A Coordinate-wise Optimization Algorithm for Sparse Inverse Covariance Selection

Ganzhao Yuan, Haoxian Tan, Wei-Shi Zheng

Abstract—Sparse inverse covariance selection is a fundamental problem for analyzing dependencies in high dimensional data. However, such a problem is difficult to solve since it is NP-hard. Existing solutions are primarily based on convex approximation and iterative hard thresholding, which only lead to sub-optimal solutions. In this work, we propose a coordinate-wise optimization algorithm to solve this problem which is guaranteed to converge to a coordinate-wise minimum point. The algorithm iteratively and greedily selects one variable or swaps two variables to identify the support set, and then solves a reduced convex optimization problem over the support set to achieve the greatest descent. As a side contribution of this paper, we propose a Newton-like algorithm to solve the reduced convex sub-problem, which is proven to always converge to the optimal solution with global linear convergence rate and local quadratic convergence rate. Finally, we demonstrate the efficacy of our method on synthetic data and real-world data sets. As a result, the proposed method consistently outperforms existing solutions in terms of accuracy.

Index Terms—Sparse Optimization, Coordinate Descent Algorithm, Inverse Covariance Selection, Nonconvex Optimization, Convex Optimization.

1 Introduction

In this paper, we mainly focus on the following nonconvex optimization problem:

$$\min_{\mathbf{X} \in \mathbb{R}^{n \times n}} f(\mathbf{X}) \triangleq \langle \mathbf{\Sigma}, \mathbf{X} \rangle - \log \det(\mathbf{X}),
s.t. \mathbf{X} \succ \mathbf{0}, \|\mathbf{X}\|_{0.\text{off}} \leq s,$$
(1)

where $\Sigma \in \mathbb{R}^{n \times n}$ is a given symmetric covariance matrix of the input data set, $\|\cdot\|_{0,\text{off}}$ counts the number of non-diagonal and non-zero elements of a square matrix, and s is a positive integer that specifies the sparsity level of the solution. $\mathbf{X} \succ 0$ means \mathbf{X} is positive definite. $\langle \cdot, \cdot \rangle$ stands for the standard inner product.

The optimization problem in (1) is known as sparse inverse covariance selection in the literature [9, 13]. It provides a good way of analyzing dependencies in high dimensional data and captures varieties of applications in computer vision and machine learning (e.g. biomedical image analysis [10], scene labeling [27], brain functional network classification [38]). The log-determinant function is introduced for maximum likelihood estimation, and the ℓ_0 norm is used to reduce over-fitting and improve the interpretability of the model. We remark that when the sparsity constraint is absent, one can set the gradient of the objective function $f(\cdot)$ to zero (i.e. $\Sigma - \mathbf{X}^{-1} = \mathbf{0}$) and output Σ^{-1} as the optimal solution.

Problem (1) is very challenging due to the introduction of the combinatorial ℓ_0 norm. Existing solutions can be categorized into

two classes: convex ℓ_1 approximation and iterative hard thresholding. Convex ℓ_1 approximation simply replaces the ℓ_0 norm by its tightest convex relaxation ℓ_1 norm. In the past decades, a plethora of approaches have been proposed to solve the ℓ_1 norm approximation problem, which include projected sub-gradient method [8], (linearized) alternating direction method [26, 34], quadratic approximation method [25, 24, 12, 11], block coordinate descent method [9, 1], Nesterov's first-order optimal method [18, 19, 6], primal-dual interior point method [17]. Despite the popularity of convex methods, they fail to control the sparsity of the solution and often lead to sub-optimal accuracy for the original non-convex problem. Recent attention has been paid to solving the original non-convex problem directly by the researchers [31, 36, 5].

Iterative hard thresholding method considers iteratively setting the small elements (in magnitude) to zero in a gradient descent manner. By using this strategy, it is able to control the sparsity of the solution directly and exactly. Due to its simplicity, it has been widely used and incorporated into the optimization framework of penalty decomposition algorithm [20] and mean doubly alternating direction method [7]. In [20], it is shown that for the general sparse optimization problem, any accumulation point of the sequence generated by the penalty decomposition algorithm always satisfies the first-order optimality condition of the problem.

Recently, A. Beck and Y. Vaisbourd present and analyze a new optimality criterion which is based on coordinate-wise optimality [3]. They show that coordinate-wise optimality is strictly stronger than the optimality criterion based on hard thresholding. They apply their algorithm to principal component analysis and show that their method consistently outperforms the well-known truncated power method [33]. Inspired by this work, we extend their method to solve sparse inverse covariance selection problem. We are also aware of the work [21] where a cyclic coordinate descent decent algorithm (combined with a randomized initialization strategy) is considered to solve the sparse inverse covariance selection problem. However, their method only addresses the ℓ_0 norm regularized optimization problem and it fails to control the sparsity

Ganzhao Yuan is with School of Data and Computer Science, Sun Yat-sen University (SYSU), China. E-mail: yuanganzhao@gmail.com

Haoxian Tan is with School of Data and Computer Science, Sun Yat-sen University (SYSU), China. E-mail: tanhx3@mail2.sysu.edu.cn

Wei-Shi Zheng is with School of Data and Computer Science, Sun Yat-sen University (SYSU), China.E-mail: zhwshi@mail.sysu.edu.cn

level of the solution.

Contributions: The contributions of this work are three-fold. (i) We propose a new coordinate-wise optimization algorithm for sparse inverse covariance selection. The algorithm iteratively and greedily selects one variable or swap two variables to identify the support set, and then solves a reduced convex optimization problem over the support set (See Section 2). (ii) An efficient Hessianfree Newton-like algorithm to solve the convex subproblem is proposed (Section 3). (iii) We provide some theoretical analysis for the proposed Coordinate-Wise Optimization Algorithm (CWOA) and the Newton-Like Optimization Algorithm (NLOA). We prove that CWOA is guaranteed to converge to a coordinatewise minimum point of the original nonconvex problem and the NLOA is guaranteed to converge to the global optimal solution of the convex subproblem with global linear rate and local quadratic convergence rate (Section 4). (iv) Extensive experiments have shown that our method consistently outperforms existing solutions in terms of accuracy (Section 5).

Notations: In this paper, boldfaced lowercase letters denote vectors and uppercase letters denote real-valued matrices. We denote $\lambda(\mathbf{X}) \in \mathbb{R}^n$ as the eigenvalues of \mathbf{X} in increasing order. All vectors are column vectors and superscript T denotes transpose. $\text{vec}(\mathbf{X}) \in \mathbb{R}^{n^2 \times 1}$ stacks the columns of the matrix $\mathbf{X} \in \mathbb{R}^{n \times n}$ into a column vector and $mat(\mathbf{x}) \in \mathbb{R}^{n \times n}$ converts $\mathbf{x} \in \mathbb{R}^{n^2 \times 1}$ into a matrix. Thus, $vec(mat(\mathbf{x})) = \mathbf{x}$ and $mat(vec(\mathbf{X})) = \mathbf{X}$. We use $\langle \mathbf{X}, \mathbf{Y} \rangle$ and $\mathbf{X} \otimes \mathbf{Y}$ to denote the Euclidean inner product and Kronecker product of X and Y, respectively. For any matrix $\mathbf{X} \in \mathbb{R}^{n \times n}$ and any $i,j \in \{1,2,...,n\}$, we denote by \mathbf{X}_{ij} the element of \mathbf{X} in i^{th} row and j^{th} column and use \mathbf{X}_k to denote the k position of vec(X). Therefore, we have $X_{ij} = X_{(j-1)\times n+i}$. We denote \mathbf{e}_i as a unit vector with a 1 in the i^{th} entry and 0 in all other entries. We use $j \in \{1, 2, ..., n^2\}$ to denote any position in a square matrix of size $n \times n$ where n is known from the context, and use row(j) and col(j) to denote the corresponding row and column for j. We denote \mathbf{E}_j is a square symmetric matrix with the entries (row(j), col(j)) and (col(j), row(j))equal 1 and 0 in all other ones. Note that when $row(j) \neq col(j)$, we have $\mathbf{E}_j = \mathbf{e}_{row(j)} \mathbf{e}_{col(j)}^T + \mathbf{e}_{col(j)} \mathbf{e}_{row(j)}^T$. Finally, for any matrices $\mathbf{X} \in \mathbb{R}^{n \times n}$ and $\mathbf{H} \in \mathbb{R}^{n^2 \times n^2}$, we define $\|\mathbf{D}\|_{\mathbf{H}}^2 \triangleq$ $\text{vec}(\mathbf{D})^T \mathbf{H} \text{vec}(\mathbf{D})$ and $\mathbf{H} \circ \mathbf{X} \triangleq \text{mat}(\mathbf{H} \text{vec}(\mathbf{X})) \in \mathbb{R}^{n \times n}$.

2 COORDINATE-WISE OPTIMIZATION ALGORITHM

This section presents our coordinate-wise optimization algorithm which is guaranteed to converge after a finite amount of iterations to a coordinate-wise minimum point [2, 3]. We denote $\bar{\mathcal{S}}(\mathbf{X})$ and $\bar{\mathcal{Z}}(\mathbf{X})$ as the index of *non-diagonally* non-zero elements and zero elements of \mathbf{X} , respectively.

First of all, we notice that when the support set S is known, problem (1) reduces to the following convex optimization problem:

$$\min_{\mathbf{X} \succ 0} f(\mathbf{X}), \ s.t. \ \mathbf{X}_Z = \mathbf{0}, \text{ with } Z \triangleq \{1, 2, ..., n^2\} \setminus S. \tag{2}$$

Our algorithm iteratively and greedily selects one variable or swaps two variables to identify the support set S, and then solves a reduced convex sub-problem in (2) to achieve the greatest descent.

We summarize our proposed method in Algorithm 1 and have a few remarks on it below.

