

# Efficient QUIC-Based Damped Inexact Iterative Reweighting for Sparse Inverse Covariance Estimation with Nonconvex Partly Smooth Regularization

Xiangyu Yang<sup>1,2</sup>, Jiarui Wei<sup>2</sup>, and Makoto Yamashita<sup>3</sup>

<sup>1</sup>School of Mathematics and Statistics, Henan University,  
Kaifeng 475000, China

<sup>2</sup>Center for Applied Mathematics of Henan Province, Henan  
University, Zhengzhou, 450046, China

<sup>3</sup>Department of Mathematical and Computing Science,  
Institute of Science Tokyo, Tokyo 152-8552, Japan

<sup>1</sup>*yangxy@henu.edu.cn*

<sup>2</sup>*jirwei@henu.edu.cn*

<sup>3</sup>*Makoto.Yamashita@comp.isct.ac.jp*

## Abstract

In this paper, we study sparse inverse covariance matrix estimation incorporating partly smooth nonconvex regularizers. To solve the resulting regularized log-determinant problem, we develop DIIR-QUIC—a novel Damped Inexact Iteratively Reweighted algorithm based on QUadratic approximate Inverse Covariance (QUIC) method. Our approach generalizes the classic iteratively reweighted  $\ell_1$  scheme through damped fixed-point updates. A key novelty of DIIR-QUIC is an inexact termination criterion for the subproblems that permits controlled inexactness in solutions to accelerate each iteration while still guaranteeing identification of the active manifold in finitely many

steps. We establish the global convergence of DIIR-QUIC and, under the Kurdyka-Łojasiewicz property, prove Q-linear convergence of the perturbed objective values and R-linear convergence of the iterates. Extensive numerical experiments on synthetic and real-world datasets demonstrate that DIIR-QUIC outperforms existing approaches in computational efficiency and estimation accuracy.

**Keywords**— Inverse covariance matrix estimation, Nonconvex regularization optimization, Smooth active manifold, Damped inexact iterative reweighting

## 1 Introduction

Estimating the inverse covariance matrix, also known as the precision matrix, is a fundamental problem in modern multivariate statistical analysis. A key motivation for the estimation is that its zero entries directly encode conditional-independence relationships among variables. This built-in sparsity makes precision estimation a versatile tool for high-dimensional problems: it underlies linear discriminant analysis in statistical learning [1], guides optimal asset allocation in portfolio optimization [2] and facilitates gene-network reconstruction in computational biology [3]. As data dimensions grow, developing effective methods that balance statistical reliability with computational scalability has become an increasingly active area of research.

Sparsity-promoting precision estimation is especially vital when the number of variables far exceeds the sample size [1]. In the setting of Gaussian Markov Random Fields (GMRF), we usually observe  $m$  independently drawn samples  $\mathbf{y}_i \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$  in  $\mathbb{R}^n$  and seek to recover the precision matrix  $\boldsymbol{\Sigma}^{-1}$ . Enforcing a sparse structure on  $\boldsymbol{\Sigma}^{-1}$  generally enhances statistical inference and interpretability in the high-dimensional regime (i.e.,  $m \ll n$ ). A zero entry in the precision matrix corresponds to a conditional independence constraint, implying no direct dependency between the associated variables. As a result, estimating a sparse precision matrix is equivalent to learning a sparse undirected graph in a GMRF.

To obtain a sparse precision matrix, one widely studied approach in high-dimensional settings is the regularized log-likelihood formulation [4, 5]. This involves maximizing the log-likelihood function over the space of positive definite matrices while incorporating a sparsity-inducing regularizer. A natural choice for enforcing sparse pattern in the solution is the  $\ell_0$  (quasi-)norm, defined as  $\|\mathbf{X}\|_0 := \sum_{ij} \mathbb{I}(X_{ij} \neq 0)$ , where  $\mathbb{I}(\cdot)$  denotes the Boolean indicator function that returns 1 when the condition is true and 0 otherwise. The  $\ell_0$ -norm is particularly appealing due to its ability to induce exact sparsity. However, its combinatorial nature renders the associated optimization problem computationally intractable, as

commented in [6]. To regain tractability, convex relaxations replace  $\|\mathbf{X}\|_0$  with its tightest convex surrogate, the  $\ell_1$ -norm. Regularizing the log-determinant function with an  $\ell_1$ -penalty ensures strong convexity properties and, through Lagrangian duality, often yields computationally efficient solutions [7, 8]. Despite these advantages, the  $\ell_1$ -penalty suffers from a well-known limitation: it tends to systematically over-penalize large coefficients, introducing bias and potentially misidentifying the true sparsity pattern. This bias generally arises due to the relaxation gap between  $\ell_1$ - and  $\ell_0$ -norm formulations [9].

To address this issue, nonconvex regularization approaches have gained considerable attention [10]. In particular, the nonconvex regularizers such as the  $\ell_p^p$ -norm ( $0 < p < 1$ ) [11], the smoothly clipped absolute deviation (SCAD) and the min-max concave penalty (MCP) [12], as well as the piecewise exponential concave approximation function and the capped  $\ell_1$  regularizers [6], have demonstrated superior numerical performance compared to  $\ell_1$ -norm regularization. These methods often yield sparser models with improved prediction accuracy. Such findings highlight the need for alternative regularization strategies that achieve a better balance between sparsity and estimation accuracy, thereby motivating further exploration of nonconvex sparsity-promoting penalties.

In this paper, we focus on the following matrix optimization problem with nonconvex sparsity-promoting regularizers:

$$\begin{aligned} \min \quad & F(\mathbf{X}) = \{f(\mathbf{X}) := \text{tr}(\mathbf{S}\mathbf{X}) - \log \det \mathbf{X}\} + \rho \left\{ \Phi(\mathbf{X}) := \sum_{ij} \phi(|X_{ij}|) \right\} \quad (\mathcal{P}) \\ \text{s.t.} \quad & \mathbf{X} \in \mathbb{S}_{++}^n, \end{aligned}$$

where  $\mathbf{S} \in \mathbb{S}_+^n$  is the empirical covariance matrix, and  $\rho > 0$  is a regularization parameter. Throughout this paper, we write  $\mathbb{S}_{++}^n$  (respectively,  $\mathbb{S}_+^n$ ) for the set of  $n$ -by- $n$  symmetric positive definite (respectively, positive semidefinite) matrices. The following assumption on the function  $\phi$  is imposed throughout.

**Assumption 1.1.** *Let  $\phi : [0, +\infty) \rightarrow [0, +\infty)$  satisfy:*

- (i)  $\phi(0) = 0$ , and  $\phi$  is concave on  $[0, +\infty)$ .
- (ii)  $\phi \in \mathcal{C}^1((0, +\infty))$ ,  $\phi'(t) \geq 0$  for all  $t > 0$ , and  $\phi'$  is nonincreasing on  $(0, +\infty)$ .
- (iii) For any  $\underline{\delta} > 0$ ,  $\phi'$  is Lipschitz on  $[\underline{\delta}, +\infty)$ .
- (iv) Let

$$\phi'(0^+) = \lim_{t \rightarrow 0^+} \phi'(t) \in (0, +\infty] \quad \text{and} \quad \phi'(+\infty) = \lim_{t \rightarrow +\infty} \phi'(t) \geq 0.$$

There exists  $t^* \in (0, +\infty]$  such that

$$\phi'(t) > 0, \quad \forall t \in [0, t^*), \quad \phi'(t) = 0, \quad \forall t \in [t^*, +\infty) \quad (\text{if } t^* < +\infty).$$

Then the inverse  $(\phi')^{-1} : (\phi'(+\infty), \phi'(0^+)) \rightarrow (0, t^*)$  defined by

$$(\phi')^{-1}(s) = \inf\{t > 0 \mid \phi'(t) \leq s\}.$$

is well-defined and continuous on its effective domain.

Under Assumption 1.1, the function  $\Phi$  covers a wide range of nonconvex surrogates for the  $\ell_0$  norm, which counts the number of nonzero entries in a matrix. Notable examples include  $\ell_p$  (quasi-)norm [13], SCAD [14] and MCP [15]. Table 3 in Appendix summarizes several representative instances.

On the algorithmic front, addressing the nonconvex optimization problem  $(\mathcal{P})$  remains a challenging task. To our knowledge, [11] was the first to address problem  $(\mathcal{P})$  with  $\Psi(\mathbf{X}) = \sum_{i \neq j} |X_{ij}|^p$  for  $0 < p < 1$  in its maximization form. The authors proposed a two-stage alternating optimization algorithm. In the first stage, they reformulated the original matrix problem into an equivalent vector problem by exploiting permutations of the matrix variables and the empirical covariance matrix. In the second stage, they applied a cyclic descent method to solve the resulting  $\ell_p^p$ -regularized least squares problem, updating each coordinate sequentially. At each iteration, their algorithm performs updates on a single row and column of the matrix, involving three components: one column vector, a scalar and a principal submatrix. However, the algorithm lacks theoretical guarantees for global convergence—precluding any convergence rate analysis—and its numerical evaluation has been confined to relatively small problem instances (typically on the order of 100 variables). Moreover, our reproduced experiments indicate that the algorithm incurs a comparatively high computational cost and generally fails to achieve satisfactorily low stationarity residuals. These limitations restrict its practical applicability to large-scale problems.

Another line of research, while not directly tackling the same problem, nonetheless offers useful insights. Notable examples include the works of Phan et al. [6] and Wei et al. [12]. Phan et al. [6] focused on a formulation involving a nonconvex loss function of the form  $f(\mathbf{X}) = \log \det \mathbf{X} + \text{tr}(\mathbf{S}\mathbf{X}^{-1})$ , combined with nonconvex sparsity-promoting penalties, specifically the difference-of-convex (DC) representations of the piecewise exponential concave approximation and the capped  $\ell_1$  regularization. The authors proposed two tailored Difference-of-Convex Algorithm (DCA) variants and established global subsequential convergence. On the other hand, Wei et al. [12] addressed related covariance estimation problems involving a loss function of the form  $f(\mathbf{X}) = \frac{1}{2}\|\Sigma - \mathbf{S}\|_F^2 - \tau \log \det \mathbf{X}$  ( $\tau > 0$  is used in their original texts), combined with the SCAD and MCP penalties. The authors applied

a standard iteratively reweighted  $\ell_1$  algorithm, with each subproblem handled by an inexact proximal gradient method and backtracking line-search. However, it is crucial to note that their statistical convergence—i.e., the consistency guarantees for the sequence of covariance estimators—and their iteration-complexity bounds for the inner subproblem solver are derived under the oracle assumption that the true support of the covariance matrix is known a priori. In addition, the global convergence of the generated iterates to a stationary point of the original nonconvex problem remains unestablished. Furthermore, the algorithms in both of these works cannot accommodate nonconvex regularizers such as the nonconvex  $\ell_p^p$  norm, since its derivative is unbounded at the origin.

We propose to estimate large-scale sparse covariance matrices by solving the log-determinant program  $(\mathcal{P})$  with more general nonconvex and nonsmooth regularizers. To this end, we develop an algorithmic framework that builds on efficient solvers for the weighted  $\ell_1$ -regularized log-determinant subproblem. Our framework generalizes the iteratively reweighted  $\ell_1$  (IR $\ell_1$ ) schemes in [16]: at each iteration, we first apply a smoothing and reweighting step to define a weighted  $\ell_1$ -regularized log-determinant subproblem and compute an intermediate solution. We then employ a damped update operator, as in [17, 18], to generate the next iterate.

While the underlying idea is simple and natural, two key challenges arise in practice. First, each weighted  $\ell_1$  subproblem can only be solved approximately, and it is unclear how much inexactness in the subproblem solution can be allowed to guarantee the global convergence of the method. Second, even if we assume the global convergence properties of the damped IR $\ell_1$  framework are guaranteed, it is also unclear whether the algorithm can identify the active manifold under inexact subproblem solutions. Smooth active manifolds are useful in nonsmooth optimization, as noted by the authors of [19], “*Once  $\mathcal{M}$  is identified, the nonsmoothness of the problem is largely irrelevant, since all future iterates lie on a smooth manifold along which  $f$  is smooth.*” Here,  $\mathcal{M}$  refers to a smooth active manifold.

In this paper, we address both issues in the damped IR $\ell_1$  framework by incorporating an inexact criterion for the subproblem solver and proving that the proposed algorithm is well-posed and converges globally. We further show that the algorithm correctly identifies the active manifold in a finite number of iterations under the proposed inexact criterion. Moreover, under the Kurdyka–Łojasiewicz (KL) property, we prove that the objective values converge at a Q-linear rate and the iterates uniquely converge at an R-linear rate. Numerical results confirm the effectiveness and efficiency of the proposed algorithm.

## 1.1 Notation and preliminaries

Throughout, let  $\mathbb{N}$ ,  $\mathbb{R}$ ,  $\mathbb{R}_+ := [0, +\infty)$  and  $\mathbb{R}_{++} := (0, +\infty)$  denote the sets of natural numbers, real numbers, nonnegative real numbers, and positive real numbers, respectively. Correspondingly,  $\mathbb{R}^{m \times n}$  denotes the Euclidean space of  $m$ -by- $n$ -dimensional real matrices, with inner product  $\langle \mathbf{X}, \mathbf{Y} \rangle = \text{tr}(\mathbf{X}^T \mathbf{Y})$ . Let  $\mathbb{S}_{>0}^n = \mathbb{R}_{++}^{n \times n} \cap \mathbb{S}^n$  denote the set of  $n$ -by- $n$ -dimensional symmetric matrices with positive entries. Meanwhile, the notation  $\mathbf{X} > \mathbf{0}$  indicates that  $\mathbf{X}$  is positive definite, while  $\mathbf{X} \geq \mathbf{0}$  means that  $\mathbf{X}$  is positive semidefinite. The Kronecker product of the matrices  $\mathbf{X}$  and  $\mathbf{Y}$  is denoted by  $\mathbf{X} \otimes \mathbf{Y}$ . We use  $\mathcal{I}(\mathbf{X}) := \{(i, j) \in [n] \times [n] \mid X_{ij} \neq 0\}$  and  $\mathcal{Z}(\mathbf{X}) := \{(i, j) \in [n] \times [n] \mid X_{ij} = 0\}$  to denote the index set of the nonzeros and zeros of  $\mathbf{X}$ , respectively. Given a matrix  $\mathbf{X} \in \mathbb{S}^n$  and a positive weight matrix  $\mathbf{W} \in \mathbb{R}_+^{n \times n}$ , we define the weighted  $\ell_1$ -norm of  $\mathbf{X}$  as  $\|\mathbf{X}\|_{1, \mathbf{W}} = \sum_{i,j=1}^n W_{ij} |X_{ij}|$ . In addition,  $\|\mathbf{X}\|_\infty := \max_{(i,j) \in [n] \times [n]} |X_{ij}|$ . For any nonempty set  $\Omega \subset \mathbb{R}^{n \times n}$ ,  $\text{rint}(\Omega)$  denotes its relative interior.

Throughout the paper we measure distances in  $\mathbb{R}^{n \times n}$  with the Frobenius norm  $\|\cdot\|_F$ . In particular:

- (i) Set-to-set distance. For any two nonempty sets  $\bar{\Omega}, \Omega \subset \mathbb{R}^{n \times n}$ , define

$$\text{dist}(\bar{\Omega}, \Omega) := \inf\{\|\mathbf{X} - \mathbf{Y}\|_F \mid \mathbf{X} \in \bar{\Omega}, \mathbf{Y} \in \Omega\}.$$

Note that here  $\text{dist}(\cdot, \cdot)$  is not a proper distance, since it may fail the triangle inequality.

- (ii) Point-to-set distance. For any  $\mathbf{X} \in \mathbb{R}^{n \times n}$  and nonempty  $\Omega \subset \mathbb{R}^{n \times n}$ , we abbreviate  $\text{dist}(\mathbf{X}, \Omega) := \text{dist}(\{\mathbf{X}\}, \Omega)$ .
- (iii) Scalar case. When  $n = 1$ , this reduces to the usual absolute-value distance  $\text{dist}(x, y) = |x - y|$ ,  $\forall x, y \in \mathbb{R}$ .

The following proposition is useful in our analysis, and its proof can be found in [20, Appendix A.3].

**Proposition 1.2** (Triangle inequality for set distances [20, Proposition 3]). *For any three sets  $\Omega_1, \Omega_2, \Omega_3 \subseteq \mathbb{R}^{n \times n}$ , it holds that*

$$\text{dist}(\Omega_1, \Omega_2) \leq \text{dist}_{\mathcal{H}}(\Omega_1 \mid \Omega_3) + \text{dist}(\Omega_3, \Omega_2), \quad (1)$$

where  $\text{dist}_{\mathcal{H}}(\Omega_1 \mid \Omega_3) := \sup_{\mathbf{X} \in \Omega_3} \inf_{\mathbf{Y} \in \Omega_1} \|\mathbf{Y} - \mathbf{X}\|_F$  refers to the Hausdorff distance between  $\Omega_1$  and  $\Omega_3$ .

Before presenting the stationarity condition of  $(\mathcal{P})$ , we first characterize the subdifferentials of  $\Phi$  in the following lemma.

**Lemma 1.3** (Subgradients [21, Definition 8.3] and their relationships [21, Theorem 8.6]). *Consider  $(\mathcal{P})$ . Let  $\mathbf{X} \in \mathbb{S}_{++}^n$ . The following holds:*

(i)  $\hat{\partial}\Phi(\mathbf{X}) = \partial\Phi(\mathbf{X}) = \partial\phi(|X_{11}|) \times \partial\phi(|X_{12}|) \times \cdots \times \partial\phi(|X_{nn}|)$  with

$$\partial\phi(|X_{ij}|) = \begin{cases} \{\phi'(|X_{ij}|)\}, & (i, j) \in \mathcal{I}(\mathbf{X}), \\ [-\phi'(0), \phi'(0)], & (i, j) \in \mathcal{Z}(\mathbf{X}) \text{ and } \lim_{s \rightarrow 0^+} \phi'(s) < +\infty, \\ \mathbb{R}, & (i, j) \in \mathcal{Z}(\mathbf{X}) \text{ and } \lim_{s \rightarrow 0^+} \phi'(s) = +\infty. \end{cases}$$

Here,  $\hat{\partial}\Phi(\mathbf{X})$  and  $\partial\Phi(\mathbf{X})$  refer to the regular, limiting (or Mordukhovich) subdifferentials of  $\Phi$  at  $\mathbf{X}$ , respectively. Both  $\partial\Phi(\mathbf{X})$  and  $\hat{\partial}\Phi(\mathbf{X})$  are closed sets, and  $\hat{\partial}\Phi(\mathbf{X})$  is convex.

(ii)  $\hat{\partial}\Phi(\mathbf{X})^\infty = \partial^\infty\Phi(\mathbf{X}) = \partial^\infty\phi(|X_{11}|) \times \partial^\infty\phi(|X_{12}|) \times \cdots \times \partial^\infty\phi(|X_{nn}|)$  with

$$\partial^\infty\phi(|X_{ij}|) = \begin{cases} \{0\}, & (i, j) \in \mathcal{I}(\mathbf{X}), \\ \{0\}, & (i, j) \in \mathcal{Z}(\mathbf{X}) \text{ and } \lim_{s \rightarrow 0^+} \phi'(s) < +\infty, \\ \mathbb{R}, & (i, j) \in \mathcal{Z}(\mathbf{X}) \text{ and } \lim_{s \rightarrow 0^+} \phi'(s) = +\infty. \end{cases}$$

Here,  $\partial^\infty\Phi(\mathbf{X})$  refers to the horizon subdifferential of  $\Phi$  at  $\mathbf{X}$ , and  $\hat{\partial}\Phi(\mathbf{X})^\infty$  refers to the horizon cone of  $\hat{\partial}\Phi(\mathbf{X})$  [21, Definition 3.3]. Furthermore,  $\partial^\infty\Phi(\mathbf{X})$  and  $\hat{\partial}\Phi(\mathbf{X})^\infty$  are closed cones, with  $\hat{\partial}\Phi(\mathbf{X})^\infty$  convex.

Consequently, by Assumption 1.1 and [21, Corollary 8.11],  $\Phi(\mathbf{X})$  is (subdifferentially) regular at  $\mathbf{X}$ . In addition, it follows from [21, Exercise 8.8] and  $f \in \mathcal{C}^1$  that  $F$  is (subdifferentially) regular for all  $\mathbf{X} \in \mathbb{S}_{++}^n$ .

In particular, given a positive weight matrix  $\mathbf{W} \in \mathbb{R}_+^{n \times n}$ , we have

$$\partial\|\mathbf{X}\|_{1,\mathbf{W}} = \left\{ \mathbf{G} \in \mathbb{S}^n \mid G_{ij} \begin{cases} = W_{ij} \text{sgn}(X_{ij}), & \text{if } X_{ij} \neq 0, \\ \in [-W_{ij}, W_{ij}], & \text{if } X_{ij} = 0. \end{cases} \right\}. \quad (2)$$

We next provide the first-order necessary optimality condition for problem  $(\mathcal{P})$ .

**Theorem 1.4** (Fermat's rule generalized [21, Theorem 10.1]). *Consider  $(\mathcal{P})$  under Assumption 1.1. If  $F$  has a local minimum at  $\mathbf{X}^* \in \mathbb{S}_{++}^n$ , then*

$$-\nabla f(\mathbf{X}^*) \in \partial\Phi(\mathbf{X}^*).$$

Indeed, it holds that

$$(i, j) \in \mathcal{I}(\mathbf{X}^*) : \quad \nabla_{ij} f(\mathbf{X}^*) + \rho \phi'(|X_{ij}^*|) \operatorname{sgn}(X_{ij}^*) = 0, \quad (3a)$$

$$(i, j) \in \mathcal{Z}(\mathbf{X}^*) : \quad \begin{cases} |\nabla_{ij} f(\mathbf{X}^*)| \leq \rho \phi'(0), & \lim_{s \rightarrow 0^+} \phi'(s) < +\infty, \\ -\nabla_{ij} f(\mathbf{X}^*) \in \mathbb{R}, & \lim_{s \rightarrow 0^+} \phi'(s) = +\infty. \end{cases} \quad (3b)$$

Then we say that a matrix  $\mathbf{X}^* \in \mathbb{S}_{++}^n$  is a stationary point of  $F$  if (3) holds.

We next recall the notion of partial smoothness, which captures the intrinsic smooth structure underlying a nonsmooth function. We begin with some elementary definitions.

**Definition 1.5.** Let  $\Omega \subset \mathbb{R}^n$  be a nonempty convex set. The subspace parallel to the set  $\Omega$ , denoted by  $\operatorname{par} \Omega$ , is defined as  $\operatorname{par} \Omega = \operatorname{aff} \Omega - \mathbf{x}$ ,  $\forall \mathbf{x} \in \Omega$ , where  $\operatorname{aff} \Omega$  is the affine span of  $\Omega$ .

**Definition 1.6** (Partly smooth function [22, Definition 2.7]). Suppose that the set  $\mathcal{M} \subset \mathbb{R}^{m \times n}$  contains the point  $\mathbf{X}$ . A function  $h : \mathbb{R}^{m \times n} \rightarrow \mathbb{R} \cup \{+\infty\}$  is said to be *partly smooth* at  $\mathbf{X}$  relative to  $\mathcal{M}$  if  $\mathcal{M}$  is a manifold around  $\mathbf{X}$  and the following four properties hold:

- (i) (**Restricted smoothness**) the restriction  $h|_{\mathcal{M}}$  is smooth at  $\mathbf{X}$ ;
- (ii) (**Regularity**) at every point close to  $\mathbf{X}$  in  $\mathcal{M}$ , the function  $h$  is regular and has a subgradient;
- (iii) (**Normals parallel to subdifferential**)  $N_{\mathcal{M}}(\mathbf{X}) \subset \operatorname{par} \partial h(\mathbf{X})$ , where  $N_{\mathcal{M}}(\mathbf{X})$  denotes the normal space to  $\mathcal{M}$  at  $\mathbf{X}$ ;
- (iv) (**Subgradient continuity**) the subdifferential map  $\partial h$  is continuous at  $\mathbf{X}$  relative to  $\mathcal{M}$ .

We say that the function  $h$  is partly smooth relative to a set  $\mathcal{M}$  if  $\mathcal{M}$  is a manifold and  $h$  is partly smooth at each point in  $\mathcal{M}$  relative to  $\mathcal{M}$ . In addition,  $\mathcal{M}$  is referred to the *active manifold* (of partial smoothness).

For a partly smooth function, the condition (iii) in Definition 1.6 reveals a “stable” property. We restate this result in the following proposition.

**Proposition 1.7** (Local normal sharpness [22, Proposition 2.10]). *If the function  $h : \mathbb{R}^{m \times n} \rightarrow \mathbb{R} \cup \{+\infty\}$  is partly smooth at the point  $\mathbf{X}_0$  relative to the set  $\mathcal{M} \subset \mathbb{R}^{m \times n}$ , then all points  $\mathbf{X} \in \mathcal{M}$  close to  $\mathbf{X}_0$  satisfy the condition  $N_{\mathcal{M}}(\mathbf{X}) = \operatorname{par} \partial h(\mathbf{X})$ .*



**Definition 1.8** (Prox-regularity [23, Definition 2.1]). A function  $h : \mathbb{R}^{m \times n} \rightarrow \mathbb{R} \cup \{+\infty\}$  is *prox-regular* at a point  $\bar{\mathbf{X}}$  for a subgradient  $\bar{\mathbf{V}} \in \partial h(\bar{\mathbf{X}})$  if  $h$  is finite at  $\bar{\mathbf{X}}$ , locally lower semi-continuous around  $\bar{\mathbf{X}}$ , and there exists  $\rho > 0$  such that

$$h(\mathbf{X}') \geq h(\mathbf{X}) + \langle \mathbf{V}, \mathbf{X}' - \mathbf{X} \rangle - \frac{\rho}{2} \|\mathbf{X}' - \mathbf{X}\|_F^2$$

whenever  $\mathbf{X}$  and  $\mathbf{X}'$  are near  $\bar{\mathbf{X}}$  with  $h(\mathbf{X})$  near  $h(\bar{\mathbf{X}})$  and  $\mathbf{V} \in \partial h(\mathbf{X})$  is near  $\bar{\mathbf{V}}$ . Furthermore,  $h$  is *prox-regular* at  $\bar{\mathbf{X}}$  if it is prox-regular at  $\bar{\mathbf{X}}$  for every  $\mathbf{V} \in \partial h(\bar{\mathbf{X}})$ .

The following proposition states the basic conditions that guarantee the algorithm identifies the active manifold in a finite number of iterations.

**Proposition 1.9** (Active manifold identification [23, Theorem 5.3]). *Let the function  $h : \mathbb{R}^{m \times n} \rightarrow \mathbb{R} \cup \{+\infty\}$  be  $\mathcal{C}^q$ -partly smooth ( $q \geq 2$ ) at the point  $\bar{\mathbf{X}}$  relative to a smooth manifold  $\mathcal{M}$ , and prox-regular there, with  $\mathbf{0} \in \text{rint } \partial h(\bar{\mathbf{X}})$ . Suppose  $\mathbf{X}^k \rightarrow \bar{\mathbf{X}}$  with  $h(\mathbf{X}^k) \rightarrow h(\bar{\mathbf{X}})$ . Then*

$$\mathbf{X}^k \in \mathcal{M} \text{ for all sufficiently large } k$$

*if and only if*

$$\text{dist}(\mathbf{0}, \partial h(\mathbf{X}^k)) \rightarrow 0.$$

The KL property is needed in our convergence analysis. For completeness, we recall its essential components below.

