
Supplement: Proximity Operator of the Matrix Perspective Function and its Applications

Joong-Ho Won

Department of Statistics
Seoul National University
wonj@stats.snu.ac.kr

A Proofs

A.1 A key lemma

Proofs of both Theorems 2 and 4 are based on the following key lemma, Lemma A.1. Recall that $\phi(x) = x_+$ for $x \in \mathbb{R}$ and $\phi^\square(\mathbf{X}) = \mathbf{P} \text{diag}(\phi(\lambda_1), \dots, \phi(\lambda_n)) \mathbf{P}^T = \mathbf{X}_+$ for $\mathbf{X} = \mathbf{P} \text{diag}(\lambda_1, \dots, \lambda_n) \mathbf{P}^T \in \mathbb{S}^n$ where \mathbf{P} satisfies $\mathbf{P}^T \mathbf{P} = \mathbf{P} \mathbf{P}^T = \mathbf{I}$. For any $\boldsymbol{\lambda} = (\lambda_1, \dots, \lambda_n)$, $\phi^{[1]}(\boldsymbol{\lambda})$ is the $n \times n$ symmetric matrix with (i, j) entry

$$\phi_{ij}^{[1]}(\boldsymbol{\lambda}) = \begin{cases} \frac{\phi(\lambda_i) + \phi(\lambda_j)}{|\lambda_i| + |\lambda_j|}, & \lambda_i \neq 0 \text{ or } \lambda_j \neq 0, \\ 0, & \lambda_i = \lambda_j = 0. \end{cases} \quad (18)$$

Also recall that $\mathbf{C}(\mu) = \bar{\mathbf{X}} - \mu \mathbf{e} \mathbf{e}^T$ so that $f(\mu) = g'(\mu) = 1 - \mathbf{e}^T \phi^\square(\mathbf{C}(\mu)) \mathbf{e}$. Lemma A.1 provides a closed-form expression of the derivative of $f(\mu)$ when it exists, in terms of the matrix function (18).

Lemma A.1. *Function f is differentiable at μ if and only if $\mathbf{e} \in \mathcal{N}(\mathbf{C}(\mu))^\perp$. In this case, the derivative is*

$$f'(\mu) = \mathbf{e}^T \mathbf{P} (\phi^{[1]}(\boldsymbol{\lambda}) \circ (\mathbf{P}^T \mathbf{e} \mathbf{e}^T \mathbf{P})) \mathbf{P}^T \mathbf{e}, \quad (\text{A.1})$$

for any $\boldsymbol{\lambda} = (\lambda_1, \dots, \lambda_n)^T$ and \mathbf{P} satisfying $\mathbf{C}(\mu) = \mathbf{P} \text{diag}(\lambda_1, \dots, \lambda_n) \mathbf{P}^T$, $\mathbf{P}^T \mathbf{P} = \mathbf{P} \mathbf{P}^T = \mathbf{I}$.

To prove this lemma, we begin by recalling the definition of directional derivatives.

Definition A.1 (Directional derivative). *For a function $F : \mathbb{R}^m \rightarrow \mathbb{R}^l$ and $\mathbf{x}, \mathbf{h} \in \mathbb{R}^m$, the directional derivative of F at \mathbf{x} along \mathbf{h} is defined and denoted by*

$$F'(\mathbf{x}; \mathbf{h}) = \lim_{t \downarrow 0} \frac{F(\mathbf{x} + t\mathbf{h}) - F(\mathbf{x})}{t}$$

if the limit exists. The F is called directionally differentiable at \mathbf{x} if $F'(\mathbf{x}; \mathbf{h})$ exists for all $\mathbf{h} \in \mathbb{R}^m$.

If F is differentiable at \mathbf{x} with Jacobian $\nabla F(\mathbf{x}) \in \mathbb{R}^{l \times m}$, then $F'(\mathbf{x}; \mathbf{h}) = \nabla F(\mathbf{x}) \mathbf{h}$.

For index sets $J, K \subset \{1, \dots, n\}$ and matrix $\mathbf{M} \in \mathbb{S}^n$, let \mathbf{M}_{JK} be the submatrix of \mathbf{M} constructed from the rows in J and the columns in K . The following lemma can be deduced from Sun and Sun [2002, Theorem 4.7]:

Lemma A.2. *Function ϕ^\square is directionally differentiable at any $\mathbf{X} \in \mathbb{S}^n$. Its directional derivative along $\mathbf{H} \in \mathbb{S}^n$ is*

$$(\phi^\square)'(\mathbf{X}; \mathbf{H}) = \mathbf{P} \begin{bmatrix} \phi_{KK}^{[1]}(\boldsymbol{\lambda}) \circ \tilde{\mathbf{H}}_{KK} & \phi_{KJ}^{[1]}(\boldsymbol{\lambda}) \circ \tilde{\mathbf{H}}_{KJ} \\ \phi_{JK}^{[1]}(\boldsymbol{\lambda}) \circ \tilde{\mathbf{H}}_{JK} & [\tilde{\mathbf{H}}_{JJ}]_+ \end{bmatrix} \mathbf{P}^T$$

where $\mathbf{X} = \mathbf{P} \text{diag}(\lambda_1, \dots, \lambda_n) \mathbf{P}^T \in \mathbb{S}^n$ with \mathbf{P} satisfying $\mathbf{P}^T \mathbf{P} = \mathbf{P} \mathbf{P}^T = \mathbf{I}$, $K = \{i \in \{1, \dots, n\} : \lambda_i \neq 0\}$, $J = \{i \in \{1, \dots, n\} : \lambda_i = 0\}$, and $\tilde{\mathbf{H}} = \mathbf{P}^T \mathbf{H} \mathbf{P}$. Furthermore, ϕ^\square is differentiable at \mathbf{X} if and only if \mathbf{X} is nonsingular, i.e., $J = \emptyset$.