• Two-stage algorithm. At each iteration of the algorithm, one or two variables of the solution are updated. At the first

Algorithm 1 CWOA: A Coordinate-wise Optimization Algorithm for Sparse Inverse Covariance Selection.

```
Input: Sparsity level s.
Output: The solution X^*.
Initialization: Set X^0 = O^{-1}, where O is a diagonal matrix
with \mathbf{O}_{ii} = \mathbf{\Sigma}_{ii}, \ \forall i \in [n]. \ \mathrm{Set} \ k = 0.
     \\ Greedy Pursuit Stage
    while \|\mathbf{X}^k\|_{0,\text{off}} < s \text{ do}
         for every j \in \bar{\mathcal{Z}}(\mathbf{X}^k) do
                 f_j = \min_{\alpha} f(\mathbf{X}^k + \theta \mathbf{E}_j), s.t. \mathbf{X}^k + \theta \mathbf{E}_j > 0
                                                                                                                          (3)
         end for
         \begin{aligned} j_k &= \arg \min\{f_j: j \in \bar{\mathcal{Z}}(\mathbf{X}^k)\} \\ \text{if } f_{j_k} &< f(\mathbf{X}^k) \text{ then} \\ \text{Solve (2) to get } \mathbf{X}^{k+1} \text{ with } S = \bar{\mathcal{S}}(\mathbf{X}^k) \cup j_k. \end{aligned}
         end if
    end while
     \\ Swap Coordinates Stage
    for every i \in \bar{\mathcal{S}}(\mathbf{X}^k) do
         for every j \in \bar{\mathcal{Z}}(\mathbf{X}^k) do
                            f_{i,j} = \min_{\alpha} f(\mathbf{X}^k - \mathbf{X}_i^k \mathbf{E}_i + \theta \mathbf{E}_j),
                                                                                                                          (4)
                                            s.t. \mathbf{X}^k - \mathbf{X}_i^k \mathbf{E}_i + \theta \mathbf{E}_i > 0
         end for
    end for
    (i_k, j_k) = \underset{k}{\operatorname{arg min}} \{ f_{i,j} : i \in \bar{\mathcal{S}}(\mathbf{X}^k), j \in \bar{\mathcal{Z}}(\mathbf{X}^k) \}
    if f_{i_k,j_k} < f(\mathbf{X}^k) then Solve (2) to get \mathbf{X}^{k+1} with S = (\bar{\mathcal{S}}(\mathbf{X}^k) \setminus i_k) \cup j_k.
          k = k + 1
    else
         set \mathbf{X}^* \leftarrow \mathbf{X}^{k+1}
         break
    end if
end while
```

greedy pursuit stage, the algorithm greedily picks one coordinate $i \in \bar{\mathcal{Z}}(\mathbf{X}^k)$ that leads to the greatest descent from $\bar{\mathcal{Z}}(\mathbf{X}^k)$. This strategy is also known as forward greedy selection in the literature [28, 37]. At the second swap coordinates stage, the algorithm enumerates all the possible pairs (i,j) with $i \in \bar{\mathcal{S}}(\mathbf{X}^k)$ and $j \in \bar{\mathcal{Z}}(\mathbf{X}^k)$ that leads to the greatest descent and changes the two coordinates from zero/non-zero to non-zero/zero. At both stages, once the support set has been updated, Algorithm 1 runs a convex subproblem procedure to solve (2) over the support set to compute a more 'compact' solution.

• One-dimensional sub-problem. The problems in (3) and (4) reduce to the following optimization problem:

$$\min_{\theta} f(\theta) \triangleq \langle \mathbf{\Sigma}, \mathbf{V} + \theta \mathbf{E}_{j} \rangle - \log \det(\mathbf{V} + \theta \mathbf{E}_{j})$$

$$s.t. \mathbf{V} + \theta \mathbf{E}_{i} \succ \mathbf{0}$$
(5)

with $\mathbf{V} = \mathbf{X}^k$ for (3) and $\mathbf{V} = \mathbf{X}^k - \mathbf{X}_i^k \cdot \mathbf{E}_i$ for (4). We now discuss how to simplify problem (5). We define $\mathbf{Y} \triangleq \mathbf{V}^{-1}$ and obtain the following equations: $\det(\mathbf{V} + \theta \mathbf{E}_j) = \det(\mathbf{V}(\mathbf{I} + \theta \mathbf{Y} \mathbf{E}_j)) = \det(\mathbf{V}) \det(\mathbf{I} + \theta \mathbf{U}_{cr}^T \mathbf{Y} \mathbf{U}_{rc}) = \det(\mathbf{V})(1 + 2\mathbf{Y}_{rc}\theta - \mathbf{O}_{rc}\theta^2)$, where $\mathbf{O}_{rc} \triangleq \mathbf{Y}_{rr} \mathbf{Y}_{cc} - \mathbf{Y}_{rc}^2 > 0$, $r \triangleq row(j)$, $c \triangleq$

 $col(j), \ \mathbf{U}_{rc} \triangleq [\mathbf{e}_r \ \mathbf{e}_c] \in \mathbb{R}^{n \times 2}, \ \mathbf{U}_{cr} \triangleq [\mathbf{e}_c \ \mathbf{e}_r] \in \mathbb{R}^{n \times 2}.$ Noticing that $-\log \det(\mathbf{V} + \theta \mathbf{E}_j) = -\log \det(\mathbf{V}) - \log(1 + 2\mathbf{Y}_{rc}\theta - \mathbf{O}_{rc}\theta^2)$ and $\langle \mathbf{\Sigma}, \mathbf{E}_j \rangle = 2\mathbf{\Sigma}_j = 2\mathbf{\Sigma}_{rc}$, we can simplify problem (5) to the following one-dimensional convex problem:

$$\min_{\theta} f(\theta) \triangleq 2\theta \mathbf{\Sigma}_{rc} - \log \left(1 + 2\mathbf{Y}_{rc}\theta - \mathbf{O}_{rc}\theta^2 \right) + C$$

s.t. $1 + 2\mathbf{Y}_{rc}\theta - \mathbf{O}_{rc}\theta^2 > 0$,

where $C = \langle \mathbf{\Sigma}, \mathbf{V} \rangle - \log \det(\mathbf{V})$ is a constant. Noting that $f(\theta)$ is differentiable, we set the gradient of $f(\theta)$ to zero, we obtain $0 = 2\mathbf{\Sigma}_{rc} + \frac{2\mathbf{O}_{rc}\theta - 2\mathbf{Y}_{rc}}{1 + 2\mathbf{Y}_{rc}\theta - \mathbf{O}_{rc}\theta^2}$. There are two solution to this equation. However, only one of these satisfies the bound constraint $1 + 2\mathbf{Y}_{rc}\theta - \mathbf{O}_{rc}\theta^2 > 0$. Thus, the optimal solution θ^* can be computed as:

$$\theta^* = \begin{cases} \mathbf{Y}_{rc}/\mathbf{O}_{rc}, & \text{if } \mathbf{\Sigma}_{rc} = 0; \\ \mathbf{Y}_{rc}/\mathbf{O}_{rc} + 1/(2\mathbf{\Sigma}_{rc}) - \\ \sqrt{\mathbf{O}_{rc}^2 + 4\mathbf{\Sigma}_{rc}^2 \mathbf{Y}_{rr} \mathbf{Y}_{cc}}/(2\mathbf{O}_{rc}\mathbf{\Sigma}_{rc}), & \text{if } \mathbf{\Sigma}_{rc} \neq 0. \end{cases}$$

• Fast matrix computation. In our algorithm, we assume that $\mathbf{Y} \triangleq \mathbf{V}^{-1}$ is available. This can be achieved by using the follow strategy. We keep a record of \mathbf{X}^{-1} in every iteration. Once the solution \mathbf{X} is changed to $\mathbf{T} = \mathbf{X} + \varpi \mathbf{E}_j$, we quickly estimate \mathbf{T}^{-1} using the well-known Sherman-Morrison-Woodbury formula $^{\mathbf{I}}$. Specifically, we rewrite \mathbf{T} as $\mathbf{T} = \mathbf{X} + \mathbf{U}_{rc} \mathrm{diag}(\varpi \mathbf{I}_2) \mathbf{U}_{cr}^T$ and apply the Sherman-Morrison-Woodbury formula with $\mathbf{A} = \mathbf{X}$, $\mathbf{P} = \mathbf{U}^{rc}$, $\mathbf{C} = \mathrm{diag}(\varpi \mathbf{I}_2)$ and $\mathbf{Q} = \mathbf{U}_{cr}^T$, leading to the following equation: $\mathbf{T}^{-1} = \mathbf{X}^{-1} - \mathbf{X}^{-1} \mathbf{U}^{rc}((\mathrm{diag}(\varpi \mathbf{I}_2))^{-1} + \mathbf{U}_{cr}^T \mathbf{X}^{-1} \mathbf{U}^{rc})^{-1} \mathbf{U}_{cr}^T \mathbf{X}^{-1}$, where $\mathrm{diag}(\mathbf{z})$ is a diagonal matrix with \mathbf{z} as the main diagonal entries and $\mathbf{I}_d \in \mathbb{R}^{d \times d}$ is an identity matrix. Finally, we obtain the following results: $\mathbf{T}^{-1} = \mathbf{X}^{-1} - \frac{\varpi}{1 - \delta \varpi^2 + 2 \mathbf{Y}_{rc} \varpi} \cdot \mathbf{E}_{rc} \mathbf{W} \mathbf{E}_{cr}^T$, where $\delta \triangleq \mathbf{Y}_{rr} \mathbf{Y}_{cc} - \mathbf{Y}_{rc}^2$, $\mathbf{E}_{rc} \triangleq [\mathbf{Y}_{:r} \mathbf{Y}_{:c}] \in \mathbb{R}^{n \times 2}$, $\mathbf{E}_{cr} \triangleq [\mathbf{Y}_{:c} \mathbf{Y}_{:r}] \in \mathbb{R}^{n \times 2}$, $\mathbf{W} \triangleq \begin{pmatrix} 1 + \varpi \mathbf{Y}_{rc} & -\varpi \mathbf{Y}_{cc} \\ -\varpi \mathbf{Y}_{rr} & 1 + \varpi \mathbf{Y}_{rc} \end{pmatrix}$, and $\mathbf{Y}_{:r} \in \mathbb{R}^n$ denotes the r-th column of \mathbf{Y} .

Remarks: (i) Algorithm 1 can be viewed as an improved version of classical greedy pursuit method for solving the sparsity-constrained inverse covariance selection problem. Given the fact that greedy pursuit methods achieve state-of-the-art performance in varieties of non-convex optimization problems (e.g. compressed sensing [28], kernel learning [14], and sensor selection [15]), our proposed method is expected to achieve state-of-the-art performance as well. (ii) Algorithm 1 is also closely related to forward-backward greedy method in the literature [37]. To obtain the greatest descent, while forward-backward strategy considers the removal step and adding step sequentially, the swapping strategy (refer to the swap coordinates stage in Algorithm 1) considers these two steps simultaneously. Thus, the swapping strategy is generally stronger than the forward-backward strategy.

3 CONVEX OPTIMIZATION OVER SUPPORT SET

After the support set has been determined, one need to solve the reduced convex sub-problem as in (2), which is equivalent to the following convex composite minimization problem:

$$\min_{\mathbf{X} \succ 0} F(\mathbf{X}) \triangleq f(\mathbf{X}) + p(\mathbf{X}),$$
with $p(\mathbf{X}) \triangleq I_{\Omega}(\mathbf{X}), \ \Omega \triangleq \{\mathbf{X} \mid \mathbf{X}_{Z} = \mathbf{0}\},$

where $Z \triangleq \{1,2,...,n^2\} \setminus S$ and I_{Ω} is an indicator function of the convex set Ω with $I_{\Omega}(\mathbf{V}) = \left\{ \begin{smallmatrix} 0, & \mathbf{V} \in \Omega \\ \infty, & \text{otherwise.} \end{smallmatrix} \right\}$. In what follows, we

1.
$$(\mathbf{A} + \mathbf{PCQ})^{-1} = \mathbf{A}^{-1} - \mathbf{A}^{-1}\mathbf{P}(\mathbf{C}^{-1} + \mathbf{Q}\mathbf{A}^{-1}\mathbf{P})^{-1}\mathbf{Q}\mathbf{A}^{-1}$$

present an efficient Newton-Like Optimization Algorithm (NLOA) to tackle this problem. This method has the good merits of greedy descent and fast convergence.