**Definition 1.10** (Desingularizing function). [24, Section 3.1.2] Let  $\wp > 0$ . We say that  $\varphi : [0, \wp] \rightarrow \mathbb{R}_+$  is a desingularizing function if (i)  $\varphi(0) = 0$ ; (ii)  $\varphi$  is continuous on  $[0, \wp]$  and of class  $C^1$  on  $(0, \wp)$ ; (iii)  $\varphi'(s) > 0$  for all  $s \in (0, \wp)$ .

**Definition 1.11** (KL property). [25, Definition 3] Let  $F : \mathbb{R}^{m \times n} \rightarrow \mathbb{R} \cup \{+\infty\}$  be proper lower semicontinuous. We say that  $F$  satisfies the Kurdyka-Łojasiewicz property at  $\bar{\mathbf{X}} \in \text{dom}(\partial F) := \{\mathbf{X} \in \mathbb{R}^{m \times n} \mid \partial F(\mathbf{X}) \neq \emptyset\}$  if there exist  $\wp > 0$ , a neighborhood  $\mathbb{U}(\bar{\mathbf{X}}, \rho)$  of  $\bar{\mathbf{X}}$ , and a concave desingularizing function  $\varphi : [0, \wp] \rightarrow \mathbb{R}_+$ , such that the Kurdyka-Łojasiewicz inequality

$$\varphi'(F(\mathbf{X}) - F(\bar{\mathbf{X}})) \text{dist}(\mathbf{0}, \partial F(\bar{\mathbf{X}})) \geq 1 \quad (4)$$

holds, for all  $\mathbf{X}$  in the strict local upper level set

$$\text{Lev}(\bar{\mathbf{X}}, \rho) := \{\mathbf{X} \in \mathbb{U}(\bar{\mathbf{X}}, \rho) \mid F(\bar{\mathbf{X}}) < F(\mathbf{X}) < F(\bar{\mathbf{X}}) + \wp\}.$$

Typical examples of desingularizing functions include those of the form  $\varphi(s) = cs^{1-\theta}$ , where  $c > 0$  and  $\theta \in [0, 1)$  is the KL exponent. If  $F$  satisfies the KL property with exponent  $\theta$  at any  $\mathbf{X} \in \text{dom}(\partial F)$ , then call  $F$  is said to be a KL function with exponent  $\theta$ . According to [26, Lemma 2.1], any proper lower semicontinuous function has the KL property with exponent  $1/2$  at all noncritical points.

## 2 Proposed Algorithm

In this section, we present Damped Inexact Iteratively Reweighted QUIC (DIIR-QUIC) to solve  $(\mathcal{P})$ . The proposed DIIR-QUIC alternates between applying the modified QUIC method to a weighted  $\ell_1$ -regularized log-determinant subproblem and updating the iterate through a damped fixed-point step that blends the subproblem solution with the previous solution estimate. Therefore, DIIR-QUIC involves three main components: (i) a smoothing technique combined with the  $\text{IR}\ell_1$  technique is employed to locally approximate the nonconvex regularizer to yield strongly-convex subproblems; (ii) an efficient subproblem solver for computational efficiency; (iii) a practical inexact termination criterion for the subproblem designed to trade off the inner and outer loop computations. An overall statement of DIIR-QUIC is formally presented in Algorithm 1.

---

**Algorithm 1** Proposed DIIR-QUIC for solving  $(\mathcal{P})$

---

- 1: Input:  $\mathbf{S} \in \mathbb{S}_{++}^n$ ,  $\rho > 0$ ,  $\alpha \in (0, 1)$  and  $\mu \in (0, 1)$
  - 2: (Initialization) Choose  $\mathbf{X}^0 \in \mathbb{S}_{++}^n$  and  $\boldsymbol{\mathcal{E}}^0 \in \mathbb{S}_{>0}^n$ . Set  $k = 0$ .
  - 3: Set  $W_{ij}^k = \phi'(|X_{ij}^k| + \mathcal{E}_{ij}^k)$ ,  $\forall (i, j) \in [n] \times [n]$
  - 4: Solve  $\mathbf{Y}^k \stackrel{(\approx)}{\leftarrow} \text{QUIC}(\rho, \mathbf{X}^k, \mathbf{W}^k, \mathbf{S})$  by invoking Algorithm 2
  - 5: Update  $\mathbf{X}^{k+1} = (1 - \alpha)\mathbf{X}^k + \alpha\mathbf{Y}^k$
  - 6: Update  $\boldsymbol{\mathcal{E}}^{k+1} = (1 - \alpha)\boldsymbol{\mathcal{E}}^k + \alpha(\mu\boldsymbol{\mathcal{E}}^k)$
  - 7: Set  $k \leftarrow k + 1$  and go to step 3
- 

To develop DIIR-QUIC, we first define a function for the objective function  $F$  in  $(\mathcal{P})$  that is nonsmooth but locally Lipschitz continuous. This construction is motivated by the extensive literature on smoothing approximation techniques designed to address the nonsmoothness of the function  $\Phi$ , which satisfies Assumption 1.1 (see, e.g., [27, 16]). To this end, we simply add perturbations to  $\Phi$  to have the following optimization problem:

$$\begin{aligned} \min \quad & F(\mathbf{X}, \boldsymbol{\mathcal{E}}) = f(\mathbf{X}) + \rho \sum_{ij} \phi(|X_{ij}| + \mathcal{E}_{ij}) \\ \text{s.t.} \quad & \mathbf{X} > \mathbf{0}, \end{aligned} \tag{5}$$

where the perturbation parameters  $\mathcal{E}_{ij} = \mathcal{E}_{ji} > 0$ ,  $\forall (i, j) \in [n] \times [n]$ . By construction and Assumption 1.1, for each entry  $X_{ij}$ , it holds that

- (i) If  $0 \leq \bar{\mathcal{E}}_{ij} \leq \mathcal{E}_{ij}$ , then  $\phi(|X_{ij}| + \bar{\mathcal{E}}_{ij}) \leq \phi(|X_{ij}| + \mathcal{E}_{ij})$ .
- (ii) For any  $\mathcal{E}_{ij} \geq 0$ ,  $0 \leq \phi(|X_{ij}| + \mathcal{E}_{ij}) - \phi(|X_{ij}|) \leq \phi(\mathcal{E}_{ij})$ .

At the  $k$ th iterate  $\mathbf{X}^k \in \mathbb{S}_{++}^n$ , we compute an intermediate point  $\mathbf{Y}^k$  by approximately solving the weighted  $\ell_1$ -regularized log-determinant subproblem:

$$\mathbf{Y}^k \approx \operatorname{argmin}_{\mathbf{Y} \succ \mathbf{0}} Q_k(\mathbf{Y}) := f(\mathbf{Y}) + \rho \|\mathbf{Y}\|_{1, \mathbf{W}^k}, \quad (\mathcal{P}_{\text{sub}})$$

where  $W_{ij}^k = \mathbf{W}(X_{ij}^k, \mathcal{E}_{ij}^k) = \phi'(|X_{ij}^k| + \mathcal{E}_{ij}^k) > 0$ ,  $\forall (i, j)$ , and the concavity of  $\phi$  underpins this construction. We accept  $\mathbf{Y}^k$  once it satisfies the inexact termination criterion:

$$\operatorname{dist}(\mathbf{0}, \partial Q_k(\mathbf{Y}^k)) \leq \tilde{\beta}_k \|\mathbf{X}^k - \mathbf{Y}^k\|_F, \quad (\mathcal{C}_{\text{inexact}}^{(k)})$$

where the scalar  $\tilde{\beta}_k \in (0, 1]$  is tied to the  $k$ th subproblem's line-search step size (see Lemma 3.3(ii)). This inexact termination condition  $(\mathcal{C}_{\text{inexact}}^{(k)})$  provides several computational and theoretical benefits: (i) it can be efficiently evaluated within QUIC; (ii) thanks to the term  $\|\mathbf{X}^k - \mathbf{Y}^k\|$ , there is no need for  $\tilde{\beta}_k$  to shrink to zero. Indeed,  $\tilde{\beta}_k$  is uniformly bounded away from zero, as shown in Lemma 3.3(ii) below; (iii) this condition allows DIIR-QUIC to identify the smooth active manifold within finitely many iterations.

With  $\mathbf{Y}^k$  at hand, we compute the new iterate as  $\mathbf{X}^{k+1} = (1 - \alpha)\mathbf{X}^k + \alpha\mathbf{Y}^k$ , and update the perturbation as  $\mathcal{E}^{k+1} = (1 - \alpha)\mathcal{E}^k + \alpha(\mu\mathcal{E}^k)$ , with a damping parameter  $\alpha \in (0, 1)$  and constant  $\mu \in (0, 1)$ . Motivated by [17], we can describe the iterative scheme of DIIR-QUIC consistently with the following relaxed fixed-point iteration:

$$\begin{bmatrix} \mathbf{X}^{k+1} \\ \mathcal{E}^{k+1} \end{bmatrix} = \mathcal{T} \left( \begin{bmatrix} \mathbf{X}^k \\ \mathcal{E}^k \end{bmatrix} \right) = \begin{bmatrix} (1 - \alpha)\mathbf{X}^k + \alpha\mathcal{S}_{\mathbf{X}}(\mathbf{X}^k, \mathcal{E}^k) \\ (1 - \alpha)\mathcal{E}^k + \alpha\mathcal{S}_{\mathcal{E}}(\mathbf{X}^k, \mathcal{E}^k) \end{bmatrix}, \quad (6)$$

where  $\mathcal{T} : \mathbb{S}_{++}^n \times \mathbb{S}_{\geq 0}^n \rightarrow \mathbb{S}_{++}^n \times \mathbb{S}_{\geq 0}^n$  denotes a relaxed mapping,  $\mathcal{S}_{\mathbf{X}}(\mathbf{X}^k, \mathcal{E}^k) = \mathbf{Y}^k \stackrel{(\approx)}{\leftarrow} \operatorname{argmin}_{\mathbf{Y} \succ \mathbf{0}} Q_k(\mathbf{Y})$  and  $\mathcal{S}_{\mathcal{E}}(\mathbf{X}^k, \mathcal{E}^k) = \mu\mathcal{E}^k$ . In view of this, (6) yields a natural characterization of stationarity for problem  $(\mathcal{P})$ , as formalized in the proposition below.

**Proposition 2.1.** *A pair  $(\mathbf{X}^*, \mathcal{E}^*) \in \mathbb{S}_{++}^n \times \mathbb{S}_{\geq 0}^n$  is a stationary point of (5) if and only if  $(\mathbf{X}^*, \mathcal{E}^*) = \mathcal{T}(\mathbf{X}^*, \mathcal{E}^*)$ . In particular, a point  $\mathbf{X}^* \in \mathbb{S}_{++}^n$  is a stationary point of  $F$  associated with  $(\mathcal{P})$  if and only if  $\mathbf{X}^* = \mathcal{S}_{\mathbf{X}}(\mathbf{X}^*, \mathbf{0})$ .*

*Proof.* Note that  $(\mathbf{X}^*, \mathcal{E}^*) = \mathcal{T}(\mathbf{X}^*, \mathcal{E}^*)$  if and only if  $\mathcal{E}^* = \mathbf{0}$ . Thus, proving this equivalence reduces to establishing the second statement.

Suppose that  $\mathbf{X}^*$  is stationary for  $(\mathcal{P})$ . Then  $-\nabla f(\mathbf{X}^*) \in \partial\Phi(\mathbf{X}^*)$  by Theorem 1.4. Consider the problem

$$\min_{\mathbf{Y} \succ \mathbf{0}} \left\{ Q_*(\mathbf{Y}) = f(\mathbf{Y}) + \rho \sum_{i,j} \mathbf{W}_{ij}^* |\mathbf{Y}_{ij}| \right\},$$

where  $\mathbf{W}_{ij}^* = \mathbf{W}(X_{ij}^*, 0) = \phi'(|X_{ij}^*|)$ . Let  $\mathbf{Y}^* = \mathcal{S}_{\mathbf{X}}(\mathbf{X}^*, \mathbf{0})$  be an inexact solution satisfying

$$\text{dist}(\mathbf{0}, \partial Q_*(\mathbf{Y}^*)) \leq \tilde{\beta}_* \|\mathbf{X}^* - \mathbf{Y}^*\|_F,$$

where  $\tilde{\beta}_* \in (0, 1]$ . Since  $\mathbf{0} \in \partial Q_*(\mathbf{X}^*)$  and  $Q_*$  is strictly convex (see Lemma 3.1),  $\mathbf{X}^*$  is its unique minimizer. Initialized at  $\mathbf{X}^*$ , any solver satisfying the above inexactness bound should return  $\mathbf{Y}^* = \mathbf{X}^*$ , proving the forward direction.

Conversely, if  $\mathbf{X}^* = \mathcal{S}_{\mathbf{X}}(\mathbf{X}^*, \mathbf{0})$ , then by the inexact termination condition, we obtain  $\text{dist}(\mathbf{0}, \partial Q_*(\mathbf{X}^*)) = 0$ . This directly implies  $\mathbf{X}^*$  satisfies the optimality conditions (3), completing the proof.  $\square$

## 2.1 Modified QUIC for subproblem solution

As shown in Algorithm 1, the proposed DIIR-QUIC method consists of approximately solving a sequence of weighted  $\ell_1$ -norm regularized log-determinant problems ( $\mathcal{P}_{\text{sub}}$ ), a subject that has garnered substantial attention over the past decade. To efficiently solve each subproblem, we employ QUIC—a proximal Newton-type method known for its computational efficiency. To contextualize its use, we give a concise overview of QUIC.

The foundation of QUIC lies within the framework of inexact successive quadratic approximation methods [28], which sequentially solve subproblems constructed from a quadratic approximation of the smooth function  $f$  of  $Q_k$ . We use  $\{\mathbf{Z}^t\}$  to denote the iterate sequence generated by QUIC for solving ( $\mathcal{P}_{\text{sub}}$ ). At a given iterate  $\mathbf{Z}^t$ , QUIC first computes an approximate Newton direction  $\mathbf{D}^t$  using a locally quadratic approximation of  $f$  at  $\mathbf{Z}^t$ , that is,

$$\mathbf{D}^t \approx \underset{\mathbf{D} \in \mathbb{S}^n}{\text{argmin}} J(\mathbf{D}; \mathbf{Z}^t) = \bar{f}(\mathbf{D}; \mathbf{Z}^t) + g(\mathbf{D}; \mathbf{Z}^t), \quad (7)$$

where  $\bar{f}(\mathbf{D}; \mathbf{Z}^t) := f(\mathbf{Z}^t) + \langle \text{vec}(\nabla f(\mathbf{Z}^t)), \text{vec}(\mathbf{D}) \rangle + \frac{1}{2} \langle \text{vec}(\mathbf{D}), \nabla^2 f(\mathbf{Z}^t) \text{vec}(\mathbf{D}) \rangle$  and  $g(\mathbf{D}; \mathbf{Z}^t) = \rho \|\mathbf{Z}^t + \mathbf{D}\|_{1, \mathbf{W}^k}$ . A key insight of QUIC is the development of a closed-form solution for computing the Newton direction associated with (7), achieved through a coordinate descent update rule. Starting with an initial guess  $\Delta^{(0)}$ , the algorithm updates  $\Delta^{(t+1)}$  via

$$\Delta^{(t+1)} = \Delta^{(t)} + \eta(\mathbf{e}_i \mathbf{e}_j^T + \mathbf{e}_j \mathbf{e}_i^T), \quad (8)$$

where  $\Delta^{(t)}$  represents a given completed updates,  $\mathbf{e}_i, \mathbf{e}_j$  are unit vectors in  $\mathbb{R}^n$ , and  $\eta \in \mathbb{R}$  is the parameter to be computed. In QUIC, this is done by solving a sequence of single-variable minimization problems:

$$\min_{\eta \in \mathbb{R}} J(\Delta^{(t)} + \eta(\mathbf{e}_i \mathbf{e}_j^T + \mathbf{e}_j \mathbf{e}_i^T); \mathbf{Z}^t). \quad (9)$$

Iterating this process under appropriate termination criteria yields the desired solution  $\mathbf{D}^t$ .

Another key insight of QUIC is a dimension-reduction strategy that leverages the optimality condition of  $(\mathcal{P}_{\text{sub}})$  to compute the Newton direction using only a subset of elements of  $\mathbf{\Delta}^{(t)}$ . The indices updated during this process are

$$\mathcal{I}_{\text{free}}(\mathbf{Z}^t) := \{(i, j) \in [n] \times [n] \mid \mathbf{W}_{ij}^k < |\nabla_{ij} f(\mathbf{Z}^t)| \text{ or } \mathbf{Z}_{ij}^t \neq 0\}, \quad (10)$$

and the remaining indices are denoted as

$$\mathcal{I}_{\text{fixed}}(\mathbf{Z}^t) := \{(i, j) \in [n] \times [n] \mid \mathbf{W}_{ij}^k \geq |\nabla_{ij} f(\mathbf{Z}^t)| \text{ and } \mathbf{Z}_{ij}^t = 0\}. \quad (11)$$

Hence, QUIC in fact solves the following optimization problem for the Newton direction:

$$\mathbf{D}^t \approx \operatorname{argmin}_{\mathbf{D}: D_{ij}=0, \forall (i,j) \in \mathcal{I}_{\text{fixed}}(\mathbf{Z}^t)} J(\mathbf{D}; \mathbf{Z}^t) = \bar{f}(\mathbf{D}; \mathbf{Z}^t) + g(\mathbf{D}; \mathbf{Z}^t). \quad (12)$$

To solve the reduced subproblem (12) by coordinate descent, we propose to use the following inexact termination criterion. We accept the current update  $\mathbf{D}^t = \mathbf{\Delta}^{(t)}$  once

$$\|\mathbf{\Delta}^{(t)} - \mathbf{\Delta}^{(t-1)}\|_1 \leq \epsilon \|\mathbf{\Delta}^{(t)}\|_1, \quad (13)$$

where  $0 < \epsilon \ll 1$  is a predetermined tolerance parameter.

*Remark 2.2.* In the original QUIC paper, the convergence analysis—both global and local—relies on the assumption that each Newton direction  $\mathbf{D}^t$  is computed exactly by solving (12). However, as noted by the authors (see [5, §5.2.1]), its actual implementation employs an iterative coordinate descent method that only approximately solves the subproblem in practice. In contrast, our proposed inexact termination criterion (13) formalizes this practical inexactness and enables a rigorous convergence analysis. Specifically, we prove that the resulting algorithm retains both global and local convergence guarantees (see Lemma 3.4(i) and Corollary 3.10(ii) below), thereby strengthening the theoretical foundation of QUIC under inexact subproblem solves and aligning it more closely with practical implementations. Moreover, as shown in Lemma 3.4(ii) below, this inexact condition ensures that the obtained solution satisfies the inexactness criterion (18), which is crucial for establishing the active manifold property of QUIC when subproblems (7) are solved approximately.

After determining  $\mathbf{D}^t$ , a step size  $\beta_t \in (0, 1]$  is chosen to ensure both feasibility (i.e.,  $\mathbf{Z}^t + \beta_t \mathbf{D}^t > \mathbf{0}$ ) and sufficient objective decrease. To this end, QUIC employs an Armijo-type line-search to identify  $\beta_t$  as the largest value in the candidate set  $\{\pi^0, \pi^1, \dots\}$ , where  $\pi \in (0, 1)$ , satisfying:

$$Q_k(\mathbf{Z}^t + \beta_t \mathbf{D}^t) \leq Q_k(\mathbf{Z}^t) + \beta_t \sigma \Delta^t, \quad (14)$$

where  $\sigma \in (0, 0.5 - \epsilon_\sigma)$  for some small  $\epsilon_\sigma > 0$ , which depends on the chosen inexactness parameter  $\epsilon$  specified in (13), and

$$\Delta^t = \langle \nabla f(\mathbf{Z}^t), \mathbf{D}^t \rangle + \rho \|\mathbf{Z}^t + \mathbf{D}^t\|_{1, \mathbf{W}^k} - \rho \|\mathbf{Z}^t\|_{1, \mathbf{W}^k}. \quad (15)$$

The updated iterate is then computed as  $\mathbf{Z}^{t+1} = \mathbf{Z}^t + \beta_t \mathbf{D}^t$ . We summarize the computational steps of QUIC in Algorithm 2.

---

**Algorithm 2** The modified QUIC for solving  $(\mathcal{P}_{\text{sub}})$

---

**Require:**  $\mathbf{S}, \mathbf{W}^k, \mathbf{Z}^0 \leftarrow \mathbf{X}^k, \beta_t \leftarrow 1, \pi \in (0, 1)$  and  $t \leftarrow 0$ .  
**while**  $\text{dist}(\mathbf{0}, \partial Q_k(\mathbf{Z}^t)) \leq \beta_t \|\mathbf{X}^k - \mathbf{Z}^t\|$  (refer to  $(\mathbf{C}_{\text{inexact}}^{(k)})$ ) is not satisfied  
**do**  
    **while** (13) is not satisfied **do**  
        Partition the variables into free and fixed sets according to (10) and (11).  
        **for**  $(i, j) \in \mathcal{I}_{\text{free}}(\mathbf{Z}^t)$  **do**  
            Solve (12) for  $\mathbf{D}^t$  by coordinate descent update (8).  
        **end for**  
    **end while**  
    **while** True **do**  
        **if**  $\mathbf{Z}^t + \beta_t \mathbf{D}^t > \mathbf{0}$  **then**  
            **if** (14) holds **then**  
                break  
            **end if**  
        **end if**  
         $\beta_t \leftarrow \pi \beta_t$   
    **end while**  
     $\mathbf{Z}^{t+1} \leftarrow \mathbf{Z}^t + \beta_t \mathbf{D}^t$  and  $t \leftarrow t + 1$ .  
**end while**  
**Return**  $\mathbf{Y}^k \leftarrow \mathbf{Z}^{t+1}$ .

---

### 3 Convergence

In this section, we show the convergence properties of the proposed Algorithm 1.

### 3.1 Well-posedness and basic properties of subproblems

We start by demonstrating that each subproblem  $(\mathcal{P}_{\text{sub}})$  is well-defined, and also present some useful properties of the subproblems solved by Algorithm 2.

**Lemma 3.1** (Well-posed subproblem). *Let Assumption 1.1 hold. Each subproblem  $(\mathcal{P}_{\text{sub}})$  is well posed in the sense that it admits a unique optimal solution whose eigenvalues are uniformly bounded away from zero and infinity.*

*Proof.* It follows from  $\mathcal{E}_{ij}^k > 0$  that  $W_{ij}^k > 0$ ,  $\forall k, \forall (i, j)$ ; hence the proof mainly follows [29, Theorem 1].  $\square$

**Proposition 3.2.** *Consider the subproblem  $(\mathcal{P}_{\text{sub}})$ . Define the sublevel set of  $Q_k$  with respect to the initialization  $\mathbf{X}^k \in \mathbb{S}_{++}^n$  as  $\text{Lev}_{Q_k}(\mathbf{X}^k) = \{\mathbf{Z} \succ \mathbf{0} \mid Q_k(\mathbf{Z}) \leq Q_k(\mathbf{X}^k)\}$ . Then the following statements hold.*

(i) *For each  $k$ , there exist constants  $\underline{c} > 0$  and  $\bar{c} < +\infty$  such that  $\underline{c} \leq l_k < u_k \leq \bar{c}$ , and the level set  $\text{Lev}_{Q_k}(\mathbf{X}^k) \subset \chi_k = \{\mathbf{Z} \succ \mathbf{0} \mid l_k \mathbf{I} \leq \mathbf{Z} \leq u_k \mathbf{I}\}$ .*

(ii) *It holds that*

$$\|\nabla f(\mathbf{X}) - \nabla f(\mathbf{Y})\| \leq l_k^{-2} \|\mathbf{X} - \mathbf{Y}\|, \quad \forall \mathbf{X}, \mathbf{Y} \in \chi_k. \quad (16)$$

*In addition,  $\nabla^2 f(\mathbf{X}) = \mathbf{X}^{-1} \otimes \mathbf{X}^{-1}$  and*

$$u_k^{-2} \mathbf{I} \leq \nabla^2 f(\mathbf{X}) \leq l_k^{-2} \mathbf{I}, \quad \forall \mathbf{X} \in \chi_k, \quad (17)$$

*where  $0 < l_k < u_k < +\infty$  are defined in statement (i).*

*Proof.* Statement (i) follows from the result in [5, Lemma 2]. Statement (ii) follows from statement (i) and [7, Lemma 6(iii)].  $\square$

We next show that the termination condition (13) implies a theoretically useful inexactness criterion, and its proof can be found in Appendix A.1.

**Lemma 3.3.** *Consider the subproblem  $(\mathcal{P}_{\text{sub}})$ . Let  $\{\mathbf{Z}^t\}_{t \geq 0}$  be the sequence generated by Algorithm 2, and let  $\{\Delta^{(t)}\}$  be the sequence generated by the update rule (8) for solving the problem associated with (12). Furthermore, let  $\mathbf{D}^{(t)}$  be an approximate solution to (7) satisfying the inexact condition (13). Then the following statements hold.*

(i) *Suppose that the initial point  $\Delta^{(0)} = \mathbf{0}$ . Then the finite quantity  $C_k := \epsilon l_k^{-2}$  satisfies the following implication for each iteration  $t$ :*

$$\|\Delta^{(t)} - \Delta^{(t-1)}\|_1 \leq \epsilon \|\Delta^{(t)}\|_1 \implies \text{dist}(\mathbf{0}, \partial J(\Delta^{(t)}; \mathbf{Z}^{(t)})) \leq C_k \|\Delta^{(t)}\|_1. \quad (18)$$

- (ii) Under (13). Suppose that  $u_k^{-2} - \epsilon l_k^{-2} > 0$ . Then there exists a constant  $\hat{\beta} > 0$  such that for all  $\beta_t \geq \hat{\beta}$ ,  $\mathbf{Z}^{t+1} = \mathbf{Z}^t + \beta_t \mathbf{D}^t$  satisfies

$$Q_k(\mathbf{Z}^t) - Q_k(\mathbf{Z}^{t+1}) \geq (\beta_t (u_k^{-2} - \epsilon l_k^{-2}) - (2l_k^2)^{-1} \beta_t^2) \|\mathbf{Z}^t - \mathbf{Z}^{t+1}\|_F^2. \quad (19)$$

Furthermore, the backtracking line-search procedure in Algorithm 2 terminates in a finite number of iterations and yields a step size  $\beta_t$  that admits the following lower bound:

$$\beta_t \geq \hat{\beta} := \min \{1, 2\pi(1 - \sigma)(u_k^{-2} - \epsilon l_k^{-2})l_k^2\}. \quad (20)$$

Consequently, it holds that

$$\lim_{t \rightarrow +\infty} \|\mathbf{Z}^{t+1} - \mathbf{Z}^t\|_F^2 = 0. \quad (21)$$

- (iii) Suppose that Algorithm 2 is run with the condition (13). Then there exist  $\bar{\epsilon}, \bar{\delta} > 0$  such that the line-search condition (14) will be satisfied with unit stepsize (i.e.,  $\beta_t = 1$ ) whenever  $\|\mathbf{Z}^t - \mathbf{Z}^*\| \leq \bar{\delta}$  and  $C_k \|\Delta^{(t)}\|_1 \leq \bar{\epsilon}$ .

