

Half-Quadratic Alternating Direction Method of Multipliers for Robust Orthogonal Tensor Approximation

Yuning Yang*

Yunlong Feng[†]

May 12, 2020

Abstract

Higher-order tensor canonical polyadic decomposition (CPD) with one or more of the latent factor matrices being columnwisely orthonormal has been well studied in recent years. However, most existing models penalize the noises, if occurring, by employing the least squares loss, which may be sensitive to non-Gaussian noise or outliers, leading to bias estimates of the latent factors. In this paper, based on the maximum a posterior estimation, we derive a robust orthogonal tensor CPD model with Cauchy loss, which is resistant to heavy-tailed noise or outliers. By exploring the half-quadratic property of the model, a new method, which is termed as half-quadratic alternating direction method of multipliers (HQ-ADMM), is proposed to solve the model. Each subproblem involved in HQ-ADMM admits a closed-form solution. Thanks to some nice properties of the Cauchy loss, we show that the whole sequence generated by the algorithm globally converges to a stationary point of the problem under consideration. Numerical experiments on synthetic and real data demonstrate the efficiency and robustness of the proposed model and algorithm.

Key words: Tensor, canonical polyadic decomposition, robust, Cauchy, HQ-ADMM

1 Introduction

A tensor is a multidimensional array. Owing to its ability to represent data with intrinsically many dimensions, tensors draw much attention from the communities of signal processing, image processing, machine learning, etc; see the surveys [7, 28, 40]. To understand the relationship behind the data tensor, decomposition tools are needed. In general, tensor decomposition aims at factorizing the data tensor into a set of lower-dimensional latent factors, where the factors can be vectors, matrices or even tensors. Among the decomposition models, tensor canonical polyadic decomposition (CPD), which factorizes a tensor into a sum of component rank-1 tensors, is one of the most important models. Tensor CPD finds applications in blind multiuser CDMA, blind source separation, and so on [40]. Different from matrix decompositions, tensor CPD is unique under quite mild conditions [28].

In some applications, one or more latent factors of the CPD are required to have orthonormal columns. For example, in linear image coding [39], one is given a set of data matrices of the same size; to explore their commonalities, one projects the matrices onto a latent lower-dimensional subspace in which the subspace can be represented by the Khatri-Rao product [28] of two columnwisely orthonormal matrices. Such a problem

*College of Mathematics and Information Science, Guangxi University, Nanning, 530004, China (yyang@gxu.edu.cn).

[†]Department of Mathematics and Statistics, State University of New York at Albany, Albany, New York 12222, USA (yifeng@albany.edu).

has been formulated as a third-order tensor CPD with two factor matrices having orthonormal columns. On the other hand, simultaneous foreground-background extraction and compression can also be formulated as a model of the same kind; this will be illustrated in Sect. 5. Other applications of CPD with orthonormal factors can be found in [8–10, 41, 44].

In reality, due to the NP-hardness of determining the tensor rank [20], and due to the presence of noise, tensor CPD model with orthonormal factors is rarely exact, and it is necessary to resort to an approximation scheme. To numerically solve the problem, one usually formulates it as an optimization problem that minimizes the Euclidean distance between the data tensor and the latent tensor over orthonormal constraints, and then applies an alternating optimization type method to solve it based on polar decomposition [5, 18, 24, 31, 35, 43, 46, 48]. Other types of methods can be found in [11, 26, 30, 37]; just to name a few.

Although the optimization model mentioned above is effective in some circumstances, note that the Euclidean distance, built upon the least squares loss that is not robust [25]. As a result, when the data tensor is contaminated by heavy-tailed noise or outliers, such least squares based models often lead to bias estimates of the true latent factors, as having been observed in practice. This drawback of the least squares based models motivates us to develop a new model that is robust to heavy-tailed noise or outliers.

In this work, from the maximum a posterior estimation, we derive a robust tensor CPD model where one or more latent factors have orthonormal columns. Such a model is based on the Cauchy loss, whose robustness comes from the redescending property of the loss function, as pointed out in robust statistics [25]. We then explore the half-quadratic property of the model, based on which, the half-quadratic alternating direction method of multipliers (HQ-ADMM) is proposed to solve the model. An advantage of HQ-ADMM is that every subproblem involved in the algorithm admits a closed-form solution. Under a very mild assumption on the parameter, HQ-ADMM is proved to globally converge to a stationary point of the problem under consideration, owing to some nice properties of the Cauchy loss. In fact, the spirit of HQ-ADMM can be extended to solving other Cauchy loss based machine learning and scientific computing problems (besides tensor problems), which will be remarked later in Sect. 3. Finally, we show via numerical experiments that the proposed model is resistant to heavy-tailed noise such as Cauchy noise, outliers, and also performs well with Gaussian noise; the proposed HQ-ADMM is observed to be efficient.