Following [29, 35, 16, 32], we develop a quadratic approximation around any solution \mathbf{X} for the objective function using second-order Taylor expansion:

$$q(\boldsymbol{\Theta}, \mathbf{X}) \triangleq f(\mathbf{X}) + \langle \boldsymbol{\Theta}, g(\mathbf{X}) \rangle + \frac{1}{2} \mathrm{vec}(\boldsymbol{\Theta})^T h(\mathbf{X}) \mathrm{vec}(\boldsymbol{\Theta}),$$

where the first-order and second-order derivatives of the objective function $f(\mathbf{X})$ can be expressed as [12]:

$$g(\mathbf{X}) = \mathbf{\Sigma} - \mathbf{X}^{-1}, \ h(\mathbf{X}) = \mathbf{X}^{-1} \otimes \mathbf{X}^{-1}.$$

Then, we keep the non-smooth function $g(\mathbf{X})$ and build a quadratic approximation for the smooth function $f(\mathbf{X})$ by:

$$d(\mathbf{X}^t) \triangleq \arg\min_{\mathbf{\Delta}} \ q(\mathbf{\Delta}; \mathbf{X}^t) + p(\mathbf{X}^t + \mathbf{\Delta}).$$
 (7)

Once the Newton direction $d(\mathbf{X}^t)$ has been computed, one can employ an Arimijo-rule based step size selection to ensure positive definiteness and sufficient descent of the next iterate. We summarize our Newton-like algorithm in Algorithm 2. Note that the initial point \mathbf{X}^0 has to be a feasible solution and the positive definiteness of all the following iterates \mathbf{X}^t will be guaranteed by the step size selection procedure (see step 7 in Algorithm 3). For notational convenience, we use the shorthand

$$f^t = f(\mathbf{X}^t), \ \mathbf{G}^t = g(\mathbf{X}^t), \ \mathbf{H}^t = h(\mathbf{X}^t), \ \mathbf{D}^t = d(\mathbf{X}^t)$$

to denote the objective value, first-order gradient, hessian matrix and the search direction at the point \mathbf{X}^t , respectively.

Algorithm 2 NLOA: Newton-Like Optimization Algorithm to Solve (2) for Optimization Over Support Set.

- 1: Input: \mathbf{X}^0 such that $\mathbf{X}^0 \succ 0$ and $\mathbf{X}_Z = \mathbf{0}$.
- 2: Output: \mathbf{X}^t
- 3: Initialize t = 0
- 4: for t=1 to T_{out} do
- 5: Solve Problem (7) by Algorithm 3 to obtain $d(\mathbf{X}^t)$.
- 6: Perform step-size search to get α^t such that:
- 7: (1) $\mathbf{X}^{t+1} = \mathbf{X}^t + \alpha^t d(\mathbf{X}^t)$ is positive definite and
- 8: (2) there is sufficient decrease in the objective.
- 9: Increment t by 1
- 10: **end for**

3.1 Computing the Search Direction

This subsection focuses on finding the search direction in (7). With the choice of $\mathbf{X}^0 \succ 0$ and $\mathbf{X}_Z^0 = \mathbf{0}$, (7) boils down to the following optimization problem:

$$\min_{\Delta} \langle \Delta, \mathbf{G}^t \rangle + \frac{1}{2} \text{vec}(\Delta)^T \mathbf{H}^t \text{vec}(\Delta), \ s.t. \ \Delta_Z = \mathbf{0}.$$
 (8)

It appears that (8) is very difficult to solve. First, it involves computing and storing an $n^2 \times n^2$ Hessian matrix \mathbf{H}^t . Second, it is a constrained optimization program with $n \times n$ variables and |Z| equality constraints.

We carefully analyze (8) and consider the following solutions. For the first issue, one can exploit the Kronecker product structure of the Hessian matrix to avoid storing it. Recall that $(\mathbf{A} \otimes \mathbf{B}) \operatorname{vec}(\mathbf{C}) = \operatorname{vec}(\mathbf{BCA}), \forall \mathbf{A}, \mathbf{B}, \mathbf{C} \in \mathbb{R}^{n \times n}$. Given any vector $\operatorname{vec}(\mathbf{D}) \in \mathbb{R}^{n^2 \times 1}$, using the fact that the Hessian

matrix can be computed as $\mathbf{H} = \mathbf{X}^{-1} \otimes \mathbf{X}^{-1}$, the Hessian-vector product can be computed efficiently as: $\mathbf{H}\text{vec}(\mathbf{D}) = (\mathbf{X}^{-1} \otimes \mathbf{X}^{-1})\text{vec}(\mathbf{D}) = \text{vec}(\mathbf{X}^{-1}\mathbf{D}\mathbf{X}^{-1})$, which only involves matrix-matrix computation. For the second issue, (8) is, in fact, a unconstrained quadratic program with $n^2 - |Z|$ variables. In order to deal with the variables indexed by Z, one can explicitly enforce the entries of Z for current solution and its corresponding gradient to $\mathbf{0}$. Therefore, the constraint $\mathbf{\Delta}_Z = \mathbf{0}$ can always be satisfied. Finally, linear conjugate gradient method can be used to solve (8).

We summarize our modified linear conjugate gradient method for computing the search direction in Algorithm 3. The algorithm involves a parameter $T_{\rm in}$ controlling the maximum number of iterations. Empirically, we found that a value of $T_{\rm in}=5$ usually leads to good overall efficiency.

Algorithm 3 A Modified Linear Conjugate Gradient to Find the Newton Direction D as in (8).

```
Input: \mathbf{Y} = (\mathbf{X}^t)^{-1}, and current gradient \mathbf{G} = g(\mathbf{X}^t), Specify the maximum iteration T \in \mathbb{N} Output: Newton direction \mathbf{D} \in \mathbb{R}^{n \times n} \mathbf{D} = \mathbf{0}, \mathbf{R} = -\mathbf{G} - \mathbf{Y}\mathbf{D}\mathbf{Y}, \mathbf{R}_Z = \mathbf{0} \mathbf{P} = \mathbf{R}, r_{old} = \langle \mathbf{R}, \mathbf{R} \rangle for p = 1 to T_{\rm in} do \mathbf{B} = \mathbf{Y}\mathbf{P}\mathbf{Y}, \ \alpha = \frac{r_{old}}{\langle \mathbf{P}, \mathbf{B} \rangle} \mathbf{D} = \mathbf{D} + \alpha \mathbf{P}, \ \mathbf{R} = \mathbf{R} - \alpha \mathbf{B} \mathbf{D}_{\mathbf{Z}} = \mathbf{0}, \ \mathbf{R}_{\mathbf{Z}} = \mathbf{0}, r_{new} = \langle \mathbf{R}, \ \mathbf{R} \rangle \mathbf{P} = \mathbf{R} + \frac{r_{\rm new}}{r_{\rm old}} \mathbf{P}, r_{old} = r_{new} end for
```

3.2 Computing the Step Size

Once the Newton direction \mathbf{D} is computed, we need to find a step size $\alpha \in (0,1]$ in order to ensure the positive definiteness of the next iterated result, i.e. $\mathbf{X} + \alpha \mathbf{D}$, so that a sufficient decrease of the objective function will be resulted. We use Armijo's rule and try step size $\alpha \in \{\eta^0, \eta^1, ...\}$ with a constant decrease rate $0 < \eta < 1$ until we find the smallest $t \in \mathbb{N}$ with $\alpha = \eta^t$ such that $\mathbf{X} + \alpha \mathbf{D}$ is (i) positive definite, and (ii) satisfies the following sufficient decrease condition [29]:

$$f(\mathbf{X}^t + \alpha^t \mathbf{D}^t) \le f(\mathbf{X}^t) + \alpha^t \omega \langle \mathbf{G}^t, \mathbf{D}^t \rangle,$$

where $0<\omega<0.5.$ In our experiments, we set $\eta=0.1$ and $\omega=0.25.$

We verify positive definiteness of the solution when we compute its Cholesky factorization (taking $\frac{1}{3}n^3$ flops). We note that the Cholesky factorization dominates the computational cost in the step-size computations. To reduce the computation cost, we can reuse the Cholesky factor in the previous iteration when evaluating the objective function (that requires the computation of $\log \det(\mathbf{X})$) and the gradient (that requires the computation of \mathbf{X}^{-1}).

4 THEORETICAL ANALYSIS

4.1 Convergence Analysis of Algorithm 1

We present the convergence results for Algorithm 1, which are analogous to the results in [3].

Proposition 1. Let \mathbf{X}^k be the sequence generated by algorithm 1. Algorithm 1 outputs a coordinate-wise minimum point \mathbf{X}^* with

$$f(\mathbf{X}^*) \leq f(\mathbf{P})$$
 for every $\mathbf{P} \in \mathcal{N}$, where $\mathcal{N} \triangleq \{\mathbf{X} \mid \mathbf{X} \succ \mathbf{0}, \|\mathbf{X}^* - \mathbf{X}\|_0 \leq 2\}$.

Proof. Note that it takes finite iterations for any convex optimization algorithm to produce an optimal solution with a given support set. Combining with the monotonicity of Algorithm 1, we conclude that the sequence of function values $f(\mathbf{X}^k)$ are monotonically decreasing and Algorithm 1 stops after a finite number of iterations. We define:

$$\mathcal{N}^{0} = \{ \mathbf{P} \mid \mathbf{P} \in \mathcal{N}, \ \bar{\mathcal{S}}(\mathbf{P}) \subseteq \bar{\mathcal{S}}(\mathbf{X}^{*}) \},$$

$$\mathcal{N}^{1} = \{ \mathbf{P} \mid \mathbf{P} \in \mathcal{N}, \ \bar{\mathcal{S}}(\mathbf{P}) = \bar{\mathcal{S}}(\mathbf{X}^{*}) \cup \{j\} \},$$

$$\mathcal{N}^{2} = \{ \mathbf{P} \mid \mathbf{P} \in \mathcal{N}, \ \bar{\mathcal{S}}(\mathbf{P}) = \bar{\mathcal{S}}(\mathbf{X}^{*})/\{i\} \cup \{j\} \},$$

for all $i \in \bar{\mathcal{S}}(\mathbf{X}^*)$, $j \in \bar{\mathcal{Z}}(\mathbf{X}^*)$. Clearly, we have $\mathcal{N} = \mathcal{N}^0 \cup \mathcal{N}^1 \cup \mathcal{N}^2$. Now we assume that point \mathbf{X}^* is generated by Algorithm 1.

For the case \mathcal{N}^0 , \mathbf{X}^* is a global optimal point generated by the convex optimization subproblem on the given support set. Therefore, $f(\mathbf{X}^*) \leq f(\mathbf{X})$ for any $\mathbf{X} \in \mathcal{N}^0$.

For the case \mathcal{N}^2 , we notice that Algorithm 1 terminates only if after the swap coordinates stage. For any $i \in \bar{\mathcal{S}}(\mathbf{X}^*)$ and $j \in \bar{\mathcal{Z}}(\mathbf{X}^*)$, we have the following inequality:

$$f_{i,j} = \min_{\theta} \{ f(\mathbf{X}^k - \mathbf{X}_i^k \mathbf{E}_i + \theta \mathbf{E}_j) \} \ge f(\mathbf{X}^*).$$

Therefore, we have that $\mathcal{N}^2=\varnothing$, which implicates that we cannot find any swap from support set and non-support set to achieve descent on the objective value. Thus, $f(\mathbf{X}^*) \leq f(\mathbf{X})$ for any $\mathbf{X} \in \mathcal{N}^2$.

