  indexed by  $\gamma$ ; for the reverse, define the *augmenting linear map*  $\gamma : \mathbb{R}^{|\gamma|} \rightarrow \mathbb{R}^p$  such that

$$(A_\gamma u)_\gamma = u, \quad (A_\gamma u)_i = 0 \text{ if } i \notin \gamma.$$

In [Obozinski *et al.*, 2011], the *overlapping group norm* is defined as the solution to

$$\|x\|_{G,*} = \min_{u_1, \dots, u_l} \left\{ \sum_{k=1}^l w_k \|u_k\| : x = \sum_{k=1}^l A_{\gamma_k} u_k \right\} \quad (2)$$

for some proper norm  $\|\cdot\|$  and nonnegative weights  $w_1, \dots, w_l$ . (A common choice is  $w_k = |\gamma_k|^{-1}$ .) Used as a penalty term or in a constraint, this norm is shown to promote *group structure*; a small subset of index sets  $\gamma_k$  are “active”, and  $x_i = 0$  whenever  $i$  is not in an active set.

We extend this concept to matrices, by defining groups implicitly through index sets  $\beta \subset \{1, \dots, p\}$ , where  $X_{\beta,\beta}$  is the submatrix of  $X$  selected by the rows and columns indicated by  $\beta$ . Let  $\mathbb{S}^p$  denote the set of  $p \times p$  matrices. We define  $A_\beta : \mathbb{S}^{|\beta|} \rightarrow \mathbb{S}^p$  such that

$$(A_\beta(U))_{\beta,\beta} = U, \quad A_\beta(U)_{i,j} = 0 \text{ if } i \notin \beta \text{ or } j \notin \beta$$

and extend the overlapping group norm as the solution

$$\|X\|_G := \begin{cases} \min_{v, U_1, \dots, U_l} w_0 \|v\|_2 + \sum_{k=1}^l w_k \|U_k\|_F \\ \text{subj. to } X = \mathbf{diag}(v) + \sum_{k=1}^l A_{\beta_k}(U_k). \end{cases} \quad (3)$$

for nonnegative weights  $w_0, \dots, w_l$ . Note that the affine constraint imposes a sparsity pattern on  $X$ ; if  $X$  does not adhere to this pattern (e.g.  $X_{ij} \neq 0$  for some  $i, j \notin \beta_k, \forall k$ ) then we define  $\|X\|_G = \infty$ .

### 2.1 Our estimator

For  $p$  random variables  $Y_1, Y_2, \dots, Y_p$ , define  $C \in \mathbb{S}^p$  and  $\hat{C} \in \mathbb{S}^p$  as the true and sample covariance matrices. The *group*

*norm regularized graphical LASSO estimator (NG-LASSO)* is the solution to

$$\begin{aligned} \min_X & -\log \det(X) + \mathbf{tr}(\hat{C}X) \\ \text{s.t. } & X = \sum_{k=1}^l A_{\beta_k}(U_k) + \mathbf{diag}(v) \\ & w_0 \|v\|_2 + \sum_{k=1}^l w_k \|U_k\|_F \leq \alpha \\ & U_k \succeq 0, \quad k = 1, \dots, l, \\ & v_i \geq 0, \quad i = 1, \dots, p. \end{aligned} \quad (4)$$

We note that the first two constraints in (4) can be equivalently written as  $\|X\|_G \leq \alpha$  with  $G$ -norm defined in (3). As defined, this constraint restricts  $X$  to be implicitly within the sparsity pattern defined by the groups  $\beta_k$ .

Problem (4) is a computationally tractable approximation of

$$\begin{aligned} \min_X & -\log \det(X) + \mathbf{tr}(\hat{C}X) \\ \text{subject to } & X \succeq 0, \quad \|X\|_G \leq \alpha \end{aligned} \quad (5)$$

a natural group norm penalized version of the G-LASSO problem. Specifically, in (5),  $\|X\|_G$  can be written in terms of smaller matrices  $W_k \in \mathbb{S}^{|\beta_k|}$  and  $z \in \mathbb{R}^p$ . If additionally the sparsity pattern is *chordal* (i.e. if the intersection graph of the groups  $\beta_k$  is a tree) then the positive semidefinite (PSD) matrix constraint  $X \succeq 0$  can be decomposed to several smaller matrix constraints, via the equivalence in the following theorem.