The rest of the paper is organized as follows. The robust tensor approximation model is formulated in Sect. 2, with some quantitative properties given. The HQ-ADMM is developed in Sect. 3; the convergence analysis of HQ-ADMM is provided in Sect. 4. Numerical results are illustrated in Sect. 5. We end this paper in Sect. 6 with conclusions.

2 Problem Formulation and the Optimization Model

Notations Vectors are written as boldface lowercase letters ($\mathbf{x}, \mathbf{y}, \dots$), matrices are denoted as italic capitals (A, B, \dots), and tensors are written as calligraphic capitals ($\mathcal{A}, \mathcal{B}, \dots$). \mathbb{R} denotes the real field. $\mathbb{R}^{m \times n}$ denotes real matrices of dimension $m \times n$ and $\mathbb{R}^{n_1 \times \dots \times n_d}$ denotes tensor space of size $n_1 \times \dots \times n_d$. The Frobenius norm, $\|\cdot\|_F$, of a matrix or a tensor, is defined to be the square root of the sum of squares of all the entries. The inner product $\langle \cdot, \cdot \rangle$ between a pair of matrices or tensors of the same size is given by the sum of entrywise product. \otimes denotes the outer product of two vectors. Other notations will be introduced whenever necessary.

Let $\mathcal{A} = (\mathcal{A}_{i_1 \dots i_d}) \in \mathbb{R}^{n_1 \times \dots \times n_d}$ be a d -th order observed data tensor. We consider the inexact CPD of \mathcal{A} , i.e., approximating \mathcal{A} by a sum of rank-1 tensors:

$$\mathcal{A} = \sum_{i=1}^R \sigma_i \bigotimes_{j=1}^d \mathbf{u}_{j,i} + \mathcal{N} \in \mathbb{R}^{n_1 \times \dots \times n_d}; \quad (2.1)$$

here $\mathbf{u}_{j,i} \in \mathbb{R}^{n_j}$, $1 \leq j \leq d$, $\bigotimes_{j=1}^d \mathbf{u}_{j,i}$ denotes the rank-1 tensor given by the outer product of $\mathbf{u}_{j,i}$'s, σ_i 's are

real scalars, $R > 0$ is a given integer, where usually R is such that $R \leq \min\{n_1, \dots, n_d\}$ for a possibly low-rank approximation, while \mathcal{N} denotes the noisy tensor.

Denote $U_j := [\mathbf{u}_{j,1}, \dots, \mathbf{u}_{j,R}] \in \mathbb{R}^{n_j \times R}$ and $\boldsymbol{\sigma} := [\sigma_1, \dots, \sigma_R] \in \mathbb{R}^R$. Then U_j 's are called the latent factor matrices of \mathcal{A} . Throughout this work, we follow [28] to write the sum of rank-1 terms as

$$\llbracket \boldsymbol{\sigma}; U_1, \dots, U_d \rrbracket := \sum_{i=1}^R \sigma_i \bigotimes_{j=1}^d \mathbf{u}_{j,i};$$

moreover, we write $\llbracket \boldsymbol{\sigma}; U_j \rrbracket := \llbracket \boldsymbol{\sigma}; U_1, \dots, U_d \rrbracket$ for short. In the sequel, we base our work on the following setup:

- One or more U_j 's are columnwise orthonormal. Without loss of generality, we assume that the last t ($1 \leq t \leq d$) matrices are columnwise orthonormal, i.e.,

$$U_j^\top U_j = I, \quad d-t+1 \leq j \leq d,$$

where I is an identity matrix of the proper size;

- The columns of the first $d-t$ matrices are normalized, i.e.,

$$\|\mathbf{u}_{j,i}\| = 1, \quad 1 \leq j \leq d-t, 1 \leq i \leq R;$$

- Entries of the noisy tensor \mathcal{N} are i.i.d..

We immediately have the following proposition.

Proposition 2.1. *There holds $\left\| \bigotimes_{j=1}^d \mathbf{u}_{j,i} \right\|_F = 1$, $1 \leq i \leq R$, and $\left\langle \bigotimes_{j=1}^d \mathbf{u}_{j,i_1}, \bigotimes_{j=1}^d \mathbf{u}_{j,i_2} \right\rangle = 0$, $i_1 \neq i_2$.*

Note that the constraints on $\mathbf{u}_{j,i}$ and U_j are all Stiefel manifolds $\text{st}(m, n) := \{P \in \mathbb{R}^{m \times n} \mid P^\top P = I\}$. Therefore, in the following, we write the constraints on $\mathbf{u}_{j,i}$ and U_j as

$$\begin{aligned} \mathbf{u}_{j,i} &\in \text{st}(n_j, 1), \quad 1 \leq j \leq d-t, 1 \leq i \leq R, \\ U_j &\in \text{st}(n_j, R), \quad d-t+1 \leq j \leq R. \end{aligned}$$