Now we can prove the lemma:

Proof of Lemma A.1. Suppose f is differentiable at μ . Then if $C(\mu)$ is nonsingular, $\mathcal{N}(C(\mu)) = \{\mathbf{0}\}$ and $e \in \mathcal{N}(C(\mu))^\perp$. If $C(\mu)$ is singular, then the two one-sided limits

$$\lim_{t \downarrow 0} \frac{f(\mu + t) - f(\mu)}{t} \quad \text{and} \quad \lim_{t \downarrow 0} \frac{f(\mu - t) - f(\mu)}{-t}$$

must coincide. The first limit is equal to

$$\begin{aligned} -e^T \left(\lim_{t \downarrow 0} \frac{\phi^\square(C(\mu + t)) - \phi^\square(C(\mu))}{t} \right) e &= -e^T \left(\lim_{t \downarrow 0} \frac{\phi^\square(C(\mu) - tee^T) - \phi^\square(C(\mu))}{t} \right) e \\ &= -e^T (\phi^\square)'(C(\mu); -ee^T) e \\ &= e^T \mathbf{P} \begin{bmatrix} \phi_{KK}^{[1]}(\boldsymbol{\lambda}) \circ (\mathbf{P}^T ee^T \mathbf{P})_{KK} & \phi_{KJ}^{[1]}(\boldsymbol{\lambda}) \circ (\mathbf{P}^T ee^T \mathbf{P})_{KJ} \\ \phi_{JK}^{[1]}(\boldsymbol{\lambda}) \circ (\mathbf{P}^T ee^T \mathbf{P})_{JK} & -[(\mathbf{P}^T ee^T \mathbf{P})_{JJ}]_+ \end{bmatrix} \mathbf{P}^T e \end{aligned}$$

by Lemma A.2, for a spectral decomposition of $C(\mu) = \mathbf{P} \text{diag}(\lambda_1, \dots, \lambda_n) \mathbf{P}^T$ satisfying the conditions of the lemma. Likewise, the second limit equals

$$e^T (\phi^\square)'(C(\mu); ee^T) e = e^T \mathbf{P} \begin{bmatrix} \phi_{KK}^{[1]}(\boldsymbol{\lambda}) \circ (\mathbf{P}^T ee^T \mathbf{P})_{KK} & \phi_{KJ}^{[1]}(\boldsymbol{\lambda}) \circ (\mathbf{P}^T ee^T \mathbf{P})_{KJ} \\ \phi_{JK}^{[1]}(\boldsymbol{\lambda}) \circ (\mathbf{P}^T ee^T \mathbf{P})_{JK} & [(\mathbf{P}^T ee^T \mathbf{P})_{JJ}]_+ \end{bmatrix} \mathbf{P}^T e.$$

Let $\mathbf{P}^T e = [\mathbf{q}_K, \mathbf{q}_J]^T = \mathbf{q}$ where $\mathbf{q}_K \in \mathbb{R}^{|K|}$ and $\mathbf{q}_J \in \mathbb{R}^{|J|}$. Note $J \neq \emptyset$ since $C(\mu)$ is singular. Then the two limits are equal if and only if $\mathbf{q}_J^T [(\mathbf{q}\mathbf{q}^T)_{JJ}]_+ \mathbf{q}_J = 0$. It is immediate to see that $(\mathbf{q}\mathbf{q}^T)_{JJ} = \mathbf{q}_J \mathbf{q}_J^T \succeq \mathbf{0}$, hence $\mathbf{q}_J^T [(\mathbf{q}\mathbf{q}^T)_{JJ}]_+ \mathbf{q}_J = \|\mathbf{q}_J\|^4$. This implies $\mathbf{q}_J = \mathbf{0}$. Finally, observe that $\mathbf{q}_J = \mathbf{P}_J^T e$ where the columns of \mathbf{P}_J span $\mathcal{N}(C(\mu))$. Thus the condition $\mathbf{q}_J = \mathbf{P}_J^T e = \mathbf{0}$ is equivalent to $e \in \mathcal{N}(C(\mu))^\perp$.

Now suppose $e \in \mathcal{N}(C(\mu))^\perp$. If $C(\mu)$ is nonsingular, then Lemma A.2 implies that f is differentiable at μ . If $C(\mu)$ is singular, then $\mathbf{P}_J^T e = \mathbf{0}$ and the two one-sided limits in the above paragraph coincide, i.e., f is differentiable at μ .

Equation (A.1) is a consequence of the coincidence of the one-sided limits, that the common limit does not depend on the order of $\lambda_1, \dots, \lambda_n$, and the definition of $\phi^{[1]}$ in equation (18). \square

A.2 Proof of Theorem 2

For a solution μ^* to the equation $f(\mu) = 0$, define a collection of matrices related to the eigenvalues $\boldsymbol{\lambda}^* = (\lambda_1^*, \dots, \lambda_n^*)^T$ of $C(\mu^*)$:

$$\mathcal{M} = \{\mathbf{M} = (m_{ij}) \in \mathbb{S}^n : m_{ij} = \phi^{[1]}(\boldsymbol{\lambda}^*) \text{ if } \lambda_i^* \neq 0 \text{ or } \lambda_j^* \neq 0; m_{ij} \in [0, 1] \text{ if } \lambda_i^* = 0 = \lambda_j^*\}.$$