Next, we show the global subsequential convergence and the manifold identification property of Algorithm 2 under the inexact condition (13).

**Lemma 3.4.** Consider the subproblem  $(\mathcal{P}_{\text{sub}})$ . Let  $\{\mathbf{Z}^t\}$  be the sequence generated by Algorithm 2, converging to  $\mathbf{Z}^*$ . The following statements hold.

- (i) (**Global subsequential convergence of QUIC under condition (13)**) Suppose that  $u_k^{-2} - \epsilon l_k^{-2} > 0$ . Then any cluster point of  $\{\mathbf{Z}^t\}$  is a stationary point of  $(\mathcal{P}_{\text{sub}})$ .
- (ii) (**Manifold identification of QUIC under condition (13)**) Given that the regularizer  $\|\cdot\|_{1, \mathbf{W}^*}$  is partly smooth at  $\mathbf{Z}^*$  relative to the manifold  $\mathcal{M}_{\text{QUIC}}(\mathbf{Z}^*) = \{\mathbf{Z} \in \chi \mid \mathcal{I}(\mathbf{Z}) \subseteq \mathcal{I}(\mathbf{Z}^*)\}$ , suppose in addition that the nondegenerate condition

$$\mathbf{0} \in \text{rint}(\partial Q_k(\mathbf{Z}^*)) \quad (22)$$

holds at  $\mathbf{Z}^*$ , and Algorithm 2 is executed under the inexactness criterion (13). Then

- (1.)  $Q_k(\mathbf{Z}^t) \rightarrow Q_k(\mathbf{Z}^*)$ .
- (2.)  $\text{dist}(\mathbf{0}, \partial Q_k(\mathbf{Z}^t)) \rightarrow 0 \implies \mathbf{Z}^t \in \mathcal{M}_{\text{QUIC}}(\mathbf{Z}^*)$  for all sufficiently large  $t$ .



Consequently, there exist constants  $\bar{\epsilon}, \bar{\delta} > 0$  such that, for sufficiently large  $t$ , the following condition holds:

$$\|\mathbf{Z}^t - \mathbf{Z}^*\| \leq \bar{\delta}, C_k \|\boldsymbol{\Delta}^{(t)}\|_1 \leq \bar{\epsilon}, \beta_t = 1 \implies \mathbf{Z}^{t+1} \in \mathcal{M}_{QUIC}(\mathbf{Z}^*).$$

*Proof.* The proof of statement (i) primarily follows from Proposition 3.2 and [28, Theorem 3.3] (applied within a monotone line-search framework (LS<sub>1</sub>) while omitting the additional regularization function (i.e., setting  $g = 0$  in their formulation)). As for statement (ii), it follows from Lemma 3.3 and [30, Lemmas 1-2 & Theorem 1].  $\square$

## 3.2 Global convergence properties

This section analyzes the convergence properties of the proposed Algorithm 1. We first show the model reduction on  $F(\mathbf{X}^k, \boldsymbol{\epsilon}^k)$  caused by  $(\mathbf{X}^{k+1}, \boldsymbol{\epsilon}^{k+1})$  with inexact subproblem solution. We now consider the outer loop of Algorithm 1. In the sequel, we define

$$l = \inf_{k \geq 0} l_k \text{ and } u = \sup_{k \geq 0} u_k, \quad (23)$$

where  $0 < \underline{c} \leq l_k < u_k \leq \bar{c} < +\infty, \forall k \in \mathbb{N}$  is stated in Proposition 3.2. Correspondingly, we define  $\chi = \{\mathbf{Z} > \mathbf{0} \mid l\mathbf{I} \leq \mathbf{Z} \leq u\mathbf{I}\}$ .

We first prove some useful results in the following proposition, and its proof can be found in Appendix A.2.

**Proposition 3.5** (Approximate sufficient descent property). *Let Assumption 1.1 hold and let  $\{(\mathbf{X}^k, \boldsymbol{\epsilon}^k)\}$  be the sequence generated by Algorithm 1 with  $\mathbf{X}^0 \in \chi$ . The following assertions hold.*

(i) Suppose  $l^2 \geq \left[ \left( 2\sqrt{\pi(1-\sigma)} \right)^{-1} + \epsilon \right] u^2$  under condition (13). The sequence  $\{F(\mathbf{X}^k, \boldsymbol{\epsilon}^k)\}$  is monotonically decreasing. Indeed, there exists some constant  $\tilde{C} > 0$  such that

$$F(\mathbf{X}^k, \boldsymbol{\epsilon}^k) - F(\mathbf{X}^{k+1}, \boldsymbol{\epsilon}^{k+1}) \geq \tilde{C} \|\mathbf{X}^k - \mathbf{X}^{k+1}\|_F^2, \quad (24)$$

where  $\tilde{C} := (\alpha^{-1} \tilde{\beta}(u^{-2} - \epsilon l^{-2}) - \frac{1}{2l^2\alpha})$  with  $\tilde{\beta} := \min \{1, 2\pi(1-\sigma)(u^{-2} - \epsilon l^{-2})l^2\}$ . Consequently,  $\lim_{k \rightarrow +\infty} F(\mathbf{X}^k, \boldsymbol{\epsilon}^k)$  exists.

(ii)  $\sum_{k=0}^{+\infty} \|\mathbf{X}^k - \mathbf{X}^{k+1}\|_F^2 < +\infty$ , indicating  $\lim_{k \rightarrow +\infty} \|\mathbf{X}^k - \mathbf{X}^{k+1}\|_F^2 = 0$  and  $\lim_{k \rightarrow +\infty} \|\mathbf{X}^k - \mathbf{Y}^k\|_F^2 = 0$ .

- (iii) There exist constants  $0 < l < u < +\infty$  such that the sequences  $\{\mathbf{X}^k\}$  and  $\{\mathbf{Y}^k\}$  are contained in a convex cone  $\chi = \{\mathbf{X} \mid l\mathbf{I} \leq \mathbf{X} \leq u\mathbf{I}\}$ . Additionally, the sequences  $\{\mathbf{X}^k\}$  and  $\{\mathbf{Y}^k\}$  are bounded. As a consequence, there exists a constant  $C_{\max} > 0$  such that  $\|\nabla f(\mathbf{X}^k)\|_{\infty} = \max_{ij} |\nabla_{ij} f(\mathbf{X}^k)| < C_{\max}$  for all  $k$ .

We next establish the global convergence of the proposed DIIR-QUIC.

**Theorem 3.6** (Global subsequential convergence). *Let Assumption 1.1 hold and let  $\{(\mathbf{X}^k, \boldsymbol{\varepsilon}^k)\}$  be the sequence generated by Algorithm 1. Then the following statements holds.*

- (i) The sequence  $\{\boldsymbol{\varepsilon}^k\}$  converges to  $\mathbf{0}$ .
- (ii) The set of cluster points  $\Omega^{\infty}$  of  $\{\mathbf{X}^k\}$  is nonempty, compact, connected, and  $\text{dist}(\mathbf{X}^k, \Omega^{\infty}) \rightarrow 0$ .
- (iii) Any cluster point  $\mathbf{X}^* \in \Omega^{\infty}$  is a stationary point of  $F$ .
- (iv) The sequence  $\{F(\mathbf{X}^k)\}$  is convergent. Moreover, the objective function  $F$  is constant on  $\Omega^{\infty}$ , and hence  $\xi := F(\mathbf{X}^*) = \lim_{k \rightarrow +\infty} F(\mathbf{X}^k, \boldsymbol{\varepsilon}^k)$ .

*Proof.* (i) This statement directly follows from the definition of Algorithm 1.

(ii) Proposition 3.5(iii) implies that  $\Omega^{\infty}$  is nonempty. On the other hand, Proposition 3.5(ii)-(iii) leads to the desired results by combining [31, Corollary 2.7 (Ostrowski)]. This establishes statement (ii).

(iii) Let  $\mathbf{X}^*$  be a cluster point of  $\{\mathbf{X}^k\}$ . Consider a subsequence  $\mathcal{S}$  such that  $\{\mathbf{X}^k\} \xrightarrow{\mathcal{S}} \mathbf{X}^*$ . It follows from  $\mathbf{Y}^k = \frac{1}{\alpha}(\mathbf{X}^{k+1} - \mathbf{X}^k) + \mathbf{X}^k$  that  $\{\mathbf{Y}^k\} \xrightarrow{\mathcal{S}} \mathbf{X}^*$  by Proposition 3.5(ii), which further implies that  $\text{sgn}(Y_{ij}^k) = \text{sgn}(X_{ij}^*), \forall (i, j) \in \mathcal{I}(\mathbf{X}^*)$ . Let  $\mathcal{R}^k \in \partial Q_k(\mathbf{Y}^k)$  be the minimum-norm subgradient [5, Definition 6] defined as follows:

$$\mathcal{R}_{ij}^k = \begin{cases} \nabla_{ij} f(\mathbf{Y}^k) + \rho W_{ij}^k, & \text{if } Y_{ij}^k > 0, \\ \nabla_{ij} f(\mathbf{Y}^k) - \rho W_{ij}^k, & \text{if } Y_{ij}^k < 0, \\ \text{sgn}(\nabla_{ij} f(\mathbf{Y}^k)) \max(|\nabla_{ij} f(\mathbf{Y}^k)| - \rho W_{ij}^k, 0), & \text{if } Y_{ij}^k = 0. \end{cases} \quad (25)$$

We then know from the first-order optimality condition of  $(\mathcal{P}_{\text{sub}})$  that

$$\mathcal{R}_{ij}^k \in \nabla_{ij} f(\mathbf{Y}^k) + \rho \mathbf{W}(X_{ij}^k, \mathcal{E}_{ij}^k) \partial |Y_{ij}^k|, \quad \forall (i, j) \in [n] \times [n]. \quad (26)$$

Consider first that  $(i, j) \in \mathcal{I}(\mathbf{X}^*)$ . By [21, Proposition 8.7] and [32, Proposition 2.1.5], we know that, for sufficiently large  $k \in \mathbb{N}$ , there exists  $\xi_{ij}^* \in \partial |X_{ij}^*|$  such that

$$0 = \nabla_{ij} f(\mathbf{X}^*) + \rho \mathbf{W}(X_{ij}^*, 0) \xi_{ij}^* \quad (27)$$

holds since  $f \in \mathcal{C}^1$ ,  $Y_{ij}^k \xrightarrow{\mathcal{S}} X_{ij}^*$ ,  $\mathcal{E}_{ij}^k \xrightarrow{\mathcal{S}} 0$  and  $\mathcal{R}_{ij}^k \xrightarrow{\mathcal{S}} 0$ ,  $\forall (i, j) \in \mathcal{I}(\mathbf{X}^*)$ . Thus, (3a) is satisfied at  $X_{ij}^*$ ,  $\forall (i, j) \in \mathcal{I}(\mathbf{X}^*)$ . On the other hand, for  $(i, j) \in \mathcal{Z}(\mathbf{X}^*)$  and  $\lim_{t \rightarrow 0^+} \phi'(t) < +\infty$ , we have from (26) that

$$0 \in \nabla_{ij} f(\mathbf{X}^*) + \rho \mathbf{W}(0, 0) \partial|0| \iff |\nabla_{ij} f(\mathbf{X}^*)| \leq \rho \phi'(0). \quad (28)$$

Hence,  $\mathbf{X}^* \in \Omega^\infty$  is stationary for  $F$ .

(iv) The convergence of  $\{F(\mathbf{X}^k)\}$  is a direct consequence of Proposition 3.5(i)-(ii) and statement (i) in this theorem. On the other hand,

$$F(\mathbf{X}^*, \mathcal{E}^*) = \lim_{k \rightarrow +\infty} f(\mathbf{X}^k) + \tilde{\Phi}(\mathbf{X}^k, \mathcal{E}^k) \stackrel{(a)}{=} f(\mathbf{X}^*) + \Phi(\mathbf{X}^*) = F(\mathbf{X}^*), \quad (29)$$

where (a) holds by Assumption 1.1 and statement (iii) and  $f \in \mathcal{C}^1$ . The proof is complete.  $\square$

The following non-degenerate condition is assumed to hold throughout.

**Assumption 3.7.** *Under Assumption 1.1, let  $\mathbf{X}^* \in \mathbb{S}_{++}^n$  be any stationary point of problem  $(\mathcal{P})$ . We assume in addition that*

(i) **(Nondegenerate condition):**

$$\mathbf{0} \in \text{rint}(\partial F(\mathbf{X}^*)) = \nabla f(\mathbf{X}^*) + \text{rint}(\partial \Phi(\mathbf{X}^*)). \quad (30)$$

(ii) **(No active-kink condition):** *Let  $\mathcal{K} = \{t > 0 \mid \phi' \text{ is not } \mathcal{C}^1 \text{ at } t\}$ . Then for each  $(i, j)$  with  $X_{ij}^* \neq 0$ ,  $|X_{ij}^*| \notin \mathcal{K}$ .*

The following lemma establishes the conditions to guarantee manifold identification in our setting, and its proof can be found in Appendix A.3.

**Lemma 3.8** ( $F$  in  $(\mathcal{P})$  admits an active manifold). *Let Assumptions 1.1-3.7 hold and let  $\{(\mathbf{X}^k, \mathcal{E}^k)\}$  be the sequence generated by Algorithm 1 converging to  $(\mathbf{X}^*, \mathcal{E}^*)$ . The following statements hold.*

(i) *Function  $\Phi$  is partly smooth at the point  $\mathbf{X}^*$  relative to the manifold  $\mathcal{M}(\mathbf{X}^*) := \{\mathbf{X} \in \mathbb{S}_{++}^n \mid \text{sgn}(\mathbf{X}_{\mathcal{I}(\mathbf{X}^*)}) = \text{sgn}(\mathbf{X}_{\mathcal{I}(\mathbf{X}^*)}^*), \mathbf{X}_{\mathcal{Z}(\mathbf{X}^*)} = \mathbf{0}\}$ , and prox-regular there. Consequently,  $F$  is also partly smooth and prox-regular at  $\mathbf{X}^*$  relative to  $\mathcal{M}(\mathbf{X}^*)$ .*

(ii) *The following implication holds:*

$$\text{dist}(\mathbf{0}, \partial F(\mathbf{X}^k)) \rightarrow 0 \iff \mathbf{X}^k \in \mathcal{M}(\mathbf{X}^*) \text{ for all sufficiently large } k. \quad (31)$$

Next, we claim that the proposed DIIR-QUIC identifies the active manifold in a finite number of iterations.

**Theorem 3.9** (Finite identification of DIIR-QUIC). *Let Assumptions 1.1-3.7 hold. Consider a point  $\mathbf{X}^*$  that satisfies the non-degeneracy condition (30). If Algorithm 1 is executed with the inexact stopping conditions  $(\mathbf{C}_{\text{inexact}}^{(k)})$ , then there exist constants  $\tilde{\tau}, \tilde{\kappa}, \tilde{\epsilon} > 0$  such that for any  $\|\mathbf{X}^k - \mathbf{X}^*\| \leq \tilde{\tau}$ ,  $\|\mathbf{X}^k - \mathbf{Y}^k\| \leq \tilde{\kappa}$ , and  $\|\boldsymbol{\mathcal{E}}^k\| \leq \tilde{\epsilon}$ , it follows that  $\mathbf{X}^{k+1} \in \mathcal{M}(\mathbf{X}^*)$ .*

*Proof.* We prove this by contradiction. Suppose the conclusion is false; then there exists a subsequence  $\mathcal{S} \subset \mathbb{N}$  such that  $\mathbf{X}^k \rightarrow \mathbf{X}^*$ ,  $\boldsymbol{\mathcal{E}}^k \rightarrow \mathbf{0}$ , and  $\mathbf{X}^k - \mathbf{Y}^k \rightarrow \mathbf{0}$  as  $k \in \mathcal{S} \rightarrow \infty$ , yet  $\mathbf{X}^{k+1} \notin \mathcal{M}$  for all  $k \in \mathcal{S}$ .

By Proposition 3.5(ii) and standard variational analysis results [21, Proposition 8.7], [32, Proposition 2.1.5], we obtain from (25) and (26) that for sufficiently large  $k \in \mathcal{S}$ , it holds that

$$\left( \mathcal{R}_{ij}^k - \nabla_{ij} f(\mathbf{Y}^k) \right) \in \rho W(X_{ij}^k, \mathcal{E}_{ij}^k) \partial |X_{ij}^k|. \quad (32)$$

Then, we derive

$$\begin{aligned} \text{dist}(\mathbf{0}, \partial F(\mathbf{X}^k, \boldsymbol{\mathcal{E}}^k)) &= \text{dist}(-\nabla f(\mathbf{X}^k), \rho \partial \Phi(\mathbf{X}^k, \boldsymbol{\mathcal{E}}^k)) \\ &= \inf_{\mathcal{L}^k \in \rho \partial \Phi(\mathbf{X}^k, \boldsymbol{\mathcal{E}}^k)} \| -\nabla f(\mathbf{X}^k) - \mathcal{L}^k \|_F \\ &\stackrel{\text{Eq. (32)}}{\leq} \| -\nabla f(\mathbf{X}^k) - (\mathcal{R}^k - \nabla f(\mathbf{Y}^k)) \|_F \\ &\leq \| \nabla f(\mathbf{X}^k) - \nabla f(\mathbf{Y}^k) \|_F + \|\mathcal{R}^k\|_F \leq l^{-2} \|\mathbf{X}^k - \mathbf{Y}^k\|_F + \|\mathcal{R}^k\|_F. \end{aligned} \quad (33)$$

On the other hand, we know that

$$\begin{aligned} &\text{dist}(\partial F(\mathbf{X}^k), \partial F(\mathbf{X}^k, \boldsymbol{\mathcal{E}}^k)) = \text{dist}(\rho \partial \Phi(\mathbf{X}^k), \rho \partial \Phi(\mathbf{X}^k, \boldsymbol{\mathcal{E}}^k)) \\ &\leq \rho \sum_{(i,j) \in \mathcal{I}(\mathbf{X}^k)} \text{dist}(\phi'(|X_{ij}^k|), \phi'(|X_{ij}^k| + \mathcal{E}_{ij}^k)) \\ &\quad + \rho \sum_{(i,j) \in \mathcal{Z}(\mathbf{X}^k)} \text{dist}(\partial \phi(0), \phi'(\mathcal{E}_{ij}^k)) \\ &\leq \rho \sum_{(i,j) \in \mathcal{I}^k} |\phi'(|X_{ij}^k|) - \phi'(|X_{ij}^k| + \mathcal{E}_{ij}^k)| + \rho \sum_{(i,j) \in \mathcal{Z}(\mathbf{X}^k)} \inf_{\xi \in \partial \phi(0)} |\xi - \phi'(\mathcal{E}_{ij}^k)| \quad (34) \\ &\leq \rho \sum_{(i,j) \in \mathcal{I}^k} L_{\phi'} \mathcal{E}_{ij}^k + \rho \sum_{(i,j) \in \mathcal{Z}(\mathbf{X}^k)} \inf_{\xi \in \partial \phi(0)} |\xi - \phi'(\mathcal{E}_{ij}^k)| \\ &= \rho L_{\phi'} \|\boldsymbol{\mathcal{E}}_{\mathcal{I}^k}^k\|_1 + \rho \sum_{(i,j) \in \mathcal{Z}(\mathbf{X}^k)} \inf_{\xi \in \partial \phi(0)} |\xi - \phi'(\mathcal{E}_{ij}^k)|. \end{aligned}$$

Therefore, combining the above bounds, we obtain

$$\begin{aligned}
& \text{dist}(\mathbf{0}, \partial F(\mathbf{X}^k)) \\
& \stackrel{\text{Eq. (1)}}{\leq} \text{dist}(\mathbf{0}, \partial F(\mathbf{X}^k, \boldsymbol{\mathcal{E}}^k)) + \sup_{\mathcal{U}^k \in \partial F(\mathbf{X}^k, \boldsymbol{\mathcal{E}}^k)} \inf_{\mathcal{V}^k \in \partial F(\mathbf{X}^k)} \|\mathcal{U}^k - \mathcal{V}^k\|_F \\
& \stackrel{\text{Eq. (33)}}{\leq} l^{-2} \|\mathbf{X}^k - \mathbf{Y}^k\|_F + \|\mathcal{R}^k\|_F + \sup_{\mathcal{U}^k \in \partial F(\mathbf{X}^k, \boldsymbol{\mathcal{E}}^k)} \inf_{\mathcal{V}^k \in \partial F(\mathbf{X}^k)} \|\mathcal{U}^k - \mathcal{V}^k\|_F \\
& \stackrel{\text{Eq. (34)}}{\leq} l^{-2} \|\mathbf{X}^k - \mathbf{Y}^k\|_F + \|\mathcal{R}^k\|_F + \rho L_{\phi'} \|\boldsymbol{\mathcal{E}}_{\mathcal{I}^k}^k\|_1 + \rho \sum_{(i,j) \in \mathcal{I}(\mathbf{X}^k)} \inf_{\xi \in \partial \phi(0)} |\xi - \phi'(\boldsymbol{\mathcal{E}}_{ij}^k)|.
\end{aligned}$$

Since  $\|\mathbf{X}^k - \mathbf{Y}^k\| \rightarrow 0$  and  $\|\mathcal{R}^k\| \rightarrow 0$  by Proposition 3.5(ii), and since Theorem 3.6(i) ensures that  $\|\boldsymbol{\mathcal{E}}_{\mathcal{I}(\mathbf{X}^k)}^k\|_1 \rightarrow 0$  and  $\inf_{\xi \in \partial \phi(0)} |\xi - \phi'(\boldsymbol{\mathcal{E}}_{ij}^k)| \rightarrow 0$ , we conclude that  $\text{dist}(\mathbf{0}, \partial F(\mathbf{X}^k)) \rightarrow 0$ , which indicates that  $\mathbf{X}^{k+1} \in \mathcal{M}$  for all sufficiently large  $k$  by Lemma 3.8. This contradicts our assumption, and therefore, the conclusion holds.  $\square$

**Corollary 3.10.** *Let Assumptions 1.1-3.7 hold. Let  $\{(\mathbf{X}^k, \boldsymbol{\mathcal{E}}^k)\}$  be the sequence generated by Algorithm 1. The following statements hold.*

- (i) (**Uniform lower bounds for nonzeros**) *For any  $(i, j) \in \mathcal{I}(\mathbf{X}^*)$ , there exists a  $\delta_k > 0$  and an index  $K \in \mathbb{N}$  such that  $X_{ij}^k \geq \delta_k > 0$  for all  $k > K$ . Indeed, it holds that*

$$|X_{ij}^k| \geq \delta_k := (\phi')^{-1} \left( \frac{C_{\max} + \mathcal{R}_{ij}^k}{\rho} \right) - \mathcal{E}_{ij}^k > 0, \forall k > K, \quad (35)$$

where  $\mathcal{R}_{ij}^k$  is defined in (25). Or equivalently, the corresponding weights for any  $(i, j) \in \mathcal{I}(\mathbf{X}^*)$  are bounded, i.e.,

$$0 < W_{ij}^k \leq \frac{C_{\max} + \mathcal{R}_{ij}^k}{\rho}, \forall k > K.$$

As a consequence, it holds for any  $(i, j) \in \mathcal{I}(\mathbf{X}^*)$  that

$$|X_{ij}^*| \geq (\phi')^{-1} \left( \frac{C_{\max}}{\rho} \right) > 0. \quad (36)$$

- (ii) (**Accelerating subproblem solution**) *There exists an index  $K \in \mathbb{N}$  such that for all  $k > K$ , Algorithm 2 solves the following optimization problem:*

$$\min_{\mathbf{Y} > \mathbf{0}} f(\mathbf{Y}) + \rho \sum_{(i,j) \in \mathcal{I}(\mathbf{X}^k)} \mathbf{W}(X_{ij}^k, 0) \phi(\text{sgn}(\mathbf{X}^k) Y_{ij}), \quad (37)$$

which has the same optimum of  $(\mathcal{P}_{\text{sub}})$ . Consequently, QUIC is equivalent to the pure Newton method.

*Proof.* (i) Let  $(i, j) \in \mathcal{I}(\mathbf{X}^*)$  and  $K \in \mathbb{N}$  be the index such that Theorem 3.9 holds. We first claim that  $\mathbf{W}(X_{ij}^k, \mathcal{E}_{ij}^k) \leq \frac{C_{\max} + \mathcal{R}_{ij}^k}{\rho}$  for all  $k > K$  with  $(i, j) \in \mathcal{I}(\mathbf{X}^*)$ . Seeking a contradiction, assume that  $\mathbf{W}(X_{ij}^k, \mathcal{E}_{ij}^k) > \frac{C_{\max} + \mathcal{R}_{ij}^k}{\rho}$  for some  $k > K$  with  $(i, j) \in \mathcal{I}(\mathbf{X}^*)$ . Now consider the first-order optimality condition of the  $k$ th subproblem  $(\mathcal{P}_{\text{sub}})$ . By (26), it holds that

$$|\nabla_{ij} f(\mathbf{Y}^k)| = |\rho \mathbf{W}(X_{ij}^k, \mathcal{E}_{ij}^k) - \mathcal{R}_{ij}^k| > \left| \rho \frac{C_{\max} + \mathcal{R}_{ij}^k}{\rho} - \mathcal{R}_{ij}^k \right| = C_{\max}, \quad (38)$$

which contradicts Proposition 3.5(iii). Therefore, we know that for any  $(i, j) \in \mathcal{I}(\mathbf{X}^*)$ ,  $\mathbf{W}(X_{ij}^k, \mathcal{E}_{ij}^k) \leq \frac{C_{\max} + \mathcal{R}_{ij}^k}{\rho}$ , which is equivalent to (35). On the other hand, since  $\mathcal{R}^k \in \partial Q_k(\mathbf{Y}^k)$ , it follows from  $(\mathbf{C}_{\text{inexact}}^{(k)})$  and Proposition 3.5(ii) that  $\mathcal{R}_{ij}^k \rightarrow 0$ . Moreover, Theorem 3.6(i)-(ii) and Theorem 3.9 guarantees that  $\mathcal{E}_{ij}^k \rightarrow 0$ . Hence, continuity of  $(\phi')^{-1}$  indicates that (35) implies (36). Moreover, Theorem 3.9, Lemma 3.3(iii), Proposition 3.5(ii), together with the arguments in [5, Theorem 16], establish statement (ii). This completes the proof.  $\square$

### 3.3 Analysis Under the KL Property

In this subsection, assuming the KL property and drawing on techniques from [33], we demonstrate that the sequence of relaxed objective values  $\{F(\mathbf{X}^k, \mathcal{E}^k)\}$  converges at a Q-linear rate. Moreover, we establish the global convergence of the iterates  $\{\mathbf{X}^k\}$  and prove that they converge at an R-linear rate.

