
High-Dimensional Tensor Regression With Oracle Properties

Anonymous Author
Anonymous Institution

Abstract

Tensor regression is a powerful tool for modeling linear relationships between multi-dimensional variables in complex data analysis. In this paper, we study a high-dimensional tensor-response tensor regression model under low-dimensional structural assumptions. Specifically, we focus on five key common structural assumptions for the tensor coefficients: element-wise sparsity, fiber-wise sparsity, slice-wise sparsity, mode-wise low-rankness, and slice-wise low-rankness. To effectively handle these structures, we employ non-convex penalties tailored to each case, yielding refined statistical convergence rates. We derive general error bounds for the resulting estimators under mild assumptions, showing that they can achieve oracle rates. Furthermore, we propose a non-convex accelerated proximal gradient algorithm to compute the estimators. Extensive numerical experiments validate the theoretical findings and show good performance on both synthetic and real-world datasets, demonstrating the efficiency and robustness of the proposed methods.

1 INTRODUCTION

Tensors, or multi-dimension arrays, generalize matrices to higher dimensions and have become fundamental in various domains of data analysis (McCullagh, 2018). They have been proven useful in numerous fields, including biostatistics (Hore et al., 2016), chemometrics (Acar et al., 2011), image signal processing (Westin et al., 1999), etc. Tensor-based methods are particularly instrumental because they preserve the multi-linear structures inherent in the data.

While applying classical methods to tensors typically requires unfolding tensors into matrices, which distorts the structure and results in information loss. Among tensor-based models, tensor regression is beneficial to uncover linear relationships between sets of high-dimensional variables and has significant applications across various fields (Liu et al., 2022). For example, in image processing, tensor regression has been successfully applied to tasks such as denoising (Zhang et al., 2023), image inpainting (Fan et al., 2017), and medical image analysis (Hua Zhou and Zhu, 2013). In multitask learning, tensor regression leverages shared information across multiple learning tasks to develop more accurate models for each task, achieving superior performance compared to learning them individually (Yang and Hospedales, 2017).

Tensor regression models frequently encounter significant challenges due to the complexity and scale of the high-order data involved, leading to pronounced issues of ill-posedness (Auddy et al., 2024). This problem is further exacerbated in high-dimensional settings where the number of model parameters substantially exceeds the number of observations (Raskutti et al., 2019). To address this issue, it is critical to impose structural assumptions on tensor models. One effective approach is to apply regularization penalties to tensor regression models to capture underlying structural data. Common penalties include sparsity (Tibshirani, 1996) or low-rankness penalties (Recht et al., 2010; Candes and Recht, 2012). Sparsity in tensors can manifest at the element-wise (Zhang et al., 2019; Raskutti et al., 2019), fiber-wise (Raskutti et al., 2019) or slice-wise level (Raskutti et al., 2019), each with different analytical implications. Raskutti et al. (2019) investigated various sparsity structures, establishing both general risk bounds and specific upper bounds in scenarios using convex optimization techniques. Similarly, low-rankness can be enforced on different matrixed forms of tensors, depending on the specific application (Farias and Li, 2017; Raskutti et al., 2019; Luo and Zhang, 2022). In their work, Raskutti et al. (2019) explored penalties for mode-wise and slice-wise low-rankness, deriving comprehensive risk and bound analysis across different scenarios through convex op-

timization. Another approach to capturing structures is via tensor decomposition, such as CANDECOMP/PARAFAC (CP) decomposition (Kiers, 2000) and Tucker decomposition (Tucker, 1966), etc. However, methods based on tensor decomposition are generally non-unique and involve solving non-convex optimization problems. To further reduce the number of estimated parameters and improve model performance, some studies consider combining both constraints and penalties (Ahmed et al., 2020; Xu, 2020).

In this paper, we investigate a high-dimensional tensor-on-tensor regression model incorporating penalty regularizers. Existing works, such as (Zhang et al., 2019; Raskutti et al., 2019), have adapted a convex approach inspired by classical regularization methods—such as Lasso (Tibshirani, 1996) and nuclear norm-based approaches (Recht et al., 2010; Candes and Recht, 2012). Despite their popularity, these convex methods introduce non-negligible biases that adversely affect estimator accuracy. Departing from these conventional approaches, we focus on non-convex regularization techniques. Specifically, we employ the smoothly clipped absolute deviation (SCAD) (Fan and Li, 2001) and the minimax concave penalty (MCP) (Zhang, 2010), which demonstrate superior performance in inducing sparsity and low-rank structures without the drawbacks of convex methods. We propose five distinct non-convex regularizers, and our results show that these estimators not only achieve faster statistical convergence rates but also enjoy oracle properties under mild assumptions. To optimize these estimators, we develop a novel accelerated proximal gradient algorithm to address the challenges of non-convex optimization. This work establishes a unifying and general framework for tensor regression, bridging gaps between empirical results and theoretical understanding in current research. We substantiate our theoretical contributions through extensive numerical experiments on both synthetic data and real-world datasets.

Notations. Scalars and matrices are denoted by standard lowercase letters and boldface capital letters, respectively. For a matrix $\mathbf{X} \in \mathbb{R}^{d_1 \times d_2}$, its singular values are denoted as $\sigma_1(\mathbf{X}) \geq \dots \geq \sigma_{\min\{d_1, d_2\}}(\mathbf{X})$. The nuclear norm is defined as $\|\mathbf{X}\|_* = \sum_{i=1}^{\min\{d_1, d_2\}} \sigma_i$, and the spectral norm is defined as $\|\mathbf{X}\|_{\text{sp}} = \sigma_1(\mathbf{X})$. We denote the i -th column of \mathbf{X} by $\mathbf{X}_{:,i}$, and the j -th row of \mathbf{X} by $\mathbf{X}_{j,:}$. The Frobenius norm of \mathbf{X} is defined as $\|\mathbf{X}\|_F = \sqrt{\sum_{i=1}^{d_2} \|\mathbf{X}_{:,i}\|_2^2}$, where $\|\mathbf{X}_{:,i}\|_2$ denotes the Euclidean norm of vector $\mathbf{X}_{:,i}$. Define $\text{vec}(\mathbf{X})$ as the column-wise vectorization of \mathbf{X} .

Tensors are denoted by boldface calligraphy letters. The order of a tensor is determined by the number

of dimensions (or modes) it has; hence, scalars, vectors, and matrices can be seen as 0th-order, 1st-order, and 2nd-order tensors, respectively. For a N -th order tensor $\mathcal{X} \in \mathbb{R}^{d_1 \times \dots \times d_N}$, where the i -th dimension is with size d_i , $i \in \{1, \dots, N\}$, the entry of \mathcal{X} at position (i_1, \dots, i_N) is denoted by $\mathcal{X}_{i_1, \dots, i_N}$. By fixing all indices of tensor \mathcal{X} except for the k -th mode, we obtain mode- (k) fibers, which are vectors. Similarly, mode- (j, k) slices are derived by fixing all indices except for j -th and k -th modes, which are matrices. For any tensor \mathcal{X} , it can be unfolded into a matrix via mode- (k) unfolding, denoted as $\mathcal{X}_{(k)} \in \mathbb{R}^{d_k \times \prod_{i \neq k} d_i}$. This unfolding arranges the mode- (k) fibers of \mathcal{X} into the columns of $\mathcal{X}_{(k)}$. Specifically, the elements of $\mathcal{X}_{(k)}$ satisfy $[\mathcal{X}_{(k)}]_{j_k, l} = \mathcal{X}_{j_1, \dots, j_k, \dots, j_N}$, where the column index $l = 1 + \sum_{s=1, s \neq k}^N (j_s - 1) \prod_{m=1, m \neq k}^{s-1} d_m$. Similarly, the tensor \mathcal{X} can be unfolded into a 3rd-order tensor via mode- (j, k) unfolding, denoted as $\mathcal{X}_{(j,k)} \in \mathbb{R}^{d_j \times d_k \times \prod_{s \neq j, k} d_s}$. This unfolding arranges the mode- (j, k) slices of \mathcal{X} into the frontal slices of $\mathcal{X}_{(j,k)}$. Specifically, the element satisfy $[\mathcal{X}_{(j,k)}]_{i_j, i_k, l} = \mathcal{X}_{i_1, \dots, i_N}$, where the slice index $l = 1 + \sum_{s=1, s \neq j, k}^N (i_s - 1) \prod_{m=1, m \neq j, k}^{s-1} d_m$. Define $\text{vec}(\mathcal{X})$ as the vectorization of tensor \mathcal{X} , where $\text{vec}(\mathcal{X}) = \text{vec}(\mathcal{X}_{(1)})$. The Frobenius norm of \mathcal{X} is defined as $\|\mathcal{X}\|_F = \sqrt{\sum_{i_1=1}^{d_1} \dots \sum_{i_N=1}^{d_N} \mathcal{X}_{i_1, \dots, i_N}^2}$.

For notational simplicity, we adopt a slight abuse of notation by using calligraphic letters to represent subspaces and index sets. The support of a set \mathcal{S} is denoted by $|\mathcal{S}|$. For a function f , $f'(\cdot)$ denotes its derivative, $\nabla f(\cdot)$ represents its gradient, and $\nabla^2 f(\cdot)$ denotes its Hessian. For functionals $f(x)$ and $g(x)$, we denote $f(x) \gtrsim g(x)$ if $f(x) \geq cg(x)$, $f(x) \lesssim g(x)$ if $f(x) \leq Cg(x)$, and $f(x) \asymp g(x)$ if $cg(x) \leq f(x) \leq Cg(x)$ for some positive constants c and C .

2 PROBLEM FORMULATION

In this section, we present a unified framework for high-dimensional tensor regression with non-convex regularization.

2.1 Tensor Regression Model

We consider the following generic observation model with tensor coefficient $\mathcal{A}^* \in \mathbb{R}^{d_1 \times \dots \times d_M}$:

$$\mathcal{Y} = \langle \mathcal{A}^*, \mathcal{X} \rangle + \mathcal{E}, \quad (1)$$

where $\mathcal{X} \in \mathbb{R}^{d_1 \times \dots \times d_N}$ is the predictor variable, $\mathcal{Y} \in \mathbb{R}^{d_{N+1} \times \dots \times d_M}$ is the response variable with $M \geq N$, and $\mathcal{E} \in \mathbb{R}^{d_{N+1} \times \dots \times d_M}$ is the noise. $\langle \cdot, \cdot \rangle$ is defined as a generalized inner product between two tensors, where its (i_{N+1}, \dots, i_M) -

th element is defined as $[\langle \mathcal{A}^*, \mathcal{X} \rangle]_{j_{N+1}, \dots, j_M} = \sum_{j_1=1}^{d_1} \dots \sum_{j_N=1}^{d_N} \mathcal{A}_{j_1, \dots, j_N}^* \mathcal{X}_{j_1, \dots, j_N, j_{N+1}, \dots, j_M}$. Specifically, when $M = N$, the output is a scalar, in which case (1) becomes a scalar-on-tensor regression model (Hua Zhou and Zhu, 2013; Gui et al., 2016); when $M > N$, the output is an $(M - N)$ -th order tensor, in which case (1) is the tensor-on-tensor regression model (Lock, 2018; Raskutti et al., 2019; Huihui Miao and Shi, 2022).

2.2 Proposed Estimators

Given a collection of n samples $\left\{ \left(\mathbf{y}^{(i)}, \mathcal{X}^{(i)} \right) \right\}_{i=1}^n$, which is assumed to be generated from the observation model (1), our goal is to estimate the unknown coefficient tensor \mathcal{A}^* by solving the following regularized optimization problem:

$$\min_{\mathcal{A} \in \mathbb{R}^{d_1 \times \dots \times d_N}} L(\mathcal{A}) + R_\lambda(\mathcal{A}), \quad (2)$$

where $L(\mathcal{A})$ is an empirical loss function defined as $L(\mathcal{A}) = \frac{1}{2n} \sum_{i=1}^n \left\| \mathbf{y}^{(i)} - \langle \mathcal{A}, \mathcal{X}^{(i)} \rangle \right\|_F^2$ and $R_\lambda(\mathcal{A})$ is a regularizer with parameter $\lambda \geq 0$. In this paper, we consider a family of decomposable non-convex penalty functions (Fan and Li, 2001; Zhang, 2010; Negahban et al., 2012), which enforces structural constraints on the regularized M-estimator $\hat{\mathcal{A}}$. We first introduce a class of typical univariate non-convex functions $p_\lambda(\cdot)$ and subsequently introduce how the regularizers $R_\lambda(\mathcal{A})$ are defined based on $p_\lambda(\cdot)$. A substantial body of literature has proposed various non-convex functions $p_\lambda(\cdot)$, with representative examples including SCAD (Fan and Li, 2001) and MCP (Zhang, 2010). These non-convex function $p_\lambda(x)$ can be decomposed as $p_\lambda(x) = \lambda|x| + q_\lambda(x)$, where $|x|$ represents the ℓ_1 penalty and $q_\lambda(x)$ is a concave function. Based on $p_\lambda(\cdot)$, we consider the following regularizers $R_\lambda(\mathcal{A})$ to enforce various structural constraints on \mathcal{A} .

2.2.1 Sparsity Penalties

A straightforward method to induce sparsity within a tensor \mathcal{A} is to encourage sparsity at element-wise level (Zhang et al., 2019; Raskutti et al., 2019), similar to the Lasso. This can be achieved by applying a non-convex penalty to each entry of the tensor. Specifically, the *element-wise sparsity penalty* is defined as:

$$R_\lambda(\mathcal{A}) = \sum_{i_1=1}^{d_1} \dots \sum_{i_M=1}^{d_M} p_\lambda(\mathcal{A}_{i_1, \dots, i_M}).$$

In many applications, sparsity manifests in structured patterns, such as entire fibers being zero (Li et al., 2015; Raskutti et al., 2019), analogous to the group

Lasso (Yuan and Lin, 2006). To capture such structured sparsity, we introduce the fiber-wise sparsity regularizer. Consider unfolding the tensor \mathcal{A} along mode- (k) , the *fiber-wise sparsity penalty* is defined as:

$$R_\lambda(\mathcal{A}) = \sum_{l=1}^{\prod_{j \neq k} d_j} p_\lambda \left(\left\| [\mathcal{A}_{(k)}]_{\cdot, l} \right\|_2 \right).$$

Beyond fiber-wise sparsity, a tensor may exhibit sparsity at the slice level, where the entire slice across two modes is zero (Raskutti et al., 2019). To promote such slice-wise sparsity, we introduce the *slice-wise sparsity penalty*:

$$R_\lambda(\mathcal{A}) = \sum_{l=1}^{\prod_{s \neq j, k} d_s} p_\lambda \left(\left\| [\mathcal{A}_{(j, k)}]_{\cdot, l} \right\|_F \right).$$

2.2.2 Low-rankness Penalties

In addition to promoting sparsity, encouraging low-rankness in tensors can be highly beneficial in various applications. Given the multiple notions of “tensor rank” (Kolda and Bader, 2009), we focus on two specific low-rank penalties, which we detail below.

One common approach to defining tensor rank is through the Tucker decomposition (Tucker, 1966). To encourage low tensor rank, we introduce the mode-wise low-rank regularizer, defined as the sum of the nuclear norms of its unfolded matrices along mode- (k) . Building upon this, the *mode-wise low-rank penalty* is defined as:

$$R_\lambda(\mathcal{A}) = \sum_{i=1}^{\min\{I_k, \prod_{j \neq k} d_j\}} p_\lambda(\sigma_i(\mathcal{A}_{(k)})),$$

where $\sigma_i(\mathcal{A}_{(k)})$ denotes the i -th singular value of the mode- (k) unfolding $\mathcal{A}_{(k)}$. This regularizer, also known as the tensor nuclear norm regularizer, encourages low rankness across all modes of the tensor by penalizing the singular values of each unfolding (Raskutti et al., 2019).

Another intuitive approach is to encourage slice-wise low-rankness (Lock, 2018; Raskutti et al., 2019). In many applications, low-rank structures are present within specific slices of a tensor. To capture and exploit this property, we introduce the *slice-wise low-rank regularizer* defined as:

$$R_\lambda(\mathcal{A}) = \sum_{l=1}^{\prod_{m \neq j, k} d_m} \sum_{s=1}^{s^{\text{all}}} p_\lambda \left(\sigma_s \left([\mathcal{A}_{(j, k)}]_{\cdot, l} \right) \right),$$

where $s^{\text{all}} = \min \left\{ d_j d_k, \prod_{l \neq j, k} d_l \right\}$.

3 MAIN THEORY

In this section, we present the primary theoretical results for the general estimator in (2). We first introduce some necessary assumptions regarding the empirical loss function and the non-convex penalty.

3.1 Technical Assumptions

Assumption 1 (Restricted strong convexity (RSC)). *For any tensor \mathcal{B} , there exists a constant $\rho^- > 0$ such that*

$$L(\mathcal{A} + \mathcal{B}) - L(\mathcal{A}) \geq \langle \nabla L(\mathcal{A}), \mathcal{B} \rangle + \frac{\rho^-}{2} \|\mathcal{B}\|_{\text{F}}^2.$$

Assumption 2 (Restricted strong smoothness (RSS)). *For any tensor \mathcal{B} , there exists a constant $\rho^+ > 0$ satisfying $\rho^+ > \rho^- > 0$ such that*

$$L(\mathcal{A} + \mathcal{B}) - L(\mathcal{A}) \leq \langle \nabla L(\mathcal{A}), \mathcal{B} \rangle + \frac{\rho^+}{2} \|\mathcal{B}\|_{\text{F}}^2.$$

Assumptions 1 and 2, which characterize the curvature behavior of the empirical loss function $L(\cdot)$, are analogous to the classical RSC and RSS conditions commonly used in the literature on simple linear regression problems (Gui et al., 2016; Elenberg et al., 2018). Following the methodology of (Candes and Tao, 2007), it can be proven that the function $L(\cdot)$ satisfies both RSC and RSS conditions with high probability.

Recall that the non-convex penalty function $R_\lambda(\cdot)$ is based on $p_\lambda(\cdot)$. We impose several conditions on the non-convex penalty function $R_\lambda(\cdot)$ in terms of functions $p_\lambda(\cdot)$ and $q_\lambda(\cdot)$.

Assumption 3. *The functions $p_\lambda(t)$ and $q_\lambda(t)$ satisfy the following conditions:*

- There exists a constant $\nu > 0$ such that the penalty function satisfies $p'_\lambda(t) = 0$ for all $t \geq \nu$;
- $q_\lambda(t)$ is symmetric, i.e., $q_\lambda(-t) = q_\lambda(t)$ for all t ;
- $q'_\lambda(t)$ is monotone and Lipschitz continuous, i.e., for $t_2 \geq t_1$, there exist two non-negative constant ζ^+, ζ^- such that $-\zeta^- \leq \frac{q'_\lambda(t_2) - q'_\lambda(t_1)}{(t_2 - t_1)} \leq -\zeta^+ \leq 0$;
- Both $q_\lambda(t)$ and $q'_\lambda(t)$ pass through the origin, i.e. $q_\lambda(0) = q'_\lambda(0) = 0$;
- There exists a positive constant λ such that $|q'_\lambda(t)| \leq \lambda$ for all t .

The third condition pertains to a curvature property that governs the degree of concavity level of $q'_\lambda(\cdot)$ and, consequently, the non-convexity level of $p_\lambda(\cdot)$. Such a condition is commonly employed in the analysis of non-convex regression problems (Wang et al., 2014; Gui et al., 2016; Fan et al., 2018).

Assumption 4. *Denote $\tilde{\mathcal{X}}$ as $\tilde{\mathcal{X}} = \left(\text{vec}(\mathcal{X}^{(1)})^\top, \dots, \text{vec}(\mathcal{X}^{(n)})^\top \right)$ the concatenation of vectorized covariates from n samples. Assume that $\tilde{\mathcal{X}}$ follows a multivariate normal distribution, i.e., $\tilde{\mathcal{X}} \sim \mathcal{N}(0, \Sigma)$, with covariance matrix $\Sigma = \text{cov}(\tilde{\mathcal{X}}, \tilde{\mathcal{X}})$. We assume that there exists constants $\kappa \geq 1$ such that the eigenvalues of covariance matrix Σ satisfy:*

$$\frac{1}{\kappa} \leq \sigma_{\min}(\Sigma) \leq \sigma_{\max}(\Sigma) \leq \kappa.$$

Assumption 4 ensures the positiveness and well-conditioning of the covariance matrix Σ , which is essential for the non-degeneracy of the parameter space. Such conditions can be verified in various statistical examples (Raskutti et al., 2019). Notably, if the covariate $\{\mathcal{X}^{(i)}\}_{i=1}^n$ are independent and identically distributed (i.i.d.), the covariance matrix Σ exhibits a block-diagonal structure. In this case, Assumption 4 simplifies to similar conditions on the covariance matrix of each sample.

3.2 Statistical Error Analysis

We now establish statistical error bounds for general decomposable regularizer $R_\lambda(\cdot)$.

Theorem 5. *Let the noise tensors $\mathcal{E}^{(i)}$ be i.i.d. with each entry following a Gaussian distribution $\mathcal{N}(0, \eta^2)$. Suppose Assumptions 1~4 hold. If $\|\mathcal{A}\|_{\text{F}} \leq R_\lambda(\mathcal{A})$, and $\lambda \geq 8\eta\sqrt{\frac{\kappa}{n}}w(\Omega)$, where $\Omega = \{\mathcal{A} \mid R_\lambda(\mathcal{A}) \leq 1\}$ represents the unit norm ball and $w(\mathcal{S}) = \mathbb{E}(\sup_{\mathcal{B} \in \mathcal{S}} \langle \mathcal{B}, \mathcal{T} \rangle)$ is the Gaussian width of the set \mathcal{S} , then there exists a constant $\gamma > 0$ such that with probability at least $1 - \exp(-\gamma w^2(\Omega))$, the estimate $\hat{\mathcal{A}}$ satisfies:*

$$\|\hat{\mathcal{A}} - \mathcal{A}^*\|_{\text{F}} \lesssim \lambda \kappa \Psi(\mathcal{A}),$$

where $\Psi(\mathcal{A}) = \sup_{\mathcal{U} \in \mathcal{A} \setminus \{0\}} \frac{R_\lambda(\mathcal{U})}{\|\mathcal{U}\|_{\text{F}}}$ is the subspace compatibility constant for subspace $\mathcal{A} \in \mathbb{R}^{I_1 \times I_2 \times \dots \times I_N}$.

Theorem 5 provides an upper bound on the estimation error between the optimal solution $\hat{\mathcal{A}}$ and the true value \mathcal{A}^* . This bound is driven by two terms: the subspace compatibility constant $\Psi(\mathcal{A})$, which captures the intrinsic complexity of the subspace \mathcal{A} , and the Gaussian width $w(\Omega)$, while quantifies the size of the unit norm ball Ω under the regularizer $R(\cdot)$. Intuitively, $\Psi(\mathcal{A})$ controls the sensitivity of the regularization to the subspace, while $w(\Omega)$ measures the richness of the subspace in terms of its capacity to fit the data.

Next, we present explicit statistical rates of convergence for different regularizers, providing deeper insights into their effectiveness under differing conditions.

Before conducting a detailed analysis of the convergence rates associated with element-wise sparsity regularizer, we first introduce a valuable theoretical benchmark—the oracle rate. The oracle rate refers to the statistical convergence rate of the oracle estimator, which serves as an idealized benchmark by assuming that the true parameter support set \mathcal{S}^* is known a priori. This assumption allows the oracle estimator to achieve the best theoretical performance. Specifically, the element-wise sparse oracle estimator is defined as $\hat{\mathcal{A}}^O = \arg \min_{\mathcal{A}: \mathcal{A}_{\overline{\mathcal{S}_1}} = 0} L(\mathcal{A})$, where $\overline{\mathcal{S}_1}$ denotes the complement of the support set \mathcal{S}_1 , i.e., $\mathcal{S}_1 = \{(i_1, \dots, i_M) \mid \mathcal{A}_{i_1, \dots, i_M}^* \neq 0\}$. Now, we have the following result.

Corollary 1 (Element-wise sparsity). *Suppose that Assumptions 1~4 hold. There exists a constant $0 < \gamma_1 < \infty$ such that the Gaussian width satisfies $w(\Omega) \leq \gamma_1$. If*

$$\rho^- > \zeta^-, \quad \lambda \asymp \sqrt{\frac{\log(d_1 d_2 \cdots d_N)}{n}},$$

and $\mathcal{A}_{i_1, \dots, i_M}^*$ satisfies the condition

$$\min_{(i_1, \dots, i_M) \in \mathcal{S}_1} |\mathcal{A}_{i_1, \dots, i_M}^*| \geq \nu,$$

the estimator $\hat{\mathcal{A}}$ satisfies

$$\|\hat{\mathcal{A}} - \mathcal{A}^*\|_{\text{F}} \lesssim \sqrt{\frac{|\mathcal{S}_1|}{n}}.$$

A similar argument applies to fiber-wise sparsity. The oracle rate refers to the statistical convergence rate of the fiber-wise sparse oracle estimator, which is assumed to know the true support set \mathcal{S}_2 in advance, where $\mathcal{S}_2 = \{i \mid \|[\mathcal{A}_{(k)}]_{\cdot, i}\|_2 \neq 0\}$. Specifically, the fiber-wise sparse oracle estimator is defined as $\hat{\mathcal{A}}^O = \arg \min_{\mathcal{A}: \mathcal{A}_{\overline{\mathcal{S}_2}} = 0} L(\mathcal{A})$.

Corollary 2 (Fiber-wise sparsity). *Suppose that Assumptions 1~4 hold. There exists a constant $0 < \gamma_2 < \infty$ such that the Gaussian width satisfies $w(\Omega) \leq \gamma_2 \sqrt{d_k}$. If*

$$\rho^- > \zeta^-, \quad \lambda \asymp \sqrt{\frac{d_k}{n}},$$

and $[\mathcal{A}_{(k)}]_{\cdot, l}$ satisfies the condition

$$\min_{l \in \mathcal{S}_2} \left\| [\mathcal{A}_{(k)}]_{\cdot, l} \right\|_2 \geq \nu,$$

the estimator $\hat{\mathcal{A}}$ satisfies

$$\|\hat{\mathcal{A}} - \mathcal{A}^*\|_{\text{F}} \lesssim \sqrt{\frac{|\mathcal{S}_2| d_k}{n}}.$$

Similarly, for slice-wise sparsity, the oracle rate refers to the statistical convergence rate of the slice-wise sparse oracle estimator, which is assumed to know the true support set \mathcal{S}_3 in advance, where $\mathcal{S}_3 = \{l \mid \|[\mathcal{A}_{(j,k)}]_{\cdot, l}\|_{\text{F}} \neq 0\}$. Specifically, the slice-wise sparse oracle estimator is defined as $\hat{\mathcal{A}}^O = \arg \min_{\mathcal{A}: \mathcal{A}_{\overline{\mathcal{S}_3}} = 0} L(\mathcal{A})$.

Corollary 3 (Slice-wise sparsity). *Suppose that Assumptions 1~4 hold. There exists a constant $0 < \gamma_3 < \infty$ such that the Gaussian width satisfies $w(\Omega) \leq \gamma_3 \sqrt{d_j d_k}$. If*

$$\rho^- > \zeta^-, \quad \lambda \asymp \sqrt{\frac{d_j d_k}{n}},$$

and $[\mathcal{A}_{(j,k)}]_{\cdot, l}$ satisfies the condition

$$\min_{l \in \mathcal{S}_3} \left\| [\mathcal{A}_{(j,k)}]_{\cdot, l} \right\|_{\text{F}} \geq \nu,$$

the estimator $\hat{\mathcal{A}}$ satisfies

$$\|\hat{\mathcal{A}} - \mathcal{A}^*\|_{\text{F}} \lesssim \sqrt{\frac{|\mathcal{S}_3| d_j d_k}{n}}.$$

In the following, we first introduce the singular value decomposition (SVD). Consider the matrix $\mathbf{X} \in \mathbb{R}^{m \times n}$, the SVD is given by $\mathbf{X} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^T$, where $\mathbf{U} \in \mathbb{R}^{m \times r}$ contains the left singular vectors, $\mathbf{V} \in \mathbb{R}^{n \times r}$ contains the right singular vectors, and $\mathbf{\Sigma} = \text{Diag}(\sigma_1(\mathbf{X}), \dots, \sigma_r(\mathbf{X})) \in \mathbb{R}^{r \times r}$ is a diagonal matrix. We further define a subspace \mathcal{F} along with its orthogonal complement \mathcal{F}^\perp as follows:

$$\begin{aligned} \mathcal{F}(\mathbf{X}) &= \{\mathbf{W} \mid \text{row}(\mathbf{W}) \subseteq \mathbf{V}, \text{col}(\mathbf{W}) \subseteq \mathbf{U}\}, \\ \mathcal{F}^\perp(\mathbf{X}) &= \{\mathbf{W} \mid \text{row}(\mathbf{W}) \perp \mathbf{V}, \text{col}(\mathbf{W}) \perp \mathbf{U}\}. \end{aligned}$$

For brevity, we adopt shorthand notations \mathcal{F} and \mathcal{F}^\perp when the dependence on \mathbf{X} is clear from context. We denote by $\Pi_{\mathcal{F}}(\cdot)$ and $\Pi_{\mathcal{F}^\perp}(\cdot)$ the projection operators onto the subspace \mathcal{F} and \mathcal{F}^\perp , respectively. Additionally, we introduce the linear operator $\mathfrak{X}(\mathcal{A}) : \mathbb{R}^{d_1 \times \cdots \times d_N} \rightarrow \mathbb{R}^{n \times d_{N+1} \times \cdots \times d_M}$ defined as $\mathfrak{X}(\mathcal{A}) = \left(\langle \mathcal{A}, \mathbf{x}^{(1)} \rangle, \dots, \langle \mathcal{A}, \mathbf{x}^{(n)} \rangle \right)^\top$ along with its adjoint operator $\mathfrak{X}^*(\mathcal{E})$. With these fundamental definitions, we present conclusions regarding tensor regression with low-rankness regularizers.