Also define the set (Bouligand subdifferential)

$$\partial_B f(\mu^*) = \left\{ \lim_{k \rightarrow \infty} f'(\mu_k) : \mu_k \rightarrow \mu^*, \mu_k \in D_f \right\}$$

where D_f denotes the set of points in which f is differentiable, so that $\partial f(\mu^*) = \text{conv } \partial_B f(\mu^*)$. The following lemma shows a representation of an element of this set in terms of \mathcal{M} :

Lemma A.3. *Suppose a spectral decomposition of $C(\mu^*)$ is $\mathbf{P}^* \text{diag}(\lambda_1^*, \dots, \lambda_n^*) \mathbf{P}^{*T}$ with $\mathbf{P}^{*T} \mathbf{P}^* = \mathbf{P}^* \mathbf{P}^{*T} = \mathbf{I}$. Then, for any $v \in \partial_B f(\mu^*)$, there exists $\mathbf{M} \in \mathcal{M}$ such that*

$$v = e^T \mathbf{P}^* (\mathbf{M} \circ (\mathbf{P}^{*T} ee^T \mathbf{P}^*)) \mathbf{P}^{*T} e.$$

Proof. By the definition of $\partial_B f(\mu^*)$, there exists a sequence $\{\mu_k\}$ such that f is differentiable at each μ_k , $\mu_k \rightarrow \mu^*$, and $f'(\mu_k) \rightarrow v$ as $k \rightarrow \infty$. Obviously $\mu_k \neq \mu$ for all k . Thus $C(\mu_k) = \bar{\mathbf{X}} - \mu_k ee^T = C(\mu) - (\mu_k - \mu) ee^T$ is a symmetric rank-1 perturbation of $C(\mu)$. Then, by Chen et al. [2003, Lemma 3.3], Rellich and Berkowitz [1969, Thm. 1], $C(\mu_k)$ has a spectral decomposition $\mathbf{P}_k \text{diag}(\lambda_{k,1}, \dots, \lambda_{k,n}) \mathbf{P}_k^T$ such that $\mathbf{P}_k \rightarrow \mathbf{P}^*$ as $k \rightarrow \infty$, by passing to a subsequence of

$\{\mu_k\}$ if necessary. Since $\lambda_{k,i} = (\mathbf{P}_k^T \mathbf{C}(\mu_k) \mathbf{P}_k)_{ii}$ and $\mathbf{C}(\mu)$ is continuous in μ , it follows that $\lim_{k \rightarrow \infty} \lambda_{k,i} = \lambda_i$ as well, for $i = 1, \dots, n$.

By Lemma A.1,

$$f'(\mu_k) = \mathbf{e}^T \mathbf{P}_k (\phi^{[1]}(\boldsymbol{\lambda}_k) \circ (\mathbf{P}_k^T \mathbf{e} \mathbf{e}^T \mathbf{P}_k)) \mathbf{P}_k^T \mathbf{e}.$$

Let

$$K = \{i \in \{1, \dots, n\} : \lambda_i^* \neq 0\}, \quad J = \{i \in \{1, \dots, n\} : \lambda_i^* = 0\}$$

and $\delta = \frac{1}{2} \min_{i \in K} |\lambda_i^*| > 0$. Then for all sufficiently large k , we have $\max_{i=1, \dots, n} |\lambda_{k,i} - \lambda_i^*| \leq \delta$. If $i \in K$ or $j \in K$, then $\lambda_{k,i} \neq 0$ or $\lambda_{k,j} \neq 0$, and

$$\phi_{ij}^{[1]}(\boldsymbol{\lambda}_k) = \frac{(\lambda_{k,i})_+ + (\lambda_{k,j})_+}{|\lambda_{k,i}| + |\lambda_{k,j}|} \rightarrow \frac{(\lambda_i^*)_+ + (\lambda_j^*)_+}{|\lambda_i^*| + |\lambda_j^*|} = \phi_{ij}^{[1]}(\boldsymbol{\lambda}^*).$$

If $i, j \in J$, then both $\lambda_{k,i}$ and $\lambda_{k,j}$ converge to 0. Since $\phi_{i,j}(\lambda_k) \in [0, 1]$ in this case, passing to a subsequence of $\{\mu_k\}$ if necessary, $\phi_{i,j}(\lambda_k)$ converges to a point $m_{ij} \in [0, 1]$. This shows that $\phi^{[1]}(\boldsymbol{\lambda}_k) \rightarrow \mathbf{M} \in \mathcal{M}$.

Finally, by the continuity of matrix multiplications, we have

$$v = \lim_{k \rightarrow \infty} f'(\mu_k) = \mathbf{e}^T \mathbf{P}^* (\mathbf{M} \circ (\mathbf{P}^{*T} \mathbf{e} \mathbf{e}^T \mathbf{P}^*)) \mathbf{P}^{*T} \mathbf{e}.$$

□

The next lemma provides a technical result useful for proving Theorem 2.

Lemma A.4. For $\mathbf{P}^* = (p_{ij})$ and $\lambda_1^*, \dots, \lambda_n^*$ in the statement of Lemma A.3, let $K_+ = \{i \in \{1, \dots, n\} : \lambda_i^* > 0\}$. Then $K_+ \neq \emptyset$ and

$$\sum_{i \in K_+} p_{ni}^2 > 0.$$

Proof. Denote the i th column of \mathbf{P} by $\mathbf{p}_i = (p_{1i}, \dots, p_{ni})^T$. Then $\phi^\square(\mathbf{C}(\mu^*)) = [\mathbf{C}(\mu^*)]_+ = \sum_{i \in K_+} \lambda_i^* \mathbf{p}_i \mathbf{p}_i^T$. From the optimality condition

$$1 = \mathbf{e}^T \phi^\square(\mathbf{C}(\mu^*)) \mathbf{e} = \sum_{i \in K_+} \lambda_i^* p_{ni}^2.$$

If $K_+ = \emptyset$ then the rightmost hand side is zero, a contradiction. That $K_+ \neq \emptyset$ and $\lambda_i^* > 0$ for all $i \in K_+$ succumbs to the fact $\sum_{i \in K_+} p_{ni}^2 > 0$. □

Now we are ready to prove the theorem.