For the case \mathcal{N}^1 , Algorithm 1 must perform greedy pursuit stage before entering the swap coordinates stage. The greedy stage terminates only if for any $j \in \bar{\mathcal{Z}}(\mathbf{X}^*)$,

$$f_j = \min_{\theta} \{ f(\mathbf{X}^k + \theta \mathbf{E}_j) \} \ge f(\mathbf{X}^*).$$

It implies that we have selected the element that leads to greatest descent as a new member of non-zero elements when $\|\mathbf{X}\|_0 \leq s$. We conclude that $f(\mathbf{X}^*) \leq f(\mathbf{X})$, for any $\mathbf{X} \in \mathcal{N}^1$.

Therefore, we finish the proof of this lemma.

4.2 Convergence Analysis of Algorithm 2

This subsection provides some convergence analysis for the proposed Newton-like optimization algorithm in Algorithm 2. We denote $\{\mathbf{X}^t\}_{t=0}^{\infty}$ as the sequence generated by the algorithm and \mathbf{X}^* as the global optimal solution set for the convex problem in (2). Throughout this subsection, we make the following assumption.

Assumption 1. The objective function $f(\mathbf{X})$ is strongly convex with the modulus σ and gradient Lipschitz continues with constant L for all \mathbf{X}^t with $t = 0, 1, 2, ..., \infty$.

Remarks: This assumption is mild and equivalent to assuming the solution is bounded since it holds that

$$\sigma \leq \lambda(\mathbf{H}^t) \leq L \Leftrightarrow \sigma \leq \lambda((\mathbf{X}^t)^{-1} \otimes (\mathbf{X}^t)^{-1}) \leq L$$

$$\Leftrightarrow \sqrt{\sigma} \leq \lambda((\mathbf{X}^t)^{-1}) \leq \sqrt{L} \Leftrightarrow 1/\sqrt{L} \leq \lambda(\mathbf{X}^t) \leq 1/\sqrt{\sigma},$$

where $\lambda(\mathbf{X}) \in \mathbb{R}^n$ as the eigenvalues of \mathbf{X} in increasing order with $\lambda_1(\mathbf{X}) \leq \lambda_2(\mathbf{X}) \leq \ldots \leq \lambda_n(\mathbf{X})$.

The following lemma characterizes the optimality of $d(\mathbf{X}^t)$. It is nearly identical to Lemma 1 in [30]. For completeness, we present the proof here.

Lemma 1. It holds that

$$\|\mathbf{D}^t\|_{\mathbf{H}^t}^2 + \langle \mathbf{G}^t, \mathbf{D}^t \rangle \le 0, \ \forall \mathbf{D}^t \ \text{with } \mathbf{X}^t + \mathbf{D}^t \in \Omega.$$
 (9)

Proof. Noticing $\mathbf{D}^t \triangleq d(\mathbf{X}^t)$ is the minimizer of (7), we have:

$$q(\mathbf{D}^t; \mathbf{X}^t) + p(\mathbf{X}^t + \mathbf{D}^t) \le q(\mathbf{Z}; \mathbf{X}^t) + p(\mathbf{X}^t + \mathbf{Z}), \ \forall \mathbf{Z}.$$

Letting $\mathbf{Z} = \alpha \mathbf{D}^t$ where α is any constant with $\alpha \in [0, 1]$, we obtain:

$$\langle \mathbf{G}^{t}, \mathbf{D}^{t} \rangle + \frac{1}{2} \| \mathbf{D}^{t} \|_{\mathbf{H}^{t}}^{2} + p(\mathbf{X}^{t} + \mathbf{D}^{t})$$

$$\leq \langle \mathbf{G}^{t}, \alpha \mathbf{D}^{t} \rangle + \frac{1}{2} \| \alpha \mathbf{D}^{t} \|_{\mathbf{H}^{t}}^{2} + p(\mathbf{X}^{t} + \alpha \mathbf{D}^{t})$$

$$\leq \langle \mathbf{G}^{t}, \alpha \mathbf{D}^{t} \rangle + \frac{1}{2} \| \alpha \mathbf{D}^{t} \|_{\mathbf{H}^{t}}^{2} + \alpha p(\mathbf{X}^{t} + \mathbf{D}^{t}) + (1 - \alpha) p(\mathbf{X}^{t}),$$

where the last inequality uses the convexity of $p(\cdot)$. Rearranging terms yields:

$$(1 - \alpha)(\langle \mathbf{G}^t, \mathbf{D}^t \rangle + p(\mathbf{X}^t + \mathbf{D}^t) - p(\mathbf{X}^t)) \le \frac{\alpha^2 - 1}{2} \|\mathbf{D}^t\|_{\mathbf{H}^t}^2$$
$$\langle \mathbf{G}^t, \mathbf{D}^t \rangle + p(\mathbf{X}^t + \mathbf{D}^t) - p(\mathbf{X}^t) \le -\frac{1 + \alpha}{2} \|\mathbf{D}^t\|_{\mathbf{H}^t}^2.$$

Since $\mathbf{X}^t \in \Omega$, $\mathbf{X}^t + \mathbf{D}^t \in \Omega$, we have $p(\mathbf{X}^t + \mathbf{D}^t) = p(\mathbf{X}^t) = 0$. Letting $\alpha = 1$, we obtain (9).

Theorem 1. (Global Convergence). We have the following results: (i) There exists a strictly positive constant $\forall t, \alpha^t \leq \min(1, 1/(\sqrt{L}C_1) - \epsilon, C_2)$ such that the positive definiteness and sufficient descent conditions (refer to step 7-8 of Algorithm 2) are satisfied. Here ϵ denotes a sufficient small positive constant. $C_1 \triangleq \lambda_n(\Sigma)/\sigma + \sigma^{-3/2}$ and $C_2 \triangleq 2\sigma(1-\omega)/L$ are some constants which are independent of the current solution \mathbf{X}^t . (ii) The sequence $f(\mathbf{X}^t)$ is non-increasing and any cluster point of the sequence \mathbf{X}^t is the global optimal solution of (6).

Proof. (i) First, we focus on the positive definiteness condition. We now bound $\lambda_n(\mathbf{D}^t)$. By Lemma 1, we have $\forall \mathbf{D}^t$ with $\mathbf{X}^t + \mathbf{D}^t \in \Omega$:

$$0 \geq \langle \mathbf{D}^{t}, \mathbf{G}^{t} \rangle + \|\mathbf{D}^{t}\|_{\mathbf{H}^{t}}^{2}$$

$$\geq -\lambda_{n}(\mathbf{D}^{t})\lambda_{n}(\mathbf{G}^{t}) + \sigma\|\mathbf{D}^{t}\|_{F}^{2}$$

$$= -\lambda_{n}(\mathbf{D}^{t})\lambda_{n}(\mathbf{\Sigma} - (\mathbf{X}^{t})^{-1}) + \sigma\|\mathbf{D}^{t}\|_{F}^{2}$$

$$\geq -\lambda_{n}(\mathbf{D}^{t}) \cdot (\lambda_{n}(\mathbf{\Sigma}) + 1/\sqrt{\sigma}) + \sigma(\lambda_{n}(\mathbf{D}^{t}))^{2}, (10)$$

where the second step uses the fact that $\lambda(\mathbf{H}^t) \geq \sigma$ and the inequality that $\langle \mathbf{A}, \mathbf{B} \rangle \geq -\lambda_n(\mathbf{A})\lambda_n(\mathbf{B})$, $\forall \mathbf{A}, \mathbf{B}$; the third step uses the definition of $\mathbf{G}^t = \mathbf{\Sigma} - (\mathbf{X}^t)^{-1}$; the last step uses the inequalities that $\lambda_n(\mathbf{\Sigma} - (\mathbf{X}^t)^{-1}) \leq \lambda_n(\mathbf{\Sigma}) + \lambda_n((\mathbf{X}^t)^{-1}) \leq \lambda_n(\mathbf{\Sigma}) + 1/\sqrt{\sigma}$. Solving the quadratic inequality in (10) gives $\lambda_n(\mathbf{D}) \leq (\lambda_n(\mathbf{\Sigma}) + 1/\sqrt{\sigma})/\sigma = \lambda_n(\mathbf{\Sigma})/\sigma + \sigma^{-3/2} \triangleq C_1$. Therefore, we have:

$$0 \prec (1/\sqrt{L} - C_1 \alpha^t) \mathbf{I} \preceq \mathbf{X}^t - \alpha^t \lambda_n(\mathbf{D}^t) \mathbf{I} \preceq \mathbf{X}^t + \alpha^t \mathbf{D}^t.$$

where the first steps uses the fact that $\alpha^t \leq 1/(\sqrt{L}C_1) - \epsilon$; the second step uses $\mathbf{X}^t \succeq (1/\sqrt{L}) \cdot \mathbf{I}$ and $\lambda_n(\mathbf{D}) \leq C_1$; the last step uses $\lambda_n(\mathbf{D}^t)\mathbf{I} \succeq -\mathbf{D}^t$.

Second, we focus on the sufficient decrease condition. For any $\alpha \in (0,1],$ we have:

$$f(\mathbf{X}^{t} + \alpha^{t}\mathbf{D}^{t}) - f(\mathbf{X}^{t})$$

$$\leq \alpha \langle \mathbf{D}^{t}, \mathbf{G}^{t} \rangle + \frac{(\alpha^{t})^{2}L}{2} \| \mathbf{D}^{t} \|_{F}^{2}$$

$$\leq \alpha^{t} \left(\langle \mathbf{D}^{t}, \mathbf{G}^{t} \rangle + \frac{\alpha^{t}L}{2\sigma} \| \mathbf{D} \|_{\mathbf{H}^{t}}^{2} \right)$$

$$\leq \alpha^{t} \left(\langle \mathbf{D}^{t}, \mathbf{G}^{t} \rangle - \frac{\alpha^{t}L}{2\sigma} \langle \mathbf{D}^{t}, \mathbf{G}^{t} \rangle \right)$$

$$= \alpha^{t} \langle \mathbf{D}^{t}, \mathbf{G}^{t} \rangle (1 - \frac{\alpha^{t}L}{2\sigma}) \leq \alpha^{t} \langle \mathbf{D}^{t}, \mathbf{G}^{t} \rangle \cdot \omega, \quad (11)$$

where the first step uses the L-Lipschitz continuity of the gradient of $f(\mathbf{X})$ that: $\forall \mathbf{X}, \mathbf{Y} \in \Omega$, $f(\mathbf{Y}) \leq f(\mathbf{X}) + \langle g(\mathbf{X}), \mathbf{Y} - \mathbf{X} \rangle + \frac{L}{2} \|\mathbf{X} - \mathbf{Y}\|_{\mathrm{F}}^2$; the second step uses the lower bound of the Hessian matrix that $\sigma \|\mathbf{D}^t\|_{\mathrm{F}}^2 \leq \|\mathbf{D}^t\|_{\mathbf{H}^t}^2$; the third step uses (9) that $\|\mathbf{D}^t\|_{\mathbf{H}^t}^2 \leq -\langle \mathbf{D}^t, \mathbf{G}^t \rangle$; the last step uses the choice that $\alpha^t \leq 2\sigma(1-\omega)/L \triangleq C_2$.