**Theorem 2.1** [Agler *et al.*, 1988; Griewank and Toint, 1984] ([Grone *et al.*, 1984] dual) *If  $X \in \mathbb{S}^p$  has chordal sparsity, corresponding to groups  $\beta_1, \dots, \beta_l$ , then*

$$X \succeq 0 \iff X = \sum_{k=1}^l A_{\beta_k}(U_k), \quad U_k \succeq 0, \quad k = 1, \dots, l.$$

In this case  $X \succeq 0$  can be decomposed into smaller matrices  $U_k \in \mathbb{S}^{|\beta_k|+}$  and a positive diagonal  $v \in \mathbb{R}_+^p$  (where  $\mathbb{S}_+^p$  and  $\mathbb{R}_+^p$  are the PSD cone and nonnegative orthant, both of order  $p$ ). Then (4) is equivalent to (5) if and only if at optimality,  $W_k = U_k$  for all  $k$ , and  $v = z$ .

## 3 Optimization

The Frank-Wolfe algorithm has regained much attention in minimizing sparse problems [Jaggi, 2013], mimicking greedy approaches yet having guaranteed optimality for convex problems. We first describe the method for a generalized vector version of problem (4)

$$\min_x \{f(x) : x \in \mathcal{D}\} \quad (6)$$

where

$$\mathcal{D} = \left\{ x = \sum_{k=1}^l w_k A_{\gamma_k} u_k : \sum_{k=1}^l w_k \|u_k\|_2 \leq \alpha, u_k \in \mathcal{C}_k \right\}. \quad (7)$$

Here, the vector variable is  $x \in \mathbb{R}^m$ ,  $f(x)$  is a differentiable convex function, and  $\mathcal{C}_1, \dots, \mathcal{C}_l$  are proper convex cones. The Frank-Wolfe algorithm for solving  $\min_x \{f(x) : x \in \mathcal{D}\}$  is described in Alg. 1. It is known that the iterates of algorithm 1 converge as  $f(x^{[t]}) - f(x^*) \leq O(1/t)$  with step size  $\eta^{[t]} = 2/(t+2)$

---

**Algorithm 1** One step of Frank-Wolfe algorithm

---

**Input:**  $x^{[t]} \in \mathcal{D}$ :  $t$ -th iteration;  $\eta$ : step size;  
1: Compute gradient  $\nabla f(x^{[t]})$   
2: Compute forward step :  $s = \arg \min_{s \in \mathcal{D}} \langle s, \nabla f(x^{[t]}) \rangle$ ;  
3: Update primal variable :  $x^{[t+1]} = (1 - \eta^{[t]})x^{[t]} + \eta^{[t]}s$   
**Output:** optimal  $x^{[t+1]}$

---

**Forward step** [Frank and Wolfe, 1956; Dunn and Harshbarger, 1978]. At each iteration, the forward step consists of  $l$  parallelizable projections on cone  $\mathcal{C}_1, \dots, \mathcal{C}_l$ . Specifically, at each forward step, we compute

$$U_j^* = \frac{\alpha}{w_j \|Z_j\|_2} Z_j, \quad U_k = 0, \quad \forall k \neq j.$$

where index  $j = \arg \max_k w_k^{-1} \|Z_k\|_2$  and  $Z_j = \mathbf{proj}_{\mathcal{C}_j}(-\nabla f(x)_{\beta_j, \beta_j})$ . Then  $s = \sum_k w_k A_{\beta_k, \beta_k}(u_k)$ . The derivations are given in appendix A.

**Gradient computation** In general, to compute the gradient  $\nabla(\log \det(X)) = X^{-1}$  requires matrix inversion, which completely negates the computational complexity gain by decomposing the PSD cone. However, of the groups  $\beta_k$  form a chordal pattern, fast inversion methods exist [Liu, 1992; Andersen *et al.*, 2013] which require at each step  $l$  inversions of matrices at most of order  $|\beta_k|$ .

Applying both techniques, Alg. 4.1 describes the procedure for one iteration to find the NG-LASSO estimator (4).

---

**Algorithm 2** One step of Frank-Wolfe algorithm for (4)

---

**Input:**  $X^{[t]} \in \mathcal{D}$   $t$ -th iteration;  $\eta := \frac{2}{t+2}$  step size;  
1: Find  $\nabla f(X) = X^{-1} + C$   
2: Find the forward direction  $U^+$ :  
 $Z_0 = \mathbf{proj}_{\mathbb{R}_+^p}(-\mathbf{diag}(\nabla f(X)))$   
 $Z_k = \mathbf{proj}_{\mathbb{S}_+^{|\beta_k|}}(-\nabla f(X)_{\beta_k, \beta_k})$   
 $j = \arg \max_k w_k^{-1} \|Z_k\|_F$   
 $U_j^+ = \frac{\alpha}{w_j \|Z_j\|_F} Z_j, \quad U_k^+ = 0, \quad \forall k \neq j$   
3: Update  $X^{[t+1]} = X^{[t]} + \frac{2}{t+2} U^+$   
**Output:**  $X^{[t+1]}$

---

The main computational bottleneck at each step is a sequence of ED of the submatrices  $\nabla f(X)_{\beta_k, \beta_k}$ , both for inverting  $X$  and for projecting on the PSD cone. For both operations, the complexity is  $O(|\beta_k|^3)$  per group. If  $|\beta_k| < p/l$  excluding the diagonal group, then the total per-iteration complexity of the proposed optimization procedure has a per-iteration complexity of  $O(p^3/l^2 + p)$ , and much smaller than  $O(p^3)$  for G-LASSO.