**Theorem 3.11.** *Suppose that Assumptions 1.1-3.7 hold and that  $F(\mathbf{X}, \mathcal{E})$  is a KL function with exponent  $1/2$ . Let  $\{(\mathbf{X}^k, \mathcal{E}^k)\}$  be the sequence generated by Algorithm 1 converging to some limit point  $(\mathbf{X}^*, \mathbf{0}) \in \Gamma^\infty$ . Then  $\{F(\mathbf{X}^k, \mathcal{E}^k)\}$  converges Q-linearly to  $F(\mathbf{X}^*, \mathbf{0})$*

*Proof.* We first note that if there exists an index  $k_0 \in \mathbb{N}$  such that  $F(\mathbf{X}^{k_0}, \mathcal{E}^{k_0}) = F(\mathbf{X}^{k_0+1}, \mathcal{E}^{k_0+1})$ , then Proposition 3.5(i) implies  $\mathbf{X}^{k_0} = \mathbf{X}^{k_0+1}$ , and hence  $\mathcal{E}^{k_0} = \mathcal{E}^{k_0+1}$ . In view of Proposition 2.1, the sequence  $\{\mathbf{X}^k\}$  then converges to a stationary point in a finite number of iterations. Hence, we assume that  $F(\mathbf{X}^k, \mathcal{E}^k) > F(\mathbf{X}^{k+1}, \mathcal{E}^{k+1})$  for all  $k \in \mathbb{N}$  in the subsequent proof.

Suppose that  $F(\mathbf{X}, \mathcal{E})$  satisfies the KL property with exponent  $1/2$ . Then, by [25, Lemma 6] (in conjunction with Theorem 3.6(ii) and (iv)), there exist constants

$\varepsilon > 0$  and  $\vartheta > 0$  such that the following holds: for every cluster point  $\hat{\mathbf{X}} \in \Omega^\infty$  and for every  $(\mathbf{X}, \boldsymbol{\varepsilon}) \in \mathbb{S}_{++}^n \times \mathbb{S}_{>0}^n$  satisfying

$$\text{dist}((\mathbf{X}, \boldsymbol{\varepsilon}), \Gamma^\infty) < \varepsilon \text{ and } F(\hat{\mathbf{X}}, \mathbf{0}) < F(\mathbf{X}, \boldsymbol{\varepsilon}) < F(\hat{\mathbf{X}}, \mathbf{0}) + \vartheta,$$

the inequality

$$\varphi'(F(\mathbf{X}, \boldsymbol{\varepsilon}) - F(\hat{\mathbf{X}}, \mathbf{0})) \cdot \text{dist}(\mathbf{0}, \partial F(\mathbf{X}, \boldsymbol{\varepsilon})) \geq 1 \quad (39)$$

holds. Here, the desingularizing function is defined by  $\varphi(t) = ct^{1/2}$  with  $c > 0$ , so that its derivative is given by  $\varphi'(t) = \frac{c}{2}t^{-1/2}$ .

Let  $(\mathbf{X}^*, \mathbf{0})$  be a cluster point of the sequence  $\{(\mathbf{X}^k, \boldsymbol{\varepsilon}^k)\}$ . By Theorem 3.6(i)-(ii) and (iv), there exists an index  $K \in \mathbb{N}$  such that for all  $k \geq K$  we have  $(\mathbf{X}^k, \boldsymbol{\varepsilon}^k) \in \{(\mathbf{X}, \boldsymbol{\varepsilon}) \in \mathbb{S}_{++}^n \times \mathbb{S}_{>0}^n \mid \text{dist}((\mathbf{X}, \boldsymbol{\varepsilon}), \Gamma^\infty) < \epsilon\} \cap \{(\mathbf{X}, \boldsymbol{\varepsilon}) \mid F(\mathbf{X}^*, \mathbf{0}) < F(\mathbf{X}^k, \boldsymbol{\varepsilon}^k) < F(\mathbf{X}^*, \mathbf{0}) + \vartheta\}$  with  $\vartheta > 0$  and . For any such  $k$ , the KL inequality yields

$$\frac{c}{2}(F(\mathbf{X}^k, \boldsymbol{\varepsilon}^k) - F(\mathbf{X}^*, \mathbf{0}))^{-1/2} \cdot \text{dist}(\mathbf{0}, \partial F(\mathbf{X}^k; \boldsymbol{\varepsilon}^k)) \geq 1. \quad (40)$$

Introducing the notion  $\tilde{\Delta}_k = F(\mathbf{X}^k, \boldsymbol{\varepsilon}^k) - F(\mathbf{X}^*, \mathbf{0})$ , (40) can be rewritten as

$$\text{dist}(\mathbf{0}, \partial F(\mathbf{X}^k; \boldsymbol{\varepsilon}^k)) \geq \frac{2}{c} \tilde{\Delta}_k^{1/2}.$$

Next, we bound the left-hand side from above in terms of the iterate differences. By virtue of (33) and using the inexactness condition  $(\mathbf{C}_{\text{inexact}}^{(k)})$ , one obtains

$$\begin{aligned} \text{dist}(\mathbf{0}, \partial F(\mathbf{X}^k, \boldsymbol{\varepsilon}^k)) &\leq l^{-2} \|\mathbf{X}^k - \mathbf{Y}^k\|_F + \|\mathcal{R}^k\|_F \\ &\leq (l^{-2} + \tilde{\beta}_k) \|\mathbf{X}^k - \mathbf{Y}^k\|_F. \end{aligned} \quad (41)$$

Then, we deduce that

$$\frac{2}{c} \tilde{\Delta}_k^{1/2} \leq \frac{l^{-2} + \tilde{\beta}_k}{\alpha} \|\mathbf{X}^k - \mathbf{X}^{k+1}\|_F.$$

On the other hand, the sufficient decrease property (see Eq. (24)) guarantees that

$$\|\mathbf{X}^k - \mathbf{X}^{k+1}\|_F^2 \leq \tilde{C}^{-1}(F(\mathbf{X}^k, \boldsymbol{\varepsilon}^k) - F(\mathbf{X}^{k+1}, \boldsymbol{\varepsilon}^{k+1})) = \tilde{C}^{-1}(\tilde{\Delta}_k - \tilde{\Delta}_{k+1}).$$

Substituting the above bound into the previous squared inequality, we obtain

$$\frac{4}{c^2} \tilde{\Delta}_k \leq \left( \frac{l^{-2} + \tilde{\beta}_k}{\alpha} \right)^2 \tilde{C}^{-1}(\tilde{\Delta}_k - \tilde{\Delta}_{k+1}).$$

Rearranging the terms, we deduce that there exists a constant  $\tilde{c} \in (0, 1)$  such that

$$\tilde{\Delta}_{k+1} \leq \tilde{c} \tilde{\Delta}_k \text{ for all sufficiently large } k.$$

This recursive inequality implies that

$$\tilde{\Delta}_k \leq (\tilde{c})^{k-K} \tilde{\Delta}_K, \quad (42)$$

i.e., the sequence  $\{\tilde{\Delta}_k\}$  converges to zero Q-linearly, which is equivalent to saying that  $\{F(\mathbf{X}^k, \boldsymbol{\varepsilon}^k)\}$  converges Q-linearly. This completes the proof.  $\square$

We next prove the convergence properties of the sequence  $\{\mathbf{X}^k\}$ .

**Theorem 3.12.** *Suppose that Assumptions 1.1-3.7 hold, and let  $\{(\mathbf{X}^k, \boldsymbol{\varepsilon}^k)\}$  be the sequence generated by Algorithm 1. The following statements hold.*

- (i) *If  $F(\mathbf{X}, \boldsymbol{\varepsilon})$  is a KL function, then  $\sum_{k=1}^{+\infty} \|\mathbf{X}^{k+1} - \mathbf{X}^k\|_F < +\infty$ . Consequently, the sequence  $\{\mathbf{X}^k\}$  converges to a stationary point of  $(\mathcal{P})$ .*
- (ii) *If, in addition,  $F(\mathbf{X}, \boldsymbol{\varepsilon})$  has the KL property with exponent  $1/2$  at  $\mathbf{X}^*$ , then  $\{\mathbf{X}^k\}$  converges at least R-linearly to  $\mathbf{X}^*$ .*

*Proof.* (i) As stated in the proof of Theorem 3.11, it suffices to consider the case in which  $F(\mathbf{X}^k, \boldsymbol{\varepsilon}^k) > F(\mathbf{X}^{k+1}, \boldsymbol{\varepsilon}^{k+1})$  for all sufficiently large  $k$ . Let  $(\mathbf{X}^*, \mathbf{0})$  be a cluster point of the sequence  $\{(\mathbf{X}^k, \boldsymbol{\varepsilon}^k)\}$ .

From Theorem 3.6(i)–(ii) and (iv), there exist  $\epsilon > 0$ ,  $\vartheta > 0$ , and an index  $K \in \mathbb{N}$  such that for all  $k \geq K$ ,  $(\mathbf{X}^k, \boldsymbol{\varepsilon}^k) \in \{(\mathbf{X}, \boldsymbol{\varepsilon}) \mid \text{dist}((\mathbf{X}, \boldsymbol{\varepsilon}), \Gamma^\infty) < \epsilon\} \cap \{(\mathbf{X}, \boldsymbol{\varepsilon}) \mid F(\mathbf{X}^*, \mathbf{0}) < F(\mathbf{X}^k, \boldsymbol{\varepsilon}^k) < F(\mathbf{X}^*, \mathbf{0}) + \vartheta\}$ . Because  $F$  is a KL function on this neighborhood, there is a *desingularizing function*  $\varphi$  such that for each  $k \geq K$ ,

$$\varphi'(F(\mathbf{X}^k, \boldsymbol{\varepsilon}^k) - F(\mathbf{X}^*, \mathbf{0})) \cdot \text{dist}(\mathbf{0}, \partial F(\mathbf{X}^k, \boldsymbol{\varepsilon}^k)) \geq 1. \quad (43)$$

Moreover, since  $\varphi$  is concave on  $[0, \vartheta]$ , for any such  $k$ , one obtains the standard KL -type inequality

$$\begin{aligned} & \varphi'(F(\mathbf{X}^k, \boldsymbol{\varepsilon}^k) - F(\mathbf{X}^*, \mathbf{0}))(F(\mathbf{X}^k, \boldsymbol{\varepsilon}^k) - F(\mathbf{X}^{k+1}, \boldsymbol{\varepsilon}^{k+1})) \\ & \leq \varphi(F(\mathbf{X}^k) - F(\mathbf{X}^*)) - \varphi(F(\mathbf{X}^{k+1}) - F(\mathbf{X}^*)). \end{aligned} \quad (44)$$

Set  $\bar{\Delta}_k := \varphi(F(\mathbf{X}^k) - F(\mathbf{X}^*))$  for each  $k$ . It follows from (43) and (44) that for each  $k \geq K$ ,

$$\begin{aligned} \bar{\Delta}_k - \bar{\Delta}_{k+1} & \geq \varphi'(F(\mathbf{X}^k, \boldsymbol{\varepsilon}^k) - F(\mathbf{X}^*, \mathbf{0}))(F(\mathbf{X}^k, \boldsymbol{\varepsilon}^k) - F(\mathbf{X}^{k+1}, \boldsymbol{\varepsilon}^{k+1})) \\ & \geq \frac{F(\mathbf{X}^k, \boldsymbol{\varepsilon}^k) - F(\mathbf{X}^{k+1}, \boldsymbol{\varepsilon}^{k+1})}{\text{dist}(\mathbf{0}, \partial F(\mathbf{X}^k, \boldsymbol{\varepsilon}^k))} \\ & \stackrel{(a)}{\geq} \frac{\tilde{C} \|\mathbf{X}^k - \mathbf{X}^{k+1}\|_F^2}{\frac{l^{-2} + \tilde{\beta}_k}{\alpha} \|\mathbf{X}^k - \mathbf{X}^{k+1}\|_F} \\ & = \frac{\alpha \tilde{C}}{l^{-2} + \tilde{\beta}_k} \|\mathbf{X}^k - \mathbf{X}^{k+1}\|_F \stackrel{(b)}{\geq} \frac{\alpha \tilde{C}}{l^{-2} + 1} \|\mathbf{X}^k - \mathbf{X}^{k+1}\|_F, \end{aligned} \quad (45)$$



where (a) follows from (24), (41) and  $\mathbf{Y}^k = \frac{1}{\alpha}(\mathbf{X}^{k+1} - \mathbf{X}^k) + \mathbf{X}^k$ , and (b) holds since  $\tilde{\beta}_k \leq 1$ .

Denote  $\omega := \frac{\alpha\tilde{C}}{l^{-2}+1} > 0$ . Summing up (45) from  $K$  to any  $\bar{k} > K$  yields

$$\sum_{k=K}^{\bar{k}} \|\mathbf{X}^k - \mathbf{X}^{k+1}\|_F \leq \frac{1}{\omega} \sum_{k=K}^{\bar{k}} (\bar{\Delta}_k - \bar{\Delta}_{k+1}) = \frac{1}{\omega} (\bar{\Delta}_K - \bar{\Delta}_{\bar{k}+1}) \leq \frac{1}{\omega} \bar{\Delta}_K.$$

Letting  $\bar{k} \rightarrow +\infty$ , we have  $\sum_{k=K}^{+\infty} \|\mathbf{X}^k - \mathbf{X}^{k+1}\|_F < +\infty$ , and hence  $\{\mathbf{X}^k\}$  is a Cauchy sequence. By Theorem 3.6(iii), its limit  $\mathbf{X}^*$  lies in  $\Omega^\infty$ . One then verifies that  $\mathbf{X}^*$  is a stationary point of  $(\mathcal{P})$ , completing the proof of statement (i).

(ii). By Theorem 3.11 (see (42)), we know that there exists an index  $K \in \mathbb{N}$  such that for all  $k \geq K$ ,

$$F(\mathbf{X}^k, \boldsymbol{\varepsilon}^k) - F(\mathbf{X}^*, \mathbf{0}) \leq (\tilde{c})^{k-K} \tilde{\Delta}_K \quad \text{for some } \tilde{c} \in (0, 1).$$

Using statement (i), we know  $\{\mathbf{X}^k\}$  converges. For any  $k \geq K$ , we have

$$\|\mathbf{X}^k - \mathbf{X}^*\|_F = \left\| \lim_{t \rightarrow \infty} \sum_{\ell=k}^t (\mathbf{X}^\ell - \mathbf{X}^{\ell+1}) \right\|_F \leq \sum_{\ell=k}^{\infty} \|\mathbf{X}^\ell - \mathbf{X}^{\ell+1}\|_F.$$

By (24), we can bound each  $\|\mathbf{X}^\ell - \mathbf{X}^{\ell+1}\|_F$  by a constant multiple of  $(F(\mathbf{X}^\ell, \boldsymbol{\varepsilon}^\ell) - F(\mathbf{X}^{\ell+1}, \boldsymbol{\varepsilon}^{\ell+1}))^{1/2}$ . Hence

$$\|\mathbf{X}^k - \mathbf{X}^*\|_F \leq \sum_{\ell=k}^{\infty} \left( \tilde{C}^{-1} (F(\mathbf{X}^\ell, \boldsymbol{\varepsilon}^\ell) - F(\mathbf{X}^{\ell+1}, \boldsymbol{\varepsilon}^{\ell+1})) \right)^{\frac{1}{2}} \leq \left( \frac{\tilde{\Delta}_K}{\tilde{C}(\tilde{c})^K} \right)^{\frac{1}{2}} \sum_{\ell=k}^{\infty} (\tilde{c})^{\frac{\ell}{2}}.$$

Since  $\sum_{\ell=k}^{\infty} (\tilde{c})^{\ell/2} = \frac{(\tilde{c})^{k/2}}{1 - \sqrt{\tilde{c}}}$ , it follows that

$$\|\mathbf{X}^k - \mathbf{X}^*\|_F \leq \left( \frac{\tilde{\Delta}_K}{\tilde{C}(\tilde{c})^K} \right)^{\frac{1}{2}} \frac{(\tilde{c})^{k/2}}{1 - \sqrt{\tilde{c}}} = \mathcal{O}((\sqrt{\tilde{c}})^k).$$

Because  $0 < \sqrt{\tilde{c}} < 1$ , this shows  $\{\mathbf{X}^k\}$  converges *R-linearly* to  $\mathbf{X}^*$ . The proof is complete.  $\square$

## 4 Numerical Experiments

This section evaluates the numerical performance of DIIR-QUIC in solving  $(\mathcal{P})$  on both synthetic and real-world datasets and investigates its convergence behavior

numerically. All experiments are conducted on a PC equipped with an AMD Ryzen 5 4600H processor (3.00 GHz base frequency), 16 GB of RAM, running 64-bit Ubuntu 22 LTS. At each outer iteration, the weighted  $\ell_1$ -regularized log-determinant subproblem ( $\mathcal{P}_{\text{sub}}$ ) is solved using a QUIC implementation<sup>1</sup> modified to enforce the inexact termination condition ( $\mathbf{C}_{\text{inexact}}^{(k)}$ ). We have publicly released the C++ implementation of DIIR-QUIC at <https://github.com/Optimizater/DIIR-QUIC>.

#### 4.1 Nonconvex $\ell_p$ -regularized log-determinant problems

The nonconvex  $\ell_p$  (quasi-)norm ( $0 < p < 1$ ) is used to promote sparsity in inverse covariance estimation [11]. As a representative instance of our broader class of sparsity-promoting formulations, we take  $\Phi(\mathbf{X}) = \|\mathbf{X}\|_p^p = \sum_{ij} |X_{ij}|^p$  and consider the following problem:

$$\begin{aligned} \min \quad & F(\mathbf{X}) = \{f(\mathbf{X}) := \text{tr}(\mathbf{S}\mathbf{X}) - \log \det \mathbf{X}\} + \rho \sum_{ij} |X_{ij}|^p \\ \text{s.t.} \quad & \mathbf{X} \in \mathbb{S}_{++}^n, \end{aligned} \quad (46)$$

where  $p \in (0, 1)$ . To our knowledge, [11] is the only existing algorithm capable of directly solving the special case of problem (46) in which the regularizer is  $\rho \sum_{i \neq j} |X_{ij}|^p$ . Accordingly, we use it as the benchmark algorithm in our numerical experiments and, following their terminology, refer to it as  $\ell_p \text{COV}$ .

For DIIR-QUIC, we set  $\alpha = 0.98$  and  $\mu = 0.1$ . In Algorithm 2, an iterative coordinate descent method is used to compute an approximate Newton direction  $\mathbf{D}^t$  in (12), under the inexactness condition (13) with  $\epsilon = 0.05$ . We terminate the outer loop when either

- (i) the iteration count reaches `MaxIter` = 3000, or
- (ii) the stationarity residual (see Theorem 1.4)

$$\max_{(i,j) \in \mathcal{I}(\mathbf{X}^k)} |S_{ij} - [(\mathbf{X}^k)^{-1}]_{ij} + \rho p |X_{ij}^k|^{p-1} \text{sgn}(X_{ij}^k)| \cdot n < \text{tol}, \quad (47)$$

is satisfied. Here,  $\text{tol} = 10^{-5}$ .

In our implementation, we treat any entry of  $\mathbf{X}^k$  whose magnitude exceeds  $10^{-8}$  as active, and these indices define  $\mathcal{I}(\mathbf{X}^k)$ . Following [34], we initialize

$$[\mathbf{X}^0]_{ij} = \begin{cases} 0, & \text{if } i \neq j \\ \frac{1}{S_{ii} + \rho}, & \text{if } i = j. \end{cases} \quad (48)$$

---

<sup>1</sup>A publicly available implementation of the original QUIC algorithm—written in C++ with a Python wrapper—can be found at <https://github.com/osdf/pyquic>.

and initialize the perturbation matrix  $\mathcal{E}^0$  by drawing an  $n \times n$  matrix  $\mathbf{W}$  with i.i.d.  $\mathcal{N}(0, 1)$ , scaling each entry by 0.5, symmetrizing via  $\widehat{\mathbf{W}} = (\mathbf{W} + \mathbf{W}^T)/2$ , and then setting  $\mathcal{E}_{ij}^0 = |\widehat{W}_{ij}|$  for all  $(i, j)$ .

For the  $\ell_p$ COV, the authors have not publicly released their MATLAB implementation. Consequently, we reimplemented the algorithm primarily based on the descriptions provided in the original paper. Since our Algorithm 2 is implemented in C++, we also developed the core computational components of  $\ell_p$ COV in C++ to ensure a relatively fair comparison. As described in [11, §IV SIMULATIONS], the authors employed a warm-start (WS) strategy with respect to the model parameter  $p$  to potentially improve numerical performance with a good initialization. This strategy has also been incorporated into our reimplementation. More precisely, this warm-start procedure generates a short sequence of exponents  $1 = p_{(0)} > p_{(1)} > \dots > p_{(K)} = p > 0$  that interpolates linearly between 1 and the target  $p \in (0, 1)$ . The number of intermediate steps  $K$  is chosen as

$$K = \begin{cases} 2, & 1 > p \geq 0.9, \\ 3, & 0.7 \leq p < 0.9, \\ 4, & 0.4 \leq p < 0.7, \\ 5, & 0.2 \leq p < 0.4, \\ 6, & 0 < p < 0.2. \end{cases}$$

Then it proceeds as follows:

- (i) **Initialization:** Set  $p_{(0)} = 1$  and adopt the GLASSO algorithm [4] to solve the resulting convex weighted  $\ell_1$ -regularized log-determinant program with initialization (48).
- (ii) **Warm-start loop:** For  $k = 1, \dots, K$ , run the  $\ell_p$ COV algorithm at exponent  $p_{(k)}$ , using the solution obtained at step  $k - 1$  as the initial estimation.

Since  $\ell_p$ COV tackles (46) with the off-diagonal regularizer  $\rho \sum_{i \neq j} |X_{ij}|^p$ , we terminate its iterations when either

- (i) the iteration count reaches `MaxIter` = 3000, or
- (ii) the stationarity residual (47) is satisfied. Here, since the regularizer excludes diagonal entries, we set  $\rho p |X_{ij}^k|^{p-1} \text{sgn}(X_{ij}^k) = 0$  for  $i = j$  and  $(i, j) \in \mathcal{I}(\mathbf{X}^k)$ .

## 4.2 Synthetic data

We generate two synthetic test cases of  $n$ -variate Gaussian data following [35], each defined by a known precision matrix:

- (i) **Tridiagonal precision:** A strongly diagonally dominant precision matrix

$$\Sigma_{ii} = 1.25, \Sigma_{i,i+1} = \Sigma_{i+1,i} = -0.5, \Sigma_{ij} = 0 \text{ otherwise.}$$

This structure induces simple chain-like dependencies among the variables and ensures positive definiteness by virtue of diagonal dominance [36].

- (ii) **Clustered precision:** A *random* structured precision matrix with  $n/100$  clusters of size 100. Each variable is conditionally dependent on approximately 10 others, with 90% of edges located within the same cluster and only 10% connecting different clusters. This yields dense intra-cluster and sparse inter-cluster connectivity.

As noted in [35], the tridiagonal case provides a simple baseline with one-dimensional dependence, while the clustered case models more realistic community-structured dependencies commonly observed in practice [37]. In particular, we follow the procedure presented in [38, Example 1] to generate the clustered precision matrix. Using the true sparse inverse covariance matrix  $\Sigma^{-1}$  generated above, we then draw  $m = n/2$  i.i.d. samples from the corresponding GMRF distribution. When  $\mathbf{S}$  is singular, we regularize it by setting  $\mathbf{S} \leftarrow \mathbf{S} + \hat{\varepsilon} \mathbf{I}_{n \times n}$ , where  $\hat{\varepsilon} \in \{10^{-8}, 10^{-7}, \dots, 10^{-4}\}$ , and repeat this update until  $\mathbf{S}$  becomes positive definite.

The regularization parameter  $\rho$  on synthetic datasets is chosen by five-fold cross-validation, as in [6]. We first construct a logarithmic grid of 10 candidate values:  $\{\rho_j = 10^{a_j}\}_{j=1}^{10}$ , where the exponents  $a_j$  are equally spaced over  $[-1, 0]$  for tridiagonal matrices and over  $[-2, 0]$  for clustered matrices. For each  $\rho_j$  and each fold  $k$ , we estimate the precision matrix on the training subset and evaluate its negative log-likelihood ( $\text{NLL}_k$ ) on the test subset:

$$\text{NLL}_k(\rho_j) = \frac{1}{2} \left( \text{tr}(\mathbf{S}_{\text{test}}^{(k)} \Sigma_{\text{train}}^{\dagger, (k)}) - \log \det \Sigma_{\text{train}}^{\dagger, (k)} \right), \quad (49)$$

where  $\mathbf{S}_{\text{test}}^{(k)}$  is the empirical covariance on fold  $k$  and  $\Sigma_{\text{train}}^{\dagger, (k)}$  is the estimated precision matrix from the training set. We then average over folds:

$$\text{NLL}_{\text{CV}} = \frac{1}{5} \sum_{k=1}^5 \text{NLL}_k(\rho_j), \quad (50)$$

and choose the  $\rho_j$  that minimizes  $\text{NLL}_{\text{CV}}$  as the optimal  $\rho^*$ .

#### 4.2.1 Empirical convergence behavior on synthetic data

We now investigate the convergence behaviors of DIIR-QUIC. All experiments run DIIR-QUIC and  $\ell_p\text{COV}$  on the same dataset. Each curve in our plots corresponds to a single run of the respective algorithm on that shared dataset. Specifically,

- (i) We plot the objective values  $F(\mathbf{X}^k, \boldsymbol{\varepsilon}^k)$  and stationarity residuals versus elapsed time for both DIIR-QUIC and  $\ell_p$ COV. In  $\ell_p$ COV, we consider the penalty  $\rho \sum_{i,j} |X_{ij}|^p$  for objective comparison. When warm-starting over a sequence of exponents  $1 = p_{(0)} > p_{(1)} > \cdots > p_{(K)} = p > 0$ , we record performance metrics only at the final exponent  $p_{(K)} = p$ . During the initial warm-up stages ( $p_{(i)}, \forall 1 \leq i < K$ ), we terminate each run when either the iteration count reaches `MaxIter` = 3000, or the change in successive objective values falls below  $10^{-4}$ , and we do not log performance metrics for these intermediate stages. Fig 1 confirms that DIIR-QUIC not only drives down the objective rapidly but also achieves high-quality first-order stationarity—effectively solving the nonconvex  $\ell_p$ -regularized log-determinant problem in far fewer iterations than  $\ell_p$ COV algorithm, across varying matrix dimensions and exponents  $p$ .
- (ii) To illustrate the Q-linear convergence of DIIR-QUIC, we plot the error  $F(\mathbf{X}^k, \boldsymbol{\varepsilon}^k) - F(\mathbf{X}^\dagger, \boldsymbol{\varepsilon}^\dagger)$  over the last few iterations, where  $\mathbf{X}^\dagger$  denotes the returned estimate of DIIR-QUIC. To avoid plotting numerical noise, we only display those  $F(\mathbf{X}^k, \boldsymbol{\varepsilon}^k) - F(\mathbf{X}^\dagger, \boldsymbol{\varepsilon}^\dagger)$  values exceeding  $10^{-8}$ . The plots are shown in Fig 2.
- (iii) We empirically validate Proposition 2.1 by checking that, once the approximate fixed-point conditions

$$\|\boldsymbol{\varepsilon}^k\|_F < 10^{-8} \quad \text{and} \quad \|\mathbf{X}^k - \mathbf{Y}^k\|_F < 10^{-5},$$

are satisfied, the stationarity residual condition (47) is generally satisfied. To illustrate this, Fig 3 plots, over the elapsed time, the scaled infinity norms

$$n \cdot \|\mathbf{X}^k - \mathbf{Y}^k\|_\infty \quad \text{and} \quad n \cdot \|\boldsymbol{\varepsilon}^k\|_\infty,$$

where we only display those  $n \cdot \|\boldsymbol{\varepsilon}^k\|_\infty$  values exceeding  $10^{-10}$ . The plots confirm that as soon as the two fixed-point thresholds are reached, the stationarity residual falls below its prescribed tolerance.