Proof of Theorem 2. Let $v \in \partial f_B(\mu^*)$. Also let J, K , and K_+ be as defined in the proof of Lemma A.3 and the statement of Lemma A.4. Define $K_- = K \setminus K_+$. Then by Lemma A.3 there exists $\mathbf{M} = (m_{ij}) \in \mathbb{S}^n$ such that

$$m_{ij} = \begin{cases} 1, & \text{if } i \in K_+, j \in K_+ \cup J, \text{ or } i \in J, j \in K_+, \\ 0, & \text{if } i \in J, j \in K_-, \text{ or } i \in K_-, j \in J \cup K_-, \\ \tau_{ij} = \frac{\lambda_i^*}{\lambda_i^* - \lambda_j^*} \in (0, 1), & \text{if } i \in K_+, j \in K_-, \text{ or } i \in K_-, j \in K_+, \\ \in [0, 1], & \text{if } i, j \in J. \end{cases}$$

and

$$v = \mathbf{e}^T \mathbf{P}^* [\mathbf{M} \circ (\mathbf{P}^{*T} \mathbf{e} \mathbf{e}^T \mathbf{P}^*)] \mathbf{P}^{*T} \mathbf{e}$$

Then,

$$\begin{aligned}
v &= \text{Tr}(e^T \mathbf{P}^* [M \circ (\mathbf{P}^{*T} e e^T \mathbf{P}^*)] \mathbf{P}^{*T} e) \\
&= \text{Tr}(\mathbf{P}^{*T} e e^T P [M \circ (\mathbf{P}^{*T} e e^T \mathbf{P}^*)]) \\
&= \text{Tr}(\mathbf{Q} [M \circ \mathbf{Q}]), \quad \text{where } \mathbf{Q} = \mathbf{P}^{*T} e e^T \mathbf{P}^* = (q_{ij}) \\
&\geq \sum_{i \in K_+} \left(\sum_{j \in K_+ \cup J} q_{ij}^2 + \sum_{j \in K_-} \tau_{ij} q_{ij}^2 \right) \quad (\text{since } m_{ij} \geq 0) \\
&\geq \left(\min_{i \in K_+, j \in K_-} \tau_{ij} \right) \sum_{i \in K_+} \sum_{j=1}^n q_{ij}^2.
\end{aligned}$$

Since $\mathbf{P}^{*T} e = (p_{n1}, \dots, p_{nn})^T$ is the last row of \mathbf{P}^* , we have $q_{ij} = p_{ni} p_{nj}$ and

$$\sum_{i \in K_+} \sum_{j=1}^n q_{ij}^2 = \sum_{i \in K_+} \sum_{j=1}^n p_{ni}^2 p_{nj}^2 = \left(\sum_{i \in K_+} p_{ni}^2 \right) \left(\sum_{j=1}^n p_{nj}^2 \right) > 0.$$

The quantity is the first pair of parentheses is positive due to Lemma A.4. The second quantity equals to $e^T \mathbf{P}^* \mathbf{P}^{*T} e = e^T e = 1$. From this and $\tau_{ij} > 0$ for all $i \in K_+$ and $j \in K_-$, it follows that $v > 0$.

Since $\partial_B f(\mu^*)$ is compact and all the elements of this set is positive, and convex combination of its elements is also positive. It follows that every element of $\partial f(\mu^*) = \text{conv } \partial_B f(\mu^*)$ is positive.

The uniqueness of solution then follows from Clarke's inverse function theorem [Clarke, 1990, Thm. 7.1.1]; existence of solution is shown in Section 2 of the main text. \square

A.3 Proof of Theorem 4

The proof of Theorem 4 also requires Lemma A.1.

Proof of Theorem 4. If f is differentiable at μ , then $\partial f(\mu) = \{f'(\mu)\}$ and the result holds by Lemma A.1. Otherwise, consider a sequence $\{\mu_k\}$ such that $\mu_k \downarrow \mu$ and f is differentiable at each μ_k . Such a sequence exists since f is Lipschitz hence almost everywhere differentiable [Rockafellar and Wets, 2009, sec. 9J]. Obviously $\mu_k > \mu$ for all k . Thus $\mathbf{C}(\mu_k) = \bar{\mathbf{X}} - \mu_k e e^T = \mathbf{C}(\mu) - (\mu_k - \mu) e e^T$ is a symmetric rank-1 perturbation of $\mathbf{C}(\mu)$. Then, by Chen et al. [2003, Lemma 3.3], Rellich and Berkowitz [1969, Thm. 1], $\mathbf{C}(\mu_k)$ has a spectral decomposition $\mathbf{P}_k \text{diag}(\lambda_{k,1}, \dots, \lambda_{k,n}) \mathbf{P}_k^T$ such that $\mathbf{P}_k \rightarrow \mathbf{P}$ as $k \rightarrow \infty$, by passing to a subsequence if necessary. Since $\lambda_{k,i} = (\mathbf{P}_k^T \mathbf{C}(\mu_k) \mathbf{P}_k)_{ii}$ and $\mathbf{C}(\mu)$ is continuous in μ , it follows that $\lim_{k \rightarrow \infty} \lambda_{k,i} = \lambda_i$ as well, for $i = 1, \dots, n$. Moreover, $\lambda_{k,i} \leq \lambda_i$ for all i [Bunch et al., 1978]. Thus if $\lambda_i = 0 = \lambda_j$, then $\lambda_{k,i}, \lambda_{k,j} \uparrow 0$, which implies that $\lim_{k \rightarrow \infty} \phi^{[1]}(\boldsymbol{\lambda}_k) = \phi^{[1]}(\boldsymbol{\lambda})$. Now since from Lemma A.1,