The oracle rate refers to the statistical convergence rate of the mode-wise low-rank oracle estimator, which

is assumed to know the true rank subspace $\mathcal{S}_4 = \{i \mid \sigma_i(\mathcal{F}(\mathcal{A}_{(k)}^*)) \neq 0\}$ in advance. Specifically, the tensor nuclear norm oracle estimator is defined as $\hat{\mathcal{A}}^O = \arg \min_{\mathcal{A}: \mathcal{A}_{\mathcal{S}_4} = 0} L(\mathcal{A})$.

Corollary 4 (Mode-wise low-rankness). *Suppose that Assumptions 1~4 hold. There exists a constant $0 < \gamma_4 < \infty$ such that the Gaussian width satisfies $w(\Omega) \leq \frac{\gamma_4}{\rho^- \sqrt{n}} \tau_k$, where $\tau_k = \left\| \Pi_{\mathcal{F}}([\mathcal{X}^*(\mathcal{E})]_{(k)}) \right\|_{\text{sp}}$. If*

$$\rho^- > \zeta^-, \quad \lambda \gtrsim \frac{1}{n} \sqrt{|\mathcal{S}_4|} \left\| [\mathcal{X}^*(\mathcal{E})]_{(k)} \right\|_{\text{sp}},$$

and $\sigma_i(\mathcal{A}_{(k)}^*)$ satisfies the condition

$$\left| \sigma_i(\mathcal{A}_{(k)}^*) \right| \geq \nu + \frac{2\sqrt{|\mathcal{S}_4|}}{n\rho^-} \left\| [\mathcal{X}^*(\mathcal{E})]_{(k)} \right\|_{\text{sp}},$$

the estimator $\hat{\mathcal{A}}$ satisfies

$$\left\| \hat{\mathcal{A}} - \mathcal{A}^* \right\|_{\text{F}} \lesssim \frac{\tau_k \sqrt{|\mathcal{S}_4|}}{n}.$$

The oracle rate refers to the statistical convergence rate of the slice-wise low-rank oracle estimator, which is assumed to know the true rank subspace $\mathcal{S}_5 = \{s \mid \sigma_s(\mathcal{F}([\mathcal{A}_{(j,k)}^*]_{\cdot, \cdot, l})) \neq 0\}$ in advance. Specifically, the slice-wise low-rank oracle estimator is defined as $\hat{\mathcal{A}}^O = \arg \min_{\mathcal{A}: \mathcal{A}_{\mathcal{S}_5} = 0} L(\mathcal{A})$.

Corollary 5 (Slice-wise low-rankness). *Suppose that Assumptions 1~4 hold. There exists a constant $0 < \gamma_5 < \infty$ such that the Gaussian width satisfies $w(\Omega) \leq \frac{\gamma_5}{\rho^- \sqrt{n}} \tau_{(j,k)}$, where $\tau_{(j,k)} = \max_l \left\| \Pi_{\mathcal{F}}([\mathcal{X}^*(\mathcal{E})]_{(j,k)})_{\cdot, \cdot, l} \right\|_{\text{sp}}$. If*

$$\rho^- > \zeta^-, \quad \lambda \gtrsim \frac{1}{n} \sqrt{|\mathcal{S}_5|} \left\| [\mathcal{X}^*(\mathcal{E})]_{(j,k)} \right\|_{\cdot, \cdot, l} \Big|_{\text{sp}},$$

and $\sigma_s([\mathcal{A}_{(j,k)}^*]_{\cdot, \cdot, l})$ satisfies the condition

$$\left| \sigma_s([\mathcal{A}_{(j,k)}^*]_{\cdot, \cdot, l}) \right| \geq \nu + \frac{2\sqrt{|\mathcal{S}_5|}}{n\rho^-} \left\| [\mathcal{X}^*(\mathcal{E})]_{(j,k)} \right\|_{\cdot, \cdot, l} \Big|_{\text{sp}},$$

the estimator $\hat{\mathcal{A}}$ satisfies

$$\left\| \hat{\mathcal{A}} - \mathcal{A}^* \right\|_{\text{F}} \lesssim \frac{\tau_{(j,k)} \sqrt{|\mathcal{S}_5|}}{n}.$$

Corollaries 1, 2, 3, 4, and 5 are direct consequences of Theorem 5. Furthermore, it is easy to see that the proposed estimators achieve a faster statistical rate of convergence—matching the oracle rate—compared to existing estimators that employ an ℓ_1 penalty (Raskutti et al., 2019) when evaluated under the Frobenius norm.

Algorithm 1: Proximal Gradient Algorithm

Require: $\eta \in (0, \frac{1}{l_p}), \delta \in (0, \frac{1}{\eta} - l_p)$;

```

1  $\mathcal{A}_0 = \mathcal{A}^* = 0$ ;
2 for  $t = 1, \dots, T$  do
3    $\mathcal{Y}_t = \mathcal{A}_t + \frac{t-1}{t+2}(\mathcal{A}_t - \mathcal{A}_{t-1})$ ;
4    $\Delta_t = \max_{s=\max(1, t-q), \dots, t} H(\mathcal{A}_t)$ ;
5   if  $H(\mathcal{Y}_t) \leq \Delta_t$  then
6      $\mathcal{V}_t = \mathcal{Y}_t$ ;
7   else
8      $\mathcal{V}_t = \mathcal{A}_t$ ;
9   end
10   $\mathcal{Z}_t = \mathcal{V}_t - \eta \nabla L(\mathcal{V}_t)$ ;
11   $\mathcal{A}_t = \text{prox}_{\eta \mathcal{R}_\lambda}(\mathcal{Z}_t)$ 

```

end

Output: \mathcal{A}_{T+1}

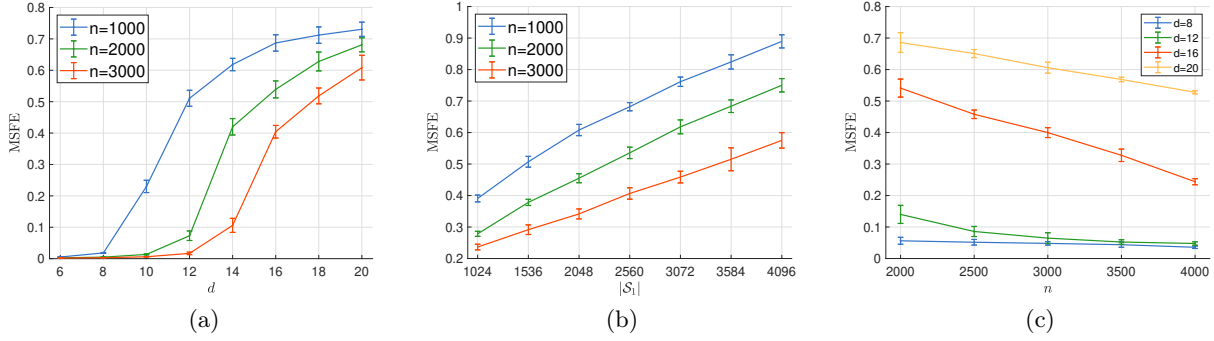
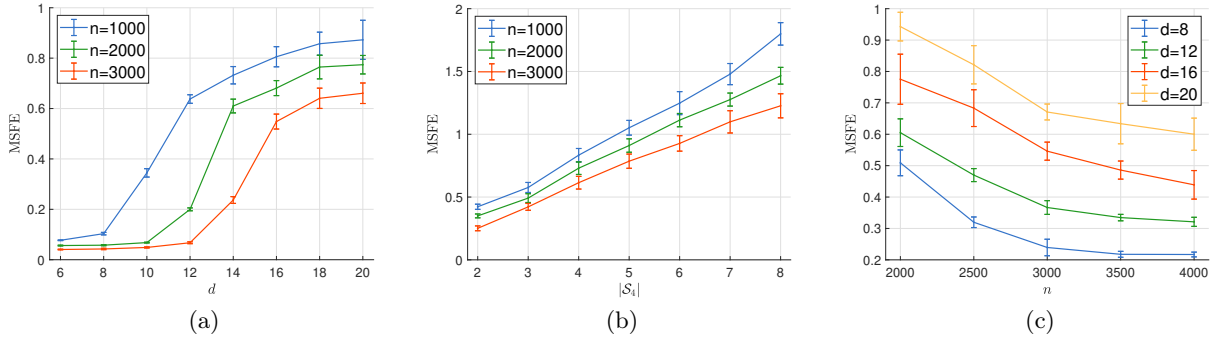
4 OPTIMIZATION ALGORITHM

In this section, we present a proximal gradient homotopy algorithm, adapted from (Yao et al., 2017), to solve the proposed estimators (2). The detailed steps are provided in Algorithm 1. We define $H(\mathcal{A}) = L(\mathcal{A}) + R_\lambda(\mathcal{A})$ for clarity and the proximal operator as $\text{prox}_{\eta R_\lambda}(v) = \arg \min_x \left\{ \frac{1}{2} \|x - v\|_2^2 + \eta R_\lambda(x) \right\}$. At each iteration t , the accelerated proximal gradient algorithm updates x_{t+1} through the following steps:

$$y_t = x_t + \theta_t(x_t - x_{t-1}),$$

$$x_{t+1} = \text{prox}_{\eta R_\lambda}(y_t - \eta \nabla L(y_t)),$$

where $\theta_t = \frac{t-1}{t+2}$. When $\theta_t = 0$, the accelerated proximal gradient reduces to the standard proximal gradient algorithm. However, when extending to non-convex problems, a significant challenge arises due to the extrapolation term \mathcal{Y}_t . The presence of \mathcal{Y}_t prevents guaranteeing a sufficient decrease of the objective function in each iteration. To address this issue, Yao et al. (2017) proposed to evaluate the objective value before the proximal step, instead of checking after it as was previously done, demonstrating that this modification guarantees a sufficient decrease under certain conditions. Moreover, to effectively implement the Algorithm 1, we need to choose an appropriate step size η . The parameter η is related to the Lipschitz constant l_p of $L(\cdot)$ (i.e., $\|\nabla L(\mathcal{X}) - \nabla L(\mathcal{Y})\|_{\text{F}} \leq l_p \|\mathcal{X} - \mathcal{Y}\|_{\text{F}}$). Choosing $\eta \leq \frac{1}{l_p}$ ensures that the gradient step does not overshoot, which is critical for the convergence of the algorithm. In line 4, q is an adjustable parameter. Setting $q = 0$ represents the simplest approach, while a larger q permits the objective function $H(\mathcal{Y}_t)$ to occasionally increase. This flexibility is inspired by the Barzilai-Borwein scheme for unconstrained smooth minimization Grippo and Sciandrone (2002). We adopt $q = 5$ in our experiments.


 Figure 1: Element-wise sparsity regularizer with the error bars of $\text{MSFE} \pm \text{standard deviation}$.

 Figure 2: Mode-wise lowrankness regularizer with the error bars of $\text{MSFE} \pm \text{standard deviation}$.

5 NUMERICAL EXPERIMENTS

In this section, we perform extensive numerical experiments to evaluate the performance of the proposed tensor regression models with various regularization techniques. For each regularizer, we take the SCAD penalty as the non-convex penalty. The tuning parameter λ and the hyperparameter associated with the SCAD penalty are selected via ten-fold cross-validation, aiming to minimize the estimation error quantified by the Mean Squared Frobenius norm Error (MSFE). Specifically, the MSFE is defined as $\frac{\|\mathcal{A}^* - \hat{\mathcal{A}}\|_F^2}{\prod_{i=1}^N d_i}$, where \mathcal{A}^* is the true parameter tensor and $\hat{\mathcal{A}}$ is the estimated one. All experiments are implemented in MATLAB, and the reported results are averaged on 100 Monte Carlo realizations to ensure statistical robustness.

5.1 Synthetic Data

In synthetic data experiments, the covariate tensors $\mathcal{X}^{(i)}$ are independently drawn from standard Gaussian ensembles. Additionally, the i.i.d. Gaussian noise terms are generated from the distribution $\mathcal{N}(0, \eta^2)$. Beyond MSFE, we also employ the Root Mean Square

Error (RMSE) between the predicted response and the true response. Specifically, the RMSE is calculated as

$$\text{RMSE} = \sqrt{\frac{1}{n} \sum_{i=1}^n \|\mathbf{y}^{(i)} - \langle \mathcal{A}, \mathbf{x}^{(i)} \rangle\|_F^2}.$$

In the subsequent experiments, we employ 3rd-order tensors $\mathcal{A} \in \mathbb{R}^{d \times d \times d}$ to evaluate the effectiveness of entry-wise sparsity, fiber-wise sparsity and mode-wise low-tensor-rankness regularizers. To better show the performance of slice-wise sparsity and low-rankness regularizers, we consider tensors $\mathcal{A} \in \mathbb{R}^{d \times d \times s}$, where the number of slices s is fixed.

Fig. 1 illustrates the performance of the estimation with the element-wise sparsity regularizer with the variation of dimensions d , element-wise sparsity s^* and the number of samples n , respectively. In Fig. 1a and Fig. 1b, three lines correspond to sample size $n = \{1000, 2000, 3000\}$. The proportion of non-zero elements $\frac{s^*}{d^3}$ is set to 0.5 in Fig. 1a and Fig. 1c. Fig. 1a demonstrates the estimation error decreases consistently with an increase in sample size n . Conversely, Fig. 1b shows that an increase in the cardinality of the support set leads to an increase in estimation error. The observed trends in Fig. 1 align with Corollary 1, thereby validating our theoretical conclusions.

Table 1: Comparisons between non-convex regularizers and convex regularizers.

Structures	Methods		Synthetic Data					Real-world Data	
			size	$ S $	η	MSFE \pm std	RMSE \pm std	MSFE \pm std	MPRE \pm std
Sparsity	Element-wise	Non-convex	$16 \times 16 \times 16$	2048	0.1	0.4042 ± 0.0201	0.0992 ± 0.0021	134.5864 ± 11.2950	8.8072 ± 0.0313
		Convex				0.6938 ± 0.0297	0.1004 ± 0.0023	144.7160 ± 14.9947	7.7498 ± 0.0457
	Fiber-wise	Non-convex	$16 \times 16 \times 16$	8	0.1	0.4406 ± 0.0157	0.0993 ± 0.0012	90.3068 ± 7.3006	4.7161 ± 0.0103
		Convex				0.7512 ± 0.0439	0.0995 ± 0.0019	102.7019 ± 9.2188	5.0330 ± 0.0118
	Slice-wise	Non-convex	$16 \times 16 \times 20$	8	0.1	0.5761 ± 0.0289	0.0997 ± 0.0027	43.8705 ± 3.0257	1.8909 ± 0.0043
		Convex				0.7201 ± 0.0314	0.1005 ± 0.0039	48.4585 ± 3.8834	1.9250 ± 0.0045
Low-rankness	Mode-wise	Non-convex	$16 \times 16 \times 16$	5	1	0.5482 ± 0.0395	0.1002 ± 0.0012	35.5536 ± 1.4889	1.0330 ± 0.0022
		Convex				1.7411 ± 0.0953	0.1096 ± 0.0020	41.2719 ± 3.5079	1.1027 ± 0.0024
	Slice-wise	Non-convex	$16 \times 16 \times 20$	5	1	0.9214 ± 0.0736	0.1004 ± 0.0010	8.9348 ± 0.7493	0.0436 ± 0.0002
		Convex				1.8261 ± 0.1066	0.1113 ± 0.0031	10.1655 ± 0.9050	0.6348 ± 0.0009

Additionally, we assess the performance of other sparsity regularizers, with detailed results presented in Appendix A.

Fig. 2 illustrates the results of the mode-wise low-rankness regularizer. In Fig. 2b, x -axis r represent the rank of the mode- (k) unfolded matrix. Fig. 2a and Fig. 2c sets the rank of the mode- (k) unfolded matrix to 5, while Fig. 2b and Fig. 2c fix the dimension of each mode to 16. Three distinct lines correspond to the estimation errors for sample size $n = \{1000, 2000, 3000\}$. These experimental results validate our theoretical findings as in Corollary 4. We also report results for the slice-wise low-rank regularizer in Appendix A.

In Table 1, we compare the performance of our proposed non-convex regularizers against traditional convex regularizers. For sparsity regularization, we set the $\eta = 0.1$, and for low-rankness regularizers, we set $\eta = 1$. We configure the tensor dimension such that tensors with slices-wise structures $\mathcal{A} \in \mathbb{R}^{d \times d \times s}$ and the others $\mathcal{A} \in \mathbb{R}^{d \times d \times d}$, where $d = 16$, $s = 20$. Depending on the tensor structure, the sparsity or the rank of the tensors varies accordingly. The results in Table 1 demonstrate that non-convex regularizers achieve lower MSFE for parameter estimation and lower RMSE for predictions compared to their convex counterparts. These empirical findings are in strong agreement with our theoretical analysis.

5.2 Real-world Datasets

To illustrate the effectiveness of the proposed methods, we conduct experiments using the ImageNet 2012 dataset (Russakovsky et al., 2015). In our regression framework, the tensor parameter is a 3rd-order color image, denoted as $\mathcal{A} \in \mathbb{R}^{64 \times 64 \times 3}$. We utilize $n = 4000$ samples in our experiments. In addition to MSFE, we evaluate performance using the Mean Prediction Relative Error (MPRE), defined as $\frac{1}{n} \sum_{i=1}^n \frac{\|\mathbf{y}^* - \hat{\mathbf{y}}^{(i)}\|_F}{\|\mathbf{y}^*\|_F}$, where $\mathbf{y}^* = \langle \mathcal{A}^*, \mathcal{X} \rangle$ represents the original tensor re-

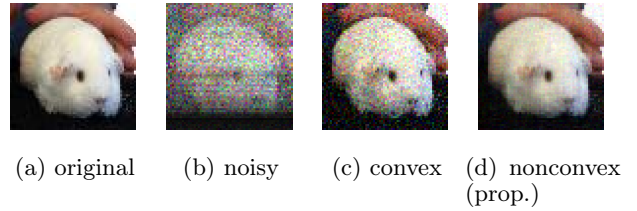


Figure 3: The original image, noisy image, and de-noised images using convex and non-convex methods.

sponse with the RGB image and $\hat{\mathbf{y}}$ is the estimated response. Fig. 3 illustrates the estimated image, and Table 1 presents the comparative results.

6 CONCLUSIONS AND FUTURE DIRECTIONS

In this paper, we propose a comprehensive framework for tensor regression estimation using non-convex penalties. Our findings demonstrate that estimators employing non-convex penalties exhibit faster convergence rates compared to those with convex penalties. Furthermore, we show that under several mild conditions, our proposed estimator possesses the oracle property. Extensive experimental results validate our theoretical claims, showcasing a close alignment between the theoretical predictions and the observed numerical performance of our estimators. Currently, we are limited to applying regularization penalties to tensor regression models. It would be desirable to derive some theoretical guarantees for alternative methods that capture structure in the tensor regression models, such as tensor decomposition; this is the aim of our future work. To conclude, our work effectively bridge the gap between practical applications and theoretical analysis of tensor-on-tensor regression with non-convex penalties. To the best of our knowledge, this is the first work to obtain the oracle statistical rate of convergence for the tensor regression problem.

References

- Evrin Acar, Daniel M Dunlavy, Tamara G Kolda, and Morten Mørup. Scalable tensor factorizations for incomplete data. *Chemometrics and Intelligent Laboratory Systems*, 106(1):41–56, 2011.
- Talal Ahmed, Haroon Raja, and Waheed U Bajwa. Tensor regression using low-rank and sparse tucker decompositions. *SIAM Journal on Mathematics of Data Science*, 2(4):944–966, 2020.
- Arnab Auddy, Dong Xia, and Ming Yuan. Tensor methods in high dimensional data analysis: Opportunities and challenges. *arXiv preprint arXiv:2405.18412*, 2024.
- Emmanuel Candes and Benjamin Recht. Exact matrix completion via convex optimization. *Communications of the ACM*, 55(6):111–119, 2012.
- Emmanuel Candes and Terence Tao. The Dantzig selector: Statistical estimation when p is much larger than n . *The Annals of Statistics*, 35(6):2313–2351, 2007.
- Ethan R Elenberg, Rajiv Khanna, Alexandros G Dimakis, and Sahand Negahban. Restricted strong convexity implies weak submodularity. *The Annals of Statistics*, 46(6B):3539–3568, 2018.
- Haiyan Fan, Yunjin Chen, Yulan Guo, Hongyan Zhang, and Gangyao Kuang. Hyperspectral image restoration using low-rank tensor recovery. *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, 10(10):4589–4604, 2017.
- Jianqing Fan and Runze Li. Variable selection via nonconcave penalized likelihood and its oracle properties. *Journal of the American Statistical Association*, 96(456):1348–1360, 2001.
- Jianqing Fan, Han Liu, Qiang Sun, and Tong Zhang. I-LAMM for sparse learning: Simultaneous control of algorithmic complexity and statistical error. *The Annals of Statistics*, 46(2):814, 2018.
- Vivek Farias and Andrew Li. Optimal recovery of tensor slices. In *Artificial Intelligence and Statistics*, pages 1394–1402. PMLR, 2017.
- Yehoram Gordon. On milman’s inequality and random subspaces which escape through a mesh in \mathbf{r}^n . In *Geometric Aspects of Functional Analysis: Israel Seminar (GAFA) 1986–87*, pages 84–106. Springer, 1988.
- Luigi Grippo and Marco Sciandrone. Nonmonotone globalization techniques for the Barzilai-Borwein gradient method. *Computational Optimization and Applications*, 23:143–169, 2002.
- Huan Gui, Jiawei Han, and Quanquan Gu. Towards faster rates and oracle property for low-rank matrix estimation. In *International Conference on Machine Learning*, pages 2300–2309. PMLR, 2016.
- Victoria Hore, Ana Vinuela, Alfonso Buil, Julian Knight, Mark I McCarthy, Kerrin Small, and Jonathan Marchini. Tensor decomposition for multiple-tissue gene expression experiments. *Nature genetics*, 48(9):1094–1100, 2016.
- Lexin Li Hua Zhou and Hongtu Zhu. Tensor regression with applications in neuroimaging data analysis. *Journal of the American Statistical Association*, 108(502):540–552, 2013.
- Bing Li Huihui Miao, Andi Wang and Jianjun Shi. Structural tensor-on-tensor regression with interaction effects and its application to a hot rolling process. *Journal of Quality Technology*, 54(5):547–560, 2022.
- Henk Kiers. Towards a standardized notation and terminology in multiway analysis. *Journal of Chemometrics*, 14:105–122, 05 2000.
- Tamara G Kolda and Brett W Bader. Tensor decompositions and applications. *SIAM Review*, 51(3):455–500, 2009.
- Yanming Li, Bin Nan, and Ji Zhu. Multivariate sparse group lasso for the multivariate multiple linear regression with an arbitrary group structure. *Biometrics*, 71(2):354–363, 2015.
- Yipeng Liu, Jiani Liu, Zhen Long, and Ce Zhu. *Tensor regression*. Springer, 2022.
- Eric F. Lock. Tensor-on-tensor regression. *Journal of Computational and Graphical Statistics*, 27(3):638–647, June 2018. ISSN 1537-2715.
- Yuetian Luo and Anru R Zhang. Tensor-on-tensor regression: Riemannian optimization, over-parameterization, statistical-computational gap, and their interplay. *arXiv preprint arXiv:2206.08756*, 2022.
- Pascal Massart. *Concentration inequalities and model selection: Ecole d’Eté de Probabilités de Saint-Flour XXXIII-2003*. Springer, 2007.
- Peter McCullagh. *Tensor methods in statistics: Monographs on statistics and applied probability*. Chapman and Hall/CRC, 2018.
- Sahand Negahban, Pradeep Ravikumar, Martin J. Wainwright, and Bin Yu. A Unified Framework for High-Dimensional Analysis of M -Estimators with Decomposable Regularizers. *Statistical Science*, 27(4):538 – 557, 2012.

Garvesh Raskutti, Ming Yuan, and Han Chen. Convex regularization for high-dimensional multiresponse tensor regression. *The Annals of Statistics*, 47(3):1554–1584, 2019.

Benjamin Recht, Maryam Fazel, and Pablo A Parrilo. Guaranteed minimum-rank solutions of linear matrix equations via nuclear norm minimization. *SIAM Review*, 52(3):471–501, 2010.

Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. Imagenet large scale visual recognition challenge. *International Journal of Computer Vision*, 115:211–252, 2015.

Robert Tibshirani. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 58(1):267–288, 1996.

Ledyard R Tucker. Some mathematical notes on three-mode factor analysis. *Psychometrika*, 31(3):279–311, 1966.

Zhaoran Wang, Han Liu, and Tong Zhang. Optimal computational and statistical rates of convergence for sparse nonconvex learning problems. *The Annals of Statistics*, 42(6):2164, 2014.

C F Westin, Stephan E Maier, Besam Khidhir, Peter Everett, Ferenc A Jolesz, and Ron Kikinis. Image processing for diffusion tensor magnetic resonance imaging. In *Medical Image Computing and Computer-Assisted Intervention–MICCAI’99: Second International Conference, Cambridge, UK, September 19–22, 1999. Proceedings 2*, pages 441–452. Springer, 1999.

Da Xu. Sparse symmetric tensor regression for functional connectivity analysis. *arXiv preprint arXiv:2010.14700*, 2020.

Yongxin Yang and Timothy Hospedales. Deep multi-task representation learning: A tensor factorisation approach. In *5th International Conference on Learning Representations*, 2017.

Quanming Yao, James T Kwok, Fei Gao, Wei Chen, and Tie-Yan Liu. Efficient inexact proximal gradient algorithm for nonconvex problems. In *Proceedings of the 26th International Joint Conference on Artificial Intelligence*, pages 3308–3314, 2017.

Ming Yuan and Yi Lin. Model selection and estimation in regression with grouped variables. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 68(1):49–67, 2006.

Cun-Hui Zhang. Nearly unbiased variable selection under minimax concave penalty. *The Annals of Statistics*, pages 894–942, 2010.

Hao Zhang, Ting-Zhu Huang, Xi-Le Zhao, Wei He, Jae Kyu Choi, and Yu-Bang Zheng. Hyperspectral image denoising: Reconciling sparse and low-tensor-ring-rank priors in the transformed domain. *IEEE Transactions on Geoscience and Remote Sensing*, 61:1–13, 2023.

Jiaqi Zhang, Yinghao Cai, Zhaoyang Wang, and Beilun Wang. Sparse and low-rank high-order tensor regression via parallel proximal method. *arXiv preprint arXiv:1911.12965*, 2019.

Checklist

The checklist follows the references. For each question, choose your answer from the three possible options: Yes, No, Not Applicable. You are encouraged to include a justification to your answer, either by referencing the appropriate section of your paper or providing a brief inline description (1-2 sentences). Please do not modify the questions. Note that the Checklist section does not count towards the page limit. Not including the checklist in the first submission won’t result in desk rejection, although in such case we will ask you to upload it during the author response period and include it in camera ready (if accepted).

In your paper, please delete this instructions block and only keep the Checklist section heading above along with the questions/answers below.

1. For all models and algorithms presented, check if you include:
 - (a) A clear description of the mathematical setting, assumptions, algorithm, and/or model. [Yes]
 - (b) An analysis of the properties and complexity (time, space, sample size) of any algorithm. [Not Applicable]
 - (c) (Optional) Anonymized source code, with specification of all dependencies, including external libraries. [Yes]
2. For any theoretical claim, check if you include:
 - (a) Statements of the full set of assumptions of all theoretical results. [Yes]
 - (b) Complete proofs of all theoretical results. [Yes]
 - (c) Clear explanations of any assumptions. [Yes]

3. For all figures and tables that present empirical results, check if you include:
 - (a) The code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL). [Yes]
 - (b) All the training details (e.g., data splits, hyperparameters, how they were chosen). [Yes]
 - (c) A clear definition of the specific measure or statistics and error bars (e.g., with respect to the random seed after running experiments multiple times). [Yes]
 - (d) A description of the computing infrastructure used. (e.g., type of GPUs, internal cluster, or cloud provider). [Yes]
4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets, check if you include:
 - (a) Citations of the creator If your work uses existing assets. [Yes]
 - (b) The license information of the assets, if applicable. [Yes]
 - (c) New assets either in the supplemental material or as a URL, if applicable. [Yes]
 - (d) Information about consent from data providers/curators. [Yes]
 - (e) Discussion of sensible content if applicable, e.g., personally identifiable information or offensive content. [Yes]
5. If you used crowdsourcing or conducted research with human subjects, check if you include:
 - (a) The full text of instructions given to participants and screenshots. [Not Applicable]
 - (b) Descriptions of potential participant risks, with links to Institutional Review Board (IRB) approvals if applicable. [Not Applicable]
 - (c) The estimated hourly wage paid to participants and the total amount spent on participant compensation. [Not Applicable]

High-Dimensional Tensor Regression With Oracle Properties: Supplementary Materials

Contents

A	ADDITIONAL NUMERICAL RESULTS	13
A.1	Synthetic Data	13
A.2	Real-world Data	15
B	Background	17
C	Proof of the Corollary 1	17
D	Proof of the Corollary 2	24
E	Proof of the Corollary 3	25
F	Proof of the Corollary 4	26
G	Proof of the Corollary 5	37
H	Proof of the Theorem 5	38

A ADDITIONAL NUMERICAL RESULTS

In this section, we provide additional results for the proposed regularizers besides the results in Section 5. In Section A.1, the performance of different regularizers based on synthetic data will be discussed, followed by the analysis of the experimental results. Moreover, the results of the regularizers for the real-world data are illustrated in Section A.2.