- (iv) To empirically confirm Theorem 3.9, we monitor the cardinality of the index set  $\mathcal{I}(\mathbf{X}^k)$  over successive iterations. In all our tests, each of these cardinalities stabilizes after a finite number of iterations, indicating that the algorithm has correctly identified the smooth active manifold. The plots are shown in Fig 4.

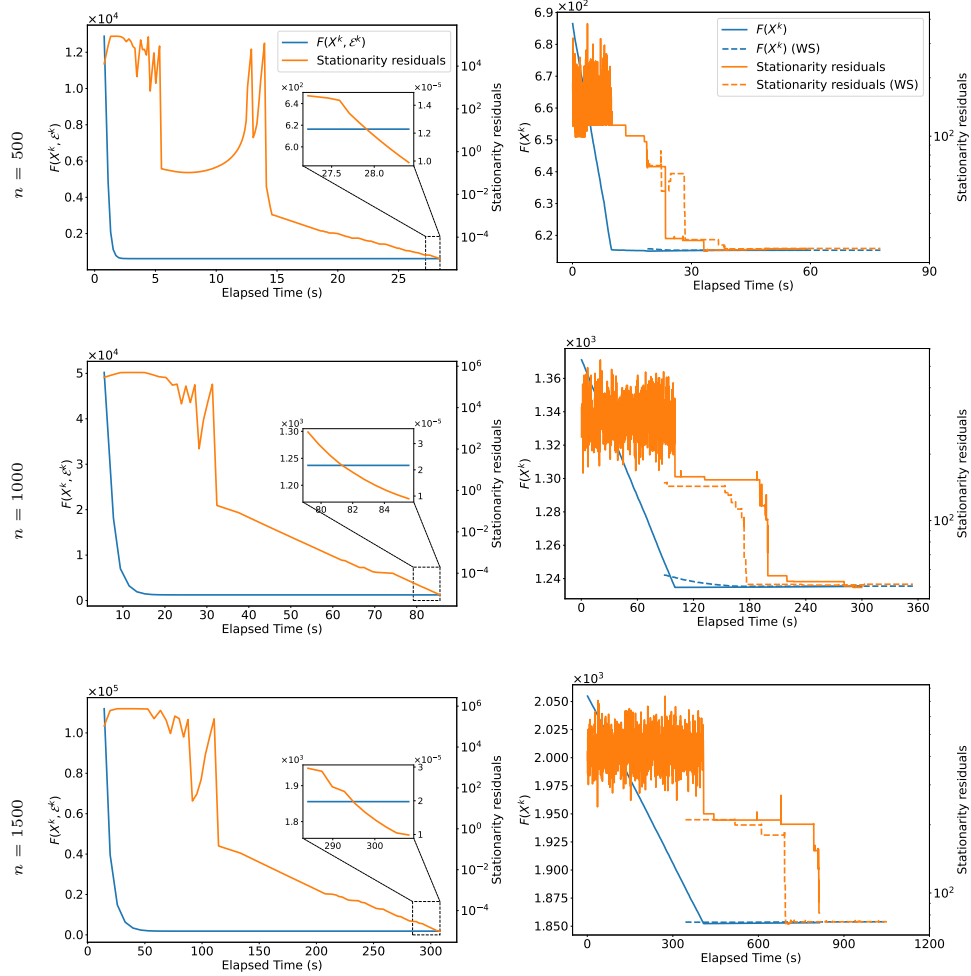


Figure 1: Convergence plots of objective value and stationarity residual (plotted on a logarithmic scale) versus elapsed time (seconds) for DIIR-QUIC (left) and  $\ell_p$ COV (right) on a tridiagonal precision matrix with  $p = 0.5$ . DIIR-QUIC not only converges faster—consistently attaining the prescribed stationarity residual when  $\ell_p$ COV often does not—but also reduces total run-time by approximately 70% across all tested dimensions.

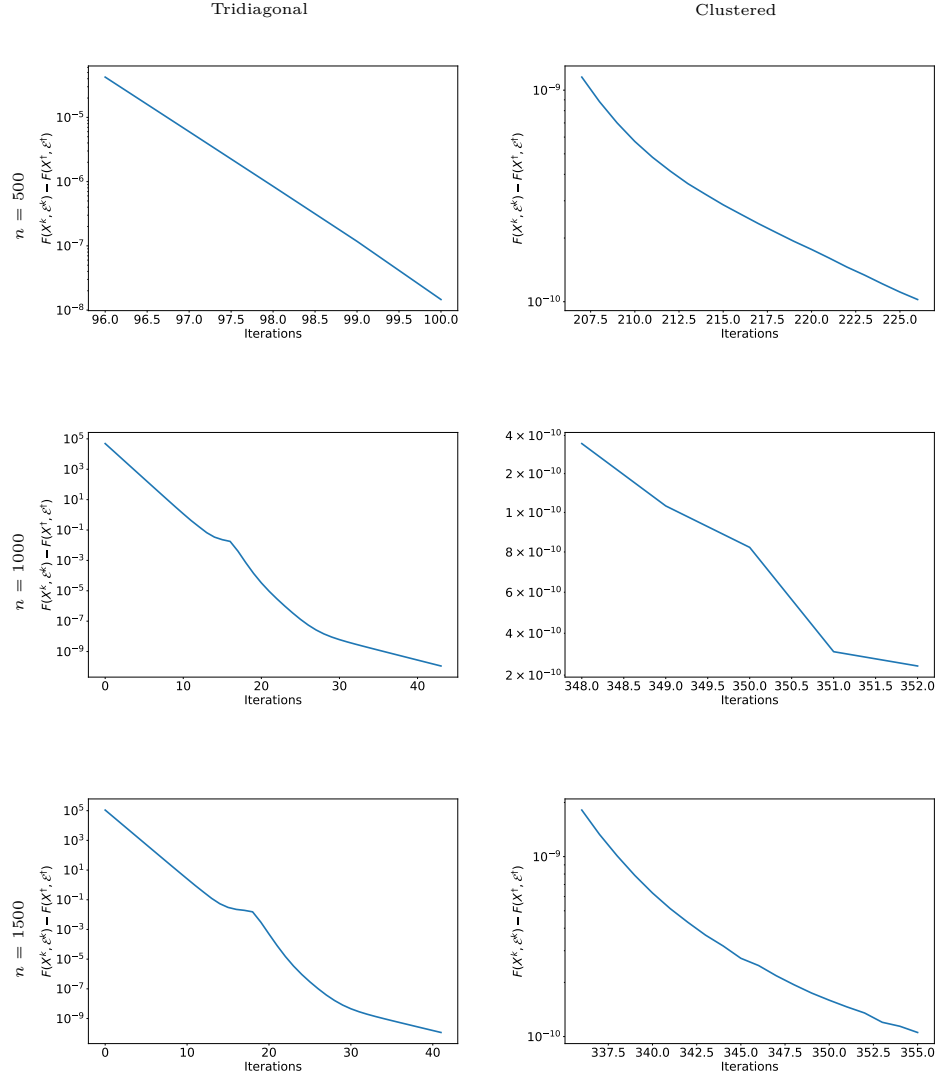


Figure 2: Q-linear convergence of the perturbed objective error  $F(\mathbf{X}^k, \boldsymbol{\varepsilon}^k) - F(\mathbf{X}^\dagger, \boldsymbol{\varepsilon}^\dagger)$  plotted versus iteration for DIIR-QUIC on tridiagonal precision matrices and clustered matrices with  $p = 0.5$ , across varying matrix dimensions. The error is displayed on a logarithmic scale, and only values above  $10^{-8}$  are shown.

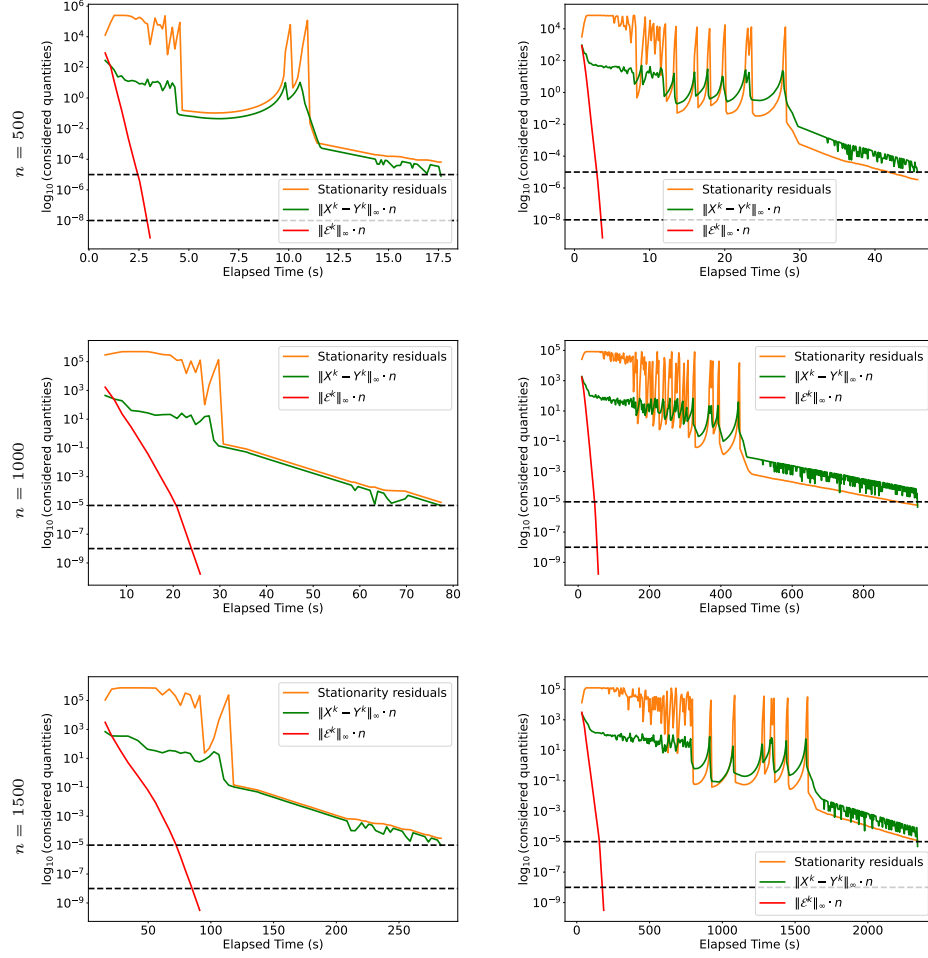


Figure 3: Fixed-point conditions and stationarity residual condition versus elapsed time (seconds) for DIIR-QUIC and  $\ell_p$ COV on tridiagonal precision matrices and cluster precision matrices with  $p = 0.5$



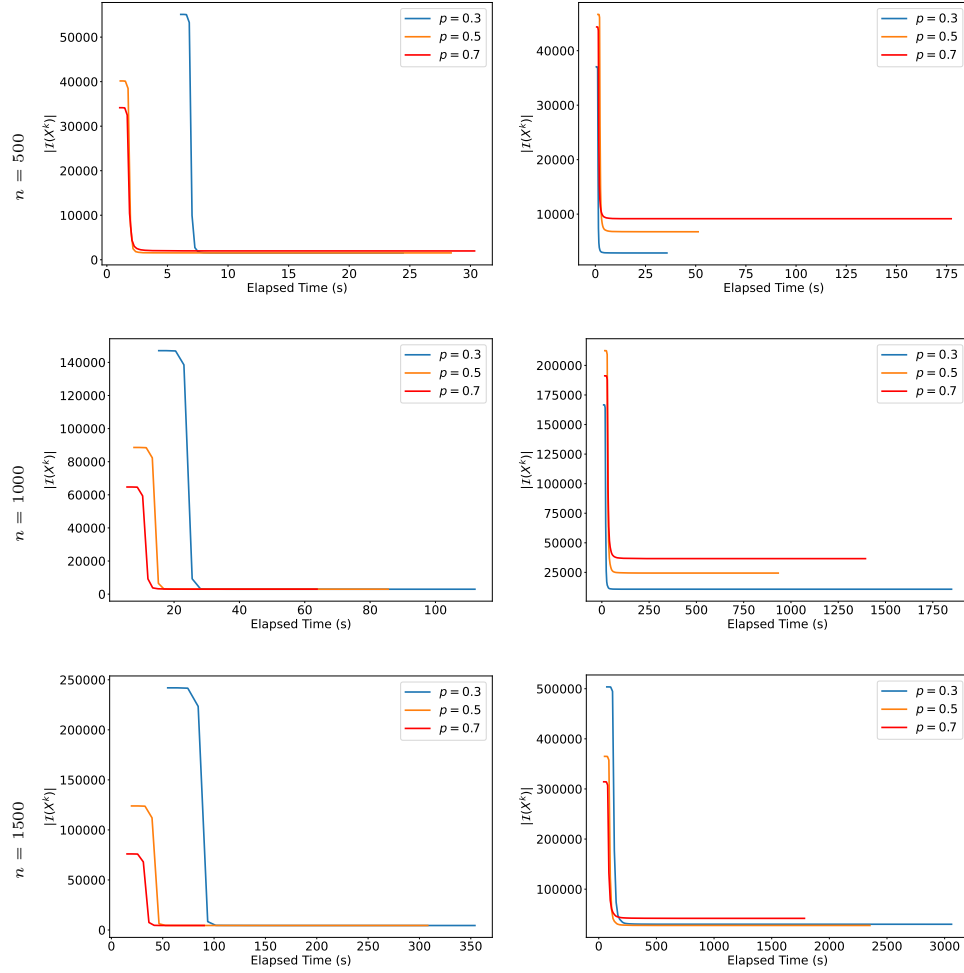


Figure 4: Plots of the cardinality of  $\mathcal{I}(\mathbf{X}^k)$  versus elapsed time for DIIR-QUIC on tridiagonal precision matrices and clustered matrices.

#### 4.2.2 Statistical properties of the estimator on synthetic data

To measure how well an algorithm recovers the true precision matrix  $\Sigma^{-1}$ , we follow [35, 11, 7] and compute the following metrics for the returned estimate  $\mathbf{X}^\dagger$ :

- (i) Normalized Kullback-Leibler (KL) loss [39]

$$\text{KL}(\mathbf{X}^\dagger) := \frac{1}{n} [\text{tr}(\Sigma \mathbf{X}^\dagger) - \log \det(\Sigma \mathbf{X}^\dagger) - n]$$

and quadratic loss ( $\text{Loss}_Q$ )

$$\text{Loss}_Q(\mathbf{X}^\dagger) := \frac{1}{n} \|\Sigma \mathbf{X}^\dagger - \mathbf{I}\|_F,$$

measure entry-wise fit between  $\Sigma \mathbf{X}^\dagger$  and the identity. For these two measures, smaller value indicates a preferred estimator.

- (ii) Matthews correlation coefficient (MCC) [40]

$$\text{MCC}(\mathbf{X}^\dagger) := \frac{\text{TP} \times \text{TN} - \text{FP} \times \text{FN}}{\sqrt{(\text{TP} + \text{FP})(\text{TP} + \text{FN})(\text{TN} + \text{FP})(\text{TN} + \text{FN})}}.$$

Here, TP, TN, FP, and FN count true positives, true negatives, false positives and false negatives in the support of  $\mathbf{X}^\dagger$ . MCC ranges from  $-1$  to  $1$ , with  $1$  denoting perfect support recovery.

- (iii) Sensitivity and specificity [7]

$$\text{Sensitivity} = \frac{\text{TP}}{\text{TP} + \text{FN}} \quad \text{and} \quad \text{Specificity} = \frac{\text{TN}}{\text{TN} + \text{FP}}.$$

Sensitivity measures the fraction of true nonzeros recovered, while specificity measures the fraction of true zeros recovered; ideal values are close to  $1$ .

- (iv)  $F_1$  score [35]:

$$F_1\text{-score} := [1 + 0.5(\text{FP} + \text{FN})/(\|\Sigma^{-1}\|_0 - \text{FN})]^{-1}.$$

Here  $\|\Sigma^{-1}\|_0$  counts the number of nonzeros of the true precision matrix  $\Sigma^{-1}$ . The  $F_1$ -score balances precision and recall on the nonzero support. A value of  $1$  indicates perfect support recovery (see, e.g., [41] for details).

To comprehensively evaluate the performance of DIIR-QUIC and  $\ell_p\text{COV}$ , we conducted 10 independent trials on synthetic datasets featuring two distinct covariance structures (tridiagonal and clustered patterns) across varying dimensions and exponents  $p$ . As shown in Table 1, the results demonstrate that DIIR-QUIC consistently outperformed competing methods in recovering the underlying sparse structure.

Table 1: Comparison of DIIR-QUIC and  $\ell_p$ COV on synthetic covariance estimation. For each method, the table reports the mean ( $\pm$  one standard deviation) over 10 independent trials under two covariance structures—tridiagonal (Tri) and clustered (Clu)—with  $p = 0.5$  at dimensions  $n = 500, 1000$ , and 1500. Boldface entries indicate the better performance.

Data	$n$	Algorithm	$F_1$	Time (s)	$\text{nnz}(\Sigma^{-1})$	$\text{nnz}(X^\dagger)$	Loss Q	KL Loss	Sensitivity	Specificity	MCC
tri	500	DIIR-QUIC	0.972 (0.004)	<b>14.615</b> (12.238)	1498.000 (0.000)	1525.600 (10.947)	0.011 (0.000)	0.028 (0.001)	0.981 (0.002)	1.000 (0.000)	0.972 (0.004)
		$\ell_p$ COV	0.999 (0.001)	59.706 (0.804)	1498.000 (0.000)	1498.400 (0.800)	<b>0.009</b> (0.000)	<b>0.015</b> (0.001)	0.999 (0.001)	1.000 (0.000)	0.999 (0.001)
		$\ell_p$ COV (WS)	<b>1.000</b> (0.000)	78.437 (1.282)	1498.000 (0.000)	1498.400 (0.800)	<b>0.009</b> (0.000)	<b>0.015</b> (0.001)	1.000 (0.000)	1.000 (0.000)	<b>1.000</b> (0.000)
	1000	DIIR-QUIC	0.984 (0.002)	<b>78.944</b> (27.952)	2998.000 (0.000)	2994.400 (4.630)	0.007 (0.000)	0.023 (0.000)	0.984 (0.001)	1.000 (0.000)	0.984 (0.002)
		$\ell_p$ COV	<b>1.000</b> (0.000)	250.391 (25.703)	2998.000 (0.000)	2998.000 (0.000)	<b>0.005</b> (0.000)	<b>0.011</b> (0.000)	1.000 (0.000)	1.000 (0.000)	<b>1.000</b> (0.000)
		$\ell_p$ COV (WS)	<b>1.000</b> (0.000)	308.795 (45.818)	2998.000 (0.000)	2998.000 (0.000)	<b>0.005</b> (0.000)	<b>0.011</b> (0.000)	1.000 (0.000)	1.000 (0.000)	<b>1.000</b> (0.000)
	1500	DIIR-QUIC	0.988 (0.001)	<b>346.117</b> (137.118)	4498.000 (0.000)	4490.200 (5.400)	0.005 (0.000)	0.021 (0.001)	0.987 (0.001)	1.000 (0.000)	0.988 (0.001)
		$\ell_p$ COV	<b>1.000</b> (0.000)	579.475 (103.963)	4498.000 (0.000)	4498.000 (0.000)	<b>0.004</b> (0.000)	0.010 (0.000)	1.000 (0.000)	1.000 (0.000)	<b>1.000</b> (0.000)
		$\ell_p$ COV (WS)	<b>1.000</b> (0.000)	840.307 (114.144)	4498.000 (0.000)	4498.000 (0.000)	<b>0.004</b> (0.000)	<b>0.009</b> (0.000)	1.000 (0.000)	1.000 (0.000)	<b>1.000</b> (0.000)
clu	500	DIIR-QUIC	<b>0.423</b> (0.005)	<b>83.094</b> (38.935)	11439.800 (136.802)	6761.000 (95.734)	<b>0.038</b> (0.000)	<b>0.252</b> (0.003)	<b>0.337</b> (0.005)	0.988 (0.000)	0.418 (0.005)
		$\ell_p$ COV	0.379 (0.007)	142.713 (38.860)	11439.800 (136.802)	3711.600 (77.386)	0.040 (0.000)	0.263 (0.004)	0.251 (0.005)	<b>0.996</b> (0.000)	0.427 (0.007)
		$\ell_p$ COV (WS)	0.403 (0.006)	142.709 (38.859)	11439.800 (136.802)	3975.400 (83.852)	0.039 (0.000)	0.256 (0.004)	0.272 (0.005)	<b>0.996</b> (0.000)	<b>0.448</b> (0.006)
	1000	DIIR-QUIC	0.568 (0.004)	<b>934.045</b> (311.597)	23021.000 (187.743)	24441.800 (264.264)	0.023 (0.000)	0.174 (0.002)	<b>0.586</b> (0.004)	0.989 (0.000)	0.558 (0.004)
		$\ell_p$ COV	0.610 (0.004)	993.646 (311.517)	23021.000 (187.743)	14357.400 (150.439)	0.023 (0.000)	0.173 (0.002)	0.495 (0.005)	<b>0.997</b> (0.000)	0.621 (0.004)
		$\ell_p$ COV (WS)	<b>0.652</b> (0.005)	993.647 (311.516)	23021.000 (187.743)	15598.600 (174.613)	<b>0.022</b> (0.000)	<b>0.160</b> (0.002)	0.547 (0.006)	<b>0.997</b> (0.000)	<b>0.658</b> (0.004)
	1500	DIIR-QUIC	0.675 (0.003)	<b>3178.596</b> (1392.902)	34577.400 (187.755)	27870.600 (229.703)	<b>0.016</b> (0.000)	0.146 (0.001)	<b>0.610</b> (0.004)	0.997 (0.000)	0.675 (0.003)
		$\ell_p$ COV	0.654 (0.003)	3238.136 (1393.000)	34577.400 (187.755)	18417.000 (170.050)	0.017 (0.000)	0.155 (0.001)	0.501 (0.003)	<b>1.000</b> (0.000)	0.684 (0.002)
		$\ell_p$ COV (WS)	<b>0.700</b> (0.003)	3238.125 (1393.003)	34577.400 (187.755)	20222.200 (174.566)	<b>0.016</b> (0.000)	<b>0.141</b> (0.001)	0.555 (0.003)	<b>1.000</b> (0.000)	<b>0.728</b> (0.002)

### 4.3 Real-world datasets

We conduct performance comparisons between DIIR-QUIC and  $\ell_p$ COV on two classes of real-world datasets: gene expression datasets [38] and stock datasets [6]. In our first experiments, we consider three benchmark datasets from [38, Examples 3, 5 and 6]: the Lymph node status data ( $n = 587$ ), the Arabidopsis thaliana data ( $n = 834$ ) and the Leukemia data ( $n = 1255$ ). Following [11], we select the regularization parameter  $\rho$  via the extended Bayesian information criterion (EBIC). For a candidate precision-matrix estimate  $\mathbf{X}^\dagger$  and sample size  $m$ , the EBIC for a Gaussian graphical model is commonly defined as

$$\text{EBIC}(\mathbf{X}^\dagger) = -m(\log \det \mathbf{X}^\dagger - \text{tr}(S\mathbf{X}^\dagger)) + \frac{\log m + 4\tilde{\gamma} \log n}{2} |\mathcal{I}(\mathbf{X}^\dagger)|_{\text{off}},$$

where  $|\mathcal{I}(\mathbf{X}^\dagger)|_{\text{off}}$  denotes the cardinality of the index set of nonzero off-diagonal entries of  $\mathbf{X}^\dagger$  and  $\tilde{\gamma} \in [0, 1]$  is a user-defined constant. We then evaluate EBIC over a logarithmic grid

$$\rho_j = 10^{a_j}, \quad a_j = -1 + \frac{j-1}{N-1}(0 - (-1)) = -1 + \frac{j-1}{N-1}, \quad j \in [N],$$

where  $N = 10$ , so that  $\rho_j$  ranges smoothly from  $10^{-1}$  to 1. Denoting the resulting EBIC values by  $\text{EBIC}(\rho_j)$ , we choose  $\rho_{\text{EBIC}} = \arg\min_{j \in [N]} \text{EBIC}(\rho_j)$ .

To evaluate both the statistical accuracy and computational efficiency of DIIR-QUIC at the chosen regularization level  $\rho_{\text{EBIC}}$ , we compute the five-fold cross-validated test-set negative log-likelihood, as defined in (50). The metric  $\text{NLL}_{\text{CV}}$  provides a cross-validated assessment of predictive performance, effectively balancing the trade-off between under- and over-fitting. The values of per-fold held-out negative log-likelihoods are summarized via box plots to visualize variability and robustness. The right column of Figure 5 plots the stationarity residuals versus elapsed time for both algorithms at  $\rho_{\text{EBIC}}$ . Under an identical wall-clock time limit (rather than an iteration cap) imposed on  $\ell_p$ COV, DIIR-QUIC converges more rapidly, demonstrating its computational efficiency in reaching a stationary point.