## 4 Numerical Simulations

Here we present simulations of sparse inverse covariance matrix. We show two simulation results of the banded sparsity and then another simulation result on the general group sparsity. Numerically, when group structure is assumed, our group structured estimator outperforms G-LASSO.

In all of the following experiments, to pick  $\alpha$  and  $\rho$ , we swept powers of two  $\rho \in \{2^{-10}, \dots, 1\}$  and  $\alpha \in \{2^{-3}, \dots, 2^{10}\}$  and then picked the best performing  $\rho$  or  $\alpha$  for each test. In all cases, the best parameter was not on the boundary.

### 4.1 Baselines

As a baseline, we solve (1); however, since group structure also reveals matrix sparsity, for fair comparison we also solve (1) restricted to the sparsity pattern induced by the groups:

$$\begin{aligned} & \underset{X}{\text{minimize}} && -\log \det(X) + \text{tr}(\hat{C}X) + \rho \|X\|_1 \\ & \text{subject to} && X \succeq 0 \\ & && X \in \mathbf{B} := \{X \mid X_{ij} = 0 \text{ if } i, j \notin \beta_k \forall k\}, \end{aligned} \quad (8)$$

which we call restricted group LASSO (RG-LASSO). We solve these baselines using the Douglas-Rachford method [Lions and Mercier, 1979; Combettes and Pesquet, 2011] for minimizing the sum of  $m$  convex functions (Alg. 4.1, also [Combettes and Pesquet, 2011], Alg 10.27), with

$$f_1(X) = -\log \det(X) + \text{tr}(\hat{C}X), \quad f_2(X) = \rho \|X\|_1$$

and  $f_3, f_4$  as indicator functions for constraints

$$f_3(X) = \begin{cases} 0 & X \succeq 0 \\ \infty & \text{else.} \end{cases}, \quad f_4(X) = \begin{cases} 0 & X \in \mathbf{B} \\ \infty & \text{else.} \end{cases}.$$

The proximal operator [Moreau, 1962] for a convex function  $f(X)$  is defined as

$$\mathbf{prox}_f(Z) = \arg \min_X f(X) + (1/2) \|X - Z\|_F^2$$

and is defined for all  $Z$ , even if  $Z$  is not in the domain of  $f$ . (This is especially useful for  $f = \log \det$  and  $Z \not\preceq 0$ .) From optimality conditions, it can be shown that

$$\mathbf{prox}_{t f_1}(Z) := V \mathbf{diag}(q) V^T, \quad 2q_i = (d_i^2) + \sqrt{d_i^2 + 4t}$$

where  $V \mathbf{diag}(d) V^T$  is the eigenvalue decomposition of  $t\hat{C} - Z$ . Similarly,  $\mathbf{prox}_{t f_2}$  is the well-known *shrinkage operator*, and  $\mathbf{prox}_{t f_3}, \mathbf{prox}_{t f_4}$  are projections on their respective constraint spaces.

---

**Algorithm 3** Douglas-Rachford for  $f = \sum_{i=1}^m f_i$ 

---

**Input:** initial  $\{Z_i\}_1^m$  in  $\mathbb{S}^p$ ,  $t_1 > 0, 0 < t_2 < 2$ .  
1: **while** not converged **do**  
2:  $X_i = \mathbf{prox}_{t_1 f_i}(Z_i)$  for  $i = 1, \dots, m$   
3:  $Y_i = 2X_i - Z_i$  for  $i = 1, \dots, m$   
4:  $Y = m^{-1} \sum_i Y_i$   
5:  $Z_i = t_2(Y - X_i)$  for  $i = 1, \dots, m$   
6: **end while**  
**Output:** Any  $X_1 = \dots = X_m$

---

### 4.2 Random sparsity

For  $X \in \mathbb{S}^p$ , we randomly select  $l$  groups  $\beta_k \subset \{1, \dots, p\}$  of size  $b$ , and assume that this is the known group structure. Additionally, select  $\sigma_G \cdot l$  "active" groups (for  $0 < \sigma_G < 1$ )—the identity of these groups are not known in training. In this