$$f'(\mu_k) = e^T \mathbf{P}_k (\phi^{[1]}(\boldsymbol{\lambda}_k) \circ (\mathbf{P}_k^T e e^T \mathbf{P}_k)) \mathbf{P}_k^T e, \quad \boldsymbol{\lambda}_k = (\lambda_{k,1}, \dots, \lambda_{k,n})^T,$$

it follows that $\lim_{k \rightarrow \infty} f'(\mu_k) = v$. From Definition 1, we see $v \in \partial f(\mu)$. \square

B Applications to proximal algorithms

B.1 Heteroskedastic scaled lasso

In the heteroskedastic scaled lasso we want to minimize

$$\ell(\boldsymbol{\Omega}, \boldsymbol{\beta}) = \phi(\boldsymbol{\Omega}, \mathbf{X}\boldsymbol{\beta} - \mathbf{y}) + \frac{1}{2\sqrt{N}} \|\boldsymbol{\Omega}\|_F + \lambda \|\boldsymbol{\beta}\|_1. \quad (\text{B.1})$$

If we define the affine map $\mathcal{K} : (\boldsymbol{\Omega}, \boldsymbol{\beta}) \mapsto (\boldsymbol{\Omega}, \mathbf{X}\boldsymbol{\beta} - \mathbf{y})$, then problem (B.1) has the form (5), where $f(\boldsymbol{\Omega}, \boldsymbol{\beta}) \equiv 0$, $g(\boldsymbol{\Omega}, \boldsymbol{\beta}) = \frac{1}{2\sqrt{N}} \|\boldsymbol{\Omega}\|_F + \lambda \|\boldsymbol{\beta}\|_1$, and $h = \phi$. The adjoint \mathcal{K}^T of the linear part of \mathcal{K}

maps $(\Theta, \zeta) \in \mathbb{S}^p \times \mathbb{R}^p$ to $(\Theta, \mathbf{X}^T \zeta)$. Thus the resulting PDHG iteration is

$$\begin{aligned}\Omega^{k+1} &= \left(1 - \frac{\tau/(2\sqrt{N})}{\max[\|\mathbf{Y}\|_F, \tau/(2\sqrt{N})]}\right) \mathbf{Y}, \quad \mathbf{Y} = \Omega^k - \tau\Theta^k, \\ \beta^{k+1} &= S_{\tau\lambda}(\beta^k - \tau\mathbf{X}^T \zeta^k), \\ \tilde{\Omega}^{k+1} &= 2\Omega^{k+1} - \Omega^k, \\ \tilde{\beta}^{k+1} &= 2\beta^{k+1} - \beta^k, \\ (\Theta^{k+1}, \zeta^{k+1}) &= \text{prox}_{\sigma\phi^*}(\Theta^k + \sigma\tilde{\Omega}^{k+1}, \zeta^k + \sigma(\mathbf{X}\tilde{\beta}^{k+1} - \mathbf{y})).\end{aligned}$$

where $S_{\tau\lambda}$ is the usual soft-thresholding operator: $[S_{\tau\lambda}(x)]_i = \min(\max(x_i - \tau\lambda, 0), x_i + \tau\lambda)$.

In order to determine the step sizes, note $\mathcal{K}^T \mathcal{K} : (\Omega, \beta) \mapsto (\Omega, \mathbf{X}^T \mathbf{X} \beta - \mathbf{X}^T \mathbf{y})$. The norm of the linear part of this affine operator equals $\max(\|\mathbf{X}^T \mathbf{X}\|_2, 1) = \max(\|\mathbf{X}\|_2^2, 1) \leq \max(\|\mathbf{X}\|_F^2, 1)$.

Setup for experiments For all combinations of (N, p) in Table 2, data matrix $\mathbf{X} \in \mathbb{R}^{N \times p}$ were generated from zero-mean independent Gaussian. Each x_i was then scaled to have norm $1/\sqrt{p}$, so that $\|\mathbf{X}\|_F = 1$. Response vector \mathbf{y} was generated by setting $\mathbf{y} = \mathbf{X}\beta + \epsilon$, where the first five components of β were independently generated from $\mathcal{N}(0, 10^2)$ and the rest set to zero; noise vector ϵ was generated from zero-mean n -variate Gaussian with covariance matrix of compound symmetry

$$\Sigma = \begin{bmatrix} 1 & \rho & \rho & \dots & \rho \\ \rho & 1 & \rho & \dots & \rho \\ \vdots & & \ddots & & \vdots \\ \rho & \rho & \rho & \dots & 1 \end{bmatrix}$$

with $\rho = 0.5$. The regularization parameter $\lambda = 0.005$. The PDHG iteration was initialized by $\Omega^0 = \mathbf{I}_N$, $\beta^0 = \mathbf{0}$, $\Theta^0 = \mathbf{0}$, and $\zeta^0 = \mathbf{0}$. The step size parameters are $\tau = 0.99$ and $\sigma = 0.99$. Convergence was declared when the relative change of the primal variables (Ω^k, β^k) was less than 10^{-6} for $p < 300$ and 10^{-5} for $p \geq 300$. The maximum number of iterations was set to 50000.