A.1 Synthetic Data

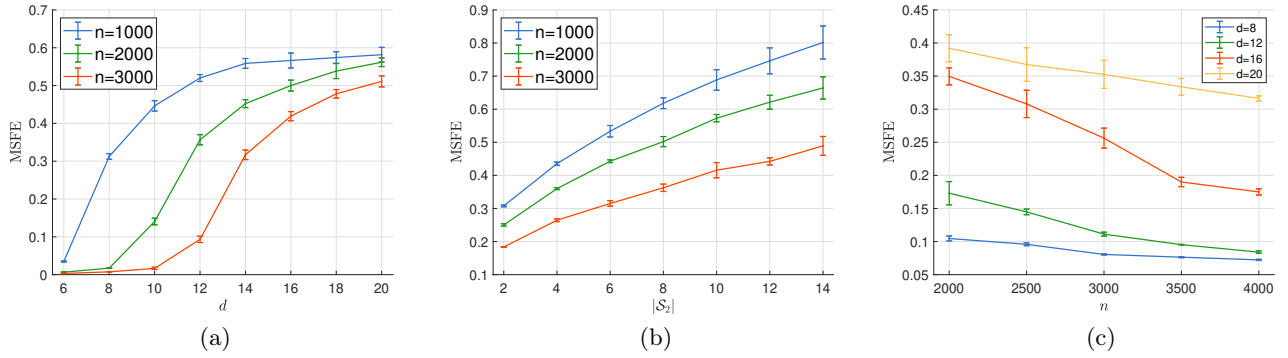


Figure 4: Fiber-wise sparsity regularizer with the error bars of $\text{MSFE} \pm \text{standard deviation}$.

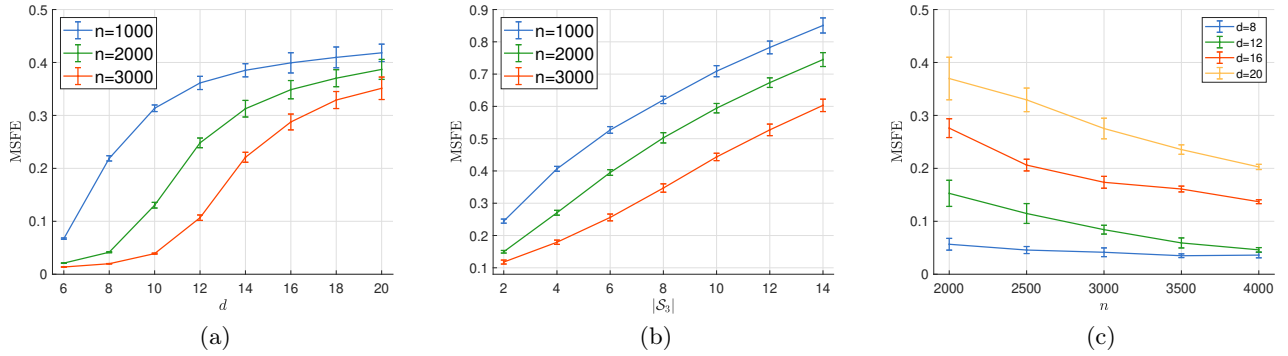
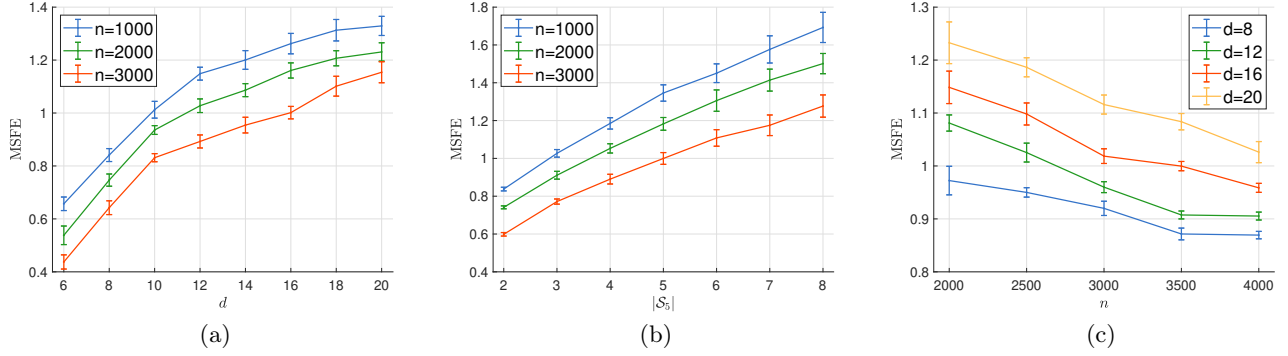


Figure 5: Slice-wise sparsity regularizer with the error bars of $\text{MSFE} \pm \text{standard deviation}$.

Figure 4 presents the results related to the fiber-wise sparsity regularizer. We display the results of the Mean Squared Frobenius norm Error (MSFE) versus the dimension d , the fiber-wise sparsity $|\mathcal{S}_2|$ and the sample size n , respectively. In these experiments, all tensors are 3rd-order tensors, with the dimension $d = 16$. In Figure 4a and Figure 4b, we fix the fiber-wise sparsity $|\mathcal{S}_2| = 4$. The three lines in different colors represent varying sample sizes $n = \{1000, 2000, 3000\}$. From Figure 4a, we observe that MSFE increases as the tensor dimension d increases. From Figure 4b, we find that increasing the sample size n decreases the MSFE. This demonstrates that larger sample sizes improve the accuracy of the tensor estimation, as expected. In Figure 4c, we see that increasing the fiber-wise-sparsity $|\mathcal{S}_2|$ leads to an increase in the estimate error. This suggests that tensors with more nonzero fibers are harder to estimate accurately. Furthermore, the standard deviation of the estimation error follows the same trend, increasing with fiber-wise sparsity.

Figure 5 presents the results of the slice-wise sparsity regularizer. In Figure 5a and Figure 5b, the number of slices is uniformly set to $s = 20$. And we set the slice-wise sparsity $|\mathcal{S}_3| = 4$ in Figure 5a and Figure 5c. We select three sample sizes while varying the dimension d or the number of non-zero slices $|\mathcal{S}_3|$. The results indicate that


 Figure 6: Slice-wise lowrankness regularizer with the error bars of $\text{MSFE} \pm \text{standard deviation}$.

the estimation error increases with increments in d or $|\mathcal{S}_3|$. The standard deviation of the MSFE also rises as the MSFE increases. Furthermore, as observed in Figure 5c, increasing the sample size reduces estimation errors when the dimension is fixed.

Figure 6 demonstrates the results of the slice-wise low-rankness regularizer. In Figure 6b, the x -axis $|\mathcal{S}_5|$ represents the rank of each slice of the tensor $\mathcal{A} \in \mathbb{R}^{d \times d \times d}$. Figure 6a and Figure 6c fix the rank of each slice to 5. The three distinct lines correspond to the estimation errors for sample size $n = \{1000, 2000, 3000\}$. From Figure 6a, we observe that with a fixed rank and sample size, the estimation error increases as the dimension d enlarges. Furthermore, Figure 6b shows that the estimation errors increase with the rank. Figure 6c demonstrates that with more samples, the estimation errors decrease.

 Table 2: The Frobenius norm $\|\hat{\mathcal{A}} - \mathcal{A}^*\|_F$ with standard variance changing the noise parameter

Structures	Methods	$d = 10 \times 10 \times 10$			$d = 20 \times 20 \times 20$		
		$\eta = 0.1$	$\eta = 1$	$\eta = 5$	$\eta = 0.1$	$\eta = 1$	$\eta = 5$
Sparsity	Element-wise	1.8909 ± 0.2004	1.9021 ± 0.2214	1.9015 ± 0.2084	19.4215 ± 2.2710	18.6899 ± 2.7556	19.2904 ± 2.3523
	Fiber-wise	1.8622 ± 0.2813	1.8461 ± 0.2783	1.8500 ± 0.2431	19.8920 ± 2.6721	19.3542 ± 2.8909	19.7628 ± 2.4745
	Slice-wise	2.0509 ± 0.3510	2.0421 ± 0.3242	1.9927 ± 0.3666	20.0062 ± 2.7153	20.4267 ± 2.7248	19.4231 ± 2.6420
Low-rankness	Mode-wise	2.6311 ± 0.3008	2.6947 ± 0.3254	2.6265 ± 0.3410	24.6129 ± 2.7010	23.8932 ± 2.7108	24.6571 ± 2.5114
	Slice-wise	3.0150 ± 0.3227	3.0045 ± 0.3754	3.0184 ± 0.3365	25.8502 ± 3.7601	25.6691 ± 3.5732	25.9134 ± 3.8282

For Corollaries 4 and 5, the error bounds contain terms involving η . However, the presence of conjugate operators, projection operators, and nuclear norm regularization complicates a direct analysis of these terms. Specifically, from equation (140), we observe that as the noise parameter η increases, the associated error term also increases, resulting in a larger overall error bound. In the table 2, we include synthetic data experiments with varying noise intensity. The results, which demonstrate the effect of changing η on the error bound, are provided in the table.

 Table 3: The Frobenius norm $\|\hat{\mathcal{A}} - \mathcal{A}^*\|_F$ with standard variance for higher dimension

Structures	Methods	3-order		4-order		5-order	
		$d = 8$	$d = 16$	$d = 8$	$d = 16$	$d = 8$	$d = 16$
Sparsity	Element-wise	1.0509 ± 0.1004	1.9021 ± 0.2214	7.9015 ± 1.2084	19.4215 ± 2.2710	58.6899 ± 7.7556	192.2904 ± 17.3523
	Fiber-wise	1.0622 ± 0.0813	1.8461 ± 0.2783	7.8500 ± 1.2431	19.8920 ± 2.6721	59.3542 ± 8.2909	190.7628 ± 16.4745
	Slice-wise	1.0909 ± 0.1510	2.0421 ± 0.3242	8.0927 ± 1.3666	20.0062 ± 2.7153	60.4267 ± 8.1248	193.4231 ± 17.6420
Low-rankness	Mode-wise	1.6311 ± 0.3008	2.6947 ± 0.3254	8.6265 ± 1.3410	34.6129 ± 3.4010	63.8932 ± 9.7108	224.6571 ± 21.5114
	Slice-wise	1.8150 ± 0.3227	2.7045 ± 0.3754	9.0184 ± 1.3365	35.8502 ± 3.7601	65.6691 ± 9.5732	245.9134 ± 22.8282

As outlined in the five corollaries of our paper, our theoretical framework is inherently generalizable to tensors of any order. Although the scope of this paper did not include experimental results for higher-order tensors, in table 3 we have conducted supplementary experiments that demonstrate promising outcomes for these cases.

Also, to explore whether the proposed methods performs robustness under a unknown structure, in table 4, we

Table 4: The Frobenius norm $\|\hat{\mathcal{A}} - \mathcal{A}^*\|_F$ with standard variance changing ground data structure of our proposed methods

Structures	Methods	Tensor Data Structures				
		element-sp	fiber-sp	slice-sp	lr-mode	lr-slice
Sparsity	Element-wise	1.0509 \pm 0.1004	1.0680 \pm 0.1027	1.1263 \pm 0.2046	1.8991 \pm 0.4002	1.9367 \pm 0.3979
	Fiber-wise	1.0931 \pm 0.1421	1.0622 \pm 0.0813	1.1305 \pm 0.2488	1.9054 \pm 0.3865	1.9274 \pm 0.4410
	Slice-wise	1.1014 \pm 0.1852	1.1226 \pm 0.2200	1.0909 \pm 0.1510	2.0221 \pm 0.4518	2.3185 \pm 0.4477
Low-rankness	Mode-wise	6.8502 \pm 1.2101	6.9333 \pm 1.2565	6.9068 \pm 1.1987	1.6311 \pm 0.3008	14.9490 \pm 1.4555
	Slice-wise	7.1481 \pm 1.3061	7.1636 \pm 1.2989	7.0701 \pm 1.3004	15.5770 \pm 1.3435	1.8150 \pm 0.3227

implement experiments on the proposed methods for each tensor structures, and the results are shown in the table.

A.2 Real-world Data

We have chosen several real-world images from the ImageNet 2012 dataset (Russakovsky et al., 2015) besides the image used in Section 5.2. We implement experiments based on different regularizers, revealing the following performance. In Figure 7, we pick four different images from the dataset, and Figure 8 shows these images with noise added. Figure 10 illustrates the images denoised with the proposed non-convex regularizers. As a comparison, the images denoised with convex methods are revealed in Figure 9.

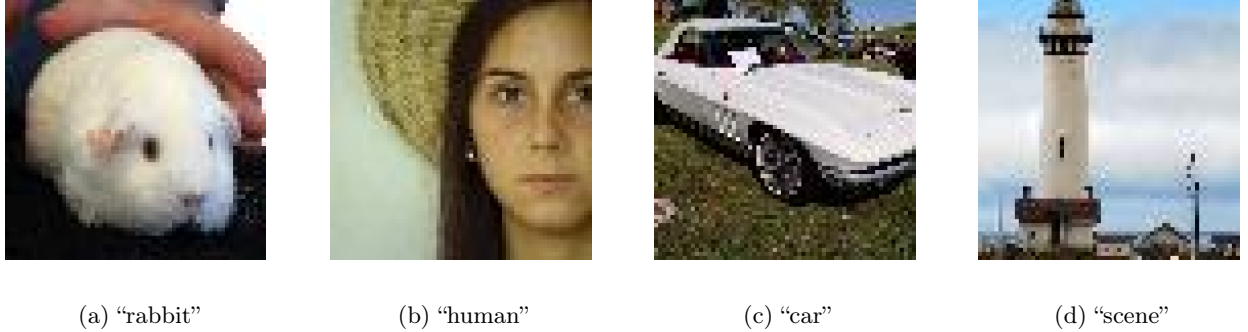


Figure 7: The original images.

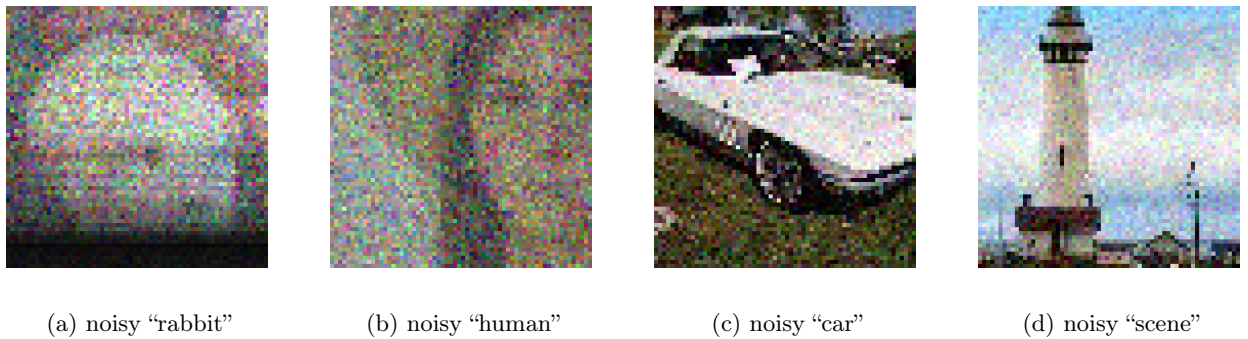


Figure 8: The noisy images.

We have also implemented additional real-world data experiments with the proposed methods. In table 5, the real data is considered the tensor to be estimate. Regarding the initialization of the covariate tensors \mathcal{A} in the

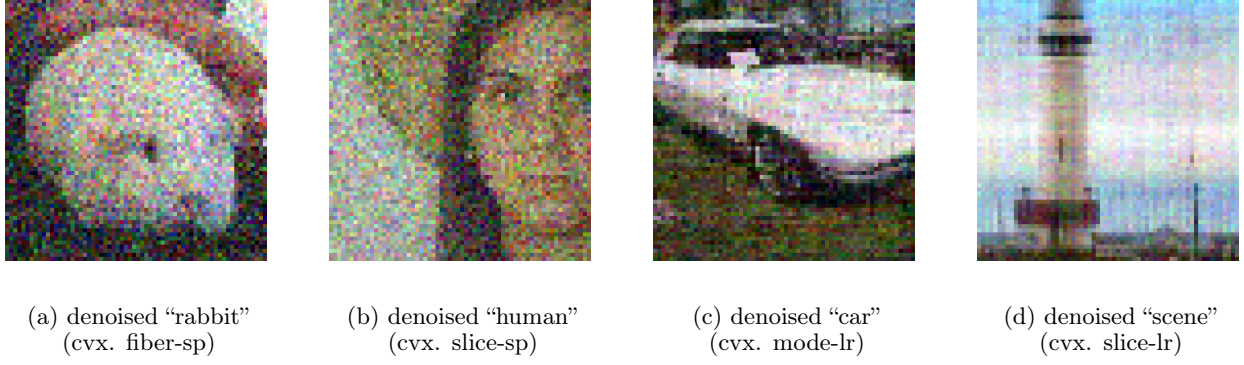


Figure 9: The denoised images using convex methods.

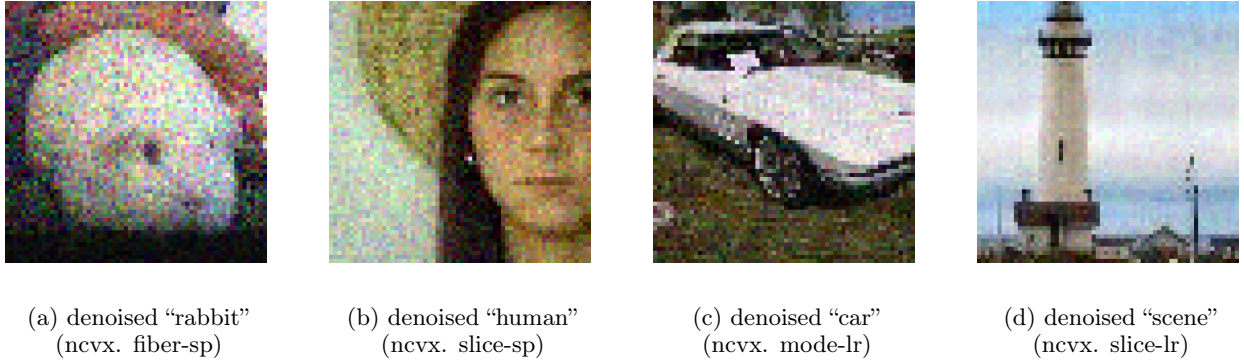


Figure 10: The denoised images using non-convex methods.

Table 5: The MSFE of the climate data(10 years) observation and alcoholic genetic predisposition data with our proposed methods

Dataset	Penalties	Sparsity		Low-rankness		
		Element	Fiber	Slice	Mode	Slice
NA-1990-2002-Monthly	SCAD	7.4556 ± 0.8235	8.0716 ± 0.9456	8.4187 ± 0.9491	11.6809 ± 2.0437	9.9750 ± 1.2203
	MCP	7.6087 ± 1.0003	7.9554 ± 0.9884	8.1305 ± 0.9050	12.6281 ± 2.1882	10.2314 ± 2.0015
	Convex	8.2502 ± 1.4887	8.7425 ± 1.7264	9.4577 ± 1.3004	14.5808 ± 3.3435	13.4508 ± 2.6688
EEG Database	SCAD	12.6865 ± 2.3544	13.8878 ± 2.2412	14.5640 ± 2.4898	18.7983 ± 5.0977	16.3202 ± 4.8331
	MCP	13.0024 ± 1.9973	14.0368 ± 2.0241	16.0431 ± 3.9314	18.4546 ± 4.8020	16.4002 ± 4.7771
	Convex	13.8001 ± 2.6764	14.5716 ± 2.1379	16.8890 ± 4.3051	19.6404 ± 5.3317	18.5783 ± 5.3854

real-data experiments, the number of covariate tensors \mathcal{A} corresponds to the sample size $n = 5000$, and the noise term \mathcal{E} are drawn independently from a Gaussian distribution with mean 0 and variance equal to $\eta = 0.01$.

The experimental data employed in this study were sourced from the University of Southern California’s Viterbi School of Engineering repository and the UCI Machine Learning Repository’s EEG Database. Specifically, the datasets can be accessed via the following links: <https://archive.ics.uci.edu/dataset/121/eeeg+database>, <https://viterbi-web.usc.edu/~liu32/data.html>.

B Background

This section provides a concise overview of the foundational concepts essential for understanding the subsequent proof.

Definition 1 (Gaussian Width). Let $\mathcal{S} \subset \mathbb{R}^{I_1 \times I_2 \times \dots \times I_M}$ be a subset of the tensor space. The Gaussian width of \mathcal{S} is defined as

$$w(\mathcal{S}) = \mathbb{E} \left(\sup_{\mathcal{B} \in \mathcal{S}} \langle \mathcal{B}, \mathcal{T} \rangle \right),$$

where $\mathcal{T} \in \mathbb{R}^{d_1 \times d_2 \times \dots \times d_M}$ is a tensor whose entry $\mathcal{T}(i_1, \dots, i_M)$ are independent standard Gaussian random variable. We refer readers to (Gordon, 1988; Negahban et al., 2012) for a comprehensive review.

Definition 2 (Subspace Compatibility Constant). For any subspace \mathcal{M} of $\mathbb{R}^{d_1 \times d_2 \times \dots \times d_M}$, the subspace compatibility constant $\Psi(\mathcal{M})$ is given by

$$\Psi(\mathcal{M}) = \sup_{\mathcal{U} \in \mathcal{M} \setminus \{0\}} \frac{R(\mathcal{U})}{\|\mathcal{U}\|_F}.$$

This constant quantifies the degree of compatibility between the regularizer $R(\mathcal{U})$ and the Frobenius norm within the subspace \mathcal{M} . It can be interpreted as a measure of the intrinsic dimensionality of \mathcal{M} .

Singular Value Decomposition. Let $\mathbf{X} \in \mathbb{R}^{m \times n}$ be a matrix with singular value decomposition (SVD) given by $\mathbf{X} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^\top$, where

- $\mathbf{U} \in \mathbb{R}^{m \times r}$ contains the left singular vectors,
- $\mathbf{V} \in \mathbb{R}^{n \times r}$ contains the right singular vectors,
- $\mathbf{\Sigma} = \text{diag}(\boldsymbol{\sigma}) \in \mathbb{R}^{r \times r}$ is a diagonal matrix with singular values $\boldsymbol{\sigma} = [\sigma_1(\mathbf{X}), \dots, \sigma_r(\mathbf{X})]$,
- $r = \text{rank}(\mathbf{X})$.

We define the subspaces \mathcal{F} and its orthogonal complement \mathcal{F}^\perp as follows:

$$\begin{aligned} \mathcal{F}(\mathbf{X}) &= \{\mathbf{W} \mid \text{row}(\mathbf{W}) \subseteq \text{span}(\mathbf{V}), \text{col}(\mathbf{W}) \subseteq \text{span}(\mathbf{U})\}, \\ \mathcal{F}^\perp(\mathbf{X}) &= \{\mathbf{W} \mid \text{row}(\mathbf{W}) \perp \text{span}(\mathbf{V}), \text{col}(\mathbf{W}) \perp \text{span}(\mathbf{U})\}. \end{aligned}$$

Here, $\text{row}(\mathbf{W})$ and $\text{col}(\mathbf{W})$ denote the row space and column space of \mathbf{W} , respectively, $\text{span}(\mathbf{V})$ denotes the subspace spanned by \mathbf{V} , and \perp denotes orthogonality with respect to the Euclidean inner product. For brevity, we will use the shorthand notations of \mathcal{F} and \mathcal{F}^\perp when the dependence on \mathbf{X} is clear from the context.

The projection operators onto the subspace \mathcal{F} and \mathcal{F}^\perp are denoted by $\Pi_{\mathcal{F}}(\cdot)$ and $\Pi_{\mathcal{F}^\perp}(\cdot)$, respectively, and are defined as:

$$\begin{aligned} \Pi_{\mathcal{F}}(\mathbf{X}) &= \mathbf{U} \mathbf{U}^\top \mathbf{X} \mathbf{V} \mathbf{V}^\top, \\ \Pi_{\mathcal{F}^\perp}(\mathbf{X}) &= (\mathbf{I} - \mathbf{U} \mathbf{U}^\top) \mathbf{X} (\mathbf{I} - \mathbf{V} \mathbf{V}^\top), \end{aligned}$$

where \mathbf{I} denotes the identity matrix with contextually appropriate dimensions.

C Proof of the Corollary 1

We begin by demonstrating that the element-wise sparsity penalty can be reformulated as the sum of the ℓ_1 penalty and a concave part. Specifically, we have

$$R_\lambda(\mathcal{A}) = \sum_{i_1=1}^{d_1} \dots \sum_{i_M=1}^{d_M} p_\lambda(\mathcal{A}_{i_1, \dots, i_M}) = \lambda \|\mathcal{A}_{i_1, \dots, i_M}\|_1 + Q_\lambda(\mathcal{A}),$$

where $\|\mathcal{A}_{i_1, \dots, i_M}\|_1 = \sum_{i_1=1}^{d_1} \dots \sum_{i_M=1}^{d_M} |\mathcal{A}_{i_1, \dots, i_M}|$ and $Q_\lambda(\mathcal{A}) = \sum_{i_1=1}^{d_1} \dots \sum_{i_M=1}^{d_M} q_\lambda(\mathcal{A}_{i_1, \dots, i_M})$. Define $\tilde{L}(\mathcal{A}) = L(\mathcal{A}) + Q_\lambda(\mathcal{A})$.

Lemma 1. Under Assumptions 1 and 2, the loss function $\tilde{L}(\mathcal{A})$ satisfies the restricted strong convexity

$$\tilde{L}(\mathcal{A}') \geq \tilde{L}(\mathcal{A}) + \langle \nabla \tilde{L}(\mathcal{A}), \mathcal{A}' - \mathcal{A} \rangle + \frac{\rho^- - \zeta^-}{2} \|\mathcal{A}' - \mathcal{A}\|_{\text{F}}^2,$$

and the restricted strong smoothness

$$\tilde{L}(\mathcal{A}') \leq \tilde{L}(\mathcal{A}) + \langle \nabla \tilde{L}(\mathcal{A}), \mathcal{A}' - \mathcal{A} \rangle + \frac{\rho^+ - \zeta^+}{2} \|\mathcal{A}' - \mathcal{A}\|_{\text{F}}^2.$$

Proof. Recall that $Q_\lambda(\mathcal{A})$ represents the concave component of the non-convex penalty $R_\lambda(\mathcal{A}_{i_1, \dots, i_N})$, implying that $-Q_\lambda(\mathcal{A})$ is convex. Specifically, $Q_\lambda(\mathcal{A})$ can be expressed as a sum over its elements $Q_\lambda(\mathcal{A}) = \sum_{i_1=1}^{d_1} \dots \sum_{i_N=1}^{d_N} q_\lambda(\mathcal{A}_{i_1, \dots, i_N})$, where $q_\lambda(\mathcal{A}_{i_1, \dots, i_N})$ satisfies the third regularity condition specified in Assumption 3. From this assumption, we have

$$-\zeta^-(\mathcal{A}'_{i_1, \dots, i_N} - \mathcal{A}_{i_1, \dots, i_N})^2 \leq (q'_\lambda(\mathcal{A}'_{i_1, \dots, i_N}) - q'_\lambda(\mathcal{A}_{i_1, \dots, i_N}))(\mathcal{A}'_{i_1, \dots, i_N} - \mathcal{A}_{i_1, \dots, i_N}) \leq -\zeta^+(\mathcal{A}'_{i_1, \dots, i_N} - \mathcal{A}_{i_1, \dots, i_N})^2. \quad (3)$$

By aggregating over all elements, we deduce that the convex function $-Q_\lambda(\mathcal{A})$ satisfies the following inequality

$$\langle (\nabla(-Q_\lambda(\mathcal{A}')) - \nabla(-Q_\lambda(\mathcal{A})))^\top, \mathcal{A}' - \mathcal{A} \rangle \leq \zeta^- \|\mathcal{A}' - \mathcal{A}\|_{\text{F}}^2, \quad (4)$$

$$\langle (\nabla(-Q_\lambda(\mathcal{A}')) - \nabla(-Q_\lambda(\mathcal{A})))^\top, \mathcal{A}' - \mathcal{A} \rangle \geq \zeta^+ \|\mathcal{A}' - \mathcal{A}\|_{\text{F}}^2. \quad (5)$$

Inequalities (4) and (5) correspond to the definitions of RSC and RSS for the function $-Q_\lambda(\mathcal{A})$, respectively. Specifically, they imply that $-Q_\lambda(\mathcal{A})$ is both ζ^- -smooth and ζ^+ -strongly convex. Consequently, we have the following

$$-Q_\lambda(\mathcal{A}') \leq -Q_\lambda(\mathcal{A}) - \langle \nabla Q_\lambda(\mathcal{A}), \mathcal{A}' - \mathcal{A} \rangle + \frac{\zeta^-}{2} \|\mathcal{A}' - \mathcal{A}\|_{\text{F}}^2, \quad (6)$$

$$-Q_\lambda(\mathcal{A}') \geq -Q_\lambda(\mathcal{A}) - \langle \nabla Q_\lambda(\mathcal{A}), \mathcal{A}' - \mathcal{A} \rangle + \frac{\zeta^+}{2} \|\mathcal{A}' - \mathcal{A}\|_{\text{F}}^2. \quad (7)$$

For the loss function $L(\mathcal{A})$, applying Taylor's theorem and the mean value theorem yields

$$L(\mathcal{A}') = L(\mathcal{A}) + \langle \nabla L(\mathcal{A}), \mathcal{A}' - \mathcal{A} \rangle + \frac{1}{2} \langle \nabla^2 L(\beta \mathcal{A}' + (1 - \beta) \mathcal{A}), (\mathcal{A}' - \mathcal{A}) \otimes (\mathcal{A}' - \mathcal{A}) \rangle, \quad (8)$$

for some $\beta \in [0, 1]$. Here, \otimes denotes the Kronecker product. Given two tensors $\mathcal{A}, \mathcal{A}' \in \mathbb{R}^{d_1 \times \dots \times d_M}$, their Kronecker product results in a tensor \mathcal{A}'' of dimension $(d_1 d_1) \times \dots \times (d_M d_M)$. Each element $\mathcal{A}''_{i_1 j_1 i_2 j_2 \dots i_M j_M}$ is defined as $\mathcal{A}_{i_1 i_2 \dots i_M} \times \mathcal{A}'_{j_1 j_2 \dots j_M}$.