In the second experiment, following [6], we evaluate portfolio performance within the Markowitz mean-variance framework. The optimal portfolio weights  $\mathbf{w} \in \mathbb{R}^p$  are obtained by solving:

$$\min_{\mathbf{w}} \mathbf{w}^T \mathbf{\Sigma} \mathbf{w} \quad \text{s.t.} \quad \mathbf{w}^T \mathbf{e} = 1, \quad (51)$$

where  $\mathbf{\Sigma}$  is the estimated covariance matrix of asset returns and  $\mathbf{e}$  is the all-ones vector. For a given weight vector  $\mathbf{w}$ , we compute on the test set  $\mathbf{X}_{\text{test}}$

$$R(\mathbf{w}) = \sum_{\mathbf{x} \in \mathbf{X}_{\text{test}}} \mathbf{w}^T \mathbf{x}, \quad \sigma(\mathbf{w}) = \sqrt{\mathbf{w}^T \mathbf{S}_{\text{test}} \mathbf{w}},$$

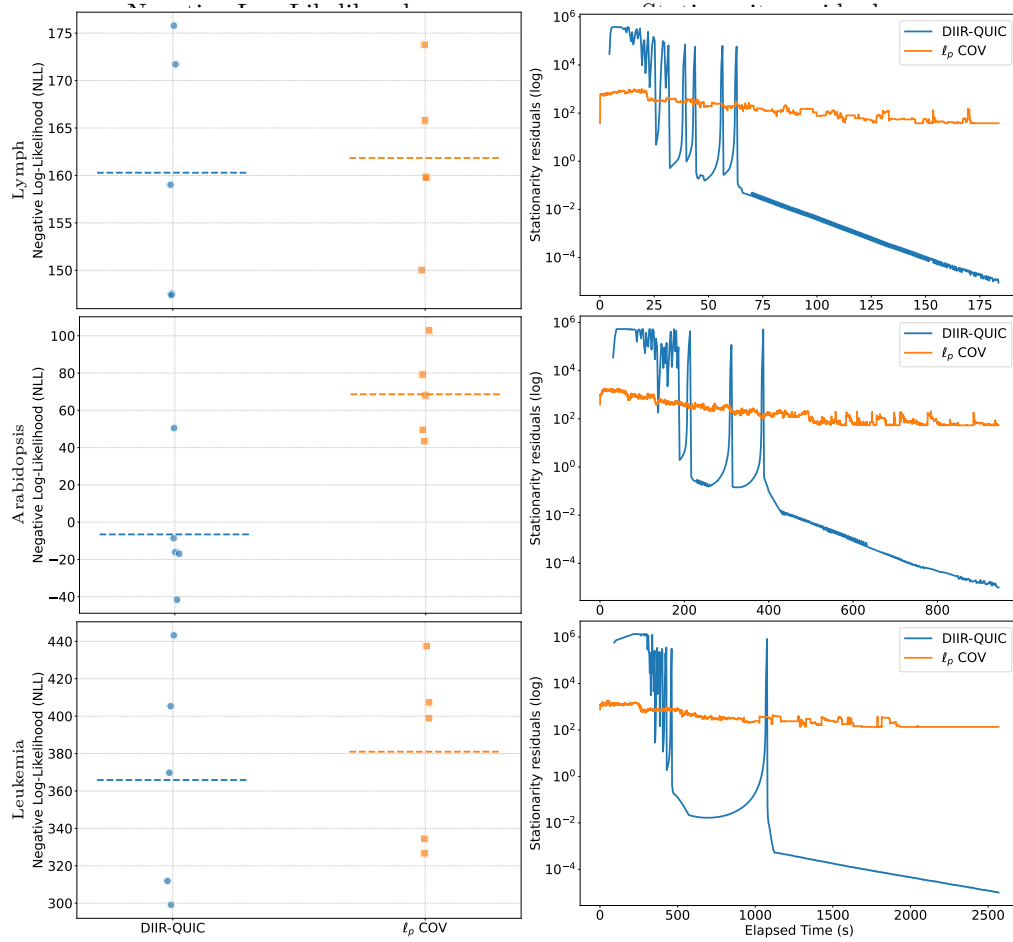


Figure 5: Comparison of DIIR-QUIC and  $\ell_p$ COV algorithms. **Left:** Five-fold cross-validated negative log-likelihood (lower is better). **Right:** Stationarity residuals versus elapsed time.

where  $\mathbf{S}_{\text{test}}$  is the sample covariance matrix of  $X_{\text{test}}$ . Portfolio performance is then measured by the Sharpe ratio

$$S(\mathbf{w}) = \frac{R(\mathbf{w})}{\sigma(\mathbf{w})}. \quad (52)$$

We analyze weekly returns of 30 Dow Jones Industrial Index components from January 6, 2020, to June 26, 2023, with data obtained from East Money Information.<sup>2</sup> Following [6], we partition the series into training (first 50 weeks), validation (next 50 weeks) and test sets. Regularization parameters are tuned to maximize the Sharpe ratio on the validation set. We compare three methods—DIIR-QUIC,  $\ell_p$ COV, and its warm-start  $\ell_p$ COV (WS)—to evaluate their impact on out-of-sample performance. Table 2 reports the results.

Table 2: Out-of-sample performance comparison of DIIR-QUIC and  $\ell_p$ COV methods on the Dow Jones weekly returns: realized return  $R(\mathbf{w})$ , realized risk  $\sigma(\mathbf{w})$ , and Sharpe ratio  $S(\mathbf{w})$ .

$p$	Algorithm	$R(\mathbf{w})$	$\sigma(\mathbf{w})$	$S(\mathbf{w})$	Time (s)	Stationarity residuals
0.3	DIIR-QUIC	0.0205	0.0263	<b>0.7815</b>	0.0037	9.65e-06
	$\ell_p$ COV	0.0191	0.0251	0.7621	0.0006	1.04e-16
	$\ell_p$ COV (WS)	0.0191	0.0251	0.7621	0.1708	1.04e-16
0.5	DIIR-QUIC	0.0231	0.0267	<b>0.8635</b>	0.0025	8.40e-06
	$\ell_p$ COV	0.0191	0.0251	0.7621	0.0007	1.04e-16
	$\ell_p$ COV (WS)	0.0191	0.0251	0.7621	0.1138	1.04e-16
0.7	DIIR-QUIC	0.0233	0.0267	<b>0.8732</b>	0.0014	3.63e-06
	$\ell_p$ COV	0.0191	0.0251	0.7621	0.0007	1.04e-16
	$\ell_p$ COV (WS)	0.0191	0.0251	0.7621	0.0577	1.04e-16

## 5 Conclusion

In this paper, we have proposed DIIR-QUIC, an inexact QUIC-based iteratively reweighting algorithm tailored for solving log-determinant optimization problems involving nonconvex partly smooth regularizers. We established that, under mild conditions, the inexactly solved subproblems are sufficient to identify the smooth

<sup>2</sup>Please refer to <https://www.eastmoney.com>

active manifold in finitely many iterations. Moreover, we proved the global convergence of the generated iterates, along with convergence rates for both the perturbed objective value sequence and the iterates under Kurdyka-Łojasiewicz property. Finally, extensive numerical experiments on synthetic and real inverse-covariance estimation tasks confirmed that DIIR-QUIC consistently delivered superior computational efficiency and estimation accuracy compared with existing methods, demonstrating its practical value for large-scale nonconvex log-determinant optimization.

## A Auxiliary proofs

For the sake of clarity, some technical proofs are deferred to this appendix.

### A.1 Proof of Lemma 3.3

*Proof.* As for statement (i), recall that  $\bar{f}$  is differentiable and

$$\partial J(\Delta^{(t)}; \mathbf{Z}^t) = \nabla \bar{f}(\Delta^{(t)}; \mathbf{Z}^t) + \partial g(\Delta^{(t)}; \mathbf{Z}^t).$$

It follows that

$$\begin{aligned} \text{dist}(\mathbf{0}, \partial J(\Delta^{(t)}; \mathbf{Z}^t)) &= \inf_{\mathbf{G}^{(t)} \in \partial g(\Delta^{(t)}; \mathbf{Z}^t)} \|\nabla \bar{f}(\Delta^{(t)}; \mathbf{Z}^t) + \mathbf{G}^{(t)}\|_F \\ &\leq \inf_{\mathbf{G}^{(t)} \in \partial g(\Delta^{(t)}; \mathbf{Z}^t)} \|\nabla \bar{f}(\Delta^{(t)}; \mathbf{Z}^t) + \mathbf{G}^{(t)}\|_1 \\ &= \inf_{\mathbf{G}^{(t)} \in \partial g(\Delta^{(t)}; \mathbf{Z}^t)} \left\{ \sum_{(i,j) \in \mathcal{I}_{\text{free}}(\mathbf{Z}^t)} |\nabla_{ij} \bar{f}(\Delta^{(t)}; \mathbf{Z}^t) + G_{ij}^{(t)}| + \right. \\ &\quad \left. \sum_{(i,j) \in \mathcal{I}_{\text{fixed}}(\mathbf{Z}^t)} |\nabla_{ij} \bar{f}(\Delta^{(t)}; \mathbf{Z}^t) + G_{ij}^{(t)}| \right\}. \end{aligned} \tag{53}$$

Define  $\varsigma_{ij}(\eta) = J(\Delta^{(t-1)} + \eta(\mathbf{e}_i \mathbf{e}_j^T + \mathbf{e}_j \mathbf{e}_i^T); \mathbf{Z}^t)$ ,  $\forall (i, j) \in \mathcal{I}_{\text{free}}(\mathbf{Z}^t)$  and let  $\eta^* \geq 0$  denote the minimizer of  $\varsigma_{ij}(\eta)$ . Then it follows from [21, Theorem 10.1] that  $0 \in \partial \varsigma_{ij}(\eta^*)$ ,  $\forall (i, j) \in \mathcal{I}_{\text{free}}(\mathbf{Z}^t)$ . By the chain rule for subdifferentials (see, e.g., [21, Theorem 10.49]), there exists  $S_{ij}^* \in \partial \psi_{ij}(\eta^*)$  with  $\psi_{ij}(\eta) = g(\Delta^{(t-1)} + \eta(\mathbf{e}_i \mathbf{e}_j^T + \mathbf{e}_j \mathbf{e}_i^T); \mathbf{Z}^t)$  such that

$$\langle \nabla \bar{f}(\Delta^{(t-1)} + \eta^*(\mathbf{e}_i \mathbf{e}_j^T + \mathbf{e}_j \mathbf{e}_i^T); \mathbf{Z}^t), \mathbf{e}_i \mathbf{e}_j^T + \mathbf{e}_j \mathbf{e}_i^T \rangle + S_{ij}^* = 0, \forall (i, j) \in \mathcal{I}_{\text{free}}(\mathbf{Z}^t).$$

On the other hand, we know the overall update is given by  $\Delta^{(t)} = \Delta^{(t-1)} + \eta^*(\mathbf{e}_i \mathbf{e}_j^T + \mathbf{e}_j \mathbf{e}_i^T), \forall (i, j) \in \mathcal{I}_{\text{free}}(\mathbf{Z}^t)$ . Then the coordinate-wise optimality conditions imply that for every updated index  $(i, j) \in \mathcal{I}_{\text{free}}(\mathbf{Z}^t)$ , it holds that

$$|\nabla_{ij} \bar{f}(\Delta^{(t)}; \mathbf{Z}^t) + G_{ij}^{(t)}| = 0.$$

Additionally, for any  $(i, j) \in \mathcal{I}_{\text{fixed}}(\mathbf{Z}^t)$ , we have  $G_{ij}^{(t-1)} = G_{ij}^{(t)}$ . Moreover, by the definition of  $\mathcal{I}_{\text{fixed}}(\mathbf{Z}^t)$ , it holds that  $Z_{ij}^t = 0$  and  $W_{ij}^k \geq |\nabla_{ij} f(\mathbf{Z}^t)|$ . These conditions, together with the expression for the subdifferential in (2) and  $\Delta_{ij}^{(t-1)} = 0$  imply that

$$\inf_{\mathbf{G}^{(t-1)} \in \partial g(\Delta^{(t-1)}; \mathbf{Z}^t)} \sum_{(i,j) \in \mathcal{I}_{\text{fixed}}(\mathbf{Z}^t)} |\nabla_{ij} \bar{f}(\Delta^{(t-1)}; \mathbf{Z}^t) + G_{ij}^{(t-1)}| = 0.$$

Combining this condition with (53) yields

$$\begin{aligned} & \text{dist}(\mathbf{0}, \partial J(\Delta^{(t)}; \mathbf{Z}^t)) \\ &= \inf_{\mathbf{G}^{(t)} \in \partial g(\Delta^{(t)}; \mathbf{Z}^t)} \sum_{(i,j) \in \mathcal{I}_{\text{fixed}}(\mathbf{Z}^t)} |\nabla_{ij} \bar{f}(\Delta^{(t)}; \mathbf{Z}^t) + G_{ij}^{(t)}| \\ &\stackrel{(a)}{\leq} \inf_{\mathbf{G}^{(t-1)} \in \partial g(\Delta^{(t-1)}; \mathbf{Z}^t)} \sum_{(i,j) \in \mathcal{I}_{\text{fixed}}(\mathbf{Z}^t)} \left\{ |\nabla_{ij} \bar{f}(\Delta^{(t)}; \mathbf{Z}^t) - \nabla_{ij} \bar{f}(\Delta^{(t-1)}; \mathbf{Z}^t)| + \right. \\ &\quad \left. |\nabla_{ij} \bar{f}(\Delta^{(t-1)}; \mathbf{Z}^t) + G_{ij}^{(t-1)}| \right\} \\ &= \sum_{(i,j) \in \mathcal{I}_{\text{fixed}}(\mathbf{Z}^t)} |\nabla_{ij} \bar{f}(\Delta^{(t)}; \mathbf{Z}^t) - \nabla_{ij} \bar{f}(\Delta^{(t-1)}; \mathbf{Z}^t)| \tag{54} \\ &\quad + \inf_{\mathbf{G}^{(t-1)} \in \partial g(\Delta^{(t-1)}; \mathbf{Z}^t)} \sum_{(i,j) \in \mathcal{I}_{\text{fixed}}(\mathbf{Z}^t)} |\nabla_{ij} \bar{f}(\Delta^{(t-1)}; \mathbf{Z}^t) + G_{ij}^{(t-1)}| \\ &= \sum_{(i,j) \in \mathcal{I}_{\text{fixed}}(\mathbf{Z}^t)} |\nabla_{ij} \bar{f}(\Delta^{(t)}; \mathbf{Z}^t) - \nabla_{ij} \bar{f}(\Delta^{(t-1)}; \mathbf{Z}^t)| \\ &\leq \|\nabla \bar{f}(\Delta^{(t)}; \mathbf{Z}^t) - \nabla \bar{f}(\Delta^{(t-1)}; \mathbf{Z}^t)\|_1 \\ &\stackrel{(b)}{\leq} l_k^{-2} \|\Delta^{(t)} - \Delta^{(t-1)}\|_1 \leq \epsilon l_k^{-2} \|\Delta^{(t)}\|_1, \end{aligned}$$

where the inequality (a) follows from the triangle inequality, and inequality (b) holds because Proposition 3.2(ii) applies to  $\bar{f}$  as well. Consequently, setting  $C_k := \epsilon l_k^{-2}$  yields the desired result.

As for statement (ii), by the standard convention in convex analysis, it holds that  $Q_k(\mathbf{Z}) = +\infty, \forall \mathbf{Z} \leq \mathbf{0}$ . Note that for any  $\mathbf{Z}, \mathbf{D}$ , scalar  $a \in [0, 1]$  and



$\mathbf{W} \in \mathbb{R}_+^{n \times n}$ , it follows from the convexity and positive homogeneity of the  $\|\cdot\|_{1, \mathbf{W}^k}$  that

$$\|\mathbf{Z} + a\mathbf{D}\|_{1, \mathbf{W}} = \|a(\mathbf{Z} + \mathbf{D}) + (1-a)\mathbf{Z}\|_{1, \mathbf{W}} \leq a\|\mathbf{Z} + \mathbf{D}\|_{1, \mathbf{W}} + (1-a)\|\mathbf{Z}\|_{1, \mathbf{W}}. \quad (55)$$

Then, we have for all  $\beta_t \in (0, 1]$  that

$$\begin{aligned} & Q_k(\mathbf{Z}^t + \beta_t \mathbf{D}^t) - Q_k(\mathbf{Z}^t) \\ &= f(\mathbf{Z}^t + \beta_t \mathbf{D}^t) - f(\mathbf{Z}^t) + \rho\|\mathbf{Z}^t + \beta_t \mathbf{D}^t\|_{1, \mathbf{W}^k} - \rho\|\mathbf{Z}^t\|_{1, \mathbf{W}^k} \\ &\stackrel{(a)}{\leq} \beta_t \langle \nabla f(\mathbf{Z}^t), \mathbf{D}^t \rangle - \beta_t \rho(\|\mathbf{Z}^t\|_{1, \mathbf{W}^k} - \|\mathbf{Z}^t + \mathbf{D}^t\|_{1, \mathbf{W}^k}) \\ &\quad + \beta_t \int_0^1 \langle \nabla f(\mathbf{Z}^t + \theta \beta_t \mathbf{D}^t) - \nabla f(\mathbf{Z}^t), \mathbf{D}^t \rangle d\theta \stackrel{(b)}{\leq} \beta_t \Delta^t + \frac{\beta_t^2}{2l_k^2} \|\mathbf{D}^t\|_F^2, \end{aligned} \quad (56)$$

where inequality (a) follows from the mean value theorem and (55), and inequality (b) holds by (15) together with the Lipschitz continuity of  $\nabla f$  as stated in (16). Then, by the strong convexity of  $J(\mathbf{D}; \mathbf{Z}^t)$ , it follows from (54) that there exists  $\bar{\mathbf{S}}^t \in \partial J(\mathbf{D}^t; \mathbf{Z}^t)$  such that

$$\begin{aligned} J(\mathbf{D}^t; \mathbf{Z}^t) - J(\mathbf{0}; \mathbf{Z}^t) &\leq \langle \bar{\mathbf{S}}^t, \mathbf{D}^t \rangle - \frac{u_k^{-2}}{2} \|\mathbf{D}^t\|_F^2 \\ &\leq \|\bar{\mathbf{S}}^t\|_F \|\mathbf{D}^t\|_F^2 - \frac{u_k^{-2}}{2} \|\mathbf{D}^t\|_F^2 \leq \left( \frac{\epsilon}{l_k^2} - \frac{1}{2u_k^2} \right) \|\mathbf{D}^t\|_F^2. \end{aligned}$$

Rearranging gives

$$\begin{aligned} & \left( \frac{\epsilon}{l_k^2} - \frac{1}{2u_k^2} \right) \|\mathbf{D}^t\|_F^2 \\ &\geq J(\mathbf{D}^t; \mathbf{Z}^t) - J(\mathbf{0}; \mathbf{Z}^t) \\ &= \langle \nabla f(\mathbf{Z}^t), \mathbf{D}^t \rangle + \frac{1}{2} \langle \text{vec}(\mathbf{D}^t), \nabla^2 f(\mathbf{Z}^t) \text{vec}(\mathbf{D}^t) \rangle + \rho\|\mathbf{Z}^t + \mathbf{D}^t\|_{1, \mathbf{W}^k} - \rho\|\mathbf{Z}^t\|_{1, \mathbf{W}^k} \\ &= \Delta^t + \frac{1}{2} \langle \text{vec}(\mathbf{D}^t), \nabla^2 f(\mathbf{Z}^t) \text{vec}(\mathbf{D}^t) \rangle \stackrel{(a)}{\geq} \Delta^t + \frac{u_k^{-2}}{2} \|\mathbf{D}^t\|_F^2, \end{aligned}$$

where inequality (a) holds by (17). Then, we have

$$\Delta^t \leq (\epsilon l_k^{-2} - u_k^{-2}) \|\mathbf{D}^t\|_F^2. \quad (57)$$

This, together with (56), leads to the desired (19).

On the other hand, combining (57) with the fact that  $\Delta^t < 0$ , we deduce from (56) that

$$\begin{aligned} Q_k(\mathbf{Z}^t + \beta_t \mathbf{D}^t) - Q_k(\mathbf{Z}^t) &\leq \beta_t \Delta^t + \frac{\beta_t^2}{2l_k^2} \|\mathbf{D}^t\|_F^2 \\ &\leq \left( \beta_t - \frac{l_k^{-2} \beta_t^2}{2(u_k^{-2} - \epsilon l_k^{-2})} \right) \Delta^t. \end{aligned} \quad (58)$$

Therefore, (14) is satisfied whenever

$$\beta_t \Delta^t - \frac{l_k^{-2}(\beta_t)^2}{2(u_k^{-2} - \epsilon l_k^{-2})} \Delta^t \leq \beta_t \sigma \Delta^t,$$

which indicates that (14) holds whenever

$$\beta_t \leq 2(1 - \sigma)(u_k^{-2} - \epsilon l_k^{-2})l_k^2.$$

Hence, we deduce that (20) holds, where  $\pi \in (0, 1)$  is introduced to mitigate potential undershooting in the backtracking procedure. In addition, rearranging (19) and letting  $t \rightarrow +\infty$ , it follows from (20) that  $\|\mathbf{Z}^{t+1} - \mathbf{Z}^t\|_F^2 \rightarrow 0$ , as desired.

(iii). For notational ease, define the function

$$\iota(s) := J(s\mathbf{D}^t; \mathbf{Z}^t), \quad \forall s \in [0, 1].$$

We first prove that

$$\iota(a) \geq \iota(1) - (1 - a)C_k \|\mathbf{D}^t\|_1^2, \quad \forall a \in (0, 1]. \quad (59)$$

Since  $g(s\mathbf{D}^t)$  is convex with respect to  $s \in [0, 1]$ , we denote its right-directional derivative at  $s = 1$  by

$$\text{Dir}^+ g(\mathbf{Z}^t + \mathbf{D}^t; \mathbf{D}^t) := \lim_{h \downarrow 0} \frac{g(\mathbf{Z}^t + (1 + h)\mathbf{D}^t; \mathbf{Z}^t) - g(\mathbf{Z}^t + \mathbf{D}^t; \mathbf{Z}^t)}{h}.$$

Thus, we obtain  $\iota'_+(1) = \frac{d}{ds} \bar{f}(s\mathbf{D}^t; \mathbf{Z}^t)|_{s=1} + \text{Dir}^+ g(\mathbf{Z}^t + \mathbf{D}^t; \mathbf{D}^t)$ . Due to the convexity of  $g$ , there exists a subgradient  $\mathbf{G}^t \in \partial g(\mathbf{Z}^t + \mathbf{D}^t; \mathbf{Z}^t)$  such that  $\text{Dir}^+ g(\mathbf{Z}^t + \mathbf{D}^t; \mathbf{D}^t) = \langle \mathbf{G}^t, \mathbf{D}^t \rangle$ . Therefore,  $\iota'_+(1) = \langle \nabla f(\mathbf{Z}^t) + \nabla^2 f(\mathbf{Z}^t)\mathbf{D}^t + \mathbf{G}^t, \mathbf{D}^t \rangle$ .

Let  $\mathbf{J}^{(t)} \in \partial J(\mathbf{D}^t; \mathbf{Z}^t)$  such that  $\iota'_+(1) = \langle \mathbf{J}^{(t)}, \mathbf{D}^t \rangle$ . Then, by applying the inexactness condition (18) together with the matrix Hölder inequality, we deduce that

$$\iota'_+(1) \leq C_k \|\mathbf{D}^t\|_1^2.$$

Exploiting the convexity of  $\iota$ , for any  $a \in (0, 1)$  we have

$$\iota(1) \leq \iota(a) + (1 - a)\iota'_+(1) \leq \iota(a) + (1 - a)C_k \|\mathbf{D}^t\|_1^2,$$

Rearranging the above inequality yields the desired (59). Moreover, by (55) and (59), we have

$$\begin{aligned} & \langle \text{vec}(\nabla f(\mathbf{Z}^t)), \text{vec}(\mathbf{D}^t) \rangle + \frac{1}{2} \langle \text{vec}(\mathbf{D}^t), \nabla^2 f(\mathbf{Z}^t) \text{vec}(\mathbf{D}^t) \rangle + \rho \|\mathbf{Z}^t + \mathbf{D}^t\|_{1, \mathbf{W}^k} \\ & \leq a \langle \text{vec}(\nabla f(\mathbf{Z}^t)), \text{vec}(\mathbf{D}^t) \rangle + \frac{a^2}{2} \langle \text{vec}(\mathbf{D}^t), \nabla^2 f(\mathbf{Z}^t) \text{vec}(\mathbf{D}^t) \rangle + \rho a \|\mathbf{Z}^t + \mathbf{D}^t\|_{1, \mathbf{W}^k} \\ & \quad + (1 - a)\rho \|\mathbf{Z}^t\|_{1, \mathbf{W}^k} + (1 - a)C_k \|\mathbf{D}^t\|_1^2. \end{aligned}$$

Rearranging gives

$$0 \geq (1-a) [\langle \text{vec}(\nabla f(\mathbf{Z}^t)), \text{vec}(\mathbf{D}^t) \rangle + \rho \|\mathbf{Z}^t + \mathbf{D}^t\|_{1, \mathbf{W}^k} - \rho \|\mathbf{Z}^t\|_{1, \mathbf{W}^k} - C_k \|\mathbf{D}^t\|_1^2] \\ + \frac{1}{2}(1-a^2) \langle \text{vec}(\mathbf{D}^t), \nabla^2 f(\mathbf{Z}^t) \text{vec}(\mathbf{D}^t) \rangle.$$

Dividing both sides of the above inequality by  $(1-a)$  and then taking  $a \uparrow 1$ , we have

$$\Delta^t \leq -\langle \text{vec}(\mathbf{D}^t), \nabla^2 f(\mathbf{Z}^t) \text{vec}(\mathbf{D}^t) \rangle + C_k \|\mathbf{D}^t\|_1^2. \quad (60)$$

Define  $\tilde{f}(\varrho) = f(\mathbf{Z}^t + \varrho \mathbf{D}^t)$ . We know that

$$|\tilde{f}''(\varrho) - \tilde{f}''(0)| = |\text{vec}(\mathbf{D}^t)^T (\nabla^2 f(\mathbf{Z}^t + t\mathbf{D}^t) - \nabla^2 f(\mathbf{Z}^t)) \text{vec}(\mathbf{D}^t)| \\ \stackrel{(a)}{\leq} \|\nabla^2 f(\mathbf{Z}^t + t\mathbf{D}^t) - \nabla^2 f(\mathbf{Z}^t)\|_F \|\mathbf{D}^t\|_F^2 \stackrel{(b)}{\leq} \varrho l_k^{-2} \|\mathbf{D}^t\|_F^3, \quad (61)$$

where (a) holds by the Cauchy-Schwartz inequality and (b) holds by the mean value theorem and Proposition 3.2(ii). Then we have that

$$\tilde{f}''(\varrho) \leq \tilde{f}''(0) + \varrho l_k^{-2} \|\mathbf{D}^t\|_F^3 = \text{vec}(\mathbf{D}^t)^T \nabla^2 f(\mathbf{Z}^t) \text{vec}(\mathbf{D}^t) + \varrho l_k^{-2} \|\mathbf{D}^t\|_F^3.$$

We next integrate both sides of the above inequality with respect to  $\varrho \in [0, 1]$  to obtain an upper bound on  $\tilde{f}'(\varrho)$ . Note that

$$\int_0^\varrho \tilde{f}''(\varrho) d\varrho \leq \int_0^\varrho (\tilde{f}''(0) + \varrho l_k^{-2} \|\mathbf{D}^t\|_F^3) d\varrho,$$

then we have

$$\tilde{f}'(\varrho) \leq \tilde{f}'(0) + \varrho \text{vec}(\mathbf{D}^t)^T \nabla^2 f(\mathbf{Z}^t) \text{vec}(\mathbf{D}^t) + \frac{1}{2} \varrho^2 l_k^{-2} \|\mathbf{D}^t\|_F^3. \quad (62)$$

We again integrate both sides of (62) with respect to  $\varrho \in [0, 1]$  to obtain an upper bound on  $\tilde{f}(\varrho)$ , and we hence have

$$\tilde{f}(\varrho) \leq \tilde{f}(0) + \varrho \langle \nabla f(\mathbf{Z}^t), \mathbf{D}^t \rangle + \frac{1}{2} \varrho^2 \text{vec}(\mathbf{D}^t)^T \nabla^2 f(\mathbf{Z}^t) \text{vec}(\mathbf{D}^t) + \frac{1}{6} \varrho^3 l_k^{-2} \|\mathbf{D}^t\|_F^3. \quad (63)$$