B.2 Gaussian joint likelihood estimation

Joint maximum likelihood estimation (MLE) of Gaussian natural parameters (Ω, η) under the variance constraints

$$\begin{aligned}\text{minimize} \quad & \ell(\Omega, \eta) = -\log \det \Omega + \text{Tr}(\Omega \mathbf{S}) - 2\bar{\mu}^T \eta + \phi(\Omega, \eta) + \frac{\epsilon}{2} \|\Omega\|_F^2 \\ \text{subject to} \quad & \mathbf{c}_i^T \Omega^{-1} \mathbf{c}_i \leq 1, \quad i = 1, \dots, m\end{aligned} \quad (\text{B.2})$$

(the ridge penalty $\frac{\epsilon}{2} \|\Omega\|_F^2$ is added to ensure existence of the solution) has the form (5) if we define

$$\begin{aligned}f(\Omega, \eta) &= 0 \\ g(\Omega, \eta) &= -\log \det \Omega + \text{Tr}(\Omega \mathbf{S}) - 2\bar{\mu}^T \eta + \frac{\epsilon}{2} \|\Omega\|_F^2 \\ h(\mathbf{Z}_0, \mathbf{Z}_1, \dots, \mathbf{Z}_m, \eta) &= \phi(\mathbf{Z}_0, \eta) + \sum_{i=1}^m \iota_{C_i}(\mathbf{Z}_i), \quad C_i = \{\Omega \in \mathbb{S}^p : \mathbf{c}_i^T \Omega^{-1} \mathbf{c}_i \leq 1\},\end{aligned}$$

and the linear map $\mathcal{K} : (\Omega, \eta) \mapsto (\Omega, \Omega, \dots, \Omega, \eta) \in \prod_{i=0}^m \mathbb{S}^p \times \mathbb{R}^p$.

Since the adjoint \mathcal{K}^T of \mathcal{K} maps $(\Theta_0, \Theta_1, \dots, \Theta_m, \zeta) \in \prod_{i=0}^m \mathbb{S}^p \times \mathbb{R}^p$ to $(\sum_{i=0}^m \Theta_i, \zeta)$, the PDHG iteration for problem (B.2) entails

$$\begin{aligned}\Omega^{k+1} &= \mathbf{prox}_{-\frac{\tau}{1+\epsilon\tau} \log \det(\cdot)} \left(\frac{1}{1+\epsilon\tau} (\Omega^k - \tau \sum_{i=0}^m \Theta_i^k - \tau \mathcal{S}) \right), \\ \eta^{k+1} &= \eta^k - \tau \zeta^k + 2\tau \bar{\mu}, \\ \tilde{\Omega}^{k+1} &= 2\Omega^{k+1} - \Omega^k, \\ \tilde{\eta}^{k+1} &= 2\eta^{k+1} - \eta^k, \\ (\Theta_0^{k+1}, \zeta^{k+1}) &= \mathbf{prox}_{\sigma\phi^*} \left(\Theta_0^k + \sigma \tilde{\Omega}^{k+1}, \zeta^k + \sigma \tilde{\eta}^{k+1} \right), \\ \Theta_i^{k+1} &= \mathbf{prox}_{\sigma\nu_{C_i}^*} \left(\Theta_i^k + \sigma \tilde{\Omega}^{k+1} \right), \quad i = 1, \dots, m.\end{aligned}$$

It is well-known that

$$\mathbf{prox}_{-\tau \log \det(\cdot)}(M) = Q \operatorname{diag} \left(\frac{\mu_1 + \sqrt{\mu_1^2 + 4\tau}}{2}, \dots, \frac{\mu_p + \sqrt{\mu_p^2 + 4\tau}}{2} \right) Q^T$$

if the eigenvalue decomposition of $M \in \mathbb{S}^p$ is $Q \operatorname{diag}(\mu_1, \dots, \mu_p) Q^T$.

It remains to compute $\mathbf{prox}_{\sigma\nu_{C_i}^*}$. The following result shows it has a closed-form expression.

Proposition B.1. *Let $S_{c,\alpha} = \{\Omega \in \mathbb{S}^p : \phi(\Omega, c) \leq \alpha\}$ where $\alpha > 0$. Then $S_{c,\alpha}$ is closed and convex. Furthermore, the projection of $Z \in \mathbb{S}^p$ onto $S_{c,\alpha}$ is*

$$P_{S_{c,\alpha}}(Z) = \left(Z - \frac{1}{2\alpha} cc^T \right)_+ + \frac{1}{2\alpha} cc^T.$$

Therefore, from the Moreau decomposition (7), for $i = 1, \dots, m$,

$$\begin{aligned}\mathbf{prox}_{\sigma\nu_{C_i}^*}(Y) &= Y - \sigma P_{S_{c_i, 1/2}}(\sigma^{-1}Y) = \sigma \left(\frac{1}{\sigma} Y - c_i c_i^T \right) - \sigma \left(\frac{1}{\sigma} Y - c_i c_i^T \right)_+ \\ &= -\sigma \left(c_i c_i^T - \frac{1}{\sigma} Y \right)_+.\end{aligned}$$

Finally, to determine the step sizes, note $\mathcal{K}^T \mathcal{K} : (\Omega, \eta) \mapsto ((m+1)\Omega, \eta)$. Hence $\|\mathcal{K}^T \mathcal{K}\|_2 = m+1$.