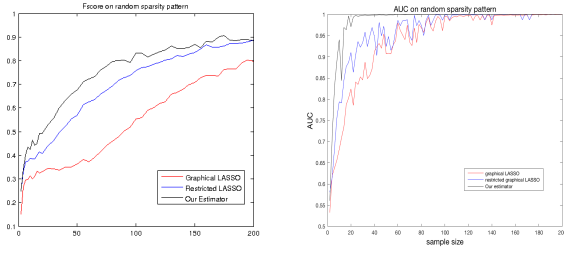


Figure 1: Random pattern sparse inverse covariance estimation for  $p = 100$ . F-measure (left) and AUC (right). Sample size range from 2 to 200. **missing x-y labels, don't need title (title should be y label). Also, change legend?**

simulation, we investigate the sample size required to recover the active groups, comparing G-LASSO, RG-LASSO, and NG-LASSO. Figure 4.2 shows the AUC as a function number of observations, where  $p = 100$ ,  $l = 100$ ,  $b = 5$  and  $\sigma_G = 0.1$ . From figure 4.2 we see that all methods recover the correct sparsity pattern given enough observations, and the sample complexity of G-LASSO is improved in RG-LASSO and even moreso with NG-LASSO.

### 4.3 Banded sparsity

For  $X \in \mathbb{S}^p$ , we assume that the true sparsity pattern consists of a nonzero diagonal and some active diagonal blocks of size  $b$ , where  $b$  is known but the true sparsity pattern is not. This gives in total  $l = p - b + 1$  candidate groups  $\beta_k = \{k, \dots, k + b\}$  for  $k = 1, \dots, l$ . Among  $l$  groups, we assume  $\sigma_G l$  groups are active (where  $0 < \sigma_G < 1$ ). Denote the set of active groups as  $I_A \subset \{\beta_1, \dots, \beta_l\}$  with  $|I_A| = \lceil \sigma_G \cdot l \rceil$ . Moreover, we simulate in-group sparsity; that is, for  $0 < \sigma_I \leq 1$ , we fix  $\Pr(X_{ij} \neq 0 | i, j \in \beta_k) = \sigma_I$ . Note that using only known information, we must assume the sparsity pattern is banded with bandwidth  $b$ .

We construct  $C$  with the true sparsity pattern, and form a sample covariance matrix  $\hat{C}$  sampling from a multivariate Gaussian with 0 mean and covariance  $C$ . The goal is to use the estimators to correctly recover the sparsity pattern of  $C$  using  $\hat{C}$  where the number of observations  $n$  is as small as possible.

Figure 2 shows a small example when  $p = 50$ ,  $n = 100$  and  $\sigma_I = 0.25$ . There are in total 90 groups in the banded sparsity pattern, where the 9 active groups (true sparsity) are in blue. We pick the estimator nonzeros by thresholding on the absolute value, choosing the threshold to, in each case, maximize  $\min\{\# \text{ true positives}, \# \text{ true negatives}\}$ . It is clear that, for this small example, G-LASSO (left) yields many spurious nonzeros. By simply restricting the sparsity pattern to  $B$ , the performance of RG-LASSO (center) already improves significantly, but NG-LASSO is still the best, since it accounts for sparsity in group selection as well.

Table 1 gives the result of a more extensive experiment, where the threshold,  $\alpha$ , and  $\gamma$  are picked to maximize AUC (Area Under the true-positive false-positive Curve), which is given for several  $p$ ,  $n$  and  $\sigma_I$ . Here, we see that NG-LASSO is comparable with RG-LASSO when  $p \approx n$ , but is consistently

better for  $p \gg n$ ; both, however, are considerably improved over G-LASSO.

Table 2 gives the per-iteration and total runtime of the three methods. In all cases, the per-iteration runtime depends only on  $p$  and  $b$ , and for larger  $p$ , NG-LASSO enjoys a much smaller per-iteration runtime. Of course, the number of iterations to convergence is also important; we notice that more iterations are usually required in all methods when  $p \gg n$ .

$p$	$n$	$\sigma_I$	$C$	$\hat{C}^\dagger$	G	RG	NG
100	10	0.1	0.43	0.44	0.50	0.57	0.65
100	10	0.25	0.39	0.40	0.48	0.58	0.64
100	100	0.1	0.59	0.60	0.89	0.88	0.88
100	100	0.25	0.52	0.69	0.83	0.80	0.85
1000	10	0.1	0.45	0.4	0.49	0.54	0.70
1000	10	0.25	0.40	0.41	0.56	0.58	0.72
1000	100	0.1	0.46	0.50	0.52	0.55	0.74
1000	100	0.25	0.40	0.55	0.60	0.62	0.82

Table 1: Best AUC scores for  $p \times p$  matrices with  $n$  samples, bandwidth  $p/10$ , and block sparsity  $\sigma_I$ . \* on boundary,  $\alpha = 2^{10}$ . G = G-LASSO. RG = RG-LASSO. NG = NG-LASSO.