Combining the positive definiteness condition, sufficient decrease condition and the fact that $\alpha \in (0,1]$, we finish the proof of the first part of this lemma.

(ii) From (11) and (9), we have:

$$\forall t, \ f(\mathbf{X}^{t+1}) - f(\mathbf{X}^{t})$$

$$\leq \alpha^{t} \omega \langle \mathbf{D}^{t}, \mathbf{G}^{t} \rangle \leq -\alpha^{t} \omega \|\mathbf{D}^{t}\|_{\mathbf{H}^{t}}^{2}$$

$$\leq -\sigma \alpha^{t} \omega \|\mathbf{D}^{t}\|_{F}^{2} = -\nu \|\mathbf{D}^{t}\|_{F}^{2}, \text{ with } \nu \triangleq \sigma \alpha^{t} \omega > 0. \quad (12)$$

Therefore, the sequence $f(\mathbf{X}^t)$ is non-increasing. Summing the inequality in (12) over i=0,1,...,t-1 and using the fact that $f(\mathbf{X}^*) \leq f(\mathbf{X}^t)$, we have:

$$\begin{aligned} f(\mathbf{X}^t) - f(\mathbf{X}^0) &\leq -\nu \sum_{i=0}^{t-1} \|\mathbf{D}^i\|_{\mathrm{F}}^2 \\ \Rightarrow & f(\mathbf{X}^*) - f(\mathbf{X}^0) \leq -\nu \sum_{i=0}^{t-1} \|\mathbf{D}^i\|_{\mathrm{F}}^2 \\ \Rightarrow & (f(\mathbf{X}^0) - f(\mathbf{X}^*))/(t\nu) \geq \min_{i=0,1,\dots,t-1} \|\mathbf{D}^i\|_{\mathrm{F}}^2. \end{aligned}$$

As $t \to \infty$, we have $\mathbf{D}^t \to 0$. We further derive the following results: $\mathbf{D}^t = \mathbf{0} \Rightarrow (\nabla q(\mathbf{D}^t))_S = \mathbf{0} \Rightarrow (\mathbf{H}^t \circ \mathbf{D}^t + \mathbf{G}^t)_S = \mathbf{0} \Rightarrow \mathbf{G}_S^t = -(\mathbf{H}^t \circ \mathbf{\Delta})_S = \mathbf{0}$. Based on the fact that $\mathbf{X}^t \succ 0$, $\mathbf{X}_Z^t = \mathbf{0}$, and $\mathbf{G}_S^t = 0$, we conclude that \mathbf{X}^t is the global optimal solution for the convex optimization problem. Therefore, any cluster point of the sequence \mathbf{X}^t is the global optimal solution.

Remarks: Due to the strongly convexity and gradient Lipschitz continuity of the objective function, there always exists a strictly positive step size α^t such that both sufficient decrease condition and positive definite condition can be satisfied for all t. This is very crucial for the global convergence of Algorithm 2.

We now prove the global linear convergence rate of Algorithm 2. The following lemma characterizes the relation between $d(\mathbf{X}^t)$ and \mathbf{X}^* for any \mathbf{X}^t .

Lemma 2. If \mathbf{X}^t is not the global optimal solution of (6), there exists a constant $\eta \in (0, \infty)$ such that $\|\mathbf{X}^t - \mathbf{X}^*\|_F \le \eta \|d(\mathbf{X}^t)\|_F$.

Proof. First, we prove that $d(\mathbf{X}^t) \neq \mathbf{0}$. This can be achieved by contradiction. Assuming that $d(\mathbf{X}^t) = \mathbf{0}$, we obtain: $\mathbf{X}^{t+1} = \mathbf{X}^t + \alpha^t d(\mathbf{X}^t)$, which implies that $\mathbf{X}^{t+1} = \mathbf{X}^t$ is the stationary point. Since $(\mathbf{6})$ is a strongly convex optimization problem, we have $\mathbf{X}^t = \mathbf{X}^*$, which contradicts with the condition that \mathbf{X}^t is not the optimal solution. Combining with the boundedness of \mathbf{X}^t and \mathbf{X}^* , we conclude that there exists a sufficiently large constant $\eta \in (0,\infty)$ such that $\|\mathbf{X}^t - \mathbf{X}^*\|_{\mathsf{F}} \leq \eta \|d(\mathbf{X}^t)\|_{\mathsf{F}}$.

Theorem 2. (Global Linear Convergence Rate). The sequence \mathbf{X}^t converges to the global optimal solution linearly.

Proof. We assume that there exists a point $\bar{\mathbf{X}}$ lying on the segment joining \mathbf{X}^{t+1} with \mathbf{X}^* . Using the fact that $\mathbf{X}^{t+1} \in \Omega$, $\mathbf{X}^* \in \Omega$ and the Mean Value Theorem, we derive the following results:

$$F(\mathbf{X}^{t+1}) - F(\mathbf{X}^*) = f(\mathbf{X}^{t+1}) - f(\mathbf{X}^*)$$

$$= \langle g(\bar{\mathbf{X}}), \mathbf{X}^{t+1} - \mathbf{X}^* \rangle$$

$$= \langle \mathbf{G}^t + \mathbf{H}^t \circ \mathbf{D}^t, \mathbf{X}^{t+1} - \mathbf{X}^* \rangle$$

$$+ \langle g(\bar{\mathbf{X}}) - \mathbf{G}^t - \mathbf{H}^t \circ \mathbf{D}^t, \mathbf{X}^{t+1} - \mathbf{X}^* \rangle. \tag{13}$$

(i) We now consider to bound the first term in (13). Since \mathbf{D}^t is the optimal solution of (7), we have

$$\mathbf{0} \in \mathbf{G}^t + \mathbf{H}^t \circ \mathbf{D}^t + \partial p(\mathbf{X}^t + \mathbf{D}^t). \tag{14}$$

Due to the convexity of $p(\cdot)$, we obtain: $\langle \mathbf{X} - \mathbf{Z}, \partial p(\mathbf{Z}) \rangle \leq p(\mathbf{X}) - p(\mathbf{Z})$. Letting $\mathbf{X} = \mathbf{X}^*$ and $\mathbf{Z} = \mathbf{X}^t + \mathbf{D}^t$, we have:

$$\langle \mathbf{X}^* - \mathbf{X}^t - \mathbf{D}^t, \ \partial p(\mathbf{X}^t + \mathbf{D}^t) \rangle \le p(\mathbf{X}^*) - p(\mathbf{X}^t + \mathbf{D}^t).$$
 (15)

Combing (14) and (15), we have: $\langle \mathbf{X}^* - \mathbf{X}^t - \mathbf{D}^t, -\mathbf{G}^t - \mathbf{H}^t \circ \mathbf{D}^t \rangle \leq p(\mathbf{X}^*) - p(\mathbf{X}^t + \mathbf{D}^t)$. Noticing $\mathbf{X}^* \in \Omega$ and $\mathbf{X}^t + \mathbf{D}^t \in \Omega$, we have $p(\mathbf{X}^*) = p(\mathbf{X}^t + \mathbf{D}^t) = 0$. We reach the following inequality:

$$\langle \mathbf{X}^t + \mathbf{D}^t - \mathbf{X}^*, \ \mathbf{G}^t + \mathbf{H}^t \circ \mathbf{D}^t \rangle \le 0.$$

Therefore, the first term in (13) is bounded by:

$$\langle \mathbf{G}^{t} + \mathbf{H}^{t} \circ \mathbf{D}^{t}, \mathbf{X}^{t+1} - \mathbf{X}^{*} \rangle$$

$$= \langle \mathbf{G}^{t} + \mathbf{H}^{t} \circ \mathbf{D}^{t}, (\alpha^{t} - 1)\mathbf{D}^{t} + \mathbf{X}^{t} + \mathbf{D}^{t} - \mathbf{X}^{*} \rangle$$

$$\leq \langle \mathbf{G}^{t} + \mathbf{H}^{t} \circ \mathbf{D}^{t}, (\alpha^{t} - 1)\mathbf{D}^{t} \rangle$$

$$\leq (\alpha^{t} - 1)\langle \mathbf{G}^{t}, \mathbf{D}^{t} \rangle + 0$$

$$\leq \frac{1 - \alpha^{t}}{\alpha^{t} \omega} (f(\mathbf{X}^{t}) - f(\mathbf{X}^{t+1})). \tag{16}$$

(ii) We now consider to bound the second term in (13). We derive the following results:

$$\langle g(\bar{\mathbf{X}}) - \mathbf{G}^{t} - \mathbf{H}^{t} \circ \mathbf{D}^{t}, \mathbf{X}^{t+1} - \mathbf{X}^{*} \rangle$$

$$\leq \|\mathbf{X}^{t+1} - \mathbf{X}^{*}\|_{F} \|(\|g(\bar{\mathbf{X}}) - \mathbf{G}^{t}\|_{F} + \|\mathbf{H}^{t} \circ \mathbf{D}^{t}\|_{F})$$

$$\leq \|\mathbf{X}^{t+1} - \mathbf{X}^{*}\|_{F} (L\|\bar{\mathbf{X}} - \mathbf{X}^{t}\|_{F} + \|\mathbf{H}^{t} \circ \mathbf{D}^{t}\|_{F})$$

$$\leq (\eta + \alpha^{t}) \|\mathbf{D}^{t}\|_{F} (L\|\bar{\mathbf{X}} - \mathbf{X}^{t}\|_{F} + \|\mathbf{H}^{t} \circ \mathbf{D}^{t}\|_{F})$$

$$\leq L(\eta + \alpha^{t}) \|\mathbf{D}^{t}\|_{F} (\|\bar{\mathbf{X}} - \mathbf{X}^{t}\|_{F} + \|\mathbf{D}^{t}\|_{F})$$

$$\leq L(\eta + \alpha^{t}) (\|\bar{\mathbf{X}} - \mathbf{X}^{*}\|_{F} + \|\mathbf{X}^{*} - \mathbf{X}^{t}\|_{F} + \|\mathbf{D}^{t}\|_{F}) \|\mathbf{D}^{t}\|_{F}$$

$$\leq L(\eta + \alpha^{t}) (\|\mathbf{X}^{t+1} - \mathbf{X}^{*}\|_{F} + \eta \|\mathbf{D}^{t}\|_{F} + \|\mathbf{D}^{t}\|_{F}) \|\mathbf{D}^{t}\|_{F}$$

$$\leq L(\eta + \alpha^{t}) ((\eta + \alpha^{t}) \|\mathbf{D}^{t}\|_{F} + \eta \|\mathbf{D}^{t}\|_{F} + \|\mathbf{D}^{t}\|_{F}) \|\mathbf{D}^{t}\|_{F}$$