Proof of Proposition B.1. Convexity and closedness of $S_{c,\alpha}$ follows from those of ϕ . The projection operator is

$$\begin{aligned}P_{S_{c,\alpha}}(Z) &= \arg \min_{\Omega \in \mathbb{S}^p} \frac{1}{2} \|\mathcal{Z} - \Omega\|_F^2 \text{ subject to } \phi(\Omega, c) \leq \alpha \\ &= \arg \min_{\Omega \in \mathbb{S}^p} \frac{1}{2} \|\mathcal{Z} - \Omega\|_F^2 \text{ subject to } \frac{1}{2} c^T \Omega^\dagger c \leq \alpha, c \in \mathcal{R}(\Omega) \\ &= \arg \min_{\Omega \in \mathbb{S}^p} \frac{1}{2} \|\mathcal{Z} - \Omega\|_F^2 \text{ subject to } \alpha - \frac{1}{2} c^T \Omega^\dagger c \geq 0, c \in \mathcal{R}(\Omega) \\ &= \arg \min_{\Omega \in \mathbb{S}^p} \frac{1}{2} \|\mathcal{Z} - \Omega\|_F^2 \text{ subject to } \Omega - \frac{1}{2\alpha} cc^T \succeq \mathbf{0} \\ &= \arg \min_{\Omega \in \mathbb{S}^p} \frac{1}{2} \left\| \mathcal{Z} - \frac{1}{2\alpha} cc^T - \left(\Omega - \frac{1}{2\alpha} cc^T \right) \right\|_F^2 \text{ subject to } \Omega - \frac{1}{2\alpha} cc^T \succeq \mathbf{0} \\ &= \left(\mathcal{Z} - \frac{1}{2\alpha} cc^T \right)_+ + \frac{1}{2\alpha} cc^T.\end{aligned}$$

The fourth equality is due to the Schur complements of

$$\begin{bmatrix} \Omega & -\frac{1}{\sqrt{2}}c \\ -\frac{1}{\sqrt{2}}c^T & \alpha \end{bmatrix} \succeq \mathbf{0}.$$

The last equality is from the fact $\arg \min_{\mathbf{X} \succeq \mathbf{0}} \frac{1}{2} \|\mathcal{Z} - \mathbf{X}\|_F^2 = \mathcal{Z}_+$. \square

Setup for experiments For all combinations of (N, p) in Table 2, data $\mathbf{x}_1, \dots, \mathbf{x}_N \in \mathbb{R}^p$ were generated from zero-mean multivariate Gaussian with covariance matrix of compound symmetry

$$\Sigma = \begin{bmatrix} 1 & \rho & \rho & \dots & \rho \\ \rho & 1 & \rho & \dots & \rho \\ \vdots & & \ddots & & \vdots \\ \rho & \rho & \rho & \dots & 1 \end{bmatrix}$$

with $\rho = 0.3$. The PDHG iteration used $\epsilon = 10/p^2$ and was initialized by

$$\begin{aligned} \Omega^0 &= (\mathbf{S} - \boldsymbol{\mu}\boldsymbol{\mu}^T + 10^{-2}\mathbf{I}_p)^{-1} \\ \eta^0 &= \Omega^0 \bar{\boldsymbol{\mu}} \\ \Theta_i^0 &= \Omega^0, \quad i = 0, 1, \dots, m \\ \zeta^0 &= \eta^0. \end{aligned}$$

The step size parameters are $\tau = 1$ and $\sigma = 1/(m+1)$. Convergence was declared when the relative change of the primal variables (Ω^k, η^k) was less than 10^{-5} . The maximum number of iterations was set to 50000.

B.3 Graphical model selection

Recall from equation (3) we want to minimize

$$-\frac{1}{N}PL(\Omega) + \lambda|\Omega|_1 = -\frac{1}{2} \sum_{i=1}^p \log \omega_{ii} + \phi(\mathcal{K}\Omega) + \lambda \sum_{i<j} |\omega_{ij}|. \quad (\text{B.3})$$

This has the form (5) if we define $f(\Omega) \equiv 0$, $g(\Omega) = -\frac{1}{2} \sum_{i=1}^p \log \omega_{ii} + \lambda \sum_{i<j} |\omega_{ij}|$, $h = \phi$, and the linear map $\mathcal{K} : \Omega \mapsto \frac{1}{N}(\mathbf{I}_N \otimes \Omega_D, \text{vec}(\Omega \mathbf{Y}^T))$. The adjoint of \mathcal{K} is

$$\mathcal{K}^T : (\mathbf{M}, \text{vec}(\mathbf{Z})) \mapsto \frac{1}{N} \sum_{i=1}^N \mathbf{M}_{ii,D} + \frac{1}{2N} (\mathbf{Z}\mathbf{Y} + \mathbf{Y}^T \mathbf{Z}^T),$$

for symmetric block matrix $\mathbf{M} = (\mathbf{M}_{ij}) \in \mathbb{S}^{Np}$ with $\mathbf{M}_{ij} = \mathbf{M}_{ji}^T \in \mathbb{R}^{p \times p}$, and $\mathbf{Z} \in \mathbb{R}^{p \times N}$. Then the PDHG iteration for problem (B.3) is

$$\begin{aligned} \Omega^{k+1} &= \text{prox}_{\tau g} \left(\Omega^k - \frac{\tau}{N} \left(\sum_{i=1}^N \Theta_{ii,D}^k + \frac{1}{2} \mathbf{Z}^k \mathbf{Y} + \frac{1}{2} \mathbf{Y}^T [\mathbf{Z}^k]^T \right) \right) \\ \tilde{\Omega}^{k+1} &= 2\Omega^{k+1} - \Omega^k \\ (\Theta^{k+1}, \text{vec}(\mathbf{Z}^{k+1})) &= \text{prox}_{\sigma \phi^*} \left(\Theta^k + \frac{\sigma}{N} (\mathbf{I}_N \otimes \tilde{\Omega}_D^{k+1}), \text{vec}(\mathbf{Z}^k + \frac{\sigma}{N} \tilde{\Omega}^{k+1} \mathbf{Y}^T) \right) \end{aligned}$$