$p$	$n$	Per Iteration			Overall		
		G	RG	NG	G	RG	NG
100	10	1.2e-2	3.1e-2	4.7e-2	1.7	8.7	4.3e1
100	100	1.9e-2	1.7e-2	3.8e-2	3.7e-1	1.7e-1	1.4
1000	10	2.7	8.1	4.8	1.9e2	8.9e2	8.5e1
1000	100	2.6	8.2	4.3	1.8e2	5.7e2	5.8e1*
2500	10	x	x	x	x	x	x
2500	100	3.8e1	1.0e2	x	2.7e3	7.3e3	x

Table 2: Runtimes in seconds for  $p \times p$  matrices with  $n$  samples, bandwidth  $p/10$ , and block sparsity  $\sigma_I = 0.1$ . G = G-LASSO. RG = RG-LASSO. NG = NG-LASSO. For  $p \leq 100$ ,  $\gamma$  and  $\alpha$  are the same as those used in Table 1. For  $p = 2500$ ,  $\gamma = 0.125$  and  $\alpha = 1$ , which was observed to work well for smaller  $p$ . 0 = unmeasurably small. \* not accurate, rerun

## 5 Financial application

We examine the performance of G-LASSO, RG-LASSO, and NG-LASSO on daily closing stock prices, obtained from Yahoo! Finance. Details on the data scraping are given in appendix B. Define  $u_i$  as the 1,005 length observation vector for stock  $i$ , and  $S_i$  as the set of indices of  $u_i$  where that stock price observation is available. Define  $V = \{1, 2, \dots, 100\}$ ,  $T = \{101, 102, \dots, 200\}$ , and  $R = \{201, 202, \dots, 200 + n\}$  as the indices of a validation, test, and train set respectively, and for all stocks  $i$ ,  $V_i = V \cap S_i$ ,  $T_i = T \cap S_i$ ,  $R_i = R \cap S_i$ .

<sup>1</sup> The sample covariance is then calculated as

$$\hat{C}_{ij} = \frac{1}{|R_i \cap R_j|} \sum_{k \in R_i \cap R_j} u_i[k] u_j[k].$$

We solve (1), (??), and (4) sweeping  $\rho$  and  $\gamma$  for powers of 2 from  $2^{-15}$  to  $2^5$ , using cross validation to pick  $\rho$  and  $\gamma$ .

<sup>1</sup>This extra detail is needed because not every day's value is provided for every stock.

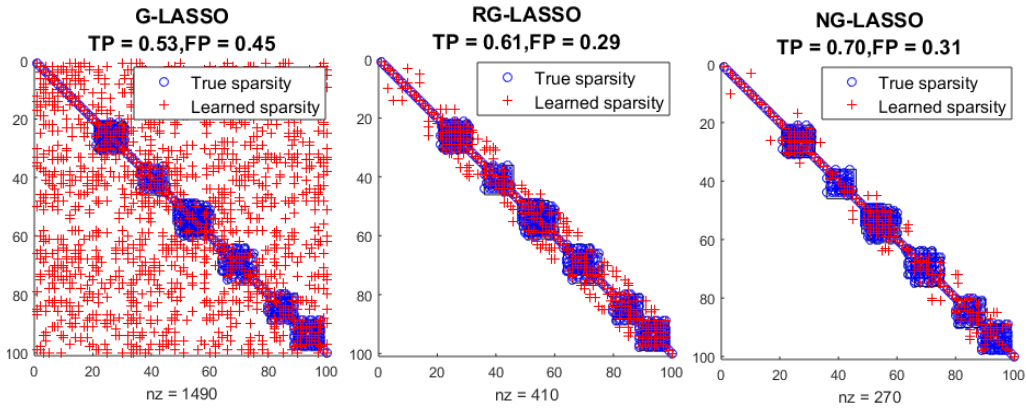


Figure 2: Banded pattern sparse inverse covariance estimation for  $p = n = 100$ . From left to right are G-LASSO (1), RG-LASSO (??) and NG-LASSO (5). TP = true positive, FP = false positive.