$$= L(\eta + \alpha^{t}) (2\eta + \alpha^{t} + 1) \|\mathbf{D}^{t}\|_{F}^{2}$$

$$\leq (L(\eta + \alpha^{t}) (2\eta + \alpha^{t} + 1) / \nu) (f(\mathbf{X}^{t}) - f(\mathbf{X}^{t+1}))$$

$$\leq (L(\eta + \alpha^{t}) (2\eta + \alpha^{t} + 1) / \sigma \alpha^{t} \omega) (f(\mathbf{X}^{t}) - f(\mathbf{X}^{t+1})),$$
(17

where the first step uses Cauchy-Schwarz inequality that $\forall \mathbf{A}, \ \mathbf{B}, \ \langle \mathbf{A}, \mathbf{B} \rangle \leq \|\mathbf{A}\|_F \|\mathbf{B}\|_F$ and the triangle inequality that $\forall \mathbf{A}, \ \mathbf{B}, \ \|\mathbf{A} + \mathbf{B}\|_F \leq \|\mathbf{A}\|_F + \|\mathbf{B}\|_F$; the second step uses the gradient Lipschitz continouity of $f(\cdot)$ that $\|g(\bar{\mathbf{X}}) - g(\mathbf{X}^t)\|_F \leq L \|\bar{\mathbf{X}} - \mathbf{X}^t\|_F$; the third step uses Lemma 2 that $\|\mathbf{X}^* - \mathbf{X}^t\|_F \leq \eta \|\mathbf{D}^t\|_F$ and $\|\mathbf{X}^{t+1} - \mathbf{X}^*\|_F = \|\mathbf{X}^t + \alpha^t\mathbf{D}^t - \mathbf{X}^*\|_F \leq \|\mathbf{X}^t - \mathbf{X}^*\|_F + \alpha^t\|\mathbf{D}^t\|_F \leq (\eta + \alpha^t)\|\mathbf{D}^t\|_F$; the fourth step uses the inequality that $\|\mathbf{H}^t \circ \mathbf{D}^t\|_F \leq L \|\mathbf{D}^t\|_F$; the fifth step uses the triangle inequality that $\|\bar{\mathbf{X}} - \mathbf{X}^t\|_F = \|\bar{\mathbf{X}} - \mathbf{X}^* + \mathbf{X}^* - \mathbf{X}^t\|_F \leq \|\bar{\mathbf{X}} - \mathbf{X}^*\|_F + \|\mathbf{X}^* - \mathbf{X}^t\|_F$; the sixth step uses the fact that $\|\bar{\mathbf{X}} - \mathbf{X}^t\|_F \leq \|\mathbf{X}^{t+1} - \mathbf{X}^*\|_F$ since $\bar{\mathbf{X}}$ is a point lying on the segment joining \mathbf{X}^{t+1} and \mathbf{X}^* ; the ninth step uses (12); the last step uses the definition of μ .

step uses the definition of μ .

We define $C^t \triangleq \frac{1-\alpha^t}{\alpha^t\omega} + \frac{L(\eta+\alpha^t)(2\eta+\alpha^t+1)}{\sigma\alpha^t\omega}$. Clearly, $C^t \leq \frac{1}{\alpha^t\omega} + \frac{2L(\eta+1)^2/\sigma}{\alpha^t\omega} \leq \frac{1+2L(\eta+1)^2/\sigma}{\min(\alpha^1,\alpha^2,\dots,\alpha^\infty)\omega} \triangleq C$. Combining (16), (17) and (13), we have the following results: $f(\mathbf{X}^{t+1}) - f(\mathbf{X}^*) \leq C(f(\mathbf{X}^t) - f(\mathbf{X}^{t+1})) = C(f(\mathbf{X}^t) - f(\mathbf{X}^*)) - C(f(\mathbf{X}^{t+1}) - f(\mathbf{X}^*))$. Finally, we obtain:

$$\frac{f(\mathbf{X}^{t+1}) - f(\mathbf{X}^*)}{f(\mathbf{X}^t) - f(\mathbf{X}^*)} \le \frac{C}{C+1}.$$

Therefore, $f(\mathbf{X}^t)$ converges to $f(\mathbf{X}^*)$ at least Q-linearly. In addition, from (12), we have: $\|\mathbf{X}^{t+1} - \mathbf{X}^t\|_{\mathrm{F}}^2 \leq (f(\mathbf{X}^t) - f(\mathbf{X}^{t+1}))/\nu$. Since $f(\mathbf{X}^{t+1}) - f(\mathbf{X}^*)$ converges to 0 at least R-linearly, this implies that $\|\mathbf{X}^{t+1} - \mathbf{X}^t\|_{\mathrm{F}}^2$ converges at least R-linearly as well. We thus complete the proof of this lemma.

Remarks: Global linear convergence rate for convex composite optimization has been extensively studied in Ref. [29]. This work extends their analysis to deal with our specific matrix optimization problem which involves an additional positive definite constraint X > 0.

Theorem 3. (Local Quadratic Convergence Rate). A full Newton step size with $\alpha^t = 1$ will be selected when \mathbf{X}^t is close enough to global optimal solution that $\|\mathbf{D}^t\|_F \leq \min(1/\sqrt{L} - \epsilon, 0.81/\sigma)$. Here ϵ denotes a sufficient small positive constant. In addition, when $\|\mathbf{X}^t - \mathbf{X}^*\|_F \leq 1/\sigma$, the sequence $\{\mathbf{X}^t\}$ converges to the global optimal solution quadratically.

Proof. (i) First, we consider the positive definiteness condition (refer to step 7 in Algorithm 2). We derive the following results:

$$\mathbf{0} \quad \prec \quad (\frac{1}{\sqrt{L}} - (\frac{1}{\sqrt{L}} - \epsilon))\mathbf{I}$$

$$\preceq \quad \mathbf{X}^t - \|\mathbf{D}^t\|_{\mathsf{F}} \cdot \mathbf{I} \preceq \mathbf{X}^t - \lambda_n(\mathbf{D}^t)\mathbf{I} \preceq \mathbf{X}^t + \mathbf{D}^t,$$

where the first step uses $\frac{1}{\sqrt{L}} - (\frac{1}{\sqrt{L}} - \epsilon) > 0$; the second step uses $\mathbf{X}^t \succeq \frac{1}{\sqrt{L}}\mathbf{I}$ and $\|\mathbf{D}^t\|_{\mathrm{F}} \leq 1/\sqrt{L} - \epsilon$; the third step uses the inequality $\lambda_n(\mathbf{D}^t) \leq \|\mathbf{D}^t\|_{\mathrm{F}}$; the last step uses $-\lambda_n(\mathbf{D}^t)\mathbf{I} \preceq \mathbf{D}^t$ and $\alpha^t = 1$. Therefore, the positive definite condition is satisfied for $\alpha^t = 1$ when

$$\|\mathbf{D}^t\|_{\mathsf{F}} \le 1/\sqrt{L} - \epsilon. \tag{18}$$

Second, we consider the sufficient decrease condition (refer to step 8 in Algorithm 2). Since $f(\mathbf{X})$ is a standard self-concordant function, it holds that (refer to Theorem 4.1.8 in [22]):

$$f(\mathbf{Y}) - f(\mathbf{X}) - \langle g(\mathbf{X}), \mathbf{Y} - \mathbf{X} \rangle \le \varphi(\|\mathbf{Y} - \mathbf{X}\|_{h(\mathbf{X})})$$

with $\varphi(t) \triangleq -t - \ln(1-t)$. Applying this inequality with $\mathbf{Y} = \mathbf{X}^{t+1}$, $\mathbf{X} = \mathbf{X}^t$ and using the update rule that $\mathbf{X}^{t+1} = \mathbf{X}^t + \alpha^t \mathbf{D}^t$, we have the following inequalities:

$$\begin{split} &f(\mathbf{X}^{t+1})\\ &\leq f(\mathbf{X}^t) + \alpha^t \langle \mathbf{G}^t, \mathbf{D}^t \rangle + \varphi(\alpha^t \| \mathbf{D}^t \|_{H^t})\\ &\leq f(\mathbf{X}^t) + \alpha^t \langle \mathbf{G}^t, \mathbf{D}^t \rangle + \frac{(\alpha^t)^2}{2} \| \mathbf{D}^t \|_{\mathbf{H}^t}^2 + (\alpha^t)^3 \| \mathbf{D}^t \|_{\mathbf{H}^t}^3\\ &\leq f(\mathbf{X}^t) + \langle \mathbf{G}^t, \mathbf{D}^t \rangle + \frac{1}{2} \| \mathbf{D}^t \|_{\mathbf{H}^t}^2 + \| \mathbf{D}^t \|_{\mathbf{H}^t}^3\\ &\leq f(\mathbf{X}^t) + \langle \mathbf{G}^t, \mathbf{D}^t \rangle - \frac{1}{2} \langle \mathbf{G}^t, \mathbf{D}^t \rangle + (-\langle \mathbf{G}^t, \mathbf{D}^t \rangle)^{3/2}\\ &= f(\mathbf{X}^t) + \langle \mathbf{G}^t, \mathbf{D}^t \rangle (\frac{1}{2} - \sqrt{-\langle \mathbf{G}^t, \mathbf{D}^t \rangle})\\ &\leq f(\mathbf{X}^t) + \omega \langle \mathbf{G}^t, \mathbf{D}^t \rangle \frac{1}{\omega}, \end{split}$$

where the second step uses the fact that $-z - \ln(1-z) \le \frac{1}{2}z^2 + z^3$ for $0 \le z \le 0.81$ (see Section 9.6 in [4]) and the inequality that $z \triangleq \alpha^t \|\mathbf{D}^t\|_{\mathbf{H}^t} \le 0.81$, where the latter is true since

$$\|\mathbf{D}^t\|_{\mathsf{F}} \le 0.81/\sigma \Rightarrow \|\mathbf{D}^t\|_{\mathbf{H}^t} \le 0.81 \tag{19}$$

and $\alpha^t \leq 1$; the third step uses the choice $\alpha^t = 1$; the fourth step uses (9); the last step uses the inequality that $\frac{1}{2} - \sqrt{-\langle \mathbf{G}^t, \ \mathbf{D}^t \rangle} \leq \frac{1}{\omega}$, which is clearly holds since $\omega < 1/2$.

Combining (18) and (19), we conclude that the full Newton step size will be achieved when $\|\mathbf{D}^t\|_{\mathrm{F}} \leq \min(1/\sqrt{L} - \epsilon, 0.81/\sigma)$.