where $\Omega^k, \tilde{\Omega}^k \in \mathbb{S}^p$, $\mathbf{Z}^k \in \mathbb{R}^{p \times N}$, and $\Theta^k = (\Theta_{ij}^k) \in \mathbb{S}^{Np}$, with $\Theta_{ij} = \Theta_{ji}^T \in \mathbb{R}^{p \times p}$. Operator $\text{prox}_{\tau g}$ has a closed form expression. For $\mathbf{W} = (w_{ij})$,

$$[\text{prox}_{\tau g}(\mathbf{W})]_{ij} = \begin{cases} \frac{1}{2}(w_{ii} + \sqrt{w_{ii}^2 + 2\tau}), & i = j, \\ S_{\tau\lambda/2}(w_{ij}), & i \neq j. \end{cases}$$

It is easy to see that $\mathcal{K}^T \mathcal{K} : \Omega \mapsto \frac{1}{N} \Omega_D + \frac{1}{2N^2} (\Omega \mathbf{Y}^T \mathbf{Y} + \mathbf{Y}^T \mathbf{Y} \Omega)$. Then $\text{vec}(\frac{1}{N} \Omega_D + \frac{1}{2N^2} [\Omega \mathbf{Y}^T \mathbf{Y} + \Omega \mathbf{Y}^T \mathbf{Y}]) = \left(\frac{1}{N} \mathbf{A} + \frac{1}{2N^2} (\mathbf{Y}^T \mathbf{Y} \otimes \mathbf{I}_p + \mathbf{I}_p \otimes \mathbf{Y}^T \mathbf{Y}) \right) \text{vec}(\Omega)$ where \mathbf{A} satis-

ties $\text{vec}(\Omega_D) = \mathbf{A} \text{vec}(\Omega)$. It follows that $\mathbf{A}^T \mathbf{A} = \mathbf{I}_{p^2}$ and $\|\mathbf{A}\|_2 = 1$. Therefore,

$$\begin{aligned} \|\mathcal{K}^T \mathcal{K}\|_2 &= \left\| \frac{1}{N} \mathbf{A} + \frac{1}{2N^2} (\mathbf{Y}^T \mathbf{Y} \otimes \mathbf{I}_p + \mathbf{I}_p \otimes \mathbf{Y} \mathbf{Y}^T) \right\|_2 \\ &\leq \frac{1}{N} \|\mathbf{A}\|_2 + \frac{1}{2N^2} \|\mathbf{Y}^T \mathbf{Y} \otimes \mathbf{I}_p\|_2 + \frac{1}{2N^2} \|\mathbf{I}_p \otimes \mathbf{Y}^T \mathbf{Y}\|_2 \\ &= \frac{1}{N} (1) + \frac{1}{2N^2} \lambda_{\max}(\mathbf{Y}^T \mathbf{Y}) + \frac{1}{2N^2} \lambda_{\max}(\mathbf{Y}^T \mathbf{Y}) \\ &= \frac{1}{N} + \frac{1}{N^2} \|\mathbf{Y}\|_2^2 \\ &\leq \frac{1}{N} + \frac{1}{N^2} \|\mathbf{Y}\|_F^2, \end{aligned}$$

which determines the step size.

Setup for experiments For all combinations of (N, p) in Table 2, data $\mathbf{y}_1, \dots, \mathbf{y}_N \in \mathbb{R}^p$ were generated from zero-mean multivariate Gaussian with precision matrix

$$\Omega = 10\mathbf{I}_p + \Xi + \Xi^T,$$

where Ξ is a $p \times p$ sparse random Gaussian matrix with 1 percent sparsity level. The regularization parameter $\lambda = 0.1$. The PDHG iteration was initialized by

$$\begin{aligned} \Omega^0 &= (\mathbf{S} + 10^{-2} \mathbf{I}_p)^{-1} \\ \Theta_i^0 &= \mathbf{I}_N \otimes \Omega_D^0 \\ \mathbf{Z}^0 &= \Omega^0 \mathbf{Y}^T. \end{aligned}$$

The step size parameters are $\tau = 2$ and $\sigma = 1/(2L_K)$ where $L_K = 1/N + \|\mathbf{Y}\|_F^2/N^2$. Convergence was declared when the relative change of the primal variable Ω^k was less than 10^{-5} . The maximum number of iterations was set to 50000. For the symmetric lasso used for comparison the implementation in the gconcord R package (<https://cran.r-project.org/web/packages/gconcord/index.html>) was used with the same input.

References

- James R Bunch, Christopher P Nielsen, and Danny C Sorensen. Rank-one modification of the symmetric eigenproblem. *Numer. Math.*, 31(1):31–48, 1978.
- Xin Chen, Houduo Qi, and Paul Tseng. Analysis of nonsmooth symmetric-matrix-valued functions with applications to semidefinite complementarity problems. *SIAM J. Optim.*, 13(4):960–985, 2003.
- Frank H Clarke. *Optimization and Nonsmooth Analysis*. Society for Industrial and Applied Mathematics, Philadelphia, USA, 1990.
- Franz Rellich and Jerome Berkowitz. *Perturbation Theory of Eigenvalue Problems*. Gordon and Breach, New York, 1969.
- R Tyrrell Rockafellar and Roger J-B Wets. *Variational Analysis*, volume 317 of *Grundlehren der mathematischen Wissenschaften*. Springer Science & Business Media, New York, 2009.
- Defeng Sun and Jie Sun. Semismooth matrix-valued functions. *Math. Oper. Res.*, 27(1):150–169, 2002.