The performance is measured as the test negative log of the maximum likelihood estimate (NLL) for precision matrix  $X$  and samples  $\{u_i\}_{i \in T}$ :

$$\text{NLL} = -\log \det(X) + \frac{X_{ij}}{|T_i \cap T_j|} \sum_{k \in T_i \cap T_j} u_i[k] u_j[k].$$

Table 5 gives the test NLL for various  $p$  and  $n$ . The benefit of NG-LASSO is most obvious when  $n/p$  is very small. However, unlike in the banded example, the RG-LASSO test NLL values are not very low. We also experiment on arbitrary groups, to confirm that it is this specific group structure that is helping us. (todo)

Table 5 gives the runtime of each experiment for the best set of parameters. Since in this application all groups are nonoverlapping (all stocks are assigned a single sector and industry) in fact RG-LASSO is fully decomposable, and can run at the same per-iteration speed as NG-LASSO. so there'd better be some advantage in performance! ... showing 1) the time to do one set of eigenvalue decompositions ( $p \times p$  for LASSO and sum of  $|\beta_k| \times |\beta_k|, k = 1, \dots, l$  for group methods), 2) the average per-iteration runtime (which can result in multiple eigenvalue decomposition sweeps if line search is used) and 3) the total runtime. Here, the performance benefits of both RG-LASSO and NG-LASSO is very clear; for large  $p$  (on order of 1000s and 10000s) it is very difficult to solve G-LASSO without any decomposition.

## A Forward step derivation

The following is the derivation the Frank-Wolfe forward step in solving (6). To compute the forward step, we reformulate into a more generalized vector optimization problem.

$$\begin{aligned} & \underset{u_k}{\text{minimize}} && \langle \nabla f(x), \sum_{k=1}^l A_{\gamma_k} u_k \rangle \\ & \text{subject to} && \sum_{k=1}^l w_k \|u_k\|_2 \leq \alpha \\ & && u_k \in \mathcal{C}_k. \end{aligned} \quad (9)$$

$p$	$n$	sectors			industry	
		G	RG	NG	RG	NG
100	10	2.2e2	2.5e2	2.7e2	7.0e2	5.5e2
500	10	4.9e4	1.8e3	1.3e3	3.7e3	3.5e3
500	100	1.3e3	9.2e2	9.3e2	4.1e3	3.1e3
1000	10	2.5e3	2.6e3	2.8e3	7.9e3	6.1e3
1000	100	x	2.9e9	2.6e97	7.1e3	7.4e3
2500	10	x	4.2e9	x	x	x
2500	100					
5000	10					
5000	100					
14910	10					
14910	100					

Table 3: Best test negative log likelihood for different methods, varying the number of stocks ( $p$ ) and observations ( $n$ ). G = G-LASSO. RG = RG-LASSO. NG = NG-LASSO.

$p$	500	1000	2500	5000	14910
$p \times p$ ED	6.2e-2	3.4e-1	5.2	4.3e1	1.1e3
S-ED	0	1.6e-2	1.4e-1	6.4e-1	1.5e1
I-ED	0	0	0	6.2e-2	1.0
G 1 it.	3.84e-1	8.0e1	x		
RG (S) 1 it.	6.07e-2	3.2e-1	x		
RG (I) 1 it.	3.65e-2	8.9e-2	x		
NG (S) 1 it.	4.46e-2	2.0e-1	x		
NG (I) 1 it.	8.02e-2	1.8e-1	x		
G all	2.0e1	7.5e2	x		
RG (S) all	3.0	1.9e2	x		
RG (I) all	2.5	5.9e1	x		
NG (S) all	2.5	2.9e1	x		
NG (I) all	4.6	2.6e1	x		

Table 4: Runtimes (in seconds) of various algorithms for different matrix sizes ( $\hat{C}$  is  $p \times p$ ). ED=eigenvalue decomposition. S = sectors. I = industries. S-ED (I-ED) = time to compute  $p_i \times p_i$  ED where  $p_1, \dots, p_l$  are the sizes of the  $l$  sector (industry) groups. G = G-LASSO. RG = RG-LASSO. NG = NG-LASSO.

Here, the vector variable is  $x \in \mathbb{R}^m$ ,  $f(x)$  is a differentiable convex function, and  $\mathcal{C}_1, \dots, \mathcal{C}_l$  are proper convex cones. As before, the parameters  $w_1, \dots, w_l > 0$  are weights. The index  $\gamma_1, \dots, \gamma_l$  define the groups. To match  $\beta_k$  with the sets  $\gamma_k$ , the equivalence is such that  $\text{vec}(A_{\beta_k, \beta_k}) = \text{vec}(A)_{\gamma_k}$ ,  $k = 1, \dots, l$ .