Under Assumptions 1 and 2, we have

$$L(\mathcal{A} + \mathcal{B}) - L(\mathcal{A}) \geq \langle \nabla L(\mathcal{A}), \mathcal{B} \rangle + \frac{\rho^-}{2} \|\mathcal{B}\|_{\text{F}}^2, \quad (9)$$

$$L(\mathcal{A} + \mathcal{B}) - L(\mathcal{A}) \leq \langle \nabla L(\mathcal{A}), \mathcal{B} \rangle + \frac{\rho^+}{2} \|\mathcal{B}\|_{\text{F}}^2. \quad (10)$$

Recall that $\tilde{L}(\mathcal{A}) = L(\mathcal{A}) + Q_\lambda(\mathcal{A})$. Thus, we obtain

$$\tilde{L}(\mathcal{A}') \geq \tilde{L}(\mathcal{A}) + \langle \nabla \tilde{L}(\mathcal{A}), \mathcal{A}' - \mathcal{A} \rangle + \frac{\rho^- - \zeta^-}{2} \|\mathcal{A}' - \mathcal{A}\|_{\text{F}}^2, \quad (11)$$

$$\tilde{L}(\mathcal{A}') \leq \tilde{L}(\mathcal{A}) + \langle \nabla \tilde{L}(\mathcal{A}), \mathcal{A}' - \mathcal{A} \rangle + \frac{\rho^+ - \zeta^+}{2} \|\mathcal{A}' - \mathcal{A}\|_{\text{F}}^2, \quad (12)$$

□

Lemma 2. Let C be a constant. Suppose there exists an integer $\tilde{s}_1 > C|\mathbf{S}_1|$, and that \mathcal{A} satisfies $\|\mathcal{A}_{\overline{\mathbf{S}_1}}\|_0 \leq \tilde{s}_1$, $\omega(\mathcal{A}) \leq \frac{\lambda}{2}$, where $\omega(\mathcal{A}) = \min_{\mathcal{H} \in \partial \|\mathcal{A}\|_1} \left\{ \left\| \nabla \tilde{L}(\mathcal{A}) + \lambda \mathcal{H} \right\|_{\max} \right\}$, and $\|\nabla L(\mathcal{A}^*)\|_{\max} \leq \lambda/8$. Under Assumptions 1 and 2, \mathcal{A} satisfies

$$\|\mathcal{A} - \mathcal{A}^*\|_F \leq \frac{21/8}{\rho^- - \zeta^-} \lambda \sqrt{|\mathbf{S}_1|}.$$

Proof. Given that $\|\mathcal{A}_{\overline{\mathbf{S}_1}}\|_0 \leq \tilde{s}_1$ and $\|\mathcal{A}_{\overline{\mathbf{S}_1}}^*\|_0 = 0$, it follows that $\|(\mathcal{A} - \mathcal{A}^*)_{\overline{\mathbf{S}_1}}\|_0 \leq \tilde{s}_1$. Based on Lemma 1, we can derive the following inequalities

$$\tilde{L}(\mathcal{A}^*) \geq \tilde{L}(\mathcal{A}) + \left\langle \nabla \tilde{L}(\mathcal{A}), \mathcal{A}^* - \mathcal{A} \right\rangle + \frac{\rho^- - \zeta^-}{2} \|\mathcal{A}^* - \mathcal{A}\|_F^2, \quad (13)$$

$$\tilde{L}(\mathcal{A}) \geq \tilde{L}(\mathcal{A}^*) + \left\langle \nabla \tilde{L}(\mathcal{A}^*), \mathcal{A} - \mathcal{A}^* \right\rangle + \frac{\rho^- - \zeta^-}{2} \|\mathcal{A} - \mathcal{A}^*\|_F^2. \quad (14)$$

Adding (13) and (14), we obtain

$$\left\langle \nabla \tilde{L}(\mathcal{A}), \mathcal{A}^* - \mathcal{A} \right\rangle \geq \left\langle \nabla \tilde{L}(\mathcal{A}^*), \mathcal{A} - \mathcal{A}^* \right\rangle + (\rho^- - \zeta^-) \|\mathcal{A}^* - \mathcal{A}\|_F^2. \quad (15)$$

Let $\mathcal{H} \in \partial \|\mathcal{A}\|_1$ denote the sub-gradient and \mathcal{S} be a set. According to the Karush-Kuhn-Tucker (KKT) condition, we have

$$\nabla \tilde{L}(\mathcal{A}) + \lambda \mathcal{H} = \mathbf{0}. \quad (16)$$

Determining the optimal solution is challenging, therefore, we introduce a measure of sub-optimality

$$\omega(\mathcal{A}) = \min_{\mathcal{H}' \in \partial \|\mathcal{A}\|_1} \max_{\mathcal{A}' \in \mathcal{S}} \left\{ \frac{1}{\|\mathcal{A} - \mathcal{A}'\|_1} \left\langle \mathcal{A} - \mathcal{A}', \nabla \tilde{L}(\mathcal{A}) + \lambda \mathcal{H}' \right\rangle \right\}. \quad (17)$$

We define our algorithm's stopping criterion as $\omega(\mathcal{A}) \leq \varepsilon$. Consequently, the sub-optimality can be expressed as

$$\omega(\mathcal{A}) = \max_{\mathcal{A}' \in \mathcal{S}} \left\{ \frac{1}{\|\mathcal{A} - \mathcal{A}'\|_1} \left\langle \mathcal{A} - \mathcal{A}', \nabla \tilde{L}(\mathcal{A}) + \lambda \mathcal{H} \right\rangle \right\}. \quad (18)$$

Adding $\lambda \langle \mathcal{A} - \mathcal{A}^*, \mathcal{H}' \rangle$ to the both sides of (15), we obtain

$$\left\langle \mathcal{A} - \mathcal{A}^*, \nabla \tilde{L}(\mathcal{A}) + \lambda \mathcal{H} \right\rangle \geq \left\langle \mathcal{A} - \mathcal{A}^*, \nabla \tilde{L}(\mathcal{A}) \right\rangle + (\rho^- - \zeta^-) \|\mathcal{A}^* - \mathcal{A}\|_F^2 + \lambda \langle \mathcal{A} - \mathcal{A}^*, \mathcal{H} \rangle. \quad (19)$$

Since $\mathcal{A}^* \in \mathcal{S}$, we have

$$\frac{1}{\|\mathcal{A} - \mathcal{A}^*\|_1} \left\langle \mathcal{A} - \mathcal{A}^*, \nabla \tilde{L}(\mathcal{A}) + \lambda \mathcal{H} \right\rangle \leq \max_{\mathcal{A}' \in \mathcal{S}} \left\{ \frac{1}{\|\mathcal{A} - \mathcal{A}'\|_1} \left\langle \mathcal{A} - \mathcal{A}', \nabla \tilde{L}(\mathcal{A}) + \lambda \mathcal{H} \right\rangle \right\} = v(\mathcal{A}). \quad (20)$$

Recall that we assume $v(\mathcal{A}) \leq \frac{\lambda}{2}$, we obtain

$$\left\langle \mathcal{A} - \mathcal{A}^*, \nabla \tilde{L}(\mathcal{A}) + \lambda \mathcal{H} \right\rangle \leq v(\mathcal{A}) \leq \frac{\lambda}{2} \|\mathcal{A} - \mathcal{A}^*\|_1. \quad (21)$$

Combining (15) and (21), we obtain

$$\frac{\lambda}{2} \|\mathcal{A} - \mathcal{A}^*\|_1 \geq \underbrace{\left\langle \mathcal{A} - \mathcal{A}^*, \nabla \tilde{L}(\mathcal{A}^*) \right\rangle}_{\text{I}} + (\rho^- - \zeta^-) \|\mathcal{A}^* - \mathcal{A}\|_F^2 + \underbrace{\lambda \langle \mathcal{A} - \mathcal{A}^*, \mathcal{H} \rangle}_{\text{II}}. \quad (22)$$

For term I, separating the support of $\mathcal{A} - \mathcal{A}^*$ into \mathcal{S}_1 and $\overline{\mathcal{S}_1}$, we have

$$\langle \mathcal{A} - \mathcal{A}^*, \nabla \tilde{L}(\mathcal{A}^*) \rangle = \langle \mathcal{A} - \mathcal{A}^*, \nabla \tilde{L}(\mathcal{A}^*) \rangle + \langle \mathcal{A} - \mathcal{A}^*, \nabla Q_\lambda(\mathcal{A}^*) \rangle \quad (23)$$

$$\geq -\|\mathcal{A} - \mathcal{A}^*\|_1 \|\nabla L(\mathcal{A}^*)\|_{\max} + \langle \mathcal{A} - \mathcal{A}^*, \nabla Q_\lambda(\mathcal{A}^*) \rangle \quad (24)$$

$$= -\|(\mathcal{A} - \mathcal{A}^*)_{\overline{\mathcal{S}_1}}\|_1 \|\nabla L(\mathcal{A}^*)\|_{\max} - \|(\mathcal{A} - \mathcal{A}^*)_{\mathcal{S}_1}\|_1 \|\nabla L(\mathcal{A}^*)\|_{\max} \quad (25)$$

$$\begin{aligned} &+ \langle (\mathcal{A} - \mathcal{A}^*)_{\mathcal{S}_1}, (\nabla Q_\lambda(\mathcal{A}^*))_{\mathcal{S}_1} \rangle + \langle (\mathcal{A} - \mathcal{A}^*)_{\overline{\mathcal{S}_1}}, (\nabla Q_\lambda(\mathcal{A}^*))_{\overline{\mathcal{S}_1}} \rangle \\ &\geq -\|(\mathcal{A} - \mathcal{A}^*)_{\overline{\mathcal{S}_1}}\|_1 \|\nabla L(\mathcal{A}^*)\|_{\max} - \|(\mathcal{A} - \mathcal{A}^*)_{\mathcal{S}_1}\|_1 \|\nabla L(\mathcal{A}^*)\|_{\max} \\ &\quad - \|(\mathcal{A} - \mathcal{A}^*)_{\mathcal{S}_1}\|_1 \|\nabla Q_\lambda(\mathcal{A}^*)\|_{\max}. \end{aligned} \quad (26)$$

For term II, separating the support of $\mathcal{A} - \mathcal{A}^*$ into \mathcal{S}_1 and $\overline{\mathcal{S}_1}$, we have

$$\lambda \langle \mathcal{A} - \mathcal{A}^*, \mathcal{H} \rangle = \lambda \langle (\mathcal{A} - \mathcal{A}^*)_{\mathcal{S}_1}, (\mathcal{H})_{\mathcal{S}_1} \rangle + \lambda \langle (\mathcal{A} - \mathcal{A}^*)_{\overline{\mathcal{S}_1}}, (\mathcal{H})_{\overline{\mathcal{S}_1}} \rangle \quad (27)$$

$$\geq -\lambda \|(\mathcal{A} - \mathcal{A}^*)_{\mathcal{S}_1}\|_1 \|(\mathcal{H})_{\mathcal{S}_1}\|_{\max} + \lambda \langle (\mathcal{A})_{\overline{\mathcal{S}_1}}, (\mathcal{H})_{\overline{\mathcal{S}_1}} \rangle \quad (28)$$

$$\geq -\lambda \|(\mathcal{A} - \mathcal{A}^*)_{\mathcal{S}_1}\|_1 + \lambda \sum_{(i_1, \dots, i_M) \in \overline{\mathcal{S}_1}} |\mathcal{A}_{i_1, \dots, i_M}| \quad (29)$$

$$= -\lambda \|(\mathcal{A} - \mathcal{A}^*)_{\mathcal{S}_1}\|_1 + \lambda \|(\mathcal{A} - \mathcal{A}^*)_{\overline{\mathcal{S}_1}}\|_1. \quad (30)$$

Thus, we obtain

$$\begin{aligned} \frac{\lambda}{2} \|\mathcal{A} - \mathcal{A}^*\|_1 &\geq -\|(\mathcal{A} - \mathcal{A}^*)_{\overline{\mathcal{S}_1}}\|_1 \|\nabla L(\mathcal{A}^*)\|_{\max} - \|(\mathcal{A} - \mathcal{A}^*)_{\mathcal{S}_1}\|_1 \|\nabla L(\mathcal{A}^*)\|_{\max} \\ &\quad - \|(\mathcal{A} - \mathcal{A}^*)_{\mathcal{S}_1}\|_1 \|\nabla Q_\lambda(\mathcal{A}^*)\|_{\max} + (\rho^- - \zeta^-) \|\mathcal{A}^* - \mathcal{A}\|_F^2 \\ &\quad - \lambda \|(\mathcal{A} - \mathcal{A}^*)_{\mathcal{S}_1}\|_1 + \lambda \|(\mathcal{A} - \mathcal{A}^*)_{\overline{\mathcal{S}_1}}\|_1. \end{aligned} \quad (31)$$

We separate the left-hand side of (31) as

$$\frac{\lambda}{2} \|\mathcal{A} - \mathcal{A}^*\|_1 = \frac{\lambda}{2} \|(\mathcal{A} - \mathcal{A}^*)_{\mathcal{S}_1}\|_1 + \frac{\lambda}{2} \|(\mathcal{A} - \mathcal{A}^*)_{\overline{\mathcal{S}_1}}\|_1. \quad (32)$$

Rearranging the terms, we obtain

$$\begin{aligned} &(\rho^- - \zeta^-) \|\mathcal{A}^* - \mathcal{A}\|_F^2 + \left(\frac{\lambda}{2} - \|\nabla L(\mathcal{A}^*)\|_{\max} \right) \|(\mathcal{A} - \mathcal{A}^*)_{\overline{\mathcal{S}_1}}\|_1 \\ &\leq \left(\frac{3\lambda}{2} + \|\nabla L(\mathcal{A}^*)\|_{\max} + \|\nabla Q_\lambda(\mathcal{A}^*)\|_{\max} \right) \|(\mathcal{A} - \mathcal{A}^*)_{\mathcal{S}_1}\|_1. \end{aligned} \quad (33)$$

Recall that $\|\nabla L(\mathcal{A}^*)\|_{\max} \leq \frac{\lambda}{8}$, we have

$$(\rho^- - \zeta^-) \|\mathcal{A}^* - \mathcal{A}\|_F^2 \leq \left(\frac{3\lambda}{2} + \|\nabla L(\mathcal{A}^*)\|_{\max} + \|\nabla Q_\lambda(\mathcal{A}^*)\|_{\max} \right) \|(\mathcal{A} - \mathcal{A}^*)_{\mathcal{S}_1}\|_1 \quad (34)$$

$$\leq \left(\frac{3\lambda}{2} + \frac{\lambda}{8} + \lambda \right) \|(\mathcal{A} - \mathcal{A}^*)_{\mathcal{S}_1}\|_1 \quad (35)$$

$$\leq \frac{21\lambda}{8} \sqrt{|\mathcal{S}_1|} \|(\mathcal{A} - \mathcal{A}^*)_{\mathcal{S}_1}\|_F \quad (36)$$

$$\leq \frac{21\lambda}{8} \sqrt{|\mathcal{S}_1|} \|\mathcal{A} - \mathcal{A}^*\|_F. \quad (37)$$

Given that $\rho^- - \zeta^- > 0$, we have

$$\|\mathcal{A} - \mathcal{A}^*\|_F \leq \frac{21/8}{\rho^- - \zeta^-} \lambda \sqrt{|\mathcal{S}_1|}. \quad (38)$$

□

Lemma 3. Consider the regularization parameter λ and assume that the derivative of the non-convex penalty satisfies $p'_\lambda(\mathcal{A}_{d_1, \dots, d_M}) = 0$ whenever $|\mathcal{A}_{d_1, \dots, d_M}| \geq \nu$ for some $\nu > 0$. Let $\mathcal{S}_1^I \cup \mathcal{S}_1^{II} = \mathcal{S}_1$. For indices $(i_1, \dots, i_M) \in \mathcal{S}_1^I \subseteq \mathcal{S}_1$, we assume $|\mathcal{A}_{d_1, \dots, d_M}^*| \geq \nu$, and for indices $(i_1, \dots, i_M) \in \mathcal{S}_1^{II} \subseteq \mathcal{S}_1$, we assume $|\mathcal{A}_{d_1, \dots, d_M}^*| \leq \nu$. Under Assumptions 1~4, we derive the following bound

$$\|\hat{\mathcal{A}} - \mathcal{A}^*\|_F \leq \frac{1}{\rho^- - \zeta^-} \|(\nabla L(\mathcal{A}^*))_{\mathcal{S}_1^I}\|_F + \frac{3}{\rho^- - \zeta^-} \lambda \sqrt{|\mathcal{S}_1^{II}|}.$$

Proof. Define the sub-gradients $\mathcal{H}^* \in \partial \|\mathcal{A}^*\|_1$ and $\hat{\mathcal{H}} \in \partial \|\hat{\mathcal{A}}\|_1$.

Note that $\hat{\mathcal{A}}$ satisfies the optimality condition that $\omega(\hat{\mathcal{A}}) \leq 0$, we have

$$\max_{\mathcal{A}' \in \mathcal{S}} \left\{ \langle \hat{\mathcal{A}} - \mathcal{A}', \nabla \tilde{L}(\hat{\mathcal{A}}) + \lambda \hat{\mathcal{H}} \rangle \right\} \leq 0. \quad (39)$$

Given that $\left\| \left(\hat{\mathcal{A}} \right)_{\overline{\mathcal{S}_1}} \right\|_0 \leq \tilde{s}_1$, since $\left\| \left(\hat{\mathcal{A}} - \mathcal{A}^* \right)_{\overline{\mathcal{S}_1}} \right\|_0 \leq \tilde{s}_1$, according to Lemma 1, we obtain

$$\tilde{L}(\hat{\mathcal{A}}) \geq \tilde{L}(\mathcal{A}^*) + \langle \nabla \tilde{L}(\mathcal{A}^*), \hat{\mathcal{A}} - \mathcal{A}^* \rangle + \frac{\rho^- - \zeta^-}{2} \|\hat{\mathcal{A}} - \mathcal{A}^*\|_F^2, \quad (40)$$

$$\tilde{L}(\mathcal{A}^*) \geq \tilde{L}(\hat{\mathcal{A}}) + \langle \nabla \tilde{L}(\hat{\mathcal{A}}), \mathcal{A}^* - \hat{\mathcal{A}} \rangle + \frac{\rho^- - \zeta^-}{2} \|\mathcal{A}^* - \hat{\mathcal{A}}\|_F^2. \quad (41)$$

By the convexity of ℓ_1 norm, we have

$$\lambda \|\hat{\mathcal{A}}\|_1 \leq \lambda \|\mathcal{A}^*\|_1 + \lambda \langle \hat{\mathcal{A}} - \mathcal{A}^*, \mathcal{H}^* \rangle, \quad (42)$$

$$\lambda \|\mathcal{A}^*\|_1 \leq \lambda \|\hat{\mathcal{A}}\|_1 + \lambda \langle \mathcal{A}^* - \hat{\mathcal{A}}, \hat{\mathcal{H}} \rangle. \quad (43)$$

Adding (40) ~ (43), we obtain

$$0 \geq \underbrace{\langle \nabla L(\mathcal{A}^*) + \nabla \mathcal{Q}_\lambda(\mathcal{A}_{i_1 \dots i_M}^*) + \lambda \mathcal{H}^*, \hat{\mathcal{A}} - \mathcal{A}^* \rangle}_{(i)} + \underbrace{\langle \nabla \tilde{L}(\hat{\mathcal{A}}) + \lambda \mathcal{H}^*, \mathcal{A}^* - \hat{\mathcal{A}} \rangle}_{(ii)} + (\rho^- - \zeta^-) \|\hat{\mathcal{A}} - \mathcal{A}^*\|_F^2. \quad (44)$$

From the optimality condition (39), we have

$$\langle \nabla \tilde{L}(\hat{\mathcal{A}}) + \lambda \hat{\mathcal{H}}, \mathcal{A}^* - \hat{\mathcal{A}} \rangle \leq \max_{\mathcal{A}' \in \mathcal{S}} \left\{ \langle \hat{\mathcal{A}} - \mathcal{A}', \nabla \tilde{L}(\hat{\mathcal{A}}) + \lambda \hat{\mathcal{H}} \rangle \right\} \leq 0, \quad (45)$$

which implies the term (ii) in (44) is non-negative. Consequently, we can arrange (44) to obtain

$$\begin{aligned} & (\rho^- - \zeta^-) \|\hat{\mathcal{A}} - \mathcal{A}^*\|_F^2 \\ & \leq \langle \nabla L(\mathcal{A}^*) + \nabla \mathcal{Q}_\lambda(\mathcal{A}^*) + \lambda_t \mathcal{H}^*, \hat{\mathcal{A}} - \mathcal{A}^* \rangle \\ & \leq \min_{\mathcal{H}^* \in \partial \|\mathcal{A}^*\|_1} \left\{ \sum_{i_1=1}^{I_1} \cdots \sum_{i_M=1}^{I_M} |(\nabla L(\mathcal{A}^*) + \nabla \mathcal{Q}_\lambda(\mathcal{A}^*) + \lambda_t \mathcal{H}^*)_{i_1 \dots i_M}| \cdot \left| (\mathcal{A}^* - \hat{\mathcal{A}})_{i_1 \dots i_M} \right| \right\}. \end{aligned} \quad (46)$$

We proceed by decomposing the summation on the right-hand side of (46) into three distinct parts

- $(i_1, \dots, i_M) \in \overline{\mathcal{S}_1}$,
- $(i_1, \dots, i_M) \in \mathcal{S}_1^I$,

- $(i_1, \dots, i_M) \in \mathcal{S}_1^{\text{II}}$.

Here, $\mathcal{S}_1^{\text{I}} = \{(i_1, \dots, i_M) \mid |\mathcal{A}_{(i_1, \dots, i_M)}| \geq \nu\}$, $\mathcal{S}_1^{\text{I}} = \{(i_1, \dots, i_M) \mid |\mathcal{A}_{(i_1, \dots, i_M)}| < \nu\}$, and $\nu > 0$ is defined in Assumption 3.

(i) For any index $(i_1, \dots, i_M) \in \overline{\mathcal{S}_1}$, the regularity condition yields

$$\nabla \mathcal{Q}_\lambda(\mathcal{A}^*)_{i_1 \dots i_M} = q'_\lambda((\mathcal{A}^*)_{i_1 \dots i_M}) = q'_\lambda(0), \quad \text{for } j \in \overline{\mathcal{S}_1}. \quad (47)$$

Assuming that $\|\nabla L(\mathcal{A}^*)\|_{\max} \leq \frac{\lambda}{8}$, it follows that

$$\max_{(i_1, \dots, i_M) \in \overline{\mathcal{S}_1}} |(\nabla L(\mathcal{A}^*))_{i_1, \dots, i_M}| \leq \|\nabla L(\mathcal{A}^*)\|_{\max} \leq \frac{\lambda}{8} \leq \lambda. \quad (48)$$

Therefore,

$$\max_{(i_1, \dots, i_M) \in \overline{\mathcal{S}_1}} |(\nabla L(\mathcal{A}^*) + \mathcal{Q}_\lambda(\mathcal{A}^*))_{i_1, \dots, i_M}| \leq \lambda. \quad (49)$$

Moreover, since $\mathcal{H}^* \in \partial \|\mathcal{A}^*\|_1$, it holds that $\lambda \mathcal{H}^*_{i_1, \dots, i_M} \in [-\lambda, \lambda]$. Consequently, for each $(i_1, \dots, i_M) \in \overline{\mathcal{S}_1}$, we can select $\mathcal{H}^*_{i_1, \dots, i_M}$ such that

$$|(\nabla L(\mathcal{A}^*) + \nabla \mathcal{Q}_\lambda(\mathcal{A}^*))_{i_1 \dots i_M} + \lambda \mathcal{H}^*_{i_1 \dots i_M}| = 0. \quad (50)$$

This implies

$$\min_{\mathcal{H}^* \in \partial \|\mathcal{A}^*\|_1} \{ |(\nabla L(\mathcal{A}^*) + \nabla \mathcal{Q}_\lambda(\mathcal{A}^*) + \lambda \mathcal{H}^*)_{i_1 \dots i_M}| \} = 0, \quad \text{for } (i_1 \dots i_M) \in \overline{\mathcal{S}_1}. \quad (51)$$

Therefore, we obtain

$$\min_{\mathcal{H}^* \in \partial \|\mathcal{A}^*\|_1} \left\{ \sum_{(i_1, \dots, i_M) \in \overline{\mathcal{S}_1}} |(\nabla L(\mathcal{A}^*) + \nabla \mathcal{Q}_\lambda(\mathcal{A}^*) + \lambda \mathcal{H}^*)_{i_1 \dots i_M}| \cdot \left| (\mathcal{A}^* - \hat{\mathcal{A}})_{i_1 \dots i_M} \right| \right\} = 0. \quad (52)$$

(ii) For indices $(i_1, \dots, i_M) \in \mathcal{S}_1^{\text{I}}$, we have $|\mathcal{A}^*_{i_1, \dots, i_M}| \geq \nu$. Given that $R(\mathcal{A}) = \lambda \|\mathcal{A}\|_1 + \mathcal{Q}_\lambda(\mathcal{A}_{i_1 \dots i_M})$, our assumption on $R(\mathcal{A})$ ensures that

$$(\nabla \mathcal{Q}_\lambda(\mathcal{A}^*) + \lambda \mathcal{H}^*)_{i_1 \dots i_M} = p'_\lambda(\mathcal{A}^*_{i_1 \dots i_M}) = 0, \quad \text{for } (i_1 \dots i_M) \in \mathcal{S}_1^{\text{I}}. \quad (53)$$

This leads to

$$\min_{\mathcal{H}^* \in \partial \|\mathcal{A}^*\|_1} \left\{ \sum_{(i_1 \dots i_M) \in \mathcal{S}_1^{\text{I}}} |(\nabla L(\mathcal{A}^*) + \nabla \mathcal{Q}_\lambda(\mathcal{A}^*) + \lambda \mathcal{H}^*)_{i_1 \dots i_M}| \cdot \left| (\mathcal{A}^* - \hat{\mathcal{A}})_{i_1 \dots i_M} \right| \right\} \quad (54)$$

$$= \sum_{(i_1 \dots i_M) \in \mathcal{S}_1^{\text{I}}} |(\nabla L(\mathcal{A}^*))_{i_1 \dots i_M}| \cdot \left| (\mathcal{A}^* - \hat{\mathcal{A}})_{i_1 \dots i_M} \right| \quad (55)$$

$$\leq \left\| (\nabla L(\mathcal{A}^*))_{\mathcal{S}_1^{\text{I}}} \right\|_{\text{F}} \cdot \left\| \mathcal{A}^* - \hat{\mathcal{A}} \right\|_{\text{F}}. \quad (56)$$

(iii) For indices $(i_1, \dots, i_M) \in \mathcal{S}_1^{\text{II}}$, we have $|\mathcal{A}^*_{i_1, \dots, i_M}| < \nu$. Given that $\|\nabla L(\mathcal{A}^*)\|_{\max} \leq \frac{\lambda}{8}$, we have

$$\max_{(i_1 \dots i_M) \in \mathcal{S}_1^{\text{II}}} |(\nabla L(\mathcal{A}^*))_{i_1 \dots i_M}| \leq |(\nabla L(\mathcal{A}^*))_{i_1 \dots i_M}|_{\max} \leq \lambda/8. \quad (57)$$