On the other hand,

$$\begin{aligned}
Q_k(\mathbf{Z}^t + \mathbf{D}^t) &= f(\mathbf{Z}^t + \mathbf{D}^t) + \rho \|\mathbf{Z}^t + \mathbf{D}^t\|_{1, \mathbf{W}^k} \\
&\stackrel{(a)}{\leq} f(\mathbf{Z}^t) + \rho \|\mathbf{Z}^t\|_{1, \mathbf{W}^k} + \langle \nabla f(\mathbf{Z}^t), \mathbf{D}^t \rangle + \rho \|\mathbf{Z}^t + \mathbf{D}^t\|_{1, \mathbf{W}^k} - \rho \|\mathbf{Z}^t\|_{1, \mathbf{W}^k} \\
&\quad + \frac{1}{2} \text{vec}(\mathbf{D}^t)^T \nabla^2 f(\mathbf{Z}^t) \text{vec}(\mathbf{D}^t) + \frac{1}{6} l_k^{-2} \|\mathbf{D}^t\|_F^3 \\
&\leq Q_k(\mathbf{Z}^t) + \Delta^t + \frac{1}{2} \text{vec}(\mathbf{D}^t)^T \nabla^2 f(\mathbf{Z}^t) \text{vec}(\mathbf{D}^t) + \frac{1}{6} l_k^{-2} \|\mathbf{D}^t\|_F^3 \\
&\stackrel{(b)}{\leq} Q_k(\mathbf{Z}^t) + \frac{\Delta^t}{2} + \frac{n \epsilon l_k^{-2}}{2} \|\mathbf{D}^t\|_F^2 + \frac{1}{6} l_k^{-2} \|\mathbf{D}^t\|_F^3, \\
&\stackrel{(c)}{\leq} Q_k(\mathbf{Z}^t) + \left( \left( \frac{1}{2} - \frac{n \epsilon l_k^{-2}}{2(u_k^{-2} - \epsilon l_k^{-2})} \right) - \frac{l_k^{-2}}{6(u_k^{-2} - \epsilon l_k^{-2})} \|\mathbf{D}^t\|_F \right) \Delta^t \\
&\stackrel{(d)}{\leq} Q_k(\mathbf{Z}^t) + \sigma \Delta^t,
\end{aligned}$$

where (a) follows from (63) with  $\varrho = 1$ , (b) follows from (60) and the fact that  $\|\mathbf{D}\|_1 \leq n \|\mathbf{D}\|_F, \forall \mathbf{D} \in \mathbb{S}^n$ , (c) follows from (57) and (d) holds provided

$$\left( \frac{1}{2} - \frac{n \epsilon l_k^{-2}}{2(u_k^{-2} - \epsilon l_k^{-2})} \right) - \frac{l_k^{-2}}{6(u_k^{-2} - \epsilon l_k^{-2})} \|\mathbf{D}^t\|_F > \sigma,$$

which is guaranteed by  $\sigma \in (0, 0.5 - \epsilon_\sigma)$ , by (21), by the boundedness of  $\beta_t$  in statement (ii) and for  $0 < \epsilon \ll 1$ . Therefore, the line-search condition (14) holds with  $\beta_t = 1$ . This completes the proof.  $\square$

## A.2 Proof of Proposition 3.5

*Proof.* (i). By the concavity of  $\phi(\cdot)$  over  $\mathbb{R}_{++}$ , we have for each  $(i, j) \in [n] \times [n]$  that

$$\phi(|X_{ij}^{k+1}| + \mathcal{E}_{ij}^k) \leq \phi(|X_{ij}^k| + \mathcal{E}_{ij}^k) + \phi'(|X_{ij}^k| + \mathcal{E}_{ij}^k)(|X_{ij}^{k+1}| - |X_{ij}^k|). \quad (64)$$

Then

$$\begin{aligned}
& F(\mathbf{X}^k, \boldsymbol{\varepsilon}^k) - F(\mathbf{X}^{k+1}, \boldsymbol{\varepsilon}^k) \\
&= f(\mathbf{X}^k) - f(\mathbf{X}^{k+1}) + \rho \sum_{ij} \phi(|X_{ij}^k| + \varepsilon_{ij}^k) - \rho \sum_{ij} \phi(|X_{ij}^{k+1}| + \varepsilon_{ij}^k) \\
&\stackrel{\text{Eq. (64)}}{\geq} f(\mathbf{X}^k) - f(\mathbf{X}^{k+1}) + \rho \sum_{ij} W_{ij}^k (|X_{ij}^k| - |X_{ij}^{k+1}|) \\
&\stackrel{(a)}{=} f(\mathbf{X}^k) + \rho \sum_{ij} W_{ij}^k |X_{ij}^k| - \left( f((1-\alpha)\mathbf{X}^k + \alpha\mathbf{Y}^k) + \rho \sum_{ij} W_{ij}^k |\alpha Y_{ij}^k + (1-\alpha)X_{ij}^k| \right) \\
&\stackrel{(b)}{\geq} \alpha \left( f(\mathbf{X}^k) - f(\mathbf{Y}^k) + \rho \sum_{ij} W_{ij}^k (|X_{ij}^k| - |Y_{ij}^k|) \right) = \alpha(Q_k(\mathbf{X}^k) - Q_k(\mathbf{Y}^k)),
\end{aligned} \tag{65}$$

where inequality (a) uses the identity  $\mathbf{X}^{k+1} = (1-\alpha)\mathbf{X}^k + \alpha\mathbf{Y}^k$  and inequality (b) holds due to the convexity of  $f$  and  $|\cdot|$ . Let  $K_k \in \mathbb{N}$  denote the number of iterations required to obtain  $\mathbf{Y}^k$  for the  $k$ th subproblem ( $\mathcal{P}_{\text{sub}}$ ). Starting from the initialization  $\mathbf{Z}^0 = \mathbf{X}^k$ , it follows that

$$\begin{aligned}
Q_k(\mathbf{X}^k) - Q_k(\mathbf{Y}^k) &= \sum_{t=0}^{K_k-1} Q_k(\mathbf{Z}^t) - Q_k(\mathbf{Z}^{t+1}) \\
&\stackrel{\text{Eq. (19)}}{\geq} \sum_{t=0}^{K_k-1} \left( \beta_t(u_k^{-2} - \epsilon l_k^{-2}) - \frac{\beta_t^2}{2l_k^2} \right) \|\mathbf{Z}^t - \mathbf{Z}^{t+1}\|_F^2 \\
&\stackrel{(a)}{\geq} \sum_{t=0}^{K_k-1} \left( \tilde{\beta}(u^{-2} - \epsilon l^{-2}) - \frac{1}{2l^2} \right) \|\mathbf{Z}^t - \mathbf{Z}^{t+1}\|_F^2 \\
&\stackrel{(b)}{\geq} \left( \tilde{\beta}(u^{-2} - \epsilon l^{-2}) - \frac{1}{2l^2} \right) \|\mathbf{X}^k - \mathbf{Y}^k\|_F^2 \\
&\stackrel{(c)}{=} \left( \frac{\tilde{\beta}(u^{-2} - \epsilon l^{-2})}{\alpha^2} - \frac{1}{2l^2\alpha^2} \right) \|\mathbf{X}^k - \mathbf{X}^{k+1}\|_F^2,
\end{aligned}$$

where (a) follows from  $\beta_t \leq 1$ , (23) and  $\tilde{\beta} := \min \{1, 2\pi(1-\sigma)(u^{-2} - \epsilon l^{-2})l^2\}$ , (b) holds by the triangle inequality, and (c) results from the identity  $\mathbf{Y}^k = \frac{1}{\alpha}(\mathbf{X}^{k+1} - \mathbf{X}^k) + \mathbf{X}^k$ . This, together with (65), gives

$$\begin{aligned}
F(\mathbf{X}^k, \boldsymbol{\varepsilon}^k) - F(\mathbf{X}^{k+1}, \boldsymbol{\varepsilon}^{k+1}) &\geq F(\mathbf{X}^k, \boldsymbol{\varepsilon}^k) - F(\mathbf{X}^{k+1}, \boldsymbol{\varepsilon}^k) \\
&\geq \left( \frac{\tilde{\beta}(u^{-2} - \epsilon l^{-2})}{\alpha} - \frac{1}{2l^2\alpha} \right) \|\mathbf{X}^k - \mathbf{X}^{k+1}\|_F^2.
\end{aligned} \tag{66}$$

(ii) Rearranging and summing up both sides of inequality (66) from  $k = 0$  to

$\bar{k} - 1$  gives

$$\begin{aligned}
& \left( \frac{\tilde{\beta}(u^{-2} - \epsilon l^{-2})}{\alpha} - \frac{1}{2l^2\alpha} \right) \sum_{k=0}^{\bar{k}-1} \|\mathbf{X}^k - \mathbf{X}^{k+1}\|_F^2 \\
& \leq \sum_{k=0}^{\bar{k}-1} F(\mathbf{X}^k, \boldsymbol{\varepsilon}^k) - F(\mathbf{X}^{k+1}, \boldsymbol{\varepsilon}^{k+1}) \\
& \leq F(\mathbf{X}^0, \boldsymbol{\varepsilon}^0) - F(\mathbf{X}^{\bar{k}}, \boldsymbol{\varepsilon}^{\bar{k}}) \stackrel{(a)}{\leq} F(\mathbf{X}^0, \boldsymbol{\varepsilon}^0) - \inf_{\mathbf{X} \in \text{Lev}_{F(\mathbf{X}, \boldsymbol{\varepsilon})}(\mathbf{X}^0, \boldsymbol{\varepsilon}^0)} F(\mathbf{X}, \boldsymbol{\varepsilon}) < +\infty,
\end{aligned}$$

where inequality (a) holds since the value  $\inf_{\mathbf{X} \in \text{Lev}_{F(\mathbf{X}, \boldsymbol{\varepsilon})}(\mathbf{X}^0, \boldsymbol{\varepsilon}^0)} F(\mathbf{X}, \boldsymbol{\varepsilon})$  is finite by [21, Theorem 1.9]. Let  $\bar{k} \rightarrow +\infty$ , and it leads us to the desired results.

(iii) It follows from  $\mathbf{X}^0 \in \chi$  and  $\{\mathbf{Y}^k\} \subset \chi$ ,  $\forall k$  (by Proposition 3.2(i)) that  $\{\mathbf{X}^k\} \subset \chi$ . This completes the proof.  $\square$

### A.3 Proof of Lemma 3.8

*Proof.* (i) It is immediate that  $\mathcal{M}(\mathbf{X}^*)$  forms a smooth manifold in a neighborhood of  $\mathbf{X}^*$ . We next claim that  $\Phi$  is partly smooth at  $\mathbf{X}^*$  relative to  $\mathcal{M}(\mathbf{X}^*)$ . To see this, observe first that for all  $\mathbf{X} \in \mathcal{M}(\mathbf{X}^*)$  we have

$$\begin{aligned}
\Phi|_{\mathcal{M}(\mathbf{X}^*)}(\mathbf{X}) &= \sum_{(i,j) \in \mathcal{I}(\mathbf{X}^*)} \phi(|X_{ij}|) + \sum_{(i,j) \in \mathcal{Z}(\mathbf{X}^*)} \phi(|X_{ij}|) \\
&= \sum_{(i,j) \in \mathcal{I}(\mathbf{X}^*)} \phi(\text{sgn}(\mathbf{X}_{ij}^*) X_{ij}),
\end{aligned} \tag{67}$$

so  $\Phi|_{\mathcal{M}(\mathbf{X}^*)}$  is smooth. Second, since  $\Phi$  is continuous by Assumption 1.1 and  $\Phi|_{\mathcal{M}(\mathbf{X}^*)}$  is smooth, we then know that  $\Phi$  is regular [21, Example 7.28] at each  $\mathbf{X} \in \mathcal{M}(\mathbf{X}^*)$  with  $\frac{\partial \Phi}{\partial X_{ij}} = \left\{ \frac{d\phi(\text{sgn}(\mathbf{X}_{ij}^*) X_{ij})}{X_{ij}} \right\}$ , for any  $(i, j) \in \mathcal{I}(\mathbf{X}^*)$  by [21, Exercise 8.8]. Next, we check the local normal sharpness condition. Routine calculation (e.g., using [21, Example 6.8]) shows that at any point  $\mathbf{X} \in \mathcal{M}(\mathbf{X}^*)$ , we have

$$\begin{aligned}
N_{\mathcal{M}(\mathbf{X}^*)}(\mathbf{X}) &= \{\mathbf{H} \in \mathbb{S}^n \mid H_{ij} = 0, \forall (i, j) \in \mathcal{I}(\mathbf{X}^*)\}, \\
\text{par } \partial \Phi(\mathbf{X}) &= \{\mathbf{H} \in \mathbb{S}^n \mid H_{ij} = 0, \forall (i, j) \in \mathcal{I}(\mathbf{X}^*)\}.
\end{aligned} \tag{68}$$

Hence the normal space is parallel to the subdifferential. Finally, the subdifferential map  $\partial \Phi$  is continuous relative to  $\mathcal{M}(\mathbf{X}^*)$  since  $\mathcal{M}(\mathbf{X}^*)$  contains  $\Phi$  to a smooth subspace by (67). On the other hand, by [21, Proposition 13.34], we also know that  $\Phi$  is prox-regular at  $\mathbf{X}^*$  relative to  $\mathcal{M}(\mathbf{X}^*)$ . Consequently, it follows from

[22, Corollary 4.7] and  $f \in \mathcal{C}^1$  that  $F$  is partly smooth at  $\mathbf{X}^*$  relative to  $\mathcal{M}(\mathbf{X}^*)$ , and prox-regular there.

(ii) Theorem 3.6(iii), together with  $f \in \mathcal{C}^1$  and Assumption 1.1, gives  $F(\mathbf{X}^k) \rightarrow F(\mathbf{X}^*)$ . Thus, the premises in Proposition 1.9 on  $F$  are satisfied. Note that

$$\begin{aligned} \text{dist}(\mathbf{0}, \partial F(\mathbf{X}^k)) \rightarrow 0 &\iff \text{dist}(-\nabla f(\mathbf{X}^k), \partial \Phi(\mathbf{X}^k)) \rightarrow 0 \\ &\stackrel{(a)}{\iff} \text{dist}(-\nabla f(\mathbf{X}^*), \partial \Phi(\mathbf{X}^*)) = 0, \end{aligned} \tag{69}$$

where (a) follows from  $\nabla f(\mathbf{X}^k) \rightarrow \nabla f(\mathbf{X}^*)$  and the non-degenerate condition (30). By Proposition 1.9, we know that  $\mathbf{X}^k \in \mathcal{M}(\mathbf{X}^*)$  for all sufficiently large  $k \in \mathbb{N}$ . This completes the proof.  $\square$

## B Tests with other nonconvex regularizers

For completeness, we provide additional experimental results with  $p = 0.3$  (Figure 6, Figure 8 and Table 4) and  $p = 0.7$  (Figure 7, Figure 9 and Table 5). Moreover, we validate the performance of DIIR-QUIC by incorporating nonconvex SCAD and MCP penalties. Synthetic precision-matrix estimation problems of size  $n \in \{500, 1000, 2000\}$  were solved under both penalty types. Figure 10 summarizes the results.

## References

- [1] D. Bertsimas, J. Lamperski, and J. Pauphilet, “Certifiably optimal sparse inverse covariance estimation,” *Mathematical Programming*, vol. 184, no. 1, pp. 491–530, 2020.
- [2] M. Gulliksson and S. Mazur, “An iterative approach to ill-conditioned optimal portfolio selection,” *Computational Economics*, vol. 56, no. 4, pp. 773–794, 2020.
- [3] T. T. Cai, H. Li, W. Liu, and J. Xie, “Covariate-adjusted precision matrix estimation with an application in genetical genomics,” *Biometrika*, vol. 100, no. 1, pp. 139–156, 2013.
- [4] J. Friedman, T. Hastie, and R. Tibshirani, “Sparse inverse covariance estimation with the graphical lasso,” *Biostatistics*, vol. 9, no. 3, pp. 432–441, 2008.

Table 3: Examples of the function  $\Phi$  with tuning parameters  $p \in (0, 1)$ ,  $\varpi > 0$ , and SCAD/MCP parameters  $(\lambda, a)$  or  $(\lambda, \gamma)$ . Here  $\phi'(|X_{ij}|)$  is the derivative at  $|X_{ij}|$ , and  $(\phi')^{-1}(s)$  is the inverse mapping defined on the strictly decreasing portion  $s \in (0, \phi'(0^+))$ .

Regularizer	$\Phi(\mathbf{X})$	$\phi'( X_{ij} )$	$\phi'(0^+)$	$\phi'(\infty)$	$(\phi')^{-1}(s)$	Flat-tail $t^*$
$\ell_p$ (quasi-norm) [13]	$\sum_{ij}  X_{ij} ^p$	$p  X_{ij} ^{p-1}$	$+\infty$	0	$\left(\frac{s}{p}\right)^{1/(p-1)}$	$\infty$
Log-sum [42]	$\sum_{ij} \log(1 + \frac{ X_{ij} }{\varpi})$	$\frac{1}{ X_{ij}  + \varpi}$	$\frac{1}{\varpi}$	0	$\frac{1}{s} - \varpi$	$\infty$
Geman [43]	$\sum_{ij} \frac{ X_{ij} }{ X_{ij}  + \varpi}$	$\frac{\varpi}{( X_{ij}  + \varpi)^2}$	$\frac{1}{\varpi}$	0	$\sqrt{\frac{\varpi}{s}} - \varpi$	$\infty$
Arctan [44]	$\sum_{ij} \arctan(\frac{ X_{ij} }{\varpi})$	$\frac{\varpi}{\varpi^2 +  X_{ij} ^2}$	$\frac{1}{\varpi}$	0	$\sqrt{\frac{\varpi}{s} - \varpi^2}$	$\infty$
Exp [45]	$\sum_{ij} (1 - e^{- X_{ij} /\varpi})$	$\frac{1}{\varpi} e^{- X_{ij} /\varpi}$	$\frac{1}{\varpi}$	0	$-\varpi \ln(\varpi s)$	$\infty$
SCAD [14]	$\sum_{ij} \phi_{\text{SCAD}}( X_{ij} )$	$\begin{cases} \lambda, &  X_{ij}  \leq \lambda, \\ \frac{a\lambda -  X_{ij} }{a-1}, & \lambda <  X_{ij}  \leq a\lambda, \\ 0, &  X_{ij}  > a\lambda, \end{cases}$	$\lambda$	0	$a\lambda - (a-1)s$	$a\lambda$
MCP [15]	$\sum_{ij} \phi_{\text{MCP}}( X_{ij} )$	$\begin{cases} \lambda - \frac{ X_{ij} }{\gamma}, &  X_{ij}  \leq \gamma\lambda, \\ 0, &  X_{ij}  > \gamma\lambda, \end{cases}$	$\lambda$	0	$\gamma(\lambda - s)$	$\gamma\lambda$



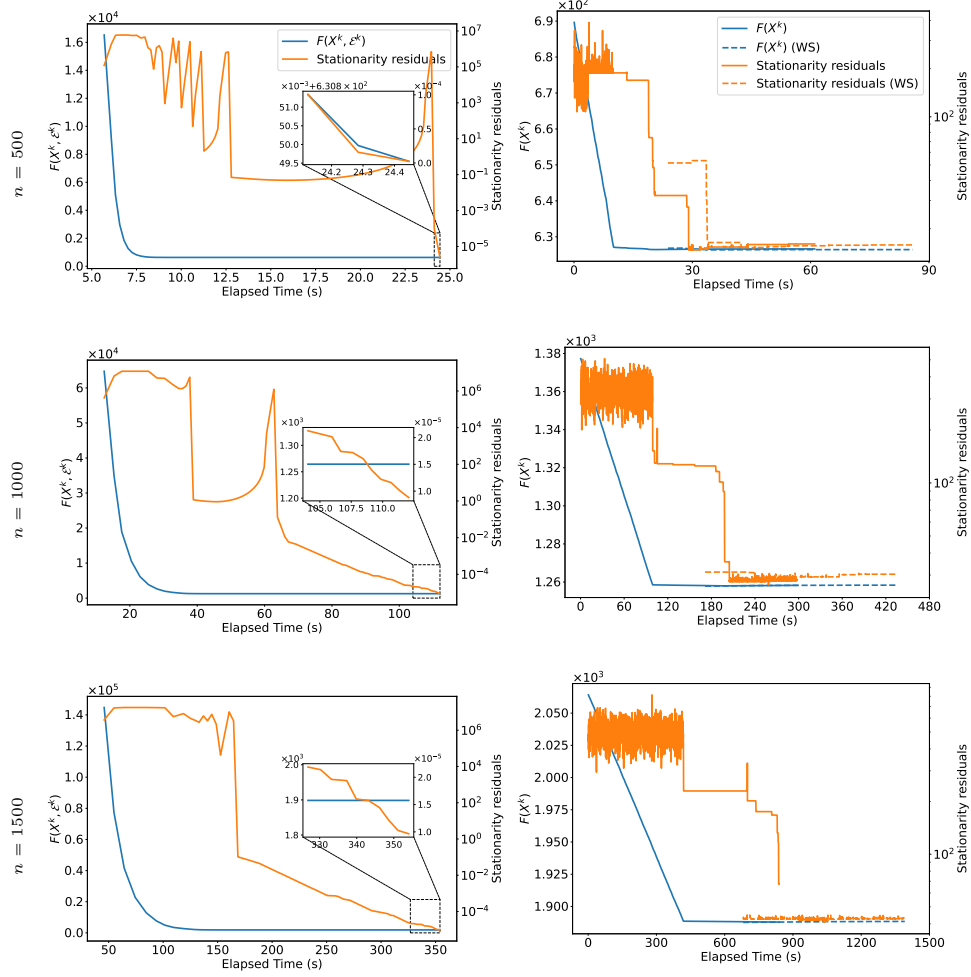


Figure 6: Convergence plots of objective value and stationarity residual (plotted on a logarithmic scale) versus elapsed time (seconds) for DIIR-QUIC and  $\ell_p$  COV on a tridiagonal precision matrix with  $p = 0.3$ .

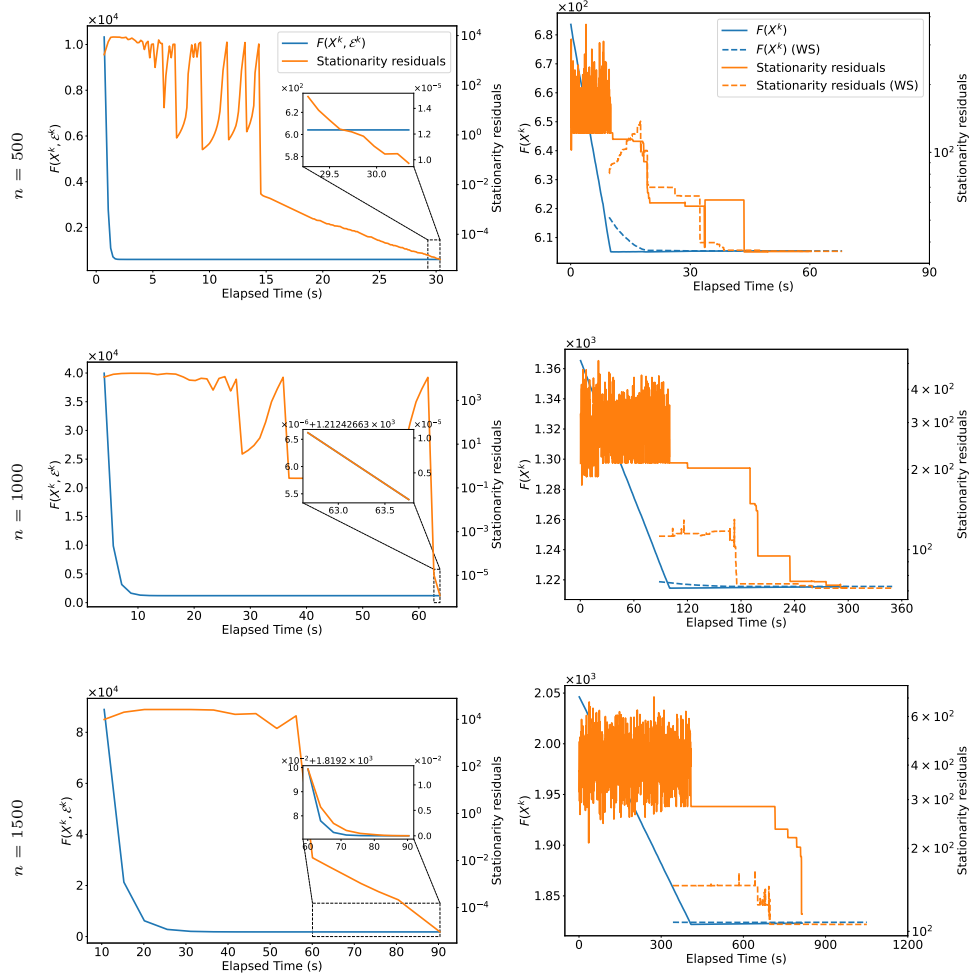


Figure 7: Convergence plots of objective value and stationarity residual (plotted on a logarithmic scale) versus elapsed time (seconds) for DIIR-QUIC and  $\ell_p$ COV on a tridiagonal precision matrix with  $p = 0.7$ .

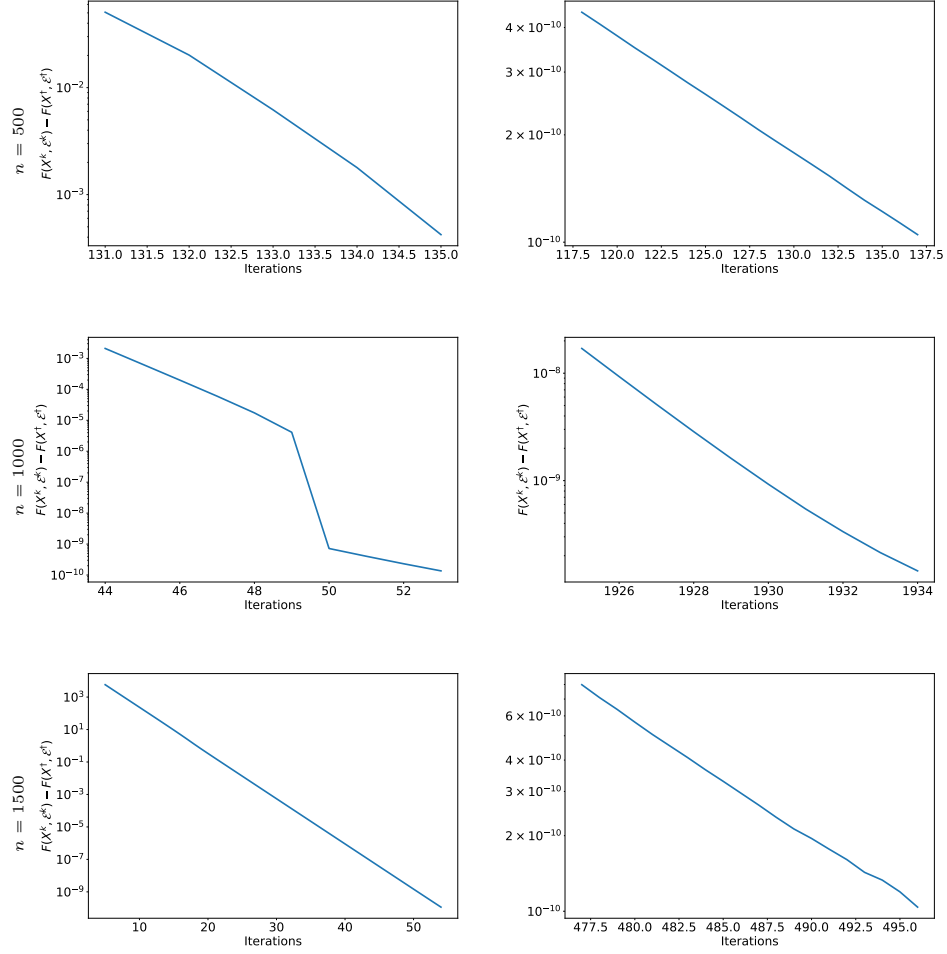


Figure 8: Q-linear convergence of the perturbed objective error  $F(\mathbf{X}^k, \mathbf{E}^k) - F(\mathbf{X}^\dagger, \mathbf{E}^\dagger)$  plotted versus iteration for DIIR-QUIC on tridiagonal precision matrices and clustered matrices with  $p = 0.3$ , across varying matrix dimensions. The error is displayed on a logarithmic scale, and only values above  $10^{-8}$  are shown.