For notational convenience, define  $c_k = \nabla f(x)_{\gamma_k}$  for  $k = 1, \dots, l$ . Then we can rewrite the forward step as

$$\begin{aligned} & \underset{x_k}{\text{maximize}} && \sum_k (-c_k)^T x_k \\ & \text{subject to} && \sum_{k=1}^l w_k \|x_k\|_2 \leq \alpha \\ & && x_k \in \mathcal{C}_k \end{aligned}$$

From Moreau's decomposition, any vector  $a$  can be written as the sum of its projection on a closed convex cone  $\mathcal{C}$  and its polar cone  $\mathcal{C}^\circ$ , of which are orthogonal. If we then expand

$$c_k^T x_k = \text{proj}_{\mathcal{C}_k}(c_k)^T x_k + \text{proj}_{\mathcal{C}_k^\circ}(c_k)^T x_k$$

then since feasible  $x_k \in \mathcal{C}_k^*$ , by definition of polar cone  $\text{proj}_{\mathcal{C}_k^\circ}(c_k)^T x_k \leq 0$ , and  $= 0$  only if  $x_k = s_k \text{proj}_{\mathcal{C}_k}(c_k)$ . This is the optimal choice of direction for  $x_k$ , since it also maximizes the first term  $\text{proj}_{\mathcal{C}_k}(c_k)^T x_k$ , and does not affect the norm constraint. If we additionally define scalars

$$a_k = \|P_{\mathcal{C}_k}(-c_k)\|_2^2, \quad b_k = w_k \|P_{\mathcal{C}_k}(-c_k)\|_2$$

then an even simpler equivalent formulation is

$$\begin{aligned} & \underset{s_k}{\text{maximize}} && a^T s \\ & \text{subject to} && b^T s \leq \alpha \\ & && s \geq 0 \end{aligned}$$

which is a linear program with a known optimal solution of

$$s_i = \begin{cases} \alpha/b_i & \text{if } i = \underset{i}{\text{argmax}} a_i/b_i \\ 0 & \text{else.} \end{cases}$$

Substituting gives the closed form solution in the text.

## B Yahoo! Finance data scraping details

Using the Yahoo! ticker downloader<sup>2</sup> we downloaded 27684 tickers for different stocks. We then used the Yahoo! finance API<sup>3</sup> to gather daily open, high, low, close, volume, and adjusted closing prices. We chose to monitor daily closing prices. We define groups as industries or sectors, as described in <https://biz.yahoo.com/p/>.

Any stock that we could not identify with an industry and sector was removed. We then prune the data to make sure it is as dense as possible; first, any day in which fewer than 14,000 stocks reported values are removed. Then, any stock with fewer than 90% of entries filled was removed. This resulted in 14,910 stocks and 1,005 daily closing prices. In total there are 9 sectors and 214 industries. with an average sector size 1656.7 and industry size 69.3. All stock vectors were then demeaned.

<sup>2</sup><https://pypi.python.org/pypi/Yahoo-ticker-downloader>

<sup>3</sup><http://chart.finance.yahoo.com/table.csv?s=TICKERNAMEHERE&a=400&b=23&c=2016&d=0&e=23&f=2017&g=d&ignore=.csv>

## References

- [Agler *et al.*, 1988] Jim Agler, William Helton, Scott McCullough, and Leiba Rodman. Positive semidefinite matrices with a given sparsity pattern. *Linear algebra and its applications*, 107:101–149, 1988.
- [Andersen *et al.*, 2013] Martin S Andersen, Joachim Dahl, and Lieven Vandenberghe. Logarithmic barriers for sparse matrix cones. *Optimization Methods and Software*, 28(3):396–423, 2013.
- [Bach *et al.*, 2011] Francis Bach, Rodolphe Jenatton, Julien Mairal, Guillaume Obozinski, et al. Convex optimization with sparsity-inducing norms. *Optimization for Machine Learning*, 5, 2011.
- [Banerjee *et al.*, 2006] Onureena Banerjee, Laurent El Ghaoui, Alexandre d'Aspremont, and Georges Natsoulis. Convex optimization techniques for fitting sparse gaussian graphical models. In *Proceedings of the 23rd international conference on Machine learning*, pages 89–96. ACM, 2006.
- [Banerjee *et al.*, 2008] Onureena Banerjee, Laurent El Ghaoui, and Alexandre d'Aspremont. Model selection through sparse maximum likelihood estimation for multivariate Gaussian or binary data. *Journal of Machine Learning Research*, 9(Mar):485–516, 2008.
- [Chandrasekaran *et al.*, 2012] Venkat Chandrasekaran, Benjamin Recht, Pablo A Parrilo, and Alan S Willsky. The convex geometry of linear inverse problems. *Foundations of Computational mathematics*, 12(6):805–849, 2012.
- [Combettes and Pesquet, 2011] Patrick L Combettes and Jean-Christophe Pesquet. Proximal splitting methods in signal processing. In *Fixed-point algorithms for inverse problems in science and engineering*, pages 185–212. Springer, 2011.
- [Dahl *et al.*, 2008] Joachim Dahl, Lieven Vandenberghe, and Vwani Roychowdhury. Covariance selection for nonchordal graphs via chordal embedding. *Optimization Methods & Software*, 23(4):501–520, 2008.
- [Danaher *et al.*, 2014] Patrick Danaher, Pei Wang, and Daniela M. Witten. The joint graphical LASSO for inverse covariance estimation across multiple classes. *Journal of the Royal Statistical Society. Series B: Statistical Methodology*, 76(2):373–397, 3 2014.
- [d'Aspremont *et al.*, 2008] Alexandre d'Aspremont, Onureena Banerjee, and Laurent El Ghaoui. First-order methods for sparse covariance selection. *SIAM Journal on Matrix Analysis and Applications*, 30(1):56–66, 2008.
- [Dunn and Harshbarger, 1978] Joseph C Dunn and S Harshbarger. Conditional gradient algorithms with open loop step size rules. *Journal of Mathematical Analysis and Applications*, 62(2):432–444, 1978.
- [Frank and Wolfe, 1956] Marguerite Frank and Philip Wolfe. An algorithm for quadratic programming. *Naval research logistics quarterly*, 3(1-2):95–110, 1956.