Meanwhile, we have

$$\max_{(i_1 \dots i_M) \in \mathcal{S}_1^\Pi} |(\nabla \mathcal{Q}_\lambda(\mathcal{A}^*))_{i_1 \dots i_M}| = \max_{(i_1 \dots i_M) \in \mathcal{S}_1^\Pi} |q'_\lambda((\mathcal{A}^*)_{i_1 \dots i_M})| \leq \max |q'_\lambda((\mathcal{A}^*)_{i_1 \dots i_M})| \leq \lambda, \quad (58)$$

Additionally, since $\mathcal{H}^* \in \partial \|\mathcal{A}^*\|_1$, it follows that $|\mathcal{H}_{i_1, \dots, i_M}^*| \leq 1$. Therefore, for each $(i_1, \dots, i_M) \in \mathcal{S}_1^\Pi$, we obtain

$$\begin{aligned} |(\nabla L(\mathcal{A}^*) + \nabla \mathcal{Q}_\lambda(\mathcal{A}^*) + \lambda \mathcal{H}^*)_{i_1 \dots i_M}| &\leq \max_{(i_1 \dots i_M) \in \mathcal{S}_1^\Pi} |(\nabla L(\mathcal{A}^*))_{i_1 \dots i_M}| + \max_{(i_1 \dots i_M) \in \mathcal{S}_1^\Pi} |(\nabla \mathcal{Q}_\lambda(\mathcal{A}^*))_{i_1 \dots i_M}| + \lambda \\ &\leq 3\lambda, \end{aligned} \quad (59)$$

which implies

$$\min_{\mathcal{H}^* \in \partial \|\mathcal{A}^*\|_1} \left\{ \sum_{(i_1 \dots i_M) \in \mathcal{S}_1^\Pi} |(\nabla L(\mathcal{A}^*) + \nabla \mathcal{Q}_\lambda(\mathcal{A}^*) + \lambda \mathcal{H}^*)_{i_1 \dots i_M}| \cdot \left| (\mathcal{A}^* - \hat{\mathcal{A}})_{i_1 \dots i_M} \right| \right\} \quad (60)$$

$$\leq 3\lambda \left| (\mathcal{A}^* - \hat{\mathcal{A}})_{i_1 \dots i_M} \right| \quad (61)$$

$$= 3\lambda \left\| (\mathcal{A}^* - \hat{\mathcal{A}})_{\mathcal{S}_1^\Pi} \right\|_{\text{F}} \quad (62)$$

$$\leq 3\lambda \sqrt{|\mathcal{S}_1|} \left\| (\mathcal{A}^* - \hat{\mathcal{A}})_{\mathcal{S}_1^\Pi} \right\|_{\text{F}} \quad (63)$$

$$\leq 3\lambda \sqrt{|\mathcal{S}_1|} \left\| \mathcal{A}^* - \hat{\mathcal{A}} \right\|_{\text{F}}. \quad (64)$$

Substituting the bounds from (52) to (64) into the right-hand side of (46), we obtain

$$\left\| \mathcal{A}^* - \hat{\mathcal{A}} \right\|_{\text{F}} \leq \frac{1}{\rho^- - \zeta^-} \left(\left\| (\nabla L(\mathcal{A}^*))_{\mathcal{S}_1^\Pi} \right\|_{\text{F}} + 3\lambda \sqrt{|\mathcal{S}_1^\Pi|} \right). \quad (65)$$

□

Lemma 4. For least-squares regression with sub-Gaussian noise, we assume that the columns of $\overline{\mathcal{X}}$ are normalized in such a way that $\max_{j \in \{1, \dots, d_1 d_2 \dots d_N\}} \|\overline{\mathcal{X}}_{\cdot j}\|_2 \leq \sqrt{n}$, where $\overline{\mathcal{X}} = \left(\text{vec}(\mathcal{X}^{(1)}), \dots, \text{vec}(\mathcal{X}^{(n)}) \right)^\top$. If $\lambda \asymp \sqrt{\frac{\log(d_1 d_2 \dots d_M)}{n}}$, then we have

$$\|\nabla L(\mathcal{A}^*)\|_{\text{F}} \lesssim \sqrt{\frac{|\mathcal{S}_1|}{n}}.$$

Proof. We begin by establishing an upper bound on the probability that the maximum entry of the gradient $\mathbb{P}(\|\nabla L(\mathcal{A})\|_{\max} \geq \frac{\lambda}{8})$, where $\nabla L(\mathcal{A}) = \frac{1}{n} \langle \overline{\mathcal{X}}, \overline{\mathcal{E}} \rangle$ and $\overline{\mathcal{E}} = \left(\text{vec}(\mathcal{E}^{(1)}), \dots, \text{vec}(\mathcal{E}^{(n)}) \right)^\top$.

For $\lambda \asymp \sqrt{\frac{\log(d_1 d_2 \dots d_M)}{n}}$, using the union bound, we obtain

$$\mathbb{P} \left(\|\nabla L(\mathcal{A})\|_{\max} \geq \frac{\lambda}{8} \right) \leq \mathbb{P} \left(\left\| \frac{1}{n} \langle \overline{\mathcal{X}}, \overline{\mathcal{E}} \rangle \right\|_{\max} \geq \frac{c \sqrt{\log d/n}}{8} \right) \quad (66)$$

$$\leq \sum_{j=1}^{d_1 d_2 \dots d_N} \mathbb{P} \left(\left| \frac{1}{n} \langle \overline{\mathcal{X}}, \overline{\mathcal{E}} \rangle \right|_{j \geq \frac{c \sqrt{\log d/n}}{8}} \right). \quad (67)$$

Let's define $\theta_k = |\langle \overline{\mathcal{X}}, \overline{\mathcal{E}} \rangle|_k$, where k is composite coordinate. Since $\overline{\mathcal{E}}_j$ is sub-Gaussian $(0, \eta^2)$, it follows that for any $t_0 > 0$,

$$\mathbb{E}(\exp\{t_0 \theta_k\} + \exp\{-t_0 \theta_k\}) \leq 2 \exp \left\{ \frac{1}{n^2} \|\overline{\mathcal{X}}_{\cdot k}\|^2 \eta^2 t_0^2 / 2 \right\}. \quad (68)$$

Taking $t_0 = \frac{tn^2}{\|\bar{\mathcal{X}}_{\cdot k}\|^2 \eta^2 t_0^2}$ yields that

$$\mathbb{P}(|\theta_k| \geq t) \leq 2 \exp \left\{ -\frac{n^2 t^2}{2 \|\bar{\mathcal{X}}_{\cdot k}\|^2 \eta^2} \right\}. \quad (69)$$

Further taking $t = \frac{\lambda}{8}$ results

$$\mathbb{P} \left(\|\nabla L(\mathcal{A})\|_{\max} \geq \frac{\lambda}{8} \right) \leq 2 (d_1 \times \dots \times d_N)^{-c^2/(128\eta^2)}. \quad (70)$$

Applying the Hanson-Wright inequality yields that

$$\begin{aligned} & \mathbb{P}(|\langle \mathcal{E}, \langle \mathcal{A}, \mathcal{E} \rangle \rangle - \mathbb{E} \langle \mathcal{E}, \langle \mathcal{A}, \mathcal{E} \rangle \rangle| > \mathbb{E} \langle \mathcal{E}, \langle \mathcal{A}, \mathcal{E} \rangle \rangle) \\ & \leq 2 \exp \left[-C \min \left\{ \frac{\mathbb{E} \langle \mathcal{E}, \langle \mathcal{A}, \mathcal{E} \rangle \rangle}{\eta^2 \|\mathcal{A}\|_F}, \left(\frac{\mathbb{E} \langle \mathcal{E}, \langle \mathcal{A}, \mathcal{E} \rangle \rangle}{\eta^2 \|\mathcal{A}\|_F} \right)^2 \right\} \right], \end{aligned} \quad (71)$$

where C is a universal constant.

Combining the above two inequalities, we have

$$\|\nabla L(\mathcal{A}^*)\|_F = \sqrt{\frac{\langle \mathcal{E}, \bar{\mathcal{X}} \rangle}{n}} \leq \sqrt{\frac{2\mathbb{E} \langle \mathcal{E}, \bar{\mathcal{X}} \rangle}{n}} \leq \sqrt{2\rho^+ \eta} \sqrt{\frac{|\mathcal{S}_1|}{n}}. \quad (72)$$

□

Building upon Lemma 1 through Lemma 4, we derive Corollary 1.

D Proof of the Corollary 2

We begin by demonstrating that the fiber-wise sparsity penalty can be reformulated as the sum of the ℓ_1 penalty and a concave part. Specifically, we have:

$$R_\lambda(\mathcal{A}) = \sum_{l=1}^{\prod_{j \neq k} d_j} p_\lambda \left(\left\| [\mathcal{A}_{(k)}]_{\cdot, l} \right\|_2 \right) = \sum_{l=1}^{\prod_{j \neq k} d_j} \lambda \left\| [\mathcal{A}_{(k)}]_{\cdot, l} \right\|_2 + Q_\lambda \left(\left\| [\mathcal{A}_{(k)}]_{\cdot, l} \right\|_2 \right),$$

where $Q_\lambda \left(\left\| [\mathcal{A}_{(k)}]_{\cdot, l} \right\|_2 \right) = \sum_{l=1}^{\prod_{j \neq k} d_j} q_\lambda \left(\left\| [\mathcal{A}_{(k)}]_{\cdot, l} \right\|_2 \right)$.

Lemma 5. Under Assumptions 1 and 2, the loss function $\tilde{L}(\mathcal{A})$ satisfies the restricted strong convexity

$$\tilde{L}(\mathcal{A}') \geq \tilde{L}(\mathcal{A}) + \left\langle \nabla \tilde{L}(\mathcal{A}), \mathcal{A}' - \mathcal{A} \right\rangle + \frac{\rho^- - \zeta^-}{2} \|\mathcal{A}' - \mathcal{A}\|_F^2,$$

and the restricted strong smoothness

$$\tilde{L}(\mathcal{A}') \leq \tilde{L}(\mathcal{A}) + \left\langle \nabla \tilde{L}(\mathcal{A}), \mathcal{A}' - \mathcal{A} \right\rangle + \frac{\rho^+ - \zeta^+}{2} \|\mathcal{A}' - \mathcal{A}\|_F^2.$$

Proof. Since the proof closely mirrors that of Lemma 1, it is omitted here for brevity. □

Lemma 6. Let C be a constant. Suppose there exists an integer $\tilde{s}_2 > C|\mathcal{S}_2|$, and that \mathcal{A} satisfies $\|\mathcal{A}_{\bar{\mathcal{S}}_2}\|_0 \leq \tilde{s}_2$, $\omega(\mathcal{A}) \leq \frac{\lambda}{2}$, where $\omega(\mathcal{A}) = \min_{\mathcal{H} \in \partial \left\| [\mathcal{A}_{(k)}]_{\cdot, l} \right\|_2} \left\{ \left\| \nabla \tilde{L}(\mathcal{A}) + \lambda \mathcal{H} \right\|_{\max} \right\}$, and $\|\nabla L(\mathcal{A}^*)\|_{\max} \leq \lambda/8$. Under Assumptions 1 and 2, \mathcal{A} satisfies

$$\|\mathcal{A} - \mathcal{A}^*\|_F \leq \frac{21/8}{\rho^- - \zeta^-} \lambda \sqrt{|\mathcal{S}_2|}.$$

Proof. We omit the proof here for brevity, as it closely mirrors that of Lemma 2. \square

Lemma 7. Consider the regularization parameter λ and assume that the derivative of the non-convex penalty satisfies $p'_\lambda \left(\left\| [\mathcal{A}_{(k)}]_{\cdot, l} \right\|_2 \right) = 0$ whenever $\left\| [\mathcal{A}_{(k)}]_{\cdot, l} \right\|_2 \geq \nu$ for some $\nu > 0$. Let $\mathcal{S}_2^I \cup \mathcal{S}_2^{II} = \mathcal{S}_2$. For indices $(i_1, \dots, i_M) \in \mathcal{S}_2^I \subseteq \mathcal{S}_2$, we assume $\min_l \left[\mathcal{A}_{(k)}^* \right]_{\cdot, l} \geq \nu$, and for indices $(i_1, \dots, i_M) \in \mathcal{S}_2^{II} \subseteq \mathcal{S}_2$, we assume $\min_l \left[\mathcal{A}_{(k)}^* \right]_{\cdot, l} \leq \nu$. Under Assumptions 1~4, we derive the following bound:

$$\left\| \hat{\mathcal{A}} - \mathcal{A}^* \right\|_F \leq \frac{1}{\rho^- - \zeta^-} \left\| (\nabla L(\mathcal{A}^*))_{\mathcal{S}_2^I} \right\|_F + \frac{3}{\rho^- - \zeta^-} \lambda \sqrt{|\mathcal{S}_2^{II}|}.$$

Proof. For brevity, we omit the proof here, as it closely resembles that of Lemma 3. \square

Lemma 8. For least-squares regression with sub-Gaussian noise, we assume that the columns of $\bar{\mathcal{X}}$ are normalized in such a way that $\max_{j \in \{1, \dots, d_1 d_2 \dots d_N\}} \left\| \bar{\mathcal{X}}_{\cdot j} \right\|_2 \leq \sqrt{n}$, where $\bar{\mathcal{X}} = \left(\text{vec}(\mathcal{X}^{(1)}), \dots, \text{vec}(\mathcal{X}^{(n)}) \right)^\top$. If $\lambda \asymp \sqrt{\frac{\log(d_k)}{n}}$, then we have

$$\left\| \nabla L(\mathcal{A}^*) \right\|_F \lesssim \sqrt{\frac{|\mathcal{S}_2|}{n}}.$$

Proof. For brevity, the proof is omitted here as it closely follows the methodology established in Lemma 4. \square

E Proof of the Corollary 3

We begin by demonstrating that the fiber-wise sparsity penalty can be reformulated as the sum of the ℓ_1 penalty and a concave part. Specifically, we have:

$$R_\lambda(\mathcal{A}) = \sum_{i=1}^{\prod_{s \neq j, k} d_s} p_\lambda \left(\left\| [\mathcal{A}_{(j, k)}]_{\cdot, l} \right\|_F \right) = \sum_{i=1}^{\prod_{s \neq j, k} d_s} \lambda \left\| [\mathcal{A}_{(j, k)}]_{\cdot, l} \right\|_F + Q_\lambda \left(\left\| [\mathcal{A}_{(j, k)}]_{\cdot, l} \right\|_F \right),$$

where $Q_\lambda \left(\left\| [\mathcal{A}_{(j, k)}]_{\cdot, l} \right\|_F \right) = \sum_{i=1}^{\prod_{s \neq j, k} d_s} q_\lambda \left(\left\| [\mathcal{A}_{(j, k)}]_{\cdot, l} \right\|_F \right)$.

Lemma 9. Under Assumptions 1 and 2, the loss function $\tilde{L}(\mathcal{A})$ satisfies the restricted strong convexity

$$\tilde{L}(\mathcal{A}') \geq \tilde{L}(\mathcal{A}) + \left\langle \nabla \tilde{L}(\mathcal{A}), \mathcal{A}' - \mathcal{A} \right\rangle + \frac{\rho^- - \zeta^-}{2} \left\| \mathcal{A}' - \mathcal{A} \right\|_F^2,$$

and the restricted strong smoothness

$$\tilde{L}(\mathcal{A}') \leq \tilde{L}(\mathcal{A}) + \left\langle \nabla \tilde{L}(\mathcal{A}), \mathcal{A}' - \mathcal{A} \right\rangle + \frac{\rho^+ - \zeta^+}{2} \left\| \mathcal{A}' - \mathcal{A} \right\|_F^2.$$

Proof. The proof can be demonstrated similarly to the proof in Lemma 1. Hence, we omit it here. \square

Lemma 10. Let C be a constant. Suppose there exists an integer $\tilde{s}_3 > C |\mathcal{S}_3|$, and that \mathcal{A} satisfies $\left\| \mathcal{A}_{\overline{\mathcal{S}_3}} \right\|_0 \leq \tilde{s}_3$, $\omega(\mathcal{A}) \leq \frac{\lambda}{2}$, where $\omega(\mathcal{A}) = \min_{\mathcal{H} \in \partial \left\| [\mathcal{A}_{(j, k)}]_{\cdot, l} \right\|_F} \left\{ \left\| \nabla \tilde{L}(\mathcal{A}) + \lambda \mathcal{H} \right\|_{\max} \right\}$, and $\left\| \nabla L(\mathcal{A}^*) \right\|_{\max} \leq \lambda/8$. Under Assumptions 1 and 2, \mathcal{A} satisfies

$$\left\| \mathcal{A} - \mathcal{A}^* \right\|_F \leq \frac{21/8}{\rho^- - \zeta^-} \lambda \sqrt{|\mathcal{S}_3|}.$$

Proof. For the sake of brevity, we omit the proof here, as it closely follows that of Lemma 2. \square

Lemma 11. Consider the regularization parameter λ and assume that the derivative of the non-convex penalty satisfies $p'_\lambda \left(\left\| [\mathcal{A}_{(j,k)}]_{\cdot, \cdot, l} \right\|_F \right) = 0$ whenever $\left\| [\mathcal{A}_{(j,k)}]_{\cdot, \cdot, l} \right\|_F \geq \nu$ for some $\nu > 0$. Let $\mathcal{S}_3^I \cup \mathcal{S}_3^{II} = \mathcal{S}_3$. For indices $(i_1, \dots, i_M) \in \mathcal{S}_3^I \subseteq \mathcal{S}_3$, we assume $\min_l \left\| [\mathcal{A}_{(j,k)}]_{\cdot, \cdot, l} \right\|_F \geq \nu$, and for indices $(i_1, \dots, i_M) \in \mathcal{S}_3^{II} \subseteq \mathcal{S}_3$, we assume $\min_l \left\| [\mathcal{A}_{(j,k)}]_{\cdot, \cdot, l} \right\|_F \leq \nu$. Under Assumptions 1~4, we derive the following bound:

$$\left\| \hat{\mathcal{A}} - \mathcal{A}^* \right\|_F \leq \frac{1}{\rho^- - \zeta^-} \left\| (\nabla L(\mathcal{A}^*))_{\mathcal{S}_3^I} \right\|_F + \frac{3}{\rho^- - \zeta^-} \lambda \sqrt{|\mathcal{S}_3^{II}|}.$$

Proof. For brevity, we omit the proof here, as it closely resembles that of Lemma 3. \square

Lemma 12. For least-squares regression with sub-Gaussian noise, we assume that the columns of $\bar{\mathcal{X}}$ are normalized in such a way that $\max_{j \in \{1, \dots, d_1 d_2 \dots d_N\}} \left\| \bar{\mathcal{X}}_{\cdot j} \right\|_2 \leq \sqrt{n}$, where $\bar{\mathcal{X}} = \left(\text{vec}(\mathcal{X}^{(1)}), \dots, \text{vec}(\mathcal{X}^{(n)}) \right)^\top$. If $\lambda \asymp \sqrt{\frac{\log(d_j d_k)}{n}}$, then we have

$$\left\| \nabla L(\mathcal{A}^*) \right\|_F \lesssim \sqrt{\frac{|\mathcal{S}_3|}{n}}.$$

Proof. For brevity, the proof is omitted here as it closely follows the methodology established in Lemma 4. \square

F Proof of the Corollary 4

The proposed mode-wise low-rankness penalty can be reformulated as the sum of a scaled norm and a concave function. Specifically, we have

$$R_\lambda(\mathcal{A}) = \sum_{i=1}^{\min\{I_k, \prod_{j \neq k} d_j\}} p_\lambda(\sigma_i(\mathcal{A}_{(k)})) = \lambda \left\| \mathcal{A}_{(k)} \right\|_* + Q_\lambda(\mathcal{A}_{(k)}),$$

where $\sigma_i(\mathcal{A}_{(k)})$ denotes the i -th singular value of the mode- (k) unfolding $\mathcal{A}_{(k)}$. For the estimation problem, we define

$$\tilde{L}(\mathcal{A}) = L(\mathcal{A}) + Q_\lambda(\mathcal{A}_{(k)}),$$

where $Q_\lambda(\mathcal{A}_{(k)}) = \sum_{i=1}^{\min\{I_k, \prod_{j \neq k} d_j\}} q_\lambda(\sigma_i(\mathcal{A}_{(k)}))$, and $\min\{I_k, \prod_{j \neq k} d_j\}$ is the number of the singular values.

Based on the restrict strongly convexity of $L(\cdot)$ in Assumption 1 and the parameter for regularity condition in Assumption 3, if $\rho^- > \zeta^-$, we have the restrict strongly convexity of $\tilde{L}(\cdot)$.

Besides, for the RSC and RSS assumption, we define the following cone of directions

$$\mathcal{C} = \{\mathcal{B} \in \mathbb{R}^{d_1 \dots d_N} \mid \|\Pi_{\mathcal{F}^\perp}(\mathcal{B})\|_* \leq 5 \|\Pi_{\mathcal{F}}(\mathcal{B})\|_*\}$$

Lemma 13. Under Assumption 1, if $\mathcal{B} \in \mathcal{C}$, we have

$$\tilde{L}(\mathcal{A} + \mathcal{B}) \geq \tilde{L}(\mathcal{A}) + \left\langle \nabla \tilde{L}(\mathcal{A}), \mathcal{B} \right\rangle + \frac{\rho^- - \zeta^-}{2} \|\mathcal{B}\|_F^2.$$

Proof. Based on Assumption 1, we have

$$L(\mathcal{A} + \mathcal{B}) \leq L(\mathcal{A}) + \left\langle \nabla L(\mathcal{A}), \mathcal{B} \right\rangle + \frac{\rho^-}{2} \|\mathcal{B}\|_F^2. \quad (73)$$

Moreover, considering the singular values of the unfolded matrices $\mathcal{A}_{(k)}$ and $\mathcal{B}_{(k)}$, we obtain

$$-\zeta^- \leq \frac{q'_\lambda(\sigma_i(\mathcal{A}_{(k)})) - q'_\lambda(\sigma_I([\mathcal{A} + \mathcal{B}]_{(k)}))}{\left[\sigma_i(\mathcal{A}_{(k)}) - \sigma_i([\mathcal{A} + \mathcal{B}]_{(k)}) \right]}, \quad (74)$$

which is similar to the proof for Lemma 1. This inequality leads to

$$\left\langle \left(-\nabla Q_\lambda(\mathcal{A}_{(k)}) \right) - \left(-\nabla Q_\lambda([\mathcal{A} + \mathcal{B}]_{(k)}) \right), \mathcal{B}_{(k)} \right\rangle \leq \zeta^- \|\mathcal{B}_{(k)}\|_F. \quad (75)$$

This inequality characterizes the strong smoothness of $-Q(\cdot)$, which is equivalent to

$$Q_\lambda([\mathcal{A} + \mathcal{B}]_{(k)}) = Q_\lambda(\mathcal{A}_{(k)} + \mathcal{B}_{(k)}) \geq Q_\lambda(\mathcal{A}_{(k)}) + \langle \nabla Q_\lambda(\mathcal{A}_{(k)}), \mathcal{B}_{(k)} \rangle - \frac{\zeta^-}{2} \|\mathcal{B}_{(k)}\|_F^2. \quad (76)$$

Noting that the Frobenius norm satisfies $\|\mathcal{B}_{(k)}\|_F^2 = \|\mathcal{B}\|_F^2$. Let $\mathcal{A}' = \mathcal{A} + \mathcal{B}$, adding (73) and (76), we have

$$\tilde{L}(\mathcal{A}') = L(\mathcal{A}') + Q_\lambda(\mathcal{A}'_{(k)}) \quad (77)$$

$$\geq L(\mathcal{A}) + Q_\lambda(\mathcal{A}_{(k)}) + \langle \nabla L(\mathcal{A}), \mathcal{B} \rangle + \langle \nabla Q_\lambda(\mathcal{A}_{(k)}), \mathcal{B}_{(k)} \rangle + \frac{\rho^- - \zeta^-}{2} \|\mathcal{B}\|_F^2 \quad (78)$$

$$= \tilde{L}(\mathcal{A}) + \langle \nabla L(\mathcal{A}), \mathcal{B} \rangle + \langle \nabla Q_\lambda(\mathcal{A}_{(k)}), \mathcal{B}_{(k)} \rangle + \frac{\rho^- - \zeta^-}{2} \|\mathcal{B}\|_F^2. \quad (79)$$

□

Lemma 14. Under Assumption 1, if $\rho^- > \zeta^-$ and the regularization parameter $\lambda \geq \frac{\|\mathfrak{X}^*(\mathcal{E})_{(k)}\|_{\text{sp}}}{2n}$, we have

$$\left\| \Pi_{\mathcal{F}^\perp}(\hat{\mathcal{A}}_{(k)} - \mathcal{A}_{(k)}^*) \right\|_* \leq 5 \left\| \Pi_{\mathcal{F}}(\hat{\mathcal{A}}_{(k)} - \mathcal{A}_{(k)}^*) \right\|_*.$$

Proof. By Lemma 13, we have

$$\tilde{L}(\hat{\mathcal{A}}) - \tilde{L}(\mathcal{A}^*) \geq \langle \nabla L(\mathcal{A}^*), \hat{\mathcal{A}} - \mathcal{A}^* \rangle + \left\langle \nabla Q(\mathcal{A}_{(k)}^*), [\hat{\mathcal{A}} - \mathcal{A}^*]_{(k)} \right\rangle. \quad (80)$$

We proceed to bound the right-hand side of inequality (80). By decomposing the inner products using the projections onto two orthogonal subspaces, we have

$$\begin{aligned} & \left\langle \nabla L(\mathcal{A}^*), \hat{\mathcal{A}} - \mathcal{A}^* \right\rangle + \left\langle \nabla Q(\mathcal{A}_{(k)}^*), [\hat{\mathcal{A}} - \mathcal{A}^*]_{(k)} \right\rangle \\ &= \left\langle [\nabla L(\mathcal{A}^*)]_{(k)} + \nabla Q(\mathcal{A}_{(k)}^*), \Pi_{\mathcal{F}}([\hat{\mathcal{A}} - \mathcal{A}^*]_{(k)}) \right\rangle + \left\langle [\nabla L(\mathcal{A}^*)]_{(k)} + \nabla Q(\mathcal{A}_{(k)}^*), \Pi_{\mathcal{F}^\perp}([\hat{\mathcal{A}} - \mathcal{A}^*]_{(k)}) \right\rangle \\ &\geq - \left\| \Pi_{\mathcal{F}}([\nabla L(\mathcal{A}^*)]_{(k)} + \nabla Q(\mathcal{A}_{(k)}^*)) \right\|_{\text{sp}} \left\| \Pi_{\mathcal{F}}([\hat{\mathcal{A}} - \mathcal{A}^*]_{(k)}) \right\|_* \\ &\quad - \left\| \Pi_{\mathcal{F}^\perp}([\nabla L(\mathcal{A}^*)]_{(k)} + \nabla Q(\mathcal{A}_{(k)}^*)) \right\|_{\text{sp}} \left\| \Pi_{\mathcal{F}^\perp}([\hat{\mathcal{A}} - \mathcal{A}^*]_{(k)}) \right\|_*. \end{aligned} \quad (81)$$

$$- \left\| \Pi_{\mathcal{F}^\perp}([\nabla L(\mathcal{A}^*)]_{(k)} + \nabla Q(\mathcal{A}_{(k)}^*)) \right\|_{\text{sp}} \left\| \Pi_{\mathcal{F}^\perp}([\hat{\mathcal{A}} - \mathcal{A}^*]_{(k)}) \right\|_*. \quad (82)$$

For (81), due to $\lambda \geq \frac{1}{2n} \left\| [\mathfrak{X}^*(\mathcal{E})]_{(k)} \right\|_{\text{sp}}$, we see that $\left\| [\nabla L(\mathcal{A}^*)]_{(k)} \right\|_{\text{sp}} \leq \lambda/2$. According to Assumption 3, we have

$$\left\| \Pi_{\mathcal{F}}([\nabla L(\mathcal{A}^*)]_{(k)} + \nabla Q(\mathcal{A}_{(k)}^*)) \right\|_{\text{sp}} \leq \frac{3}{2} \lambda. \quad (83)$$

For (82), since $\Pi_{\mathcal{F}^\perp}(\mathcal{A}_{(k)}^*) = 0$, we obtain

$$\left\| \Pi_{\mathcal{F}^\perp}([\nabla L(\mathcal{A}^*)]_{(k)} + \nabla Q(\mathcal{A}_{(k)}^*)) \right\|_{\text{sp}} \leq \frac{1}{2} \lambda. \quad (84)$$

Combine (83) and (84), we have

$$\left\langle \nabla L(\mathcal{A}^*), \hat{\mathcal{A}} - \mathcal{A}^* \right\rangle + \left\langle \nabla Q(\mathcal{A}_{(k)}^*), [\hat{\mathcal{A}} - \mathcal{A}^*]_{(k)} \right\rangle \quad (85)$$

$$\geq -\frac{3}{2} \lambda \left\| \Pi_{\mathcal{F}}([\hat{\mathcal{A}} - \mathcal{A}^*]_{(k)}) \right\|_* - \frac{1}{2} \lambda \left\| \Pi_{\mathcal{F}^\perp}([\hat{\mathcal{A}} - \mathcal{A}^*]_{(k)}) \right\|_* \quad (86)$$

Moreover, noting that $\lambda \|\hat{\mathcal{A}}\|_* - \lambda \|\mathcal{A}^*\|_* \geq -\lambda \left\| \Pi_{\mathcal{F}} \left([\hat{\mathcal{A}} - \mathcal{A}^*]_{(k)} \right) \right\| + \lambda \left\| \Pi_{\mathcal{F}^\perp} \left([\hat{\mathcal{A}} - \mathcal{A}^*]_{(k)} \right) \right\|_*$, and combining with (80), we obtain

$$\left\langle \nabla L(\mathcal{A}^*) + \nabla Q(\mathcal{A}_{(k)}^*), \hat{\mathcal{A}} - \mathcal{A}^* \right\rangle + \lambda \|\hat{\mathcal{A}}\|_* - \lambda \|\mathcal{A}^*\|_* \quad (87)$$

$$\geq -\frac{5}{2}\lambda \left\| \Pi_{\mathcal{F}} \left([\hat{\mathcal{A}} - \mathcal{A}^*]_{(k)} \right) \right\|_* + \frac{1}{2}\lambda \left\| \Pi_{\mathcal{F}^\perp} \left([\hat{\mathcal{A}} - \mathcal{A}^*]_{(k)} \right) \right\|_*. \quad (88)$$

Since $\hat{\mathcal{A}}$ is the global minimizer of the general estimator (2) and given that $\rho^- > \zeta^-$, it follows that

$$\tilde{L}(\hat{\mathcal{A}}) + \lambda \|\hat{\mathcal{A}}\|_* - \tilde{L}(\mathcal{A}^*) - \lambda \|\mathcal{A}^*\|_* \leq 0. \quad (89)$$

Substituting (80) and (89) into (88), we obtain

$$\frac{1}{2}\lambda \left\| \Pi_{\mathcal{F}^\perp} \left([\hat{\mathcal{A}} - \mathcal{A}^*]_{(k)} \right) \right\|_* \leq \frac{5}{2}\lambda \left\| \Pi_{\mathcal{F}} \left([\hat{\mathcal{A}} - \mathcal{A}^*]_{(k)} \right) \right\|_*. \quad (90)$$

Since $\lambda > 0$, we obtain

$$\left\| \Pi_{\mathcal{F}^\perp} \left([\hat{\mathcal{A}} - \mathcal{A}^*]_{(k)} \right) \right\|_* \leq 5 \left\| \Pi_{\mathcal{F}} \left([\hat{\mathcal{A}} - \mathcal{A}^*]_{(k)} \right) \right\|_*. \quad (91)$$

□

Lemma 15. *Considering the mode-wise low-rankness regularizer, under Assumptions 1~4, for the estimated parameter tensor $\hat{\mathcal{A}}$ and the true parameter tensor \mathcal{A}^* , we have*

$$\|\hat{\mathcal{A}} - \mathcal{A}^*\|_{\text{F}} \leq \frac{1}{(\rho^- - \zeta^-)} \left[\sqrt{|\mathcal{S}_4^{\text{I}}|} \left\| \Pi_{\mathcal{F}} \left([\nabla L(\mathcal{A}^*)]_{(k)} \right) \right\|_{\text{sp}} + 3\lambda \sqrt{|\mathcal{S}_4^{\text{II}}|} \right].$$

where \mathcal{S}_4^{I} and $\mathcal{S}_4^{\text{II}}$ are subsets of the support set of \mathcal{S}_4 . The set \mathcal{S}_4^{I} include all indices $i \in \mathcal{S}_4^{\text{I}}$ which satisfy $\sigma_i(\mathcal{A}_{(k)}^*) \geq \nu$, and $\mathcal{S}_4^{\text{II}}$ includes all indices with $\sigma_i(\mathcal{A}_{(k)}^*) < \nu$.