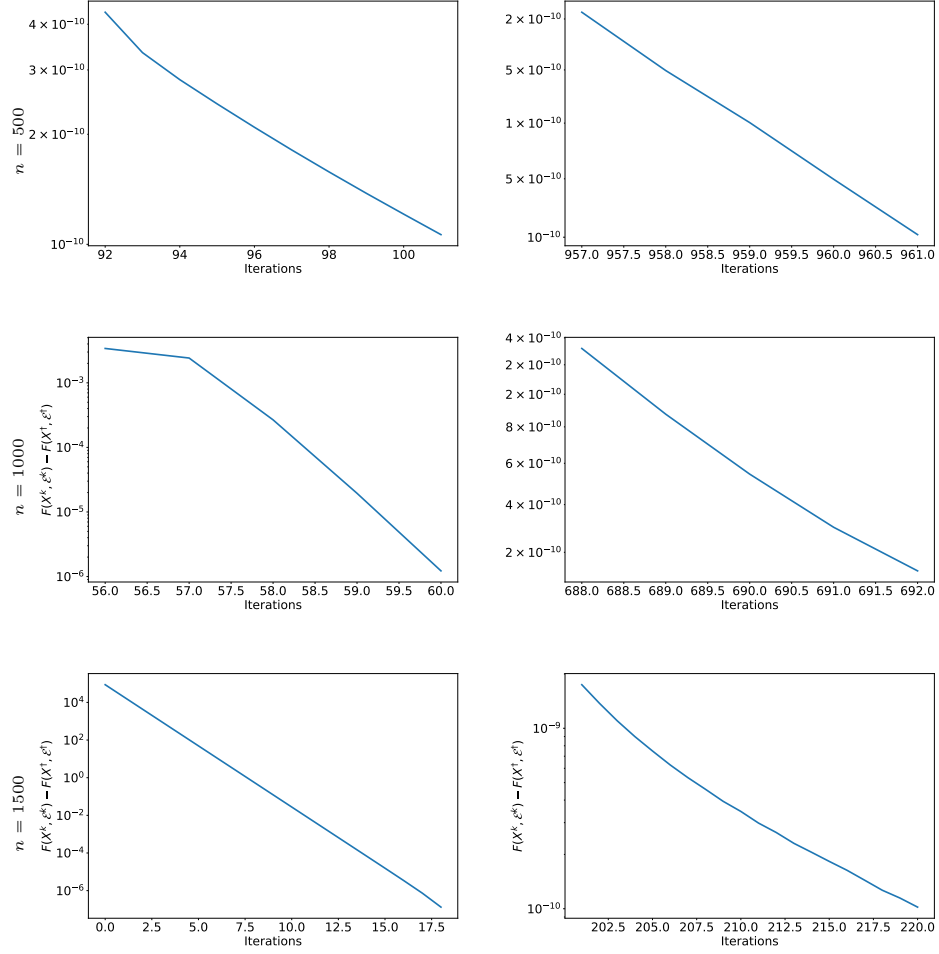


Figure 9: Q-linear convergence of the perturbed objective error  $F(\mathbf{X}^k, \mathbf{E}^k) - F(\mathbf{X}^*, \mathbf{E}^*)$  plotted versus iteration for DIIR-QUIC on tridiagonal precision matrices and clustered matrices with  $p = 0.7$ , across varying matrix dimensions. The error is displayed on a logarithmic scale, and only values above  $10^{-8}$  are shown.

Table 4: Comparison of DIIR-QUIC and  $\ell_p$ COV on synthetic covariance estimation. For each method, the table reports the mean ( $\pm$  one standard deviation) over 10 independent trials under two covariance structures—tridiagonal (Tri) and clustered (Clu)—with  $p = 0.3$  at dimensions  $n = 500, 1000$ , and 1500. Boldface entries indicate the better performance.

Data	$n$	Algorithm	$F_1$	Time (s)	$\text{nnz}(\Sigma^{-1})$	$\text{nnz}(X^\dagger)$	Loss Q	KL Loss	Sensitivity	Specificity	MCC
tri	500	DIIR-QUIC	0.959 (0.002)	<b>20.296</b> ( <b>6.954</b> )	1498.000 (0.000)	1513.200 (7.111)	0.011 (0.000)	0.027 (0.001)	0.964 (0.001)	1.000 (0.000)	0.959 (0.002)
		$\ell_p$ COV	0.995 (0.001)	59.292 (1.507)	1498.000 (0.000)	1492.800 (2.857)	0.009 (0.000)	0.013 (0.001)	0.994 (0.002)	1.000 (0.000)	0.995 (0.001)
		$\ell_p$ COV (WS)	<b>0.999</b> ( <b>0.001</b> )	78.470 (2.660)	1498.000 (0.000)	1495.800 (1.887)	<b>0.008</b> ( <b>0.000</b> )	<b>0.012</b> ( <b>0.001</b> )	<b>0.998</b> ( <b>0.001</b> )	1.000 (0.000)	<b>0.999</b> ( <b>0.001</b> )
	1000	DIIR-QUIC	0.971 (0.001)	<b>112.508</b> ( <b>35.052</b> )	2998.000 (0.000)	2976.600 (5.219)	0.007 (0.000)	0.022 (0.001)	0.968 (0.001)	1.000 (0.000)	0.971 (0.001)
		$\ell_p$ COV	<b>1.000</b> ( <b>0.000</b> )	248.449 (24.727)	2998.000 (0.000)	2997.600 (0.800)	<b>0.004</b> ( <b>0.000</b> )	<b>0.007</b> ( <b>0.000</b> )	1.000 (0.000)	1.000 (0.000)	1.000 (0.000)
		$\ell_p$ COV (WS)	<b>1.000</b> ( <b>0.000</b> )	381.302 (53.667)	2998.000 (0.000)	2997.800 (0.600)	<b>0.004</b> ( <b>0.000</b> )	<b>0.007</b> ( <b>0.000</b> )	1.000 (0.000)	1.000 (0.000)	1.000 (0.000)
	1500	DIIR-QUIC	0.973 (0.001)	<b>397.025</b> ( <b>63.827</b> )	4498.000 (0.000)	4462.200 (9.775)	0.005 (0.000)	0.020 (0.000)	0.970 (0.001)	1.000 (0.000)	0.973 (0.001)
		$\ell_p$ COV	<b>1.000</b> ( <b>0.000</b> )	592.189 (119.818)	4498.000 (0.000)	4498.000 (0.000)	<b>0.003</b> ( <b>0.000</b> )	0.007 (0.000)	1.000 (0.000)	1.000 (0.000)	1.000 (0.000)
		$\ell_p$ COV (WS)	<b>1.000</b> ( <b>0.000</b> )	1113.915 (172.059)	4498.000 (0.000)	4498.000 (0.000)	<b>0.003</b> ( <b>0.000</b> )	<b>0.006</b> ( <b>0.000</b> )	1.000 (0.000)	1.000 (0.000)	1.000 (0.000)
clu	500	DIIR-QUIC	<b>0.333</b> ( <b>0.009</b> )	<b>49.281</b> ( <b>22.177</b> )	11432.800 (158.622)	2974.200 (71.181)	<b>0.039</b> ( <b>0.000</b> )	<b>0.279</b> ( <b>0.004</b> )	<b>0.210</b> ( <b>0.007</b> )	0.998 (0.000)	<b>0.400</b> ( <b>0.009</b> )
		$\ell_p$ COV	0.146 (0.003)	108.808 (22.124)	11432.800 (158.622)	905.200 (21.264)	0.041 (0.000)	0.327 (0.005)	0.079 (0.002)	1.000 (0.000)	0.273 (0.004)
		$\ell_p$ COV (WS)	0.155 (0.005)	108.811 (22.124)	11432.800 (158.622)	970.000 (26.803)	0.041 (0.000)	0.324 (0.005)	0.084 (0.003)	1.000 (0.000)	0.282 (0.005)
	1000	DIIR-QUIC	<b>0.539</b> ( <b>0.006</b> )	<b>613.051</b> ( <b>274.081</b> )	23093.400 (193.779)	10933.800 (113.900)	<b>0.024</b> ( <b>0.000</b> )	<b>0.208</b> ( <b>0.003</b> )	<b>0.397</b> ( <b>0.006</b> )	0.998 (0.000)	<b>0.571</b> ( <b>0.005</b> )
		$\ell_p$ COV	0.287 (0.003)	672.529 (274.001)	23093.400 (193.779)	3915.000 (44.553)	0.027 (0.000)	0.275 (0.003)	0.168 (0.002)	1.000 (0.000)	0.403 (0.002)
		$\ell_p$ COV (WS)	0.324 (0.006)	672.537 (273.988)	23093.400 (193.779)	4513.600 (83.258)	0.026 (0.000)	0.262 (0.004)	0.194 (0.004)	1.000 (0.000)	0.434 (0.005)
	1500	DIIR-QUIC	<b>0.682</b> ( <b>0.003</b> )	<b>4185.551</b> ( <b>1598.464</b> )	34701.000 (233.142)	30211.200 (277.186)	<b>0.017</b> ( <b>0.000</b> )	<b>0.140</b> ( <b>0.001</b> )	<b>0.638</b> ( <b>0.004</b> )	0.996 (0.000)	0.679 (0.003)
		$\ell_p$ COV	0.559 (0.003)	4245.262 (1598.553)	34701.000 (233.142)	13793.800 (149.019)	0.019 (0.000)	0.182 (0.002)	0.391 (0.003)	1.000 (0.000)	0.617 (0.002)
		$\ell_p$ COV (WS)	0.643 (0.002)	4245.306 (1598.517)	34701.000 (233.142)	16760.200 (152.054)	<b>0.017</b> ( <b>0.000</b> )	0.153 (0.002)	0.477 (0.003)	1.000 (0.000)	<b>0.683</b> ( <b>0.002</b> )

Table 5: Comparison of DIIR-QUIC and  $\ell_p$ COV on synthetic covariance estimation. For each method, the table reports the mean ( $\pm$  one standard deviation) over 10 independent trials under two covariance structures—tridiagonal (Tri) and clustered (Clu)—with  $p = 0.7$  at dimensions  $n = 500, 1000$ , and 1500. Boldface entries indicate the better performance.

Data	$n$	Algorithm	$F_1$	Time (s)	$\text{nnz}(\Sigma^{-1})$	$\text{nnz}(X^\dagger)$	Loss Q	KL Loss	Sensitivity	Specificity	MCC
tri	500	DIIR-QUIC	0.874 (0.005)	<b>43.538</b> ( <b>27.415</b> )	1498.000 (0.000)	1925.000 (18.852)	0.012 (0.000)	0.032 (0.001)	0.998 (0.001)	0.998 (0.000)	0.880 (0.004)
		$\ell_p$ COV	<b>0.944</b> ( <b>0.003</b> )	58.648 (1.675)	1498.000 (0.000)	1674.400 (10.575)	<b>0.010</b> (0.000)	<b>0.020</b> (0.001)	<b>1.000</b> (0.000)	<b>0.999</b> (0.000)	<b>0.946</b> ( <b>0.003</b> )
		$\ell_p$ COV (WS)	0.943 (0.003)	66.769 (3.191)	1498.000 (0.000)	1677.600 (11.586)	<b>0.010</b> (0.000)	<b>0.020</b> (0.001)	<b>1.000</b> (0.000)	<b>0.999</b> (0.000)	0.945 (0.003)
	1000	DIIR-QUIC	0.995 (0.001)	<b>82.146</b> ( <b>32.319</b> )	2998.000 (0.000)	3018.800 (5.075)	0.007 (0.000)	0.025 (0.000)	0.998 (0.000)	<b>1.000</b> (0.000)	0.995 (0.001)
		$\ell_p$ COV	<b>1.000</b> ( <b>0.000</b> )	252.932 (27.301)	2998.000 (0.000)	3000.800 (1.833)	<b>0.005</b> (0.000)	<b>0.013</b> (0.000)	<b>1.000</b> (0.000)	<b>1.000</b> (0.000)	<b>1.000</b> (0.000)
		$\ell_p$ COV (WS)	0.999 (0.000)	303.844 (50.265)	2998.000 (0.000)	3001.200 (2.400)	<b>0.005</b> (0.000)	<b>0.013</b> (0.000)	<b>1.000</b> (0.000)	<b>1.000</b> (0.000)	0.999 (0.000)
	1500	DIIR-QUIC	0.998 (0.000)	<b>128.492</b> ( <b>28.726</b> )	4498.000 (0.000)	4502.600 (1.562)	0.005 (0.000)	0.024 (0.000)	0.999 (0.000)	<b>1.000</b> (0.000)	0.998 (0.000)
		$\ell_p$ COV	<b>1.000</b> ( <b>0.000</b> )	588.834 (93.311)	4498.000 (0.000)	4498.000 (0.000)	<b>0.004</b> (0.000)	<b>0.012</b> (0.000)	<b>1.000</b> (0.000)	<b>1.000</b> (0.000)	<b>1.000</b> (0.000)
		$\ell_p$ COV (WS)	<b>1.000</b> ( <b>0.000</b> )	856.249 (128.172)	4498.000 (0.000)	4498.000 (0.000)	<b>0.004</b> (0.000)	<b>0.012</b> (0.000)	<b>1.000</b> (0.000)	<b>1.000</b> (0.000)	<b>1.000</b> (0.000)
clu	500	DIIR-QUIC	0.428 (0.005)	<b>75.779</b> ( <b>26.243</b> )	11522.200 (111.905)	9139.600 (128.473)	<b>0.036</b> (0.000)	<b>0.247</b> (0.001)	<b>0.384</b> (0.006)	0.980 (0.000)	0.407 (0.005)
		$\ell_p$ COV	0.430 (0.007)	135.312 (26.093)	11522.200 (111.905)	7171.600 (106.714)	0.039 (0.000)	0.251 (0.001)	0.349 (0.007)	<b>0.987</b> (0.000)	0.422 (0.006)
		$\ell_p$ COV (WS)	<b>0.440</b> ( <b>0.006</b> )	135.311 (26.093)	11522.200 (111.905)	7362.600 (105.864)	0.039 (0.000)	0.248 (0.001)	0.361 (0.007)	<b>0.987</b> (0.000)	<b>0.431</b> ( <b>0.006</b> )
	1000	DIIR-QUIC	0.503 (0.005)	<b>1222.015</b> ( <b>210.478</b> )	23026.600 (212.091)	36873.600 (438.393)	<b>0.022</b> (0.000)	0.173 (0.002)	<b>0.654</b> (0.007)	0.978 (0.000)	0.503 (0.005)
		$\ell_p$ COV	0.558 (0.005)	1281.443 (210.630)	23026.600 (212.091)	28733.200 (342.888)	0.023 (0.000)	0.169 (0.002)	0.627 (0.007)	<b>0.985</b> (0.000)	0.550 (0.005)
		$\ell_p$ COV (WS)	<b>0.569</b> ( <b>0.005</b> )	1281.438 (210.626)	23026.600 (212.091)	29595.600 (353.259)	0.023 (0.000)	<b>0.166</b> (0.002)	0.650 (0.006)	<b>0.985</b> (0.000)	<b>0.562</b> ( <b>0.005</b> )
	1500	DIIR-QUIC	0.611 (0.006)	<b>5759.561</b> ( <b>3675.491</b> )	34713.600 (256.051)	42335.800 (313.913)	<b>0.015</b> (0.000)	0.145 (0.002)	<b>0.678</b> (0.008)	0.992 (0.000)	0.607 (0.007)
		$\ell_p$ COV	0.663 (0.006)	5818.995 (3675.481)	34713.600 (256.051)	32665.200 (258.192)	0.016 (0.000)	0.140 (0.002)	0.644 (0.007)	<b>0.995</b> (0.000)	0.659 (0.007)
		$\ell_p$ COV (WS)	<b>0.678</b> ( <b>0.007</b> )	5818.996 (3675.463)	34713.600 (256.051)	33525.200 (260.962)	<b>0.015</b> (0.000)	<b>0.136</b> (0.002)	0.667 (0.008)	<b>0.995</b> (0.000)	<b>0.674</b> ( <b>0.007</b> )

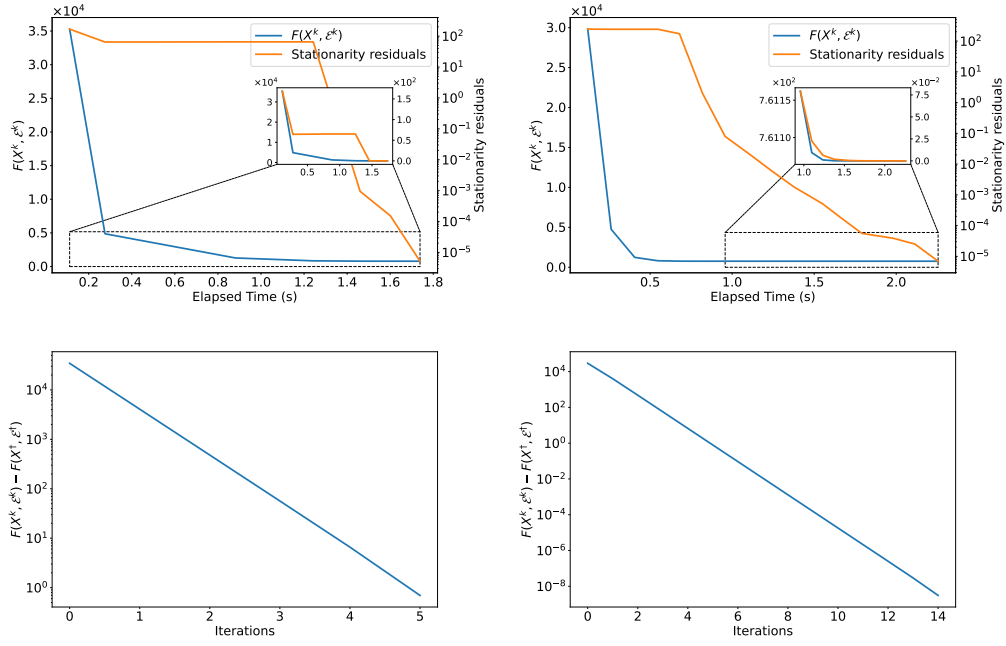


Figure 10: Convergence behavior of DIIR-QUIC on a tridiagonal synthetic precision matrix ( $n = 500$ ) with SCAD and MCP penalties. **Top row:** Evaluation of the penalized objective value and stationarity residuals versus elapsed time. **Bottom row:** Q-linear convergence of the objective error in log scale.

- [5] C.-J. Hsieh, M. A. Sustik, I. S. Dhillon, P. Ravikumar, *et al.*, “QUIC: quadratic approximation for sparse inverse covariance estimation,” *Journal of Machine Learning Research*, vol. 15, no. 1, pp. 2911–2947, 2014.
- [6] D. N. Phan, H. A. Le Thi, and T. P. Dinh, “Sparse covariance matrix estimation by DCA-based algorithms,” *Neural Computation*, vol. 29, no. 11, pp. 3040–3077, 2017.
- [7] T. Nakagaki, M. Fukuda, S. Kim, and M. Yamashita, “A dual spectral projected gradient method for log-determinant semidefinite problems,” *Computational Optimization and Applications*, vol. 76, pp. 33–68, 2020.
- [8] C. Wang, “On how to solve large-scale log-determinant optimization problems,” *Computational Optimization and Applications*, vol. 64, pp. 489–511, 2016.
- [9] D. Bertsimas, A. King, and R. Mazumder, “Best subset selection via a modern optimization lens,” *The Annals of Statistics*, vol. 44, no. 2, pp. 813–852, 2016.
- [10] J. Fan, H. Liu, Q. Sun, and T. Zhang, “I-lamm for sparse learning: Simultaneous control of algorithmic complexity and statistical error,” *Annals of Statistics*, vol. 46, no. 2, p. 814, 2018.
- [11] G. Marjanovic and V. Solo, “On  $\ell_q$  optimization and sparse inverse covariance selection,” *IEEE Transactions on Signal Processing*, vol. 62, no. 7, pp. 1644–1654, 2014.
- [12] Q. Wei and Z. Zhao, “Large covariance matrix estimation with oracle statistical rate via majorization-minimization,” *IEEE Transactions on Signal Processing*, vol. 71, pp. 3328–3342, 2023.
- [13] L. E. Frank and J. H. Friedman, “A statistical view of some chemometrics regression tools,” *Technometrics*, vol. 35, no. 2, pp. 109–135, 1993.
- [14] J. Fan and R. Li, “Variable selection via nonconcave penalized likelihood and its oracle properties,” *Journal of the American Statistical Association*, vol. 96, no. 456, pp. 1348–1360, 2001.
- [15] C.-H. Zhang, “Nearly unbiased variable selection under minimax concave penalty,” *The Annals of Statistics*, vol. 38, no. 2, pp. 894–942, 2010.
- [16] Z. Lu, “Iterative reweighted minimization methods for  $\ell_p$  regularized unconstrained nonlinear programming,” *Mathematical Programming*, vol. 147, no. 1, pp. 277–307, 2014.



- [17] L. Bai, Y. Hu, H. Wang, and X. Yang, “Avoiding strict saddle points of nonconvex regularized problems,” *arXiv preprint arXiv:2401.09274*, 2024.
- [18] D. Davis and D. Drusvyatskiy, “Proximal methods avoid active strict saddles of weakly convex functions,” *Foundations of Computational Mathematics*, vol. 22, no. 2, pp. 561–606, 2022.
- [19] D. Davis, D. Drusvyatskiy, and L. Jiang, “Active manifolds, stratifications, and convergence to local minima in nonsmooth optimization,” *Foundations of Computational Mathematics*, pp. 1–83, 2025.
- [20] A. Flinth, F. de Gournay, and P. Weiss, “Grid is good. adaptive refinement algorithms for off-the-grid total variation minimization,” *Open Journal of Mathematical Optimization*, vol. 6, pp. 1–27, 2025.
- [21] R. T. Rockafellar and R. J.-B. Wets, *Variational Analysis*, vol. 317. Heidelberg, Germany: Springer Berlin, 2009.
- [22] A. S. Lewis, “Active sets, nonsmoothness, and sensitivity,” *SIAM Journal on Optimization*, vol. 13, no. 3, pp. 702–725, 2002.
- [23] W. L. Hare and A. S. Lewis, “Identifying active constraints via partial smoothness and prox-regularity,” *Journal of Convex Analysis*, vol. 11, no. 2, pp. 251–266, 2004.
- [24] G. Garrigos, *Descent dynamical systems and algorithms for tame optimization and multi-objective problems*. PhD thesis, Université de Montpellier; Universidad Tecnica Federico Santa Maria, 2015.
- [25] J. Bolte, S. Sabach, and M. Teboulle, “Proximal alternating linearized minimization for nonconvex and nonsmooth problems,” *Mathematical Programming*, vol. 146, no. 1, pp. 459–494, 2014.
- [26] H. Attouch, J. Bolte, P. Redont, and A. Soubeyran, “Proximal alternating minimization and projection methods for nonconvex problems: An approach based on the Kurdyka-Łojasiewicz inequality,” *Mathematics of Operations Research*, vol. 35, no. 2, pp. 438–457, 2010.
- [27] X. Chen, L. Niu, and Y. Yuan, “Optimality conditions and a smoothing trust region newton method for nonlipschitz optimization,” *SIAM Journal on Optimization*, vol. 23, no. 3, pp. 1528–1552, 2013.
- [28] T. Liu and A. Takeda, “An inexact successive quadratic approximation method for a class of difference-of-convex optimization problems,” *Computational Optimization and Applications*, vol. 82, no. 1, pp. 141–173, 2022.

- [29] O. Banerjee, L. El Ghaoui, and A. d’Aspremont, “Model selection through sparse maximum likelihood estimation for multivariate Gaussian or binary data,” *The Journal of Machine Learning Research*, vol. 9, pp. 485–516, 2008.
- [30] C.-P. Lee, “Accelerating inexact successive quadratic approximation for regularized optimization through manifold identification,” *Mathematical Programming*, vol. 201, no. 1, pp. 599–633, 2023.
- [31] H. H. Bauschke, M. N. Dao, and W. M. Moursi, “On Fejér monotone sequences and nonexpansive mappings,” *arXiv preprint arXiv:1507.05585*, 2015.
- [32] F. H. Clarke, *Optimization and nonsmooth analysis*. Philadelphia, PA: SIAM, 1990.
- [33] Y. Wu, S. Pan, and X. Yang, “A regularized newton method for  $\ell_q$ -norm composite optimization problems,” *SIAM Journal on Optimization*, vol. 33, no. 3, pp. 1676–1706, 2023.
- [34] S. Wu, “Stochastic optimization methods for structure learning in Gaussian graphical models and the general log-determinant optimization,” *Doctoral Dissertation*, 2020.
- [35] A. Eftekhari, L. Gaedke-Merzhäuser, D. Pasadakis, M. Bollhöfer, S. Scheidegger, and O. Schenk, “Algorithm 1042: Sparse precision matrix estimation with squic,” *ACM Transactions on Mathematical Software*, vol. 50, no. 2, pp. 1–18, 2024.
- [36] M. Bollhöfer, A. Eftekhari, S. Scheidegger, and O. Schenk, “Large-scale sparse inverse covariance matrix estimation,” *SIAM Journal on Scientific Computing*, vol. 41, no. 1, pp. A380–A401, 2019.
- [37] D. Pasadakis, M. Bollhöfer, and O. Schenk, “Sparse quadratic approximation for graph learning,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 45, no. 9, pp. 11256–11269, 2023.
- [38] L. Li and K.-C. Toh, “An inexact interior point method for  $\ell_1$ -regularized sparse covariance selection,” *Mathematical Programming Computation*, vol. 2, pp. 291–315, 2010.
- [39] M. Yuan and Y. Lin, “Model selection and estimation in the Gaussian graphical model,” *Biometrika*, vol. 94, no. 1, pp. 19–35, 2007.
- [40] P. Baldi, S. Brunak, Y. Chauvin, C. A. Andersen, and H. Nielsen, “Assessing the accuracy of prediction algorithms for classification: an overview,” *Bioinformatics*, vol. 16, no. 5, pp. 412–424, 2000.

- [41] H. Dalianis and H. Dalianis, “Evaluation metrics and evaluation,” *Clinical Text Mining: Secondary Use of Electronic Patient Records*, pp. 45–53, 2018.
- [42] M. Fazel, H. Hindi, and S. P. Boyd, “Log-det heuristic for matrix rank minimization with applications to hankel and euclidean distance matrices,” in *Proceedings of the 2003 IEEE American Control Conference, 2003*, vol. 3, pp. 2156–2162, 2003.
- [43] D. Geman and C. Yang, “Nonlinear image recovery with half-quadratic regularization,” *IEEE Transactions on Image Processing*, vol. 4, no. 7, pp. 932–946, 1995.
- [44] Y. Wang and L. Zhu, “Variable selection and parameter estimation with the atan regularization method,” *Journal of Probability and Statistics*, vol. 2016, no. 1, p. 6495417, 2016.
- [45] C. Gao, N. Wang, Q. Yu, and Z. Zhang, “A feasible nonconvex relaxation approach to feature selection,” in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 25, pp. 356–361, 2011.