- [Friedman *et al.*, 2008] Jerome Friedman, Trevor Hastie, and Robert Tibshirani. Sparse inverse covariance estimation with the graphical LASSO. *Biostatistics*, 9(3):432–441, 2008.
- [Griewank and Toint, 1984] Andreas Griewank and Ph L Toint. On the existence of convex decompositions of partially separable functions. *Mathematical Programming*, 28(1):25–49, 1984.
- [Grone *et al.*, 1984] Robert Grone, Charles R Johnson, Eduardo M Sá, and Henry Wolkowicz. Positive definite completions of partial hermitian matrices. *Linear algebra and its applications*, 58:109–124, 1984.
- [Hosseini and Lee, 2016] Seyed Mohammad Javad Hosseini and Su-In Lee. Learning sparse gaussian graphical models with overlapping blocks. In *Advances in Neural Information Processing Systems*, pages 3801–3809, 2016.
- [Hsieh *et al.*, 2013] Cho-Jui Hsieh, Mátyás A Sustik, Inderjit S Dhillon, Pradeep K Ravikumar, and Russell Poldrack. BIG & QUIC: Sparse inverse covariance estimation for a million variables. In *Advances in Neural Information Processing Systems*, pages 3165–3173, 2013.
- [Jaggi, 2013] Martin Jaggi. Revisiting Frank-Wolfe: Projection-free sparse convex optimization. In *ICML (1)*, pages 427–435, 2013.
- [Lions and Mercier, 1979] Pierre-Louis Lions and Bertrand Mercier. Splitting algorithms for the sum of two nonlinear operators. *SIAM Journal on Numerical Analysis*, 16(6):964–979, 1979.
- [Liu, 1992] Joseph WH Liu. The multifrontal method for sparse matrix solution: Theory and practice. *SIAM review*, 34(1):82–109, 1992.
- [Mazumder and Hastie, 2012] Rahul Mazumder and Trevor Hastie. Exact covariance thresholding into connected components for large-scale graphical lasso. *Journal of Machine Learning Research*, 13(Mar):781–794, 2012.
- [Meinshausen and Bühlmann, 2006] Nicolai Meinshausen and Peter Bühlmann. High-dimensional graphs and variable selection with the LASSO. *The annals of statistics*, pages 1436–1462, 2006.
- [Moreau, 1962] Jean-Jacques Moreau. Fonctions convexes duales et points proximaux dans un espace hilbertien. *CR Acad. Sci. Paris Sér. A Math*, 255:2897–2899, 1962.
- [Negahban *et al.*, 2009] Sahand Negahban, Bin Yu, Martin J Wainwright, and Pradeep K Ravikumar. A unified framework for high-dimensional analysis of  $m$ -estimators with decomposable regularizers. In *Advances in Neural Information Processing Systems*, pages 1348–1356, 2009.
- [Obozinski *et al.*, 2011] Guillaume Obozinski, Laurent Jacob, and Jean-Philippe Vert. Group LASSO with overlaps: the latent group LASSO approach. *arXiv preprint arXiv:1110.0413*, 2011.
- [Rolfs *et al.*, 2012] Benjamin Rolfs, Bala Rajaratnam, Dominique Guillot, Ian Wong, and Arian Maleki. Iterative thresholding algorithm for sparse inverse covariance estimation. In *Advances in Neural Information Processing Systems*, pages 1574–1582, 2012.
- [Scheinberg and Rish, 2009] Katya Scheinberg and Irina Rish. Sinco-a greedy coordinate ascent method for sparse inverse covariance selection problem. *preprint*, 2009.
- [Yuan and Lin, 2007] Ming Yuan and Yi Lin. Model selection and estimation in the gaussian graphical model. *Biometrika*, 94(1):19–35, 2007.
- [Yuan, 2009] XiaoMing Yuan. Alternating direction methods for sparse covariance selection. *preprint*, 2(1), 2009.