Proof. Since $\|\cdot\|_*$ is convex, we have

$$\lambda \|\hat{\mathcal{A}}_{(k)}\|_* \geq \lambda \|\mathcal{A}_{(k)}^*\|_* + \lambda \left\langle [\hat{\mathcal{A}} - \mathcal{A}^*]_{(k)}, \mathcal{H}^* \right\rangle. \quad (92)$$

$$\lambda \|\mathcal{A}_{(k)}^*\|_* \geq \lambda \|\hat{\mathcal{A}}_{(k)}\|_* + \lambda \left\langle [\mathcal{A}^* - \hat{\mathcal{A}}]_{(k)}, \hat{\mathcal{H}} \right\rangle. \quad (93)$$

where $\mathcal{H}^* \in \partial \|\mathcal{A}_{(k)}^*\|_*$ and $\hat{\mathcal{H}} \in \partial \|\hat{\mathcal{A}}_{(k)}\|_*$. From (92) and (93), we have

$$\lambda \|\hat{\mathcal{A}}_{(k)}\|_* + \lambda \|\mathcal{A}_{(k)}^*\|_* \geq \lambda \|\mathcal{A}_{(k)}^*\|_* + \lambda \|\hat{\mathcal{A}}_{(k)}\|_* + \lambda \left\langle [\hat{\mathcal{A}} - \mathcal{A}^*]_{(k)}, \mathcal{H}^* \right\rangle + \lambda \left\langle [\mathcal{A}^* - \hat{\mathcal{A}}]_{(k)}, \hat{\mathcal{H}} \right\rangle. \quad (94)$$

This equals to

$$0 \geq \left(\lambda \left\langle [\hat{\mathcal{A}} - \mathcal{A}^*]_{(k)}, \mathcal{H}^* \right\rangle + \lambda \left\langle [\mathcal{A}^* - \hat{\mathcal{A}}]_{(k)}, \hat{\mathcal{H}} \right\rangle \right). \quad (95)$$

Moreover, according to Lemma 13, we have

$$\tilde{L}(\hat{\mathcal{A}}) \geq \tilde{L}(\mathcal{A}^*) + \left\langle \nabla L(\mathcal{A}^*), \hat{\mathcal{A}} - \mathcal{A}^* \right\rangle + \left\langle \nabla Q_\lambda(\mathcal{A}_{(k)}^*), [\hat{\mathcal{A}} - \mathcal{A}^*]_{(k)} \right\rangle + \frac{\rho^- - \zeta^-}{2} \|\hat{\mathcal{A}} - \mathcal{A}^*\|_{\text{F}}^2. \quad (96)$$

$$\tilde{L}(\mathcal{A}^*) \geq \tilde{L}(\hat{\mathcal{A}}) + \langle \nabla L(\hat{\mathcal{A}}), \mathcal{A}^* - \hat{\mathcal{A}} \rangle + \left\langle \nabla Q_\lambda(\hat{\mathcal{A}}_{(k)}), [\mathcal{A}^* - \hat{\mathcal{A}}]_{(k)} \right\rangle + \frac{\rho^- - \zeta^-}{2} \|\mathcal{A}^* - \hat{\mathcal{A}}\|_F^2. \quad (97)$$

Summing (95), (96), and (97), we have

$$\begin{aligned} 0 &\geq \langle \nabla L(\mathcal{A}^*), \hat{\mathcal{A}} - \mathcal{A}^* \rangle + \langle \nabla L(\hat{\mathcal{A}}), \mathcal{A}^* - \hat{\mathcal{A}} \rangle + (\rho^- - \zeta^-) \|\hat{\mathcal{A}} - \mathcal{A}^*\|_F^2 \\ &\quad + \left(\left\langle \nabla Q_\lambda(\mathcal{A}_{(k)}^*) + \lambda \mathcal{H}^*, [\hat{\mathcal{A}} - \mathcal{A}^*]_{(k)} \right\rangle + \left\langle \nabla Q_\lambda(\hat{\mathcal{A}}_{(k)}) + \lambda \hat{\mathcal{H}}, [\mathcal{A}^* - \hat{\mathcal{A}}]_{(k)} \right\rangle \right). \end{aligned} \quad (98)$$

Since $\hat{\mathcal{A}}$ is the solution to the estimation problem and $\hat{\mathcal{A}}$ satisfies the optimality condition, for any $\mathcal{A}' \in \mathbb{R}^{d_1 \times \dots \times d_N}$, it holds that

$$\max_{\mathcal{A}'} \left\{ \langle \nabla L(\hat{\mathcal{A}}), \hat{\mathcal{A}} - \mathcal{A}' \rangle + \left\langle \nabla Q_\lambda(\hat{\mathcal{A}}_{(k)}) + \lambda \hat{\mathcal{H}}, [\hat{\mathcal{A}} - \mathcal{A}']_{(k)} \right\rangle \right\} \leq 0. \quad (99)$$

which implies

$$\langle \nabla L(\hat{\mathcal{A}}), \mathcal{A}^* - \hat{\mathcal{A}} \rangle + \left\langle \nabla Q_\lambda(\hat{\mathcal{A}}_{(k)}) + \lambda \hat{\mathcal{H}}, [\mathcal{A}^* - \hat{\mathcal{A}}]_{(k)} \right\rangle \geq 0. \quad (100)$$

Since $\langle \nabla L(\mathcal{A}^*), \mathcal{A}^* - \hat{\mathcal{A}} \rangle = \left\langle [\nabla L(\mathcal{A}^*)]_{(k)}, [\mathcal{A}^* - \hat{\mathcal{A}}]_{(k)} \right\rangle$, we have

$$\begin{aligned} (\rho^- - \zeta^-) \|\hat{\mathcal{A}} - \mathcal{A}^*\|_F^2 &\leq \left[\left\langle [\nabla L(\mathcal{A}^*)]_{(k)}, [\mathcal{A}^* - \hat{\mathcal{A}}]_{(k)} \right\rangle + \left\langle \nabla Q_\lambda(\mathcal{A}_{(k)}^*) + \lambda \mathcal{H}^*, [\mathcal{A}^* - \hat{\mathcal{A}}]_{(k)} \right\rangle \right] \\ &\leq \left\langle \Pi_{\mathcal{F}^\perp} \left[[\nabla L(\mathcal{A}^*)]_{(k)} + \nabla Q_\lambda(\mathcal{A}_{(k)}^*) + \lambda \mathcal{H}^* \right], [\mathcal{A}^* - \hat{\mathcal{A}}]_{(k)} \right\rangle \\ &\quad + \left\langle \Pi_{\mathcal{F}} \left[[\nabla L(\mathcal{A}^*)]_{(k)} + \nabla Q_\lambda(\mathcal{A}_{(k)}^*) + \lambda \mathcal{H}^* \right], [\mathcal{A}^* - \hat{\mathcal{A}}]_{(k)} \right\rangle. \end{aligned} \quad (101)$$

We have defined $\sigma_i(\mathcal{A}_{(k)}^*)$ as the i -th singular value of tensor \mathcal{A}^* . With regard to the magnitudes of the singular values of \mathcal{A}^* , we can decompose (101) into three parts:

- $i \in \mathcal{S}_4^I$ that $\sigma_i(\mathcal{A}_{(k)}^*) \geq \nu$,
- $i \in \mathcal{S}_4^{II}$ that $\nu \geq \sigma_i(\mathcal{A}_{(k)}^*) > 0$,
- $i \in \mathcal{S}_4^c$ that $\sigma_i(\mathcal{A}_{(k)}^*) = 0$.

(i) For $i \in \mathcal{S}_4^I$ that $\sigma_i(\mathcal{A}_{(k)}^*) \geq \nu$, define a subspace of \mathcal{F} associated with \mathcal{S}_4^I as follows

$$\mathcal{F}_{\mathcal{S}_4^I}(\mathbf{U}^*, \mathbf{V}^*) := \{ \mathbf{W} \mid \text{row}(\mathbf{W}) \subset \mathbf{V}_I^*, \text{col}(\mathbf{W}) \subset \mathbf{U}_I^* \}, \quad (102)$$

where \mathbf{V}_I^* and \mathbf{U}_I^* is the matrix with the i -th row of \mathbf{V}_I^* and \mathbf{U}_I^* with $i \in \mathcal{S}_4^I$.

Recall that $R_\lambda(\mathcal{A}_{(k)}^*) = \lambda \|\mathcal{A}_{(k)}^*\|_* + Q_\lambda(\mathcal{A}_{(k)}^*)$, we have

$$\nabla R_\lambda(\mathcal{A}_{(k)}^*) = \nabla Q_\lambda(\mathcal{A}_{(k)}^*) + \lambda_k (\mathbf{U}_I^* \mathbf{V}_I^{*\top} + \mathbf{Z}_I^*), \quad (103)$$

where $\mathbf{Z}_I^* = -\lambda^{-1} \Pi_{\mathcal{F}_{S_4^I}} \left([\nabla L(\mathcal{A}^*)]_{(k)} \right)$. Since $\|\mathbf{Z}_I^*\| \leq 1$ and $\mathbf{Z}_I^* \in \mathcal{F}_{S_4^I}$, which satisfies the condition of \mathbf{W}^* to be sub-gradient of $\|\mathcal{A}_{(k)}^*\|$. Projecting $R_\lambda \left(\mathcal{A}_{(k)}^* \right)$ into the subspace $\mathcal{F}_{S_4^I}$, we have

$$\Pi_{\mathcal{F}_{S_4^I}} \left(\nabla R_\lambda \left(\mathcal{A}_{(k)}^* \right) \right) = \Pi_{\mathcal{F}_{S_4^I}} \left(\nabla Q_\lambda \left(\mathcal{A}_{(k)}^* \right) + \lambda \mathbf{U}_I^* \mathbf{V}_I^{*\top} + \lambda \mathbf{Z}_I^* \right) \quad (104)$$

$$= \mathbf{U}_I^* q'_\lambda \left(\Sigma_I^* \right) \mathbf{V}_I^{*\top} + \lambda \mathbf{U}_I^* \mathbf{V}_I^{*\top} \quad (105)$$

$$= \mathbf{U}_I^* [q'_\lambda \left(\Sigma_I^* \right) + \lambda \mathbf{I}_I] \mathbf{V}_I^{*\top}, \quad (106)$$

where \mathbf{I}_I is an identity matrix with the size $\min\{d_k, \Pi_{j \neq k} d_j\}$ and $(q'_\lambda \left(\Sigma_I^* \right) + \lambda \mathbf{I}_I)$ is a diagonal matrix that for $i \notin S_4^I$, the i -th element on the diagonal equals 0, i.e. $[q'_\lambda \left(\Sigma_I^* \right) + \lambda \mathbf{I}_I]_{ii} = 0$, and for all $i \in S_4^I$, we have

$$[q'_\lambda \left(\Sigma_I^* \right) + \lambda \mathbf{I}_I]_{ii} = q'_\lambda \left(\sigma_i \left(\mathcal{A}_{(k)}^* \right) \right) + \lambda = p'_\lambda \left(\sigma_i \left(\mathcal{A}_{(k)}^* \right) \right) = 0. \quad (107)$$

The last equality is derived from fact that $i \in S_4^I$ satisfies Assumption 3, $p'_\lambda(t) = 0$. Therefore, we have $q'_\lambda \left(\Sigma_I^* \right) + \lambda \mathbf{I}_I = 0$, which indicates that $\Pi_{\mathcal{F}_{S_4^I}} \left(\nabla R_\lambda \left(\mathcal{A}_{(k)}^* \right) \right) = 0$. For $\mathcal{H}^* = \mathbf{U}_I^* \mathbf{V}_I^{*\top} + \mathbf{Z}_I^* \in \partial \|\mathcal{A}_{(k)}^*\|_*$, we have

$$\left\langle \Pi_{\mathcal{F}_{S_4^I}} \left[[\nabla L(\mathcal{A}^*)]_{(k)} + \lambda \mathcal{H}^* + \nabla Q_\lambda \left(\mathcal{A}_{(k)}^* \right) \right], [\mathcal{A}^* - \hat{\mathcal{A}}]_{(k)} \right\rangle \quad (108)$$

$$= \left\langle \Pi_{\mathcal{F}_{S_4^I}} \left[[\nabla L(\mathcal{A}^*)]_{(k)} + \nabla R_\lambda \left(\mathcal{A}_{(k)}^* \right) \right], [\mathcal{A}^* - \hat{\mathcal{A}}]_{(k)} \right\rangle \quad (109)$$

$$= \left\langle \Pi_{\mathcal{F}_{S_4^I}} \left([\nabla L(\mathcal{A}^*)]_{(k)} \right), \Pi_{\mathcal{F}_{S_4^I}} \left([\mathcal{A}^* - \hat{\mathcal{A}}]_{(k)} \right) \right\rangle \quad (110)$$

$$\leq \left\| \Pi_{\mathcal{F}_{S_4^I}} \left([\nabla L(\mathcal{A}^*)]_{(k)} \right) \right\|_{\text{sp}} \cdot \left\| \Pi_{\mathcal{F}_{S_4^I}} \left([\mathcal{A}^* - \hat{\mathcal{A}}]_{(k)} \right) \right\|_* \quad (111)$$

where the last inequality is derived from the Hölder inequality. For $\left\| \Pi_{\mathcal{F}_{S_4^I}} \left([\mathcal{A}^* - \hat{\mathcal{A}}]_{(k)} \right) \right\|_*$, from the properties of projection on to the subspace $\mathcal{F}_{S_4^I}$, we have

$$\left\| \Pi_{\mathcal{F}_{S_4^I}} \left([\mathcal{A}^* - \hat{\mathcal{A}}]_{(k)} \right) \right\|_* \leq \sqrt{|S_4^I|} \left\| \Pi_{\mathcal{F}_{S_4^I}} \left([\mathcal{A}^* - \hat{\mathcal{A}}]_{(k)} \right) \right\|_F \quad (112)$$

$$\leq \sqrt{|S_4^I|} \left\| [\mathcal{A}^* - \hat{\mathcal{A}}]_{(k)} \right\|_F = \sqrt{|S_4^I|} \left\| \mathcal{A}^* - \hat{\mathcal{A}} \right\|_F. \quad (113)$$

We obtain the second inequality from that the rank of the matrix $\Pi_{\mathcal{F}_{S_4^I}} \left([\mathcal{A}^* - \hat{\mathcal{A}}]_{(k)} \right) \leq |S_4^I|$.

Thus, we have

$$\begin{aligned} & \left\langle \Pi_{\mathcal{F}_{S_4^I}} \left[[\nabla L(\mathcal{A}^*)]_{(k)} + \nabla Q_\lambda \left(\mathcal{A}_{(k)}^* \right) + \lambda \mathcal{H}^* \right], [\mathcal{A}^* - \hat{\mathcal{A}}]_{(k)} \right\rangle \\ & \leq \sqrt{|S_4^I|} \left\| \Pi_{\mathcal{F}_{S_4^I}} \left([\mathcal{A}^* - \hat{\mathcal{A}}]_{(k)} \right) \right\|_{\text{sp}} \cdot \left\| \mathcal{A}^* - \hat{\mathcal{A}} \right\|_F. \end{aligned} \quad (114)$$

(ii) For $i \in S_4^{\text{II}}$, $\nu \geq \sigma_i \left(\mathcal{A}_{(k)}^* \right) > 0$, define a subspace of \mathcal{F} associated with S_4^{II} as follows

$$\mathcal{F}_{S_4^{\text{II}}}(\mathbf{U}^*, \mathbf{V}^*) := \{\mathbf{W} \mid \text{row}(\mathbf{W}) \subset \mathbf{V}_{\text{II}}^*, \text{col}(\mathbf{W}) \subset \mathbf{U}_{\text{II}}^*\}.$$

where \mathbf{V}_{II}^* and \mathbf{U}_{II}^* is the matrix with the i -th row of \mathbf{U}^* and \mathbf{V}^* with $i \in S_4^{\text{II}}$. Obviously, for all \mathbf{W} , the following decomposition holds

$$\Pi_{\mathcal{F}}(\mathbf{W}) = \Pi_{\mathcal{F}_{S_4^I}}(\mathbf{W}) + \Pi_{\mathcal{F}_{S_4^{\text{II}}}}(\mathbf{W}).$$

In addition, since \mathbf{U}^* , \mathbf{V}^* are unitary matrices, for subspace $\mathcal{F}_{\mathcal{S}_4^I}$ and $\mathcal{F}_{\mathcal{S}_4^{II}}$, we have the complementary subspace $\mathcal{F}_{\mathcal{S}_4^I}^\perp$, $\mathcal{F}_{\mathcal{S}_4^{II}}^\perp$, thus we have

$$\mathcal{F}_{\mathcal{S}_4^I} \subset \mathcal{F}_{\mathcal{S}_4^I}^\perp, \text{ and } \mathcal{F}_{\mathcal{S}_4^{II}} \subset \mathcal{F}_{\mathcal{S}_4^{II}}^\perp.$$

Similar to analysis in (i) on \mathcal{S}_4^I , we have

$$\Pi_{\mathcal{F}_{\mathcal{S}_4^{II}}} \left(\nabla Q_\lambda \left(\mathcal{A}_{(k)}^* \right) \right) = \mathbf{U}_\Pi^* q'_\lambda \left(\boldsymbol{\Sigma}_\Pi^* \right) \mathbf{V}_\Pi^{*\top}. \quad (115)$$

where $q'_\lambda \left(\boldsymbol{\Sigma}_\Pi^* \right)$ is a diagonal matrix that $[q'_\lambda \left(\boldsymbol{\Sigma}_\Pi^* \right)]_{ii} = 0$ for $i \notin \mathcal{S}_4^{II}$, and for all $i \in \mathcal{S}_4^{II}$,

$$[q'_\lambda \left(\boldsymbol{\Sigma}_\Pi^* \right)]_{ii} = \left[q'_\lambda \left(\sigma_i \left(\mathcal{A}_{(k)}^* \right) \right) \right]_{ii} \leq \lambda. \quad (116)$$

Since $\sigma_i \left(\mathcal{A}_{(k)}^* \right) \leq \nu$, and $q_\lambda \left(\cdot \right)$ satisfies the regularity Assumption 3, $|q'_\lambda \left(t \right)| \leq \lambda$. Therefore

$$\left\| \Pi_{\mathcal{F}_{\mathcal{S}_4^{II}}} \left(\nabla Q_\lambda \left(\mathcal{A}_{(k)}^* \right) \right) \right\|_{\text{sp}} = \max_{i \in \mathcal{S}_4^{II}} [q'_\lambda \left(\boldsymbol{\Sigma}_\Pi^* \right)]_{ii} \leq \lambda. \quad (117)$$

Meanwhile, because of the fact that $\mathcal{F}_{\mathcal{S}_4^{II}} \subset \mathcal{F}_{\mathcal{S}_4}$, we have

$$\left\| \Pi_{\mathcal{F}_{\mathcal{S}_4^{II}}} \left(\lambda \mathcal{H}^* \right) \right\|_{\text{sp}} \leq \left\| \Pi_{\mathcal{F}} \left(\lambda \mathbf{U}_\Pi^* \mathbf{V}_\Pi^{*\top} \right) \right\|_{\text{sp}}. \quad (118)$$

Since $\left\| \mathbf{U}^* \mathbf{V}^{*\top} \right\|_{\text{sp}} = 1$, we have

$$\left\| \Pi_{\mathcal{F}} \left(\lambda \mathbf{U}_\Pi^* \mathbf{V}_\Pi^{*\top} \right) \right\|_{\text{sp}} = \lambda. \quad (119)$$

Thus, from (118) and (119), we have

$$\left\| \Pi_{\mathcal{F}_{\mathcal{S}_4^{II}}} \left(\lambda \mathcal{H}^* \right) \right\|_{\text{sp}} \leq \lambda. \quad (120)$$

Additionally, due to the fact that $\left\| \Pi_{\mathcal{F}_{\mathcal{S}_4^{II}}} \left([\nabla L \left(\mathcal{A}^* \right)]_{(k)} \right) \right\|_{\text{sp}} \leq \left\| [\nabla L \left(\mathcal{A}^* \right)]_{(k)} \right\|_{\text{sp}} \leq \lambda$, which indicates that

$$\begin{aligned} & \left\langle \Pi_{\mathcal{F}_{\mathcal{S}_4^{II}}} \left([\nabla L \left(\mathcal{A}^* \right)]_{(k)} + \nabla Q_\lambda \left(\mathcal{A}_{(k)}^* \right) + \lambda \mathcal{H}^* \right), \left[\mathcal{A}^* - \hat{\mathcal{A}} \right]_{(k)} \right\rangle \\ &= \left\langle \Pi_{\mathcal{F}_{\mathcal{S}_4^{II}}} \left([\nabla L \left(\mathcal{A}^* \right)]_{(k)} \right), \left[\mathcal{A}^* - \hat{\mathcal{A}} \right]_{(k)} \right\rangle + \left\langle \Pi_{\mathcal{F}_{\mathcal{S}_4^{II}}} \left(\nabla Q_\lambda \left(\mathcal{A}_{(k)}^* \right) \right), \left[\mathcal{A}^* - \hat{\mathcal{A}} \right]_{(k)} \right\rangle + \left\langle \Pi_{\mathcal{F}_{\mathcal{S}_4^{II}}} \left(\lambda \mathcal{H}^* \right), \left[\mathcal{A}^* - \hat{\mathcal{A}} \right]_{(k)} \right\rangle \\ &\leq \left(\left\| \Pi_{\mathcal{F}_{\mathcal{S}_4^{II}}} \left([\nabla L \left(\mathcal{A}^* \right)]_{(k)} \right) \right\|_{\text{sp}} + \left\| \Pi_{\mathcal{F}_{\mathcal{S}_4^{II}}} \left(\nabla Q_\lambda \left(\mathcal{A}_{(k)}^* \right) \right) \right\|_{\text{sp}} + \left\| \Pi_{\mathcal{F}_{\mathcal{S}_4^{II}}} \left(\lambda \mathcal{H}^* \right) \right\|_{\text{sp}} \right) \cdot \left\| \Pi_{\mathcal{F}_{\mathcal{S}_4^{II}}} \left(\left[\mathcal{A}^* - \hat{\mathcal{A}} \right]_{(k)} \right) \right\|_*, \end{aligned}$$

where the last inequality is derived from the Hölder inequality. Since we have obtained the bound for each term, as in (119) and (120), we have

$$\begin{aligned} \left\langle \Pi_{\mathcal{F}_{\mathcal{S}_4^{II}}} \left([\nabla L \left(\mathcal{A}^* \right)]_{(k)} + \nabla Q_\lambda \left(\mathcal{A}_{(k)}^* \right) + \lambda \mathcal{H}^* \right), \left[\mathcal{A}^* - \hat{\mathcal{A}} \right]_{(k)} \right\rangle &\leq 3\lambda \left\| \Pi_{\mathcal{F}_{\mathcal{S}_4^{II}}} \left(\left[\mathcal{A}^* - \hat{\mathcal{A}} \right]_{(k)} \right) \right\|_* \\ &\leq 3\lambda \sqrt{|\mathcal{S}_4^{II}|} \left\| \left[\mathcal{A}^* - \hat{\mathcal{A}} \right]_{(k)} \right\|_{\text{F}} \\ &= 3\lambda \sqrt{|\mathcal{S}_4^{II}|} \left\| \mathcal{A}^* - \hat{\mathcal{A}} \right\|_{\text{F}}. \end{aligned} \quad (121)$$

where the second inequality utilizes the fact that $\text{rank} \left(\Pi_{\mathcal{F}_{\mathcal{S}_4^{II}}} \left(\left[\mathcal{A}^* - \hat{\mathcal{A}} \right]_{(k)} \right) \right) \leq |\mathcal{S}_4^{II}|$.

(iii) For $i \in \mathcal{S}_4^c$, which correspond to the projector $\Pi_{\mathcal{F}^\perp}$ since $\sigma_i \left(\Pi_{\mathcal{F}^\perp} \left(\mathcal{A}_{(k)}^* \right) \right) = 0$.