(ii) We now prove the second part of this theorem. Recall that when $f(\cdot)$ is a standard self-concordant function, it holds that (refer to Lemma 1 in [23]):

$$||g(\mathbf{Y}) - g(\mathbf{X}) - h(\mathbf{X})(\mathbf{Y} - \mathbf{X})||_{h(\mathbf{X})} \le \frac{r^2}{1 - r}$$
 (20)

for all $r \triangleq \|\mathbf{X} - \mathbf{Y}\|_{h(\mathbf{X})} < 1$. We define the generalized proximal operator as:

$$\operatorname{prox}_p^{\mathbf{N}}(\mathbf{X}) \triangleq \arg\min_{\mathbf{Y}} \ \tfrac{1}{2} \|\mathbf{Y} - \mathbf{X}\|_{\mathbf{N}}^2 + p(\mathbf{Y}).$$

Thus, \mathbf{D}^t can be represented as:

$$\begin{split} \mathbf{D}^t &= & \arg\min_{\mathbf{\Delta}} \ \langle \mathbf{\Delta}, \mathbf{G}^t \rangle + \frac{1}{2} \|\mathbf{\Delta}\|_{\mathbf{H}^t}^2 + p(\mathbf{X}^t + \mathbf{\Delta}) \\ &= & \arg\min_{\mathbf{\Delta}} \ \frac{1}{2} \|\mathbf{\Delta} + (\mathbf{H}^t)^{-1} \circ \mathbf{G}^t\|_{\mathbf{H}^t}^2 + p(\mathbf{X}^t + \mathbf{\Delta}) \\ &= & \operatorname{prox}_p^{\mathbf{H}^t} (\mathbf{X}^t - (\mathbf{H}^t)^{-1} \circ \mathbf{G}^t) - \mathbf{X}^t. \end{split}$$

We derive the following inequalities:

$$\begin{split} &\|\mathbf{X}^{t+1} - \mathbf{X}^*\|_{\mathbf{H}^t} \\ &= \|\mathbf{X}^t + \mathbf{D}^t - \mathbf{X}^*\|_{\mathbf{H}^t} \\ &= \|\operatorname{prox}_p^{\mathbf{H}^t}(\mathbf{X}^t - (\mathbf{H}^t)^{-1} \circ \mathbf{G}^t) - \mathbf{X}^*\|_{\mathbf{H}^t} \\ &= \|\operatorname{prox}_p^{\mathbf{H}^t}(\mathbf{X}^t - (\mathbf{H}^t)^{-1} \circ \mathbf{G}^t) - \mathbf{X}^*\|_{\mathbf{H}^t} \\ &= \|\operatorname{prox}_p^{\mathbf{H}^t}(\mathbf{X}^* - (\mathbf{H}^t)^{-1} \circ g(\mathbf{X}^*))\|_{\mathbf{H}^t} \\ &\leq \|\mathbf{X}^t - \mathbf{X}^* + (\mathbf{H}^t)^{-1} \circ (\mathbf{G}^* - \mathbf{G}^t)\|_{\mathbf{H}^t} \\ &= \|(\mathbf{H}^t)^{-1} \circ (\mathbf{H}^t \circ (\mathbf{X}^t - \mathbf{X}^* + (\mathbf{H}^t)^{-1}(\mathbf{G}^* - \mathbf{G}^t)))\|_{\mathbf{H}^t} \\ &= \|(\mathbf{H}^t)^{-1} \circ (\mathbf{H}^t \circ (\mathbf{X}^t - \mathbf{X}^*) + (\mathbf{G}^* - \mathbf{G}^t))\|_{\mathbf{H}^t} \\ &\leq \|(\mathbf{H}^t)^{-1} \circ \mathbf{I}\|_{\mathbf{H}^t} \cdot \|(\mathbf{H}^t \circ (\mathbf{X}^t - \mathbf{X}^*) + (\mathbf{G}^* - \mathbf{G}^t))\|_{\mathbf{H}^t} \\ &\leq \frac{1}{\sigma} \cdot \|\mathbf{H}^t(\mathbf{X}^t - \mathbf{X}^*) - \mathbf{G}^t + \mathbf{G}^*\|_{\mathbf{H}^t} \\ &\leq \frac{1}{\sigma} \|\mathbf{H}^t(\mathbf{X}^t - \mathbf{X}^*)\|_{\mathbf{H}^t}, \end{split}$$

where the third step uses the fact that $\operatorname{prox}_p^{\mathbf{H}^t}(\mathbf{X}^* - (\mathbf{H}^t)^{-1} \circ g(\mathbf{X}^*))|_{\mathbf{H}^t} = \mathbf{X}^*$; the fourth step uses the fact that the generalized proximal mappings are firmly non-expansive in the generalized vector norm; the seventh step uses the Cauchy-Schwarz inequality; the eighth step uses the fact that $\|(\mathbf{H}^t)^{-1} \circ \mathbf{I}\|_{\mathbf{H}^t} \leq \frac{1}{\sigma}$ with \mathbf{I} being an identity matrix of dimension n; the last step uses (20).

In particular, when $\|\mathbf{X} - \mathbf{Y}\|_{h(\mathbf{X})} < 1 \iff \|\mathbf{X} - \mathbf{Y}\|_{F} < 1/\sigma$, we have:

$$\begin{split} & \|\mathbf{X}^{t+1} - \mathbf{X}^*\|_{\mathbf{H}^t} \leq \frac{1}{\sigma} \|\mathbf{X}^t - \mathbf{X}^*\|_{\mathbf{H}^t}^2 \\ \Rightarrow & \|\mathbf{X}^{t+1} - \mathbf{X}^*\|_{\mathbf{F}} \leq \frac{L}{\sigma^2} \|\mathbf{X}^t - \mathbf{X}^*\|_{\mathbf{F}}^2. \end{split}$$

In other words, NLOA converges to the global optimal solution \mathbf{X}^* with asymptotic quadratic convergence rate.

Remarks: We are also aware of the work of [12], which shows the quadratic convergence rate for their Newton-like method. However, their work focus on a different convex ℓ_1 norm regularized sparse inverse covariance selection problem and their results are not applicable to our problem. In addition, their work is based on the assumption that the objective is Hessian Lipschitz continuous, while ours is based on the self-concordant analysis [22, 4] of the objective function.

5 EXPERIMENTS

This section demonstrates the performance of the proposed Coordinated-Wise Optimization Algorithm (CWOA) on synthetic and real-world data sets. All codes are implemented in MATLAB on an Intel 3.20GHz CPU with 8 GB RAM. Some Matlab code can be found in the authors' research webpages.

- Data sets. Three types of data sets are considered in our experiments. (i) Gaussian random data sets. We generate a data matrix $\mathbf{Z} \in \mathbb{R}^{m \times n}$ sampled from a standard normal distribution. The parameter m is fixed to 500. The covariance matrix Σ is computed by $\mathbf{\Sigma} = 1/(m-1)\sum_{i=1}^{m}(\mathbf{z_i} - \boldsymbol{\mu})(\mathbf{z_i} - \boldsymbol{\mu})^T$, where $\mathbf{z_i}$ denotes i^{th} column of \mathbf{Z} and $\boldsymbol{\mu} = \frac{1}{m}\sum_{i=1}^{m}\mathbf{z_i} \in \mathbb{R}^{n\times 1}$. We consider different values for $n \in \{500, 1000, 1500, 2000\}$ and denote the data sets as 'Gaussian-Random-n'. (ii) Sparsestructured data sets. We generate the synthetic data in a similar manner as described in [20]. Roughly speaking, we first generate a true inverse covariance matrix \mathbf{X}^* which only contains p nonzero entries with fixing p = 500. Then we inject Gaussian noise to \mathbf{X}^* to obtain a noisy covariance matrix Σ . We consider different values for $n \in \{500, 1000, 1500, 2000\}$ and denote the data sets as 'Sparse-Structure-n'. (iii) Real-world data sets. We use four well-known real-world data sets {'isolet', 'mnist', 'usps', 'w1a'} in our experiments, all of which can be download in the LIBSVM website 2 . The size of the data sets are 7797×617 , 10000×780 , 9298×256 and 49749×300 , respectively. We construct Σ from the data sets using the same strategy as in Gaussian random data
- Compared methods. We compared the following methods. (a) QUIC applies a Newton-like method [12] to solve the convex ℓ_1 regularized problem ³. Since this method cannot control the sparsity of the solution, we solve the convex problem where the regulation parameter is swept over $2^{\{-10,-9,\ldots,10\}}$. Finally, the solution that leads to smallest objective value after a hard thresholding projection (which reduces to setting the small values of the solution in magnitude to 0) is selected. We use the default stopping criterion for QUIC-L1. (b) ADMM directly applies alternating direction method of multipliers to solve the non-convex ℓ_0 norm problem in(1). (c) Penalty Decomposition Algorithm (PDA) [20] decomposes the ℓ_0 norm problem into a sequence of penalty subproblems which are solved by a block coordinate descent method 4. (d) CWOA is proposed in this paper to solve the original ℓ_0 norm problem. The four methods above are denoted as QUIC-L1, ADMM-L0, PDA-L0, and CWOA-L0, respectively. We vary the parameter s with the range $\{30, 70, 110, 150, 190, 230, 270\}$.

We remark that this paper pays more attention to the solution quality of the non-convex optimization problem in (1). When ℓ_1 convex relaxation is considered, the resulting problem is strongly convex and existing convex methods will exactly lead to the same unique solution. Therefore, we only select QUIC as the representative of convex methods for comparision.

• Quantitative Comparisons. We demonstrate the accuracy of all methods by comparing their objective values in Figure 1. We also report their responding computational time in Figure 2.

Several conclusions can be drawn. (i) CWOA-L0 consistently outperforms existing state-of-the-art approaches in all data sets in term of accuracy. In addition, with increasing sparsity level s, the gap between our method and others becomes larger in some

- $2.\ https://www.csie.ntu.edu.tw/{\sim}cjlin/libsvmtools/datasets$
- $3. \ Code: \ http://www.cs.utexas.edu/{\sim} sustik/QUIC$
- 4. Code: http://people.math.sfu.ca/~zhaosong

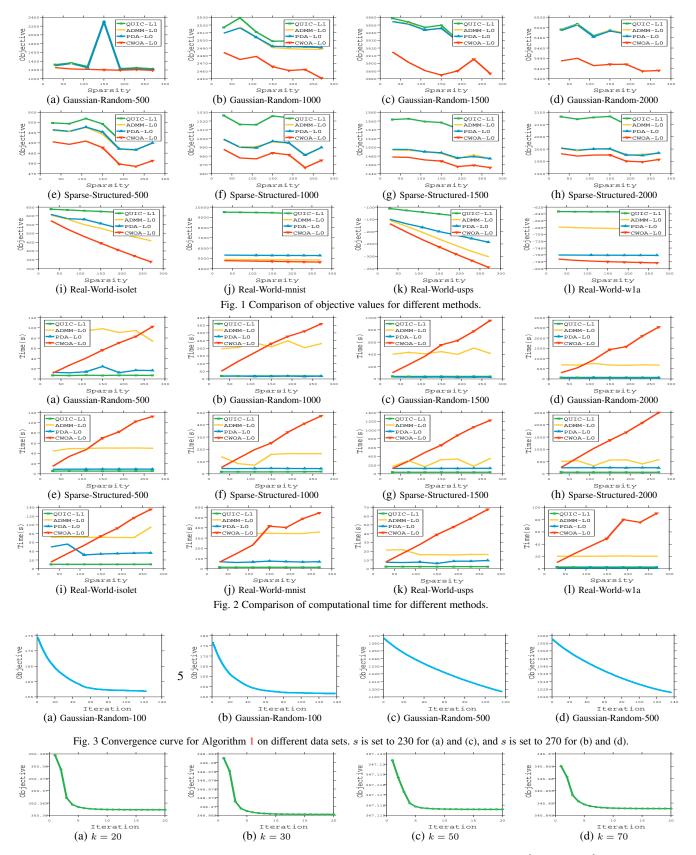


Fig. 4 Convergence curve for Algorithm 2 on Real-World-mnist data set with s=150 for different iterations $t=\{20,30,50,70\}$ of Algorithm 1.

data sets. (ii) The solutions generated by QUIC-L1 and PDA-L0 may not always satisfy the positive definite constraint and incur much larger objective values. In addition, ADM-L0 seems

to present the second best results in terms of accuracy. (iii) While the computational time of the other methods are insensitive to the change in sparsity level s, the computational time of our method

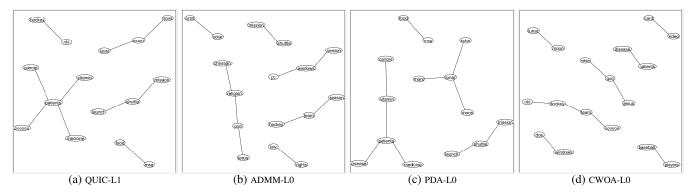


Fig. 5 A practical example on '20newsgroups' data set. The objective values generated by the four methods are 717.04, 715.88, 717.45, and 714.76, respectively.

scales linearly with the sparsity level. This is expected since CWOA-L0 needs to solve the reduced convex problem for at least s times. (iv) The proposed method is slower than the compared method when the s is large. However, the computational time pays off since our method consistently achieves lower objective values.