Based on Assumption 3, $q_\lambda(0) = q'_\lambda(0) = 0$. We have that $\nabla Q_\lambda \left(\mathcal{A}_{(k)}^* \right) = \mathbf{U}_c^* q'_\lambda(\Sigma_c^*) \mathbf{V}_c^{*\top}$, where $\Sigma_c^* \in \mathbb{R}^{r \times r}$ is a diagonal matrix and $r = \min\{d_k, \Pi_{j \neq k} d_j\}$. Now we have

$$\Pi_{\mathcal{F}^\perp} \left(\nabla Q_\lambda \left(\mathcal{A}_{(k)}^* \right) \right) = (\mathbf{I}_c - \mathbf{U}_c^* \mathbf{U}_c^{*\top}) \mathbf{U}_c^* q'_\lambda(\Sigma_c^*) \mathbf{V}_c^{*\top} (\mathbf{I}_c - \mathbf{V}_c^* \mathbf{V}_c^{*\top}) \quad (122)$$

$$= (\mathbf{U}_c^* - \mathbf{U}_c^*) q'_\lambda(\Sigma_c^*) (\mathbf{V}_c^{*\top} - \mathbf{V}_c^{*\top}) \quad (123)$$

$$= 0, \quad (124)$$

where \mathbf{I}_c is the identity matrix. Meanwhile, since

$$\left\| \Pi_{\mathcal{F}^\perp} \left([\nabla L(\mathcal{A}^*)]_{(k)} \right) \right\|_{\text{sp}} \leq \left\| [\nabla L(\mathcal{A}^*)]_{(k)} \right\|_{\text{sp}} = \frac{\left\| [\mathfrak{X}^*(\mathcal{E})]_{(k)} \right\|_{\text{sp}}}{n} \leq \lambda. \quad (125)$$

For $\mathbf{Z}_c^* = -\lambda^{-1} \Pi_{\mathcal{F}^\perp} \left([\nabla L(\mathcal{A}^*)]_{(k)} \right)$ and $\mathcal{H}^* = \mathbf{U}_c^* \mathbf{V}_c^{*\top} + \mathbf{Z}_c^* \in \partial \left\| \mathcal{A}_{(k)}^* \right\|_*$, we have

$$\Pi_{\mathcal{F}^\perp} \left[[\nabla L(\mathcal{A}^*)]_{(k)} + \lambda \mathcal{H}^* \right] = \Pi_{\mathcal{F}^\perp} \left([\nabla L(\mathcal{A}^*)]_{(k)} \right) + \lambda \mathbf{Z}_c^* = 0, \quad (126)$$

which implies that

$$\left\langle \Pi_{\mathcal{F}^\perp} \left([\nabla L(\mathcal{A}^*)]_{(k)} + \lambda \mathcal{H}^* + \nabla Q_\lambda \left(\mathcal{A}_{(k)}^* \right) \right), [\mathcal{A}^* - \hat{\mathcal{A}}]_{(k)} \right\rangle = \left\langle 0, [\mathcal{A}^* - \hat{\mathcal{A}}]_{(k)} \right\rangle = 0. \quad (127)$$

Adding (114), (121) and (127), which indicate that

$$(\rho^- - \zeta^-) \left\| \hat{\mathcal{A}} - \mathcal{A}^* \right\|_{\text{F}} \quad (128)$$

$$\leq \left\langle \Pi_{\mathcal{F}} \left([\nabla L(\mathcal{A}^*)]_{(k)} + \nabla Q_\lambda \left(\mathcal{A}_{(k)}^* \right) + \lambda \mathcal{H}^* \right), [\mathcal{A}^* - \hat{\mathcal{A}}]_{(k)} \right\rangle \quad (129)$$

$$\leq \sqrt{|\mathcal{S}_4^{\text{I}}|} \left\| \Pi_{\mathcal{S}_4^{\text{I}}} \left([\nabla L(\mathcal{A}^*)]_{(k)} \right) \right\|_{\text{sp}} \cdot \left\| \mathcal{A}^* - \hat{\mathcal{A}} \right\|_{\text{F}} + 3\lambda \sqrt{|\mathcal{S}_4^{\text{II}}|} \left\| \mathcal{A}^* - \hat{\mathcal{A}} \right\|_{\text{F}} \quad (130)$$

$$= \left\| \mathcal{A}^* - \hat{\mathcal{A}} \right\|_{\text{F}} \sqrt{|\mathcal{S}_4^{\text{I}}|} \left\| \Pi_{\mathcal{S}_4^{\text{I}}} \left([\nabla L(\mathcal{A}^*)]_{(k)} \right) \right\|_{\text{sp}} + 3\lambda \sqrt{|\mathcal{S}_4^{\text{II}}|}. \quad (131)$$

Thus, we have

$$\left\| \hat{\mathcal{A}} - \mathcal{A}^* \right\|_{\text{F}} \leq \frac{1}{\rho^- - \zeta^-} \left[\sqrt{|\mathcal{S}_4^{\text{I}}|} \left\| \Pi_{\mathcal{S}_4^{\text{I}}} \left([\nabla L(\mathcal{A}^*)]_{(k)} \right) \right\|_{\text{sp}} + 3\lambda \sqrt{|\mathcal{S}_4^{\text{II}}|} \right].$$

□

Lemma 16. Suppose $\mathcal{A}^* \in \mathbb{R}^{d_1 \times \dots \times d_N}$ with rank of each mode- (k) unfolding $|\mathcal{S}_4|$. Then the error bound between the oracle estimator $\hat{\mathcal{A}}^O$ and the true \mathcal{A}^* satisfies

$$\left\| \hat{\mathcal{A}}^O - \mathcal{A}^* \right\|_{\text{F}} = \left\| [\hat{\mathcal{A}}^O - \mathcal{A}^*]_{(k)} \right\|_{\text{F}} \leq \frac{2\sqrt{|\mathcal{S}_4|} \left\| \Pi_{\mathcal{F}} \left([\nabla L(\mathcal{A}^*)]_{(k)} \right) \right\|_{\text{sp}}}{\rho^-}. \quad (132)$$

Proof. Let $\mathcal{B}' = \hat{\mathcal{A}}^O - \mathcal{A}^*$. According to the general estimator (2) and the definition of the adjoint operator $\mathfrak{X}(\cdot)$, we can express the difference in loss as follows

$$L(\hat{\mathcal{A}}^O) - L(\mathcal{A}^*) = \frac{1}{2n} \sum_{i=1}^n \left[\mathcal{Y}^{(i)} - \mathfrak{X}^{(i)}(\mathcal{A}^* + \mathcal{B}') \right]^2 - \frac{1}{2n} \sum_{i=1}^n \left[\mathcal{Y}^{(i)} - \mathfrak{X}^{(i)}(\mathcal{A}^*) \right]^2 \quad (133)$$

$$= \frac{1}{2n} \sum_{i=1}^n \left[\mathcal{E}^{(i)} - \mathfrak{X}^{(i)}(\mathcal{B}') \right]^2 - \frac{1}{2n} \sum_{i=1}^n \mathcal{E}^{(i)} \quad (134)$$

$$= \frac{1}{2n} \left\| [\mathfrak{X}(\mathcal{B}')]_{(k)} \right\|_{\text{sp}}^2 - \frac{1}{n} \langle \mathfrak{X}^*(\mathcal{E}), \mathcal{B}' \rangle. \quad (135)$$

Given that $\hat{\mathcal{A}}^O$ minimizes $L(\cdot)$ over the subspace \mathcal{F} and $\mathcal{A}_{(k)}^* \in \mathcal{F}$, we have

$$L(\hat{\mathcal{A}}^O) - L(\mathcal{A}^*) \leq 0. \quad (136)$$

Thus, it follows that

$$\frac{1}{2n} \left\| [\mathfrak{X}(\mathcal{B}')]_{(k)} \right\|_{\text{sp}}^2 \leq \frac{1}{n} \langle \mathfrak{X}^*(\mathcal{E}), \mathcal{B}' \rangle. \quad (137)$$

By the RSC condition 1, we know that

$$L(\mathcal{A} + \mathcal{B}) - L(\mathcal{A}) \geq \langle \nabla L(\mathcal{A}), \mathcal{B} \rangle + \frac{\rho^-}{2} \|\mathcal{B}\|_{\text{F}}^2. \quad (138)$$

Applying this to \mathcal{B}' ,

$$\begin{aligned} \frac{\rho^-}{2} \|\mathcal{B}'\|_{\text{F}}^2 &\leq L(\mathcal{B}') - L(\mathcal{A}^*) - \langle \nabla L(\mathcal{A}^*), \mathcal{B}' \rangle \\ &= \frac{1}{2n} \left\| [\mathfrak{X}(\mathcal{B}')]_{(k)} \right\|_{\text{sp}}^2 - \frac{1}{n} \langle \mathfrak{X}^*(\mathcal{E}), \mathcal{B}' \rangle - \langle \nabla L(\mathcal{A}^*), \mathcal{B}' \rangle. \end{aligned} \quad (139)$$

Substituting (137) into (139) gives

$$\frac{\rho^-}{2} \|\mathcal{B}'\|_{\text{F}}^2 \leq \frac{1}{2n} \left\| [\mathfrak{X}(\mathcal{B}')]_{(k)} \right\|_{\text{sp}}^2 \leq \frac{1}{n} \langle \mathfrak{X}^*(\mathcal{E}), \mathcal{B}' \rangle. \quad (140)$$

Therefore, we have

$$\|\mathcal{B}'\|_{\text{F}}^2 \leq \frac{2 \langle \Pi_{\mathcal{F}}([\mathfrak{X}^*(\mathcal{E})]_{(k)}), \mathcal{B}' \rangle}{n\rho^-} \leq \frac{2 \left\| \Pi_{\mathcal{F}}([\mathfrak{X}^*(\mathcal{E})]_{(k)}) \right\|_{\text{sp}} \cdot \|\mathcal{B}'\|_*}{n\rho^-}. \quad (141)$$

Using the fact that $\text{rank}(\mathcal{B}') = |\mathcal{S}_4|$, we have

$$\left\| \mathcal{B}'_{(k)} \right\|_* \leq \sqrt{|\mathcal{S}_4|} \left\| \mathcal{B}'_{(k)} \right\|_{\text{F}}. \quad (142)$$

Thus, it follows that

$$\left\| \mathcal{B}'_{(k)} \right\|_{\text{F}}^2 \leq \frac{2\sqrt{|\mathcal{S}_4|} \left\| \Pi_{\mathcal{F}}([\mathfrak{X}^*(\mathcal{E})]_{(k)}) \right\|_{\text{sp}} \cdot \left\| \mathcal{B}'_{(k)} \right\|_{\text{F}}^2}{n\rho^-}. \quad (143)$$

Recalling that $\nabla L(\mathcal{A}^*) = -\frac{\mathfrak{X}^*(\mathcal{E})}{n}$, we conclude

$$\left\| \mathcal{B}'_{(k)} \right\|_{\text{F}} \leq \frac{2\sqrt{|\mathcal{S}_4|} \left\| \Pi_{\mathcal{F}}([\mathfrak{X}^*(\mathcal{E})]_{(k)}) \right\|_{\text{sp}}}{n\rho^-} = \frac{2\sqrt{|\mathcal{S}_4|} \left\| \Pi_{\mathcal{F}}([\nabla L(\mathcal{A}^*)]_{(k)}) \right\|_{\text{sp}}}{\rho^-}. \quad (144)$$

Thus, since $\left\| \mathcal{B}'_{(k)} \right\|_{\text{F}} = \|\mathcal{B}'\|_{\text{F}}$, we have the desired error bound

$$\left\| \hat{\mathcal{A}}^O - \mathcal{A}^* \right\|_{\text{F}} = \|\mathcal{B}'\|_{\text{F}} \leq \frac{2\sqrt{|\mathcal{S}_4|} \left\| \Pi_{\mathcal{F}}([\nabla L(\mathcal{A}^*)]_{(k)}) \right\|_{\text{sp}}}{\rho^-}. \quad (145)$$

□

Proof of Corollary 4

Proof. Suppose $\hat{\mathcal{H}} \in \partial \left\| \left(\hat{\mathcal{A}}_{(k)} \right) \right\|_*$, since $\hat{\mathcal{A}}$ satisfies the optimality condition, for any $\mathcal{A}' \in \mathbb{R}^{d_1 \times \dots \times d_N}$, it holds that

$$\max_{\mathcal{A}'} \left\{ \left\langle \nabla L \left(\hat{\mathcal{A}} \right), \hat{\mathcal{A}} - \mathcal{A}' \right\rangle + \left\langle \nabla Q_\lambda \left(\hat{\mathcal{A}}_{(k)} \right) + \lambda \hat{\mathcal{H}}, \left[\hat{\mathcal{A}} - \mathcal{A}' \right]_{(k)} \right\rangle \right\} \leq 0. \quad (146)$$

In the following, we will show some $\hat{\mathcal{H}}^O \in \partial \left\| \hat{\mathcal{A}}_{(k)}^O \right\|_*$ satisfy that

$$\max_{\mathcal{A}'} \left\{ \left\langle \left[\nabla L \left(\hat{\mathcal{A}}^O \right) \right]_{(k)} + \nabla Q_\lambda \left(\hat{\mathcal{A}}_{(k)}^O \right) + \lambda \hat{\mathcal{H}}^O, \left[\hat{\mathcal{A}}^O - \mathcal{A}' \right]_{(k)} \right\rangle \right\} \leq 0. \quad (147)$$

Recall that $\tilde{L}(\mathcal{A}) = L(\mathcal{A}) + Q_\lambda(\mathcal{A}_{(k)})$. Projecting the components of the inner product of the LHS in (147) into two complementary spaces \mathcal{F} and \mathcal{F}^\perp , we have the following decomposition

$$\begin{aligned} & \left\langle \left[\nabla L \left(\hat{\mathcal{A}}^O \right) \right]_{(k)} + \nabla Q_\lambda \left(\hat{\mathcal{A}}_{(k)}^O \right) + \lambda \hat{\mathcal{H}}^O, \left[\hat{\mathcal{A}}^O - \mathcal{A}' \right]_{(k)} \right\rangle \\ &= \underbrace{\left\langle \left[\nabla L \left(\hat{\mathcal{A}}^O \right) \right]_{(k)} + \nabla Q_\lambda \left(\hat{\mathcal{A}}_{(k)}^O \right) + \lambda \hat{\mathcal{H}}^O, \Pi_{\mathcal{F}} \left(\left[\hat{\mathcal{A}}^O - \mathcal{A}' \right]_{(k)} \right) \right\rangle}_{P_1} \\ &+ \underbrace{\left\langle \left[\nabla L \left(\hat{\mathcal{A}}^O \right) \right]_{(k)} + \nabla Q_\lambda \left(\hat{\mathcal{A}}_{(k)}^O \right) + \lambda \hat{\mathcal{H}}^O, \Pi_{\mathcal{F}^\perp} \left(\left[\hat{\mathcal{A}}^O - \mathcal{A}' \right]_{(k)} \right) \right\rangle}_{P_2}. \end{aligned} \quad (148)$$

For Term P_1 .

By applying Weyl's inequality for singular values, we obtain

$$\max_l \left| \sigma_l \left(\mathcal{A}_{(k)}^* \right) - \sigma_l \left(\hat{\mathcal{A}}_{(k)}^O \right) \right| \leq \left\| \mathcal{A}_{(k)}^* - \hat{\mathcal{A}}_{(k)}^O \right\|_{\text{sp}}. \quad (149)$$

Further, from the properties of the Frobenius norm, we have

$$\left\| \mathcal{A}^* - \hat{\mathcal{A}}^O \right\|_{\text{F}} = \left\| \left[\mathcal{A}^* - \hat{\mathcal{A}}^O \right]_{(k)} \right\|_{\text{F}}. \quad (150)$$

From Lemma 16, the estimation error $\mathcal{A}^* - \hat{\mathcal{A}}^O$ yields

$$\max_l \left| \sigma_l \left(\mathcal{A}_{(k)}^* \right) - \sigma_l \left(\hat{\mathcal{A}}_{(k)}^O \right) \right| \leq \frac{2\sqrt{|\mathcal{S}_4|} \left\| [\mathcal{X}^*(\mathcal{E})]_{(k)} \right\|_{\text{sp}}}{n\rho^-}, \quad (151)$$

where $|\mathcal{S}_4|$ denotes the rank of the unfolded matrix $\mathcal{A}_{(k)}^*$. Utilizing the weak condition of the singular values, we find

$$\min_{i \in \mathcal{S}_4} \left| \sigma_i \left(\mathcal{A}_{(k)}^* \right) \right| \geq \nu + \frac{2\sqrt{|\mathcal{S}_4|}}{n\rho^-} \left\| [\mathcal{X}^*(\mathcal{E})]_{(k)} \right\|_{\text{sp}}. \quad (152)$$

Applying the triangle inequality, we derive

$$\begin{aligned} \min_{i \in \mathcal{S}_4} \left| \sigma_i \left(\hat{\mathcal{A}}_{(k)}^O \right) \right| &= \min_{i \in \mathcal{S}_4} \left| \sigma_i \left(\hat{\mathcal{A}}_{(k)}^O \right) - \sigma_i \left(\mathcal{A}_{(k)}^* \right) + \sigma_i \left(\mathcal{A}_{(k)}^* \right) \right| \\ &\geq -\max_{i \in \mathcal{S}_4} \left| \sigma_i \left(\hat{\mathcal{A}}_{(k)}^O \right) - \sigma_i \left(\mathcal{A}_{(k)}^* \right) \right| + \min_{i \in \mathcal{S}_4} \left| \sigma_i \left(\mathcal{A}_{(k)}^* \right) \right| \\ &\geq -\frac{2\sqrt{|\mathcal{S}_4|}}{n\rho^-} \left\| [\mathcal{X}^*(\mathcal{E})]_{(k)} \right\|_{\text{sp}} + \nu + \frac{2\sqrt{|\mathcal{S}_4|}}{n\rho^-} \left\| [\mathcal{X}^*(\mathcal{E})]_{(k)} \right\|_{\text{sp}} \\ &= \nu. \end{aligned} \quad (153)$$

Considering the definition of oracle estimator, $\hat{\mathcal{A}}^O \in \mathcal{F}$, which implies the tensor rank of each mode-(k) unfolding rank $\left(\hat{\mathcal{A}}_{(k)}^O\right) = |\mathcal{S}_4|$. And we have the singular value decomposition $\hat{\mathcal{A}}_{(k)}^O = \mathbf{U}^* \hat{\Sigma}^O \mathbf{V}^{*\top}$. Since $R_\lambda(\mathcal{A}_{(k)}) = \lambda \|\mathcal{A}_{(k)}\|_* + Q_\lambda(\mathcal{A}_{(k)})$, for $\hat{\mathbf{Z}}^O \in \mathcal{F}^\perp$, $\|\hat{\mathbf{Z}}^O\|_{\text{sp}} \leq 1$, and $\left(\hat{\Sigma}^O\right)_{\mathcal{S}_4} \in \mathbb{R}^{|\mathcal{S}_4| \times |\mathcal{S}_4|}$ is a diagonal matrix where $\Pi_{\mathcal{F}}\left(q'_\lambda\left(\hat{\Sigma}^O\right)\right) = q'_\lambda\left(\left(\hat{\Sigma}^O\right)_{\mathcal{S}_4}\right)$. Based on the definition of $\nabla Q_\lambda(\cdot)$ and $\partial \|\cdot\|_*$, we have

$$\begin{aligned} \Pi_{\mathcal{F}}\left(\nabla R_\lambda\left(\hat{\mathcal{A}}_{(k)}^O\right)\right) &= \Pi_{\mathcal{F}}\left(Q_\lambda\left(\hat{\mathcal{A}}_{(k)}^O\right)\right) + \lambda \partial \left\|\hat{\mathcal{A}}_{(k)}^O\right\|_* \\ &= \Pi_{\mathcal{F}}\left(\mathbf{U}^* q'_\lambda\left(\hat{\Sigma}^O\right) \mathbf{V}^{*\top} + \lambda \mathbf{U}^* \mathbf{V}^{*\top} + \lambda \hat{\mathbf{Z}}^O\right) \\ &= \mathbf{U}^* \left(q'_\lambda\left(\left(\hat{\Sigma}^O\right)_{\mathcal{S}_4}\right) + \lambda \mathbf{I}_{\mathcal{S}_4}\right) \mathbf{V}^{*\top}, \end{aligned} \quad (154)$$

where the second equality in (154) is to simply project each component into the subspace \mathcal{F} . $\mathbf{I}_{\mathcal{S}_4}$ is the identity matrix and $\mathbf{I}_{\mathcal{S}_4} \in \mathbb{R}^{|\mathcal{S}_4| \times |\mathcal{S}_4|}$. Since $p_\lambda(t) = q_\lambda(t) + \lambda|t|$, we have $p'_\lambda(t) = q'_\lambda(t) + \lambda t$ for all $t > 0$. Consider the diagonal matrix $q'_\lambda\left(\left(\hat{\Sigma}^O\right)_{\mathcal{S}_4}\right) + \lambda \mathbf{I}_{\mathcal{S}_4}$, we have the i -th ($i \in \mathcal{S}_4$) element on the diagonal that

$$\left[q'_\lambda\left(\left(\hat{\Sigma}^O\right)_{\mathcal{S}_4}\right) + \lambda \mathbf{I}_{\mathcal{S}_4}\right]_{ii} = q'_\lambda\left(\sigma_i\left(\hat{\mathcal{A}}_{(k)}^O\right)\right) + \lambda = p'_\lambda\left(\sigma_i\left(\hat{\mathcal{A}}_{(k)}^O\right)\right). \quad (155)$$

Since $p_\lambda(\cdot)$ satisfies the regularity condition (iii) in Assumption 3 that $p'_\lambda(t) = 0$ for all $t \geq \nu$, we have $p'_\lambda\left(\sigma_i\left(\hat{\mathcal{A}}_{(k)}^O\right)\right) = 0$ for $i \in \mathcal{S}_4$, due to the fact that $\sigma_i\left(\hat{\mathcal{A}}_{(k)}^O\right) \geq \nu > 0$.

Therefore, the diagonal matrix $q'_\lambda\left(\left(\hat{\Sigma}^O\right)_{\mathcal{S}_4}\right) + \lambda \mathbf{I}_{\mathcal{S}_4} = 0$, substituting which in to (154) yields

$$\Pi_{\mathcal{F}}\left(\nabla R_\lambda\left(\hat{\mathcal{A}}_{(k)}^O\right)\right) = 0. \quad (156)$$

Since $\hat{\mathcal{A}}^O$ is the estimator over \mathcal{F} , we have the optimality condition that for any $\mathcal{A}' \in \mathbb{R}^{d_1 \times \dots \times d_N}$, it holds that

$$\max_{\mathcal{A}'} \left\langle \left[\nabla L\left(\hat{\mathcal{A}}^O\right)\right]_{(k)}, \Pi_{\mathcal{F}}\left(\left[\hat{\mathcal{A}}^O - \mathcal{A}'\right]_{(k)}\right) \right\rangle \leq 0. \quad (157)$$

Substitute (156) and (157) into P_1 , for all $\hat{\mathcal{H}}^O \in \partial \left\|\hat{\mathcal{A}}_{(k)}^O\right\|_*$ we have

$$\begin{aligned} &\max_{\mathcal{A}'} \left\langle \left[\nabla L\left(\hat{\mathcal{A}}^O\right)\right]_{(k)} + \nabla Q_\lambda\left(\hat{\mathcal{A}}_{(k)}^O\right) + \lambda \hat{\mathcal{H}}^O, \Pi_{\mathcal{F}}\left(\left[\hat{\mathcal{A}}^O - \mathcal{A}'\right]_{(k)}\right) \right\rangle \\ &= \max_{\mathcal{A}'} \left\langle \left[\nabla L\left(\hat{\mathcal{A}}^O\right)\right]_{(k)}, \Pi_{\mathcal{F}}\left(\left[\hat{\mathcal{A}}^O - \mathcal{A}'\right]_{(k)}\right) \right\rangle + \max_{\mathcal{A}'} \left\langle \Pi_{\mathcal{F}}\left(\nabla R_\lambda\left(\hat{\mathcal{A}}_{(k)}^O\right)\right), \Pi_{\mathcal{F}}\left(\left[\hat{\mathcal{A}}^O - \mathcal{A}'\right]_{(k)}\right) \right\rangle \\ &\leq 0. \end{aligned} \quad (158)$$

For Term P_2 .

By definition of $\nabla Q_\lambda(\cdot)$, and the regularity condition (v) in Assumption 3, we do the decomposition that $\nabla Q_\lambda\left(\hat{\mathcal{A}}_{(k)}^O\right) = \mathbf{U}^* q'_\lambda\left(\hat{\Sigma}^O\right) \mathbf{V}^{*\top}$, where $\hat{\Sigma}^O$ is diagonal matrix. Projecting $\nabla Q_\lambda\left(\hat{\mathcal{A}}_{(k)}^O\right)$ into \mathcal{F}^\perp yields that

$$\begin{aligned} \Pi_{\mathcal{F}^\perp}\left(\nabla Q_\lambda\left(\hat{\mathcal{A}}_{(k)}^O\right)\right) &= (\mathbf{I}_{\mathcal{S}_4} - \mathbf{U}^* \mathbf{U}^{*\top}) \mathbf{U}^* q'_\lambda\left(\hat{\Sigma}^O\right) \mathbf{V}^{*\top} (\mathbf{I}_{\mathcal{S}_4} - \mathbf{V}^* \mathbf{V}^{*\top}) \\ &= (\mathbf{U}^* - \mathbf{U}^*) q'_\lambda\left(\left(\hat{\Sigma}^O\right)_{\mathcal{S}_4}\right) (\mathbf{V}^{*\top} - \mathbf{V}^{*\top}) \\ &= 0. \end{aligned} \quad (159)$$

Therefore,

$$P_2 = \left\langle \Pi_{\mathcal{F}^\perp} \left(\left[\nabla L \left(\hat{\mathcal{A}}^O \right) \right]_{(k)} + \lambda \hat{\mathcal{H}}^O \right), \Pi_{\mathcal{F}^\perp} \left(\left[\hat{\mathcal{A}}^O - \mathcal{A}' \right]_{(k)} \right) \right\rangle. \quad (160)$$

Moreover, with the triangle inequality, we have

$$\begin{aligned} \left\| \left[\nabla L \left(\hat{\mathcal{A}}^O \right) \right]_{(k)} \right\|_{\text{sp}} &\leq \left\| \left[\nabla L \left(\mathcal{A}^* \right) \right]_{(k)} \right\|_{\text{sp}} + \left\| \left[\nabla L \left(\mathcal{A}^* \right) \right]_{(k)} - \left[\nabla L \left(\hat{\mathcal{A}}^O \right) \right]_{(k)} \right\|_{\text{sp}} \\ &\leq \left\| \left[\nabla L \left(\mathcal{A}^* \right) \right]_{(k)} \right\|_{\text{sp}} + \left\| \left[\nabla L \left(\mathcal{A}^* \right) \right]_{(k)} - \left[\nabla L \left(\hat{\mathcal{A}}^O \right) \right]_{(k)} \right\|_{\text{F}}, \end{aligned} \quad (161)$$

where the second inequality comes from the fact that

$$\left\| \left[\nabla L \left(\mathcal{A}^* \right) \right]_{(k)} - \left[\nabla L \left(\hat{\mathcal{A}}^O \right) \right]_{(k)} \right\|_{\text{sp}} \leq \left\| \left[\nabla L \left(\mathcal{A}^* \right) \right]_{(k)} - \left[\nabla L \left(\hat{\mathcal{A}}^O \right) \right]_{(k)} \right\|_{\text{F}}. \quad (162)$$

From Restricted Strong Smoothness in Assumption 2 where $\left\| \nabla L \left(\mathcal{A} \right) - \nabla L \left(\mathcal{A} + \mathcal{B}' \right) \right\|_{\text{F}} \leq \rho^+ \left\| \mathcal{B}' \right\|_{\text{F}}$, we can substitute it into (161), and we have

$$\left\| \left[\nabla L \left(\hat{\mathcal{A}}^O \right) \right]_{(k)} \right\|_{\text{sp}} \leq \left\| \left[\nabla L \left(\mathcal{A}^* \right) \right]_{(k)} \right\|_{\text{sp}} + \rho^+ \left\| \mathcal{A}^* - \hat{\mathcal{A}}^O \right\|_{\text{F}}. \quad (163)$$

Since $\Pi_{\mathcal{F}^\perp} \left(\mathcal{B}' \right) = 0$, it is evident that $\mathcal{B}' \in \mathcal{C}$. Substitute (132) from Lemma 16 into (163), from the choice of λ , we have

$$\left\| \Pi_{\mathcal{F}^\perp} \left(\left[\nabla L \left(\hat{\mathcal{A}}^O \right) \right]_{(k)} \right) \right\|_{\text{sp}} \leq \left\| \left[\nabla L \left(\hat{\mathcal{A}}^O \right) \right]_{(k)} \right\|_{\text{sp}} \quad (164)$$

$$\leq \left\| \left[\nabla L \left(\mathcal{A}^* \right) \right]_{(k)} \right\|_{\text{sp}} + \frac{2\sqrt{|\mathcal{S}_4|}\rho^+}{n\rho^-} \left\| [\mathfrak{X}^* \left(\mathcal{E} \right)]_{(k)} \right\|_{\text{sp}} \quad (165)$$

$$\leq \lambda. \quad (166)$$

By setting $\hat{\mathbf{Z}}^O = -\lambda^{-1} \Pi_{\mathcal{F}^\perp} \left(\left[\nabla L \left(\hat{\mathcal{A}}^O \right) \right]_{(k)} \right)$, such that $\hat{\mathcal{H}}^O = \mathbf{U}^* \mathbf{V}^{*\top} + \hat{\mathbf{Z}}^O \in \partial \left\| \hat{\mathcal{A}}_{(k)}^O \right\|_*$, since $\hat{\mathbf{Z}}^O$ satisfies the condition $\hat{\mathbf{Z}}^O \in \mathcal{F}^\perp$. $\left\| \hat{\mathbf{Z}}^O \right\|_{\text{sp}} \leq 1$, we have

$$\Pi_{\mathcal{F}^\perp} \left(\left[\nabla L \left(\hat{\mathcal{A}}^O \right) \right]_{(k)} + \lambda \hat{\mathcal{H}}^O \right) = 0, \quad (167)$$

which implies that

$$P_2 = \left\langle 0, \Pi_{\mathcal{F}^\perp} \left(\left[0 - \mathcal{A}' \right]_{(k)} \right) \right\rangle = 0. \quad (168)$$

Substituting (158) and (168) into (148), we obtain (147) that

$$\max_{\mathcal{A}'} \left\langle \left[\nabla L \left(\hat{\mathcal{A}}^O \right) \right]_{(k)} + Q_\lambda \left(\hat{\mathcal{A}}_{(k)}^O \right) + \lambda \hat{\mathcal{H}}^O, \left[\hat{\mathcal{A}}^O - \mathcal{A}' \right]_{(k)} \right\rangle \leq 0. \quad (169)$$

Now we are going to prove that $\hat{\mathcal{A}}^O = \hat{\mathcal{A}}$ and the error bound between $\hat{\mathcal{A}}^O$ and \mathcal{A}^* .

Similar to the proof of Lemma 15, since $\|\cdot\|_*$ is convex, and applying Lemma 13, we have

$$\begin{aligned} 0 &\geq \left\langle \nabla L \left(\hat{\mathcal{A}} \right), \hat{\mathcal{A}}^O - \hat{\mathcal{A}} \right\rangle + \left\langle \nabla Q_\lambda \left(\hat{\mathcal{A}}_{(k)} \right) + \lambda \hat{\mathcal{H}}, \left[\hat{\mathcal{A}}^O - \hat{\mathcal{A}} \right]_{(k)} \right\rangle \\ &\quad + \left\langle \nabla L \left(\hat{\mathcal{A}}^O \right), \hat{\mathcal{A}} - \hat{\mathcal{A}}^O \right\rangle + \left\langle \nabla Q_\lambda \left(\hat{\mathcal{A}}_{(k)}^O \right) + \lambda \hat{\mathcal{H}}^O, \left[\hat{\mathcal{A}} - \hat{\mathcal{A}}^O \right]_{(k)} \right\rangle + (\rho^- - \zeta^-) \left\| \hat{\mathcal{A}}^O - \hat{\mathcal{A}} \right\|_{\text{F}}^2. \end{aligned} \quad (170)$$

From (146), we have

$$\begin{aligned} & \left\langle \nabla L(\hat{\mathcal{A}}), \hat{\mathcal{A}} - \hat{\mathcal{A}}^O \right\rangle + \left\langle \nabla Q_\lambda(\hat{\mathcal{A}}_{(k)}) + \lambda \hat{\mathcal{H}}, [\hat{\mathcal{A}} - \hat{\mathcal{A}}^O]_{(k)} \right\rangle \\ & \leq \max_{\mathcal{A}'} \left\{ \left\langle \nabla L(\hat{\mathcal{A}}), \hat{\mathcal{A}} - \mathcal{A}' \right\rangle + \left\langle \nabla Q_\lambda(\hat{\mathcal{A}}_{(k)}) + \lambda \hat{\mathcal{H}}, [\hat{\mathcal{A}} - \mathcal{A}']_{(k)} \right\rangle \right\} \leq 0. \end{aligned} \quad (171)$$

Therefore, in (170),

$$\left\langle \nabla L(\hat{\mathcal{A}}), \hat{\mathcal{A}}^O - \hat{\mathcal{A}} \right\rangle + \left\langle \nabla Q_\lambda(\hat{\mathcal{A}}_{(k)}) + \lambda \hat{\mathcal{H}}, [\hat{\mathcal{A}}^O - \hat{\mathcal{A}}]_{(k)} \right\rangle \geq 0. \quad (172)$$

From (147), we have

$$\begin{aligned} & \left\langle \nabla L(\hat{\mathcal{A}}^O), \hat{\mathcal{A}}^O - \hat{\mathcal{A}} \right\rangle + \left\langle \nabla Q_\lambda(\hat{\mathcal{A}}_{(k)}^O) + \lambda \hat{\mathcal{H}}, [\hat{\mathcal{A}}^O - \hat{\mathcal{A}}]_{(k)} \right\rangle \\ & \leq \max_{\mathcal{A}'} \left\{ \left\langle \nabla L(\hat{\mathcal{A}}^O), \hat{\mathcal{A}}^O - \mathcal{A}' \right\rangle + \left\langle \nabla Q_\lambda(\hat{\mathcal{A}}_{(k)}^O) + \lambda \hat{\mathcal{H}}, [\hat{\mathcal{A}}^O - \mathcal{A}']_{(k)} \right\rangle \right\} \leq 0. \end{aligned} \quad (173)$$

Therefore, in (170),

$$\left\langle \nabla L(\hat{\mathcal{A}}^O), \hat{\mathcal{A}} - \hat{\mathcal{A}}^O \right\rangle + \left\langle \nabla Q_\lambda(\hat{\mathcal{A}}_{(k)}^O) + \lambda \hat{\mathcal{H}}^O, [\hat{\mathcal{A}} - \hat{\mathcal{A}}^O]_{(k)} \right\rangle \geq 0. \quad (174)$$

Substituting (170) and (171) into (173) such that

$$(\rho^- - \zeta^-) \left\| \hat{\mathcal{A}}^O - \hat{\mathcal{A}} \right\|_F^2 \geq 0. \quad (175)$$

Since $\rho^- > \zeta^-$, the inequation holds only if

$$\hat{\mathcal{A}}^O = \hat{\mathcal{A}}.$$

And by Lemma 16, we obtain the statistical oracle bound for the penalty

$$\left\| \hat{\mathcal{A}}^O - \mathcal{A}^* \right\|_F = \left\| [\hat{\mathcal{A}} - \mathcal{A}^*]_{(k)} \right\|_F \leq \frac{2\sqrt{|\mathcal{S}_4|} \left\| \Pi_{\mathcal{F}}([\mathcal{X}^*(\mathcal{E})]_{(k)}) \right\|_{\text{sp}}}{n\rho^-}. \quad (176)$$

Furthermore, we can have

$$\left\| \hat{\mathcal{A}}^O - \mathcal{A}^* \right\|_F = \frac{2\sqrt{|\mathcal{S}_4|} \tau_k}{n\rho^-}, \quad (177)$$

where $\tau_k = \left\| \Pi_{\mathcal{F}}([\mathcal{X}^*(\mathcal{E})]_{(k)}) \right\|_{\text{sp}}$, which completes the proof. \square

G Proof of the Corollary 5

Recall that the proposed slice-wise lowrankness penalty can be reformulated as the sum of the ℓ_1 penalty and a concave part. Specifically, we have:

$$R_\lambda(\mathcal{A}) = \sum_{l=1}^{\Pi_{m \neq j, k} d_m} \sum_{s=1}^{s^{\text{all}}} p_\lambda \left(\sigma_s \left([\mathcal{A}_{(j, k)}]_{\cdot, \cdot, l} \right) \right) = \sum_{l=1}^{\Pi_{m \neq j, k} d_m} \left[\lambda \left\| [\mathcal{A}_{(j, k)}]_{\cdot, \cdot, l} \right\|_* + Q_\lambda \left([\mathcal{A}_{(j, k)}]_{\cdot, \cdot, l} \right) \right],$$

where $s^{\text{all}} = \min\{d_j d_k, \prod_{l \neq j,k} d_l\}$, $[\mathcal{A}_{(j,k)}]_{\cdot,\cdot,l}$ denotes the l -th slice of the mode- (j,k) unfolding $\mathcal{A}_{(j,k)}$ and $\sigma_s([\mathcal{A}_{(j,k)}]_{\cdot,\cdot,l})$ denotes the s -th singular value of the slice. For the estimation problem, we define

$$\tilde{L}(\mathcal{A}) = L(\mathcal{A}) + \sum_{l=1}^{\prod_{m \neq j,k} d_m} Q_\lambda([\mathcal{A}_{(j,k)}]_{\cdot,\cdot,l}),$$

where $Q_\lambda([\mathcal{A}_{(j,k)}]_{\cdot,\cdot,l}) = \sum_{s=1}^{s^{\text{all}}} q_\lambda(\sigma_s([\mathcal{A}_{(j,k)}]_{\cdot,\cdot,l}))$.

Based on Lemma 13, for slice-wise lowrankness regularizer, we can similarly prove the following lemmas

Lemma 17. Under Assumption 1, $\rho^- > \zeta^-$, and the regularization parameter $\lambda \geq \frac{\|[\mathfrak{X}^*(\mathcal{E})]_{(j,k)}\|_{\text{sp}}}{2n}$, we have

$$\left\| \Pi_{\mathcal{F}^\perp} \left([\hat{\mathcal{A}}_{(j,k)}]_{\cdot,\cdot,l} - [\mathcal{A}_{(j,k)}^*]_{\cdot,\cdot,l} \right) \right\|_* \leq 5 \left\| \Pi_{\mathcal{F}} \left([\hat{\mathcal{A}}_{(j,k)}]_{\cdot,\cdot,l} - [\mathcal{A}_{(j,k)}^*]_{\cdot,\cdot,l} \right) \right\|_*.$$

Proof. Similar to the proof of Lemma 14, we can prove the Lemma 17 □

From Lemma 13 and Lemma 17, we can prove the following general deterministic bound.

Lemma 18. For the estimated parameter tensor $\hat{\mathcal{A}}$ and the true parameter tensor \mathcal{A}^* , we have

$$\|\hat{\mathcal{A}} - \mathcal{A}^*\|_{\text{F}} \leq \frac{1}{(\rho^- - \zeta^-)} \sqrt{\sum_{l=1}^{\prod_{m \neq j,k} d_m} \left[\sqrt{|\mathcal{S}_5^{\text{I}}|} \left\| \Pi_{\mathcal{F}} \left([\mathfrak{X}^*(\mathcal{E})]_{(j,k)} \right)_{\cdot,\cdot,l} \right\|_{\text{sp}} + 3\lambda \sqrt{|\mathcal{S}_5^{\text{II}}|} \right]^2}.$$

Proof. □

Similar to the proof for Lemma 16, we can derive the error bound for the slice-wise low-rankness regularizer.

Lemma 19. Suppose $\mathcal{A}^* \in \mathbb{R}^{d_1 \times \dots \times d_N}$ with rank of each slices $|\mathcal{S}_5|$. Then the error bound between the oracle estimator $\hat{\mathcal{A}}^O$ and the true \mathcal{A}^* satisfies

$$\|\hat{\mathcal{A}}^O - \mathcal{A}^*\|_{\text{F}} = \sqrt{\sum_{l=1}^{\prod_{m \neq j,k} d_m} \left\| [\hat{\mathcal{A}}^O - \mathcal{A}^*]_{(j,k)} \right\|_{\text{F}}^2} \lesssim \frac{2\sqrt{|\mathcal{S}_5|} \left\| \Pi_{\mathcal{F}} \left([\mathfrak{X}^*(\mathcal{E})]_{(j,k)} \right)_{\cdot,\cdot,l} \right\|_{\text{sp}}}{n\rho^-}.$$

Proof. With Lemma 19, we can also obtain that $\hat{\mathcal{A}}^O = \hat{\mathcal{A}}$. Similarly, we can prove the Corollary 5. □

H Proof of the Theorem 5

Lemma 20 (Slepian's Lemma). Let $\{G_s, s \in \mathcal{S}\}$ and $\{H_s, s \in \mathcal{S}\}$ be two centered Gaussian processes defined over the same index set \mathcal{S} . Suppose that both processes are almost surely bounded. For any $s, t \in \mathcal{S}$, if

$$\mathbb{E}(G_s - G_t)^2 \leq \mathbb{E}(H_s - H_t)^2,$$

then it follows that

$$\mathbb{E}[\sup_{s \in \mathcal{S}} G_s] \leq \mathbb{E}[\sup_{s \in \mathcal{S}} H_s].$$

Furthermore, if the second moments of both processes are equal, i.e., $\mathbb{E}(G_s^2) = \mathbb{E}(H_s^2)$ for all $s \in \mathcal{S}$, then for all $x > 0$, the following inequality holds:

$$\mathbb{P} \left\{ \sup_{s \in \mathcal{S}} G_s > x \right\} \leq \mathbb{P} \left\{ \sup_{s \in \mathcal{S}} H_s > x \right\}.$$

Lemma 21. Let Z denote a centralized χ^2 -distributed random variable with m degrees of freedom. For all $x \geq 0$, the following tail bounds hold:

$$\begin{aligned}\mathbb{P}[Z - m \geq 2\sqrt{mx} + 2x] &\leq \exp(-x), \quad \text{and} \\ \mathbb{P}[Z - m \leq -2\sqrt{mx}] &\leq \exp(-x).\end{aligned}$$

Lemma 22 (Theorem 3.8 in Massart (2007)). Let $\alpha \sim \mathbb{N}(0, I_{d \times d})$ denote a d -dimensional Gaussian random variable. Consider a function $F : \mathbb{R}^d \rightarrow \mathbb{R}$ that satisfies a Lipschitz condition, i.e.,

$$|F(x) - F(y)| \leq L\|x - y\|_2 \quad \text{for all } x, y \in \mathbb{R}^d,$$

where L is the Lipschitz constant. Then for all $t > 0$, we have the following concentration inequality for the deviation of $F(\alpha)$ from its expectation:

$$\mathbb{P}[|F(\alpha) - \mathbb{E}[F(\alpha)]| \geq t] \leq 2 \exp\left(-\frac{t^2}{2L^2}\right).$$

Lemma 23. Let \mathcal{B} be an tensor with each entry following a Gaussian distribution. Suppose that $\lambda \geq 4\eta\sqrt{\frac{\kappa}{n}}\mathbb{E}[R_\lambda^*(\mathcal{B})]$, where $R_\lambda^*(\mathcal{B})$ is the dual norm of $R_\lambda(\mathcal{B})$. Then, with probability at least $1 - \exp(-\mathbb{E}[R_\lambda^*(\mathcal{B})]^2)$, we have

$$\lambda \geq 4R_\lambda^*\left(\frac{1}{n} \sum_{i=1}^n \mathcal{E}^{(i)} \otimes \mathcal{X}^{(i)}\right).$$

Proof. Recall that we have set

$$\lambda \geq 4\eta\sqrt{\frac{\kappa}{n}}\mathbb{E}[R_\lambda^*(\mathcal{B})].$$

First we show that $\lambda \geq 2\eta\sqrt{\frac{\kappa}{n}}R_\lambda^*(\mathcal{B})$ with high probability using concentration of Lipschitz functions for Gaussian random variables, due to Lemma 22. First we prove that $f(\mathcal{B}) = R_\lambda^*(\mathcal{B}) = \sup_{\mathcal{A}: R_\lambda(\mathcal{A}) \leq 1} \langle \mathcal{B}, \mathcal{A} \rangle$ is a 1-Lipschitz function in terms of \mathcal{B} . In particular note that:

$$f(\mathcal{B}) - f(\mathcal{B}') = \sup_{\mathcal{A}: R_\lambda(\mathcal{A}) \leq 1} \langle \mathcal{B}, \mathcal{A} \rangle - \sup_{\mathcal{A}: R_\lambda(\mathcal{A}) \leq 1} \langle \mathcal{B}', \mathcal{A} \rangle.$$

Let $\tilde{\mathcal{A}} := \arg \max_{\mathcal{A}: R_\lambda(\mathcal{A}) \leq 1} \langle \mathcal{B}, \mathcal{A} \rangle$. Then

$$\sup_{\mathcal{A}: R_\lambda(\mathcal{A}) \leq 1} \langle \mathcal{B}, \mathcal{A} \rangle - \sup_{\mathcal{A}: R_\lambda(\mathcal{A}) \leq 1} \langle \mathcal{B}', \mathcal{A} \rangle \leq \langle \mathcal{B}, \tilde{\mathcal{A}} \rangle - \langle \mathcal{B}', \tilde{\mathcal{A}} \rangle \leq \sup_{\mathcal{A}: R_\lambda(\mathcal{A}) \leq 1} \langle \mathcal{B} - \mathcal{B}', \mathcal{A} \rangle \leq \|\mathcal{B} - \mathcal{B}'\|_F,$$

Therefore, applying Lemma 22, we obtain

$$\mathbb{P}\left\{\left|\sup_{\mathcal{A}: R_\lambda(\mathcal{A}) \leq 1} \langle \mathcal{B}, \mathcal{A} \rangle - \mathbb{E}\left[\sup_{\mathcal{A}: R_\lambda(\mathcal{A}) \leq 1} \langle \mathcal{B}, \mathcal{A} \rangle\right]\right| > w(\Omega)\right\} \leq 2 \exp\left(-\frac{1}{2}w^2[\Omega]\right).$$

Therefore,

$$\lambda \geq 2\eta\sqrt{\frac{\kappa}{n}}R_\lambda^*(\mathcal{B})$$

with probability at least $1 - 2 \exp\left(-\frac{w^2(\Omega)}{2}\right)$.

To complete the proof, we use a Gaussian comparison inequality between the supremum of the process $\sqrt{\frac{\kappa}{n}}\langle \mathcal{B}, \mathcal{A} \rangle$ and $\sum_{i=1}^n \langle \mathcal{E}^{(i)} \otimes \mathcal{X}^{(i)}, \mathcal{A} \rangle / n$ over the set Ω . Recall that:

$$R_\lambda^*\left(\frac{1}{n} \sum_{i=1}^n \mathcal{E}^{(i)} \otimes \mathcal{X}^{(i)}\right) = \sup_{\mathcal{A} \in \Omega} \left\langle \mathcal{A}, \frac{1}{n} \sum_{i=1}^n \mathcal{E}^{(i)} \otimes \mathcal{X}^{(i)} \right\rangle.$$

Recall that each $\mathcal{E}^{(i)} \in \mathbb{R}^{d_{M+1} + d_{M+2} + \dots + d_N}$ with each entry having variance η^2 and $\text{vec}(\mathcal{X}) \in \mathbb{R}^{nd_1 d_2 \dots d_M}$ is a Gaussian vector covariance Σ . First we condition on all the $\mathcal{E}^{(i)}$'s, which are independent of the $\mathcal{X}^{(i)}$'s. Further

let $\{u^{(i)} : i = 1, \dots, n\}$ be i.i.d. standard normal Gaussian tensors where $u^{(i)} \in \mathbb{R}^{d_d \times d_2 \times \dots \times d_N}$. First we condition on all the $\mathcal{E}^{(i)}$'s, which are independent of the $\mathcal{X}^{(i)}$'s. Using a standard Gaussian comparison inequality, if we condition on the $\mathcal{E}^{(i)}$'s we have

$$\mathbb{P} \left\{ \sup_{\mathcal{A}: R_\lambda(\mathcal{A}) \leq 1} \frac{1}{n} \sum_{i=1}^n \langle \mathcal{E}^{(i)} \otimes \mathcal{X}^{(i)}, \mathcal{A} \rangle > x \right\} \leq \mathbb{P} \left\{ \sup_{\mathcal{A}: R_\lambda(\mathcal{A}) \leq 1} \frac{1}{n} \sum_{i=1}^n \langle \mathcal{E}^{(i)} \otimes u^{(i)}, \mathcal{A} \rangle > \frac{x}{\sqrt{\kappa}} \right\},$$

since $\text{cov}(\tilde{\mathcal{X}}) = \Sigma \preceq \kappa \mathbf{I}$.

Now we define the $u_j \in \mathbb{R}^n$ as the standard random vector where $1 \leq j \leq d_M$ and $u_j = (u_j^{(1)}, \dots, u_j^{(n)})$. Conditioning on the $u^{(i)}$'s and dealing with the randomness in the $\mathcal{E}^{(i)}$'s,

$$\frac{1}{n} \sum_{i=1}^n \langle \mathcal{E}^{(i)} \otimes u^{(i)}, \mathcal{A} \rangle \leq \max_{1 \leq j \leq d_M} \frac{\|u_j\|_{l^2}}{\sqrt{n}} \frac{\eta}{\sqrt{n}} \langle \mathcal{B}, \mathcal{A} \rangle,$$

where $\mathcal{B} \in \mathbb{R}^{d_1 \times d_2 \times \dots \times d_N}$ is an i.i.d. standard normal tensor, since the $\mathcal{E}^{(i)}$'s are i.i.d. standard normal. Now we upper bound

$$\max_{1 \leq j \leq d_M} \frac{\|u_j\|_{l^2}}{\sqrt{n}},$$

using standard χ^2 tail bounds. Since $\|u_j\|_{l^2}^2/n$ is a χ^2 random variable with n degrees of freedom, for each j ,

$$\mathbb{P} \{ \|u_j\|_{l^2}^2/n \geq 4 \} \leq \exp(-n),$$

using χ^2 tail bounds.

Now taking the union bound over d_M ,

$$\mathbb{P} \left\{ \max_{1 \leq j \leq d_M} \frac{\|u_j\|_{l^2}}{\sqrt{n}} \geq 2 \right\} \leq \exp(\log(d_M) - n),$$

and provided $n \geq 2 \log(d_M)$, it follows that with probability greater than $1 - \exp(-n/2)$, $\max_{1 \leq j \leq d_M} \|u_j\|_{l^2} \leq 2$. Therefore, with probability at least $1 - \exp(-n/2)$,

$$\frac{1}{n} \sum_{i=1}^n \langle \mathcal{E}^{(i)} \otimes u^{(i)}, \mathcal{A} \rangle \leq \frac{2\eta}{\sqrt{n}} \langle \mathcal{B}, \mathcal{A} \rangle.$$

Now we apply Lemma 20 to complete the proof,

$$\mathbb{P} \left\{ \sup_{R_\lambda(\mathcal{A}) \leq 1} \frac{1}{n} \sum_{i=1}^n \langle \mathcal{E}^{(i)} \otimes u^{(i)}, \mathcal{A} \rangle > x \right\} \leq \mathbb{P} \left\{ \sup_{R_\lambda(\mathcal{A}) \leq 1} \frac{2\eta}{\sqrt{n}} \langle \mathcal{B}, \mathcal{A} \rangle > x \right\},$$

for all $x > 0$. Substituting x by $x/\sqrt{\kappa}$ means that

$$\mathbb{P} \left\{ R_\lambda^* \left(\frac{1}{n} \sum_{i=1}^n \langle \mathcal{E}^{(i)} \otimes u^{(i)}, \mathcal{A} \rangle \right) > x \right\} \leq \mathbb{P} \left\{ 2\eta \sqrt{\frac{\kappa}{n}} R_\lambda^*(\mathcal{B}) > x \right\},$$

for any $x > 0$. This completes the proof.

In light of Lemma 11, for the remainder of the proof, we can condition on the event that

$$\lambda \geq 2R_\lambda^* \left(\frac{1}{n} \sum_{i=1}^n \mathcal{E}^{(i)} \otimes \mathcal{X}^{(i)} \right).$$

Under this event,

$$\frac{1}{2n} \sum_{i=1}^n \|\mathcal{X}^{(i)}, \Delta\|_F^2 \leq \frac{1}{2} R_\lambda(\Delta) + \lambda(R_\lambda(\Delta_0) - R_\lambda(\Delta^\perp)) \leq \left(\frac{3}{2}\right) \mathcal{R}(\Delta_0) - \left(c_R - \frac{1}{\eta_R}\right) \mathcal{R}(\Delta^\perp),$$

where Δ^\perp represents a projection into the complementary space of Δ and Δ_0 represents a projection into the zero space of Δ .

Since

$$\frac{1}{2n} \sum_{i=1}^n \|\langle \Delta, \mathbf{x}^{(i)} \rangle\|_{\mathbb{F}}^2 > 0,$$

we get

$$R_\lambda(\Delta^\perp) \leq 3R_\lambda(\Delta_0).$$

Here we define the cone

$$\mathcal{C} = \left\{ \Delta \mid R_\lambda(\Delta^\perp) \leq 3R_\lambda(\Delta_0) \right\}$$

and know that $\Delta \in \mathcal{C}$. Hence

$$\frac{1}{2n} \sum_{i=1}^n \|\langle \mathbf{x}^{(i)}, \Delta \rangle\|_{\mathbb{F}}^2 \leq \frac{3}{2} \lambda R_\lambda(\Delta_0) \leq \frac{3}{2} (s(\mathcal{A})\lambda)^{-1/4} \|\Delta\|_{\mathbb{F}}.$$

Let us define the following set:

$$C(\delta_{\mathbb{F}}) := \{ \Delta \in \mathbb{R}^{d_d \times d_2 \times \dots \times d_N} \mid R_\lambda(\Delta^\perp) \leq 3R_\lambda(\Delta_0) \}.$$

For convenience, in the remainder of this proof let

$$\delta_n := 3(s(\mathcal{A})\lambda)^{-1/4}.$$

Further, let us define the event:

$$E(\delta_{\mathbb{F}}) := \{ \Delta \in C(\delta_{\mathbb{F}}) \mid \|\Delta\|_{\mathbb{F}} \geq \sqrt{\kappa} \delta_{\mathbb{F}} \}.$$

We claim that it suffices to show that $E(\delta_{\mathbb{F}})$ holds with probability at least $1 - \exp(-2c)$ for some constant $c > 0$. In particular, given an arbitrary non-zero $\Delta \in C(\delta_{\mathbb{F}})$, consider the re-scaled tensor

$$\tilde{\Delta} = \kappa \eta \frac{\Delta}{\|\Delta\|_{\mathbb{F}}}.$$

Since $\Delta \in C(\delta_{\mathbb{F}})$ and $C(\delta_{\mathbb{F}})$ is star-shaped, we have $\tilde{\Delta} \in C(\Delta_{\mathbb{F}})$ and $\|\tilde{\Delta}\|_{\mathbb{F}} = \frac{\kappa \Delta_{\mathbb{F}}}{\|\Delta\|_{\mathbb{F}}}$ by construction. Consequently, it is sufficient to prove that $E(\delta_{\mathbb{F}})$ holds with high probability. \square

Lemma 24. *Suppose that Assumptions 1 and 2 hold. Assume that for any $c > 0$, there exists an n such that $\sqrt{s}\lambda < c$. Then there exists a $\tilde{c} > 0$ such that*

$$\mathbb{P}(E(\delta_{\mathbb{F}})) \geq 1 - \exp(-2\tilde{c}).$$

Proof. We define the random variable

$$Z_n(C(\delta_{\mathbb{F}})) = \sup_{\Delta \in C(\delta_{\mathbb{F}})} \left\{ \frac{\kappa^4 \delta_{\mathbb{F}}^2}{4} - \frac{1}{n} \sum_{i=1}^n \|\langle \Delta, \mathbf{x}^{(i)} \rangle\|_{\mathbb{F}}^2 \right\},$$

then it suffices to show that

$$Z_n(C(\delta_{\mathbb{F}})) \leq \frac{\kappa^4 \delta_{\mathbb{F}}^2}{2}.$$

With a slight abuse of notation, it follows that

$$\|\Delta\|_{\mathbb{F}}^2 = \frac{1}{n} \sum_{i=1}^n \frac{d_N}{d_M} \sum_{m=1}^{d_N/d_M} \langle \tilde{\Delta}_m, \text{vec}(\mathbf{x}^{(i)}) \rangle^2,$$

where $\tilde{\Delta}_m \in \mathbb{R}^{d_M}$ and clearly $\text{vec}(\mathcal{X}^{(i)}) \in \mathbb{R}^{d_M}$. In order to complete the proof we make use of a truncation argument. For a constant $\tau > 0$ to be chosen later, consider the truncated function

$$\phi_\tau \left(\tilde{\Delta}_m, \text{vec}(\mathcal{X}) \right)$$

and define

$$\Delta_{m,\tau}(\mathcal{X}) = \text{sign} \left(\langle \tilde{\Delta}_m, \text{vec}(\mathcal{X}) \rangle \right) \sqrt{\phi_\tau \left(\langle \tilde{\Delta}_m, \text{vec}(\mathcal{X}) \rangle \right)}.$$

Further let

$$\Delta_r(\mathcal{X}) = (\Delta_{1,\tau}(\mathcal{X}), \Delta_{2,\tau}(\mathcal{X}), \dots, \Delta_{d_N/d_M,\tau}(\mathcal{X})).$$

The remainder of the proof consists of showing that for a suitable τ ,

$$\|\Delta_r\|_{\mathbb{F}}^2 \geq \frac{3}{4} \|\Delta\|_{\mathbb{F}}^2,$$

for all $\Delta \in C(\Delta_{\mathbb{F}})$ and

$$P \left\{ Z_n \geq \frac{\kappa^2 \delta_{\mathbb{F}}^2}{4} \right\} \leq c_1 \exp(-c_2 n \delta_{\mathbb{F}}^2),$$

where $Z_n := \sup_{\Delta \in C(\delta_{\mathbb{F}})} \|\Delta\|_{\mathbb{F}}^2 - \|\Delta_r\|_{\mathbb{F}}^2$. By definition

$$\|\Delta_m\|_{\mathbb{F}}^2 - \|\Delta_r\|_{\mathbb{F}}^2 \leq \mathbb{E} \left[\left(\langle \tilde{\Delta}_m, \text{vec}(\mathcal{X}) \rangle \right)^2 \mathbf{1} \left[\langle \tilde{\Delta}_m, \text{vec}(\mathcal{X}) \rangle \geq \tau \right] \right] \leq \frac{\mathbb{E}[\langle \tilde{\Delta}_m, \text{vec}(\mathcal{X}) \rangle]^4}{\tau^2},$$

where the last inequality follows from the Cauchy-Schwarz inequality and the final inequality follows from Markov's inequality. Since $\langle \tilde{\Delta}_m, \text{vec}(\mathcal{X}) \rangle$ is a Gaussian random variable,

$$E \left[\langle \tilde{\Delta}_m, \text{vec}(\mathcal{X}) \rangle \right] \leq 3E[\|\tilde{\Delta}_m\|_{\mathbb{F}}].$$

Therefore $\|\Delta\|_{\mathbb{F}}^2 \geq \frac{3}{4} \|\Delta\|_{\mathbb{F}}^2$.

From [Raskutti et al. \(2019\)](#), we see the high probability bound on Z_n by first upper bounding $E[Z_n]$. Thus, we obtain

$$\left\| \hat{\mathcal{A}} - \mathcal{A}^* \right\|_{\mathbb{F}} \leq 3\lambda\kappa\Psi(\mathcal{A}),$$

This completes the proof. □