• Convergence Behavior. First of all, we demonstrate the convergence behavior of Algorithm 1 with different sparsity level s on different data sets in Figure 3. Since the methods {QUIC-L1, ADMM-L0, PDA-L0} may violate positive definite constraint or the sparsity constraint and our method always generates feasible solutions ($\mathbf{X} \succ 0$ and $\|\mathbf{X}\|_{0,\text{off}} \leq s$) in *all* iterations, we do not compare the objective values for different iterations of the algorithms.

We make two important observations from these results. (i) The objective value decrease monotonically. This is because Algorithm 1 is a greedy descent algorithm. (ii)We observe from Figure 3 that Algorithm 1 terminates at iteration $\{122,138,116,136\}$. Therefore, Algorithm 1 spent $\{230,270,230,270\}/2$ iterations and $\{122,138,116,136\}-\{230,270,230,270\}/2=\{7,3,1,1\}$ iterations to perform the greedy pursuit stage and swap coordinates stage, respectively.

Secondly, we demonstrate the convergence behavior of Algorithm 2 on Real-World-mnist data set with different s for different iterations $t=\{20,30,50,70\}$ of Algorithm 1 in Figure 4. We make two important observations from these results. (i) The objective value decreases monotonically. (ii) The objective values stabilize after the 10th iteration, which means that our algorithm has converged, and the decrease of the objective value is negligible after the 10th iteration. This implies that one may use a looser stopping criterion without sacrificing accuracy.

• A Practical Example. We consider different methods to solve the sparse inverse covariance selection problem on a processed version of the 20 newsgroups data set ⁵. We expect to obtain a small relation graph with 10 edges, thus, k is set to 20 in our experiments. Two conclusions can be drawn from Figure 5. (i) Our method achieves the lowest objective value for solving the sparse inverse covariance selection problem. (ii) The proposed method CWOA-L0 is observed to output strong relation patterns in this example.

6 CONCLUSIONS

In this paper, we have developed an effective and efficient coordinate-wise optimization algorithm for solving the non-convex ℓ_0 norm sparse inverse covariance selection problem. The

algorithm is guaranteed to converge to a desirable coordinatewise minimum point. Extensive experiments have shown that the proposed method *consistently* outperforms existing methods.

REFERENCES

- [1] 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 (JMLR)*, 9:485–516, 2008.
- [2] Amir Beck and Yonina C Eldar. Sparsity constrained nonlinear optimization: Optimality conditions and algorithms. SIAM Journal on Optimization (SIOPT), 23(3):1480–1509, 2013.
- [3] Amir Beck and Yakov Vaisbourd. The sparse principal component analysis problem: Optimality conditions and algorithms. *Journal of Optimization Theory and Applications*, 170(1):119–143, 2016.
- [4] Stephen Boyd and Lieven Vandenberghe. Convex Optimization. Cambridge University Press, New York, NY, USA, 2004.
- [5] Emmanuel J Candes, Michael B Wakin, and Stephen P Boyd. Enhancing sparsity by reweighted ℓ_1 minimization. Journal of Fourier Analysis and Applications, 14(5-6):877–905, 2008.
- [6] Alexandre d'Aspremont, Onureena Banerjee, and Laurent El Ghaoui. First-order methods for sparse covariance selection. SIAM Journal on Matrix Analysis and Applications (SIMAX), 30(1):56–66, 2008.
- [7] Bin Dong and Yong Zhang. An efficient algorithm for ℓ_0 minimization in wavelet frame based image restoration. *Journal of Scientific Computing*, 54(2-3):350–368, 2013.
- [8] John C. Duchi, Stephen Gould, and Daphne Koller. Projected subgradient methods for learning sparse gaussians. In *Con*ference in Uncertainty in Artificial Intelligence (UAI), pages 145–152, 2008.
- [9] Jerome Friedman, Trevor Hastie, and Robert Tibshirani. Sparse inverse covariance estimation with the graphical lasso. *Biostatistics*, 9(3):432–441, 2008.
- [10] Jean Honorio, Dimitris Samaras, Nikos Paragios, Rita Goldstein, and Luis E Ortiz. Sparse and locally constant gaussian graphical models. In *Advances in Neural Information Processing Systems*, pages 745–753, 2009.
- [11] Cho-Jui Hsieh, Inderjit S. Dhillon, Pradeep Ravikumar, and Arindam Banerjee. A divide-and-conquer method for sparse

- inverse covariance estimation. In *Neural Information Processing Systems (NIPS)*, pages 2339–2347, 2012.
- [12] Cho-Jui Hsieh, Mátyás A. Sustik, Inderjit S. Dhillon, and Pradeep Ravikumar. QUIC: quadratic approximation for sparse inverse covariance estimation. *Journal of Machine Learning Research (JMLR)*, 15(1):2911–2947, 2014.
- [13] Michael Irwin Jordan. Learning in graphical models, volume 89. Springer Science & Business Media, 1998.
- [14] S. Sathiya Keerthi, Olivier Chapelle, and Dennis DeCoste. Building support vector machines with reduced classifier complexity. *Journal of Machine Learning Research*, 7:1493– 1515, 2006.
- [15] Andreas Krause and Carlos Guestrin. Near-optimal observation selection using submodular functions. In *AAAI Conference on Artificial Intelligence (AAAI)*, pages 1650–1654, 2007.
- [16] Jason D. Lee, Yuekai Sun, and Michael A. Saunders. Proximal newton-type methods for minimizing composite functions. SIAM Journal on Optimization (SIOPT), 24(3):1420–1443, 2014.
- [17] Lu Li and Kim-Chuan Toh. An inexact interior point method for 11-regularized sparse covariance selection. *Mathematical Programming Computation*, 2(3):291–315, 2010.
- [18] Zhaosong Lu. Smooth optimization approach for sparse covariance selection. SIAM Journal on Optimization (SIOPT), 19(4):1807–1827, 2009.
- [19] Zhaosong Lu. Adaptive first-order methods for general sparse inverse covariance selection. SIAM Journal on Matrix Analysis and Applications (SIMAX), 31(4):2000–2016, 2010.
- [20] Zhaosong Lu and Yong Zhang. Sparse approximation via penalty decomposition methods. *SIAM Journal on Optimization (SIOPT)*, 23(4):2448–2478, 2013.
- [21] Goran Marjanovic and Alfred O Hero. $l_{-}\{0\}$ sparse inverse covariance estimation. *IEEE Transactions on Signal Processing*, 63(12):3218–3231, 2015.
- [22] Y. E. Nesterov. *Introductory lectures on convex optimization: a basic course*, volume 87 of *Applied Optimization*. Kluwer Academic Publishers, 2003.
- [23] Yurii Nesterov. Towards non-symmetric conic optimization. *Optimization Methods and Software*, 27(4-5):893–917, 2012
- [24] Peder A. Olsen, Figen Öztoprak, Jorge Nocedal, and Steven J. Rennie. Newton-like methods for sparse inverse covariance estimation. In *Neural Information Processing Systems (NIPS)*, pages 764–772, 2012.
- [25] Figen Oztoprak, Jorge Nocedal, Steven Rennie, and Peder A Olsen. Newton-like methods for sparse inverse covariance estimation. In *Neural Information Processing Systems* (*NIPS*), pages 755–763, 2012.
- [26] Katya Scheinberg, Shiqian Ma, and Donald Goldfarb. Sparse inverse covariance selection via alternating linearization methods. In *Neural Information Processing Systems (NIPS)*, pages 2101–2109, 2010.
- [27] Nasim Souly and Mubarak Shah. Scene labeling using sparse precision matrix. In *IEEE Conference on Computer Vision* and Pattern Recognition (CVPR), pages 3650–3658, 2016.
- [28] Joel A Tropp and Anna C Gilbert. Signal recovery from random measurements via orthogonal matching pursuit. *IEEE Transactions on Information Theory*, 53(12):4655–4666, 2007.

- [29] Paul Tseng and Sangwoon Yun. A coordinate gradient descent method for nonsmooth separable minimization. *Mathematical Programming*, 117(1):387–423, 2009.
- [30] Paul Tseng and Sangwoon Yun. A coordinate gradient descent method for nonsmooth separable minimization. *Mathematical Programming*, 117(1-2):387–423, 2009.
- [31] Ganzhao Yuan and Bernard Ghanem. $\ell_0 tv$: A new method for image restoration in the presence of impulse noise. In Computer Vision and Pattern Recognition (CVPR), pages 5369–5377, 2015.
- [32] Ganzhao Yuan, Yin Yang, Zhenjie Zhang, and Zhifeng Hao. Convex optimization for linear query processing under approximate differential privacy. In ACM International Conference on Knowledge Discovery and Data Mining (SIGKDD), pages 2005–2014, 2016.
- [33] Xiao-Tong Yuan and Tong Zhang. Truncated power method for sparse eigenvalue problems. *Journal of Machine Learning Research*, 14(Apr):899–925, 2013.
- [34] Xiaoming Yuan. Alternating direction method for covariance selection models. *Journal of Scientific Computing*, 51(2):261–273, 2012.
- [35] Sangwoon Yun, Paul Tseng, and Kim-Chuan Toh. A block coordinate gradient descent method for regularized convex separable optimization and covariance selection. *Mathematical Programming*, 129(2):331–355, 2011.
- [36] Tong Zhang. Analysis of multi-stage convex relaxation for sparse regularization. *The Journal of Machine Learning Research (JMLR)*, 11:1081–1107, 2010.
- [37] Tong Zhang. Adaptive forward-backward greedy algorithm for learning sparse representations. *IEEE Transactions on Information Theory*, 57(7):4689–4708, 2011.
- [38] Luping Zhou, Lei Wang, and Philip Ogunbona. Discriminative sparse inverse covariance matrix: Application in brain functional network classification. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3097–3104, 2014.