In the presence of the noisy term \mathcal{N} , it is natural to deal with (2.1) via solving the following optimization problem [5, 18, 43, 46, 48]:

$$\min_{\boldsymbol{\sigma}, \mathbf{u}_{j,i} \in \text{st}(n_j, 1), U_j \in \text{st}(n_j, R)} \|\mathcal{A} - \llbracket \boldsymbol{\sigma}; U_j \rrbracket\|_F^2 = \sum_{i_1=1, \dots, i_d=1}^{n_1, \dots, n_d} \left(\mathcal{A}_{i_1 \dots i_d} - \llbracket \boldsymbol{\sigma}; U_j \rrbracket_{i_1 \dots i_d} \right)^2. \quad (2.2)$$

From a statistical estimation viewpoint, the above model is built upon the least squares loss $\ell_2(t) := t^2/2$, i.e., it employs the $\ell_2(\cdot)$ loss to deal with noise. However, it is commonly known that the estimators induced by the least squares loss are sensitive to heavy-tailed noise or outliers; in other words, by using the model (2.2), one assumes that every entry of \mathcal{N} obeys the standard Gaussian distribution by default.

Derivation of our model In real-world applications, data may be contaminated by heavy-tailed noise, and even outliers/impulsive noise. A typical non-Gaussian and heavy-tailed noise is the Cauchy noise, whose probability density function is given by

$$P_{\text{Cauchy}}(t) \propto \frac{1}{1 + (t-c)^2/\delta^2},$$

where $\delta > 0$ is the scale parameter and c is the location parameter. By assuming the symmetry of the noise, we let $c = 0$ in the above function.

We derive our model from the maximum a posterior (MAP) estimation by assuming that \mathcal{N} obeys the Cauchy distribution whose density function is given above. To this end, denote respectively the indicator function $\mathbf{1}_C(\cdot)$ and the characteristic function $\iota_C(\cdot)$ of a closed set C as follows

$$\begin{aligned}\mathbf{1}_C(\mathbf{x}) &= 1, \text{ if } \mathbf{x} \in C; \mathbf{1}_C(\mathbf{x}) = 0, \text{ if } \mathbf{x} \notin C, \\ \iota_C(\mathbf{x}) &= 0, \text{ if } \mathbf{x} \in C; \iota_C(\mathbf{x}) = +\infty, \text{ if } \mathbf{x} \notin C.\end{aligned}$$

From the constraints on $\mathbf{u}_{j,i}$ and U_j , it is natural to impose a uniform prior belief distributional assumption on $\{\mathbf{u}_{j,i}, U_j\}$ as follows

$$P(\llbracket \boldsymbol{\sigma}; U_j \rrbracket) \propto \prod_{j=1}^{d-t} \prod_{i=1}^R \mathbf{1}_{\text{st}(n_j,1)}(\mathbf{u}_{j,i}) \cdot \prod_{j=d-t+1}^d \mathbf{1}_{\text{st}(n_j,R)}(U_j). \quad (2.3)$$

On the other hand, in the presence of Cauchy noise, the probability of the observed data tensor \mathcal{A} conditioned on $\llbracket \boldsymbol{\sigma}; U_j \rrbracket$ is given by

$$P(\mathcal{A}_{i_1 \dots i_d} \mid \llbracket \boldsymbol{\sigma}; U_j \rrbracket_{i_1 \dots i_d}) \propto \frac{1}{1 + \left(\llbracket \boldsymbol{\sigma}; U_j \rrbracket_{i_1 \dots i_d} - \mathcal{A}_{i_1 \dots i_d} \right)^2 / \delta^2}, \quad 1 \leq i_j \leq n_j, \quad 1 \leq j \leq d. \quad (2.4)$$

With (2.3) and (2.4) at hand, using Bayes's rule, the MAP estimation is given by

$$\begin{aligned}\{\boldsymbol{\sigma}^*, U_j^*\} &= \arg \max P(\llbracket \boldsymbol{\sigma}; U_j \rrbracket \mid \mathcal{A}) \\ &= \arg \max \frac{P(\mathcal{A} \mid \llbracket \boldsymbol{\sigma}; U_j \rrbracket) \cdot P(\llbracket \boldsymbol{\sigma}; U_j \rrbracket)}{P(\mathcal{A})} \\ &= \arg \max \prod_{i_1=1, \dots, i_d=1}^{n_1, \dots, n_d} P(\mathcal{A}_{i_1 \dots i_d} \mid \llbracket \boldsymbol{\sigma}; U_j \rrbracket_{i_1 \dots i_d}) \cdot P(\llbracket \boldsymbol{\sigma}; U_j \rrbracket) \\ &\stackrel{t \leftarrow -\log(t)}{=} \arg \min \sum_{i_1=1, \dots, i_d=1}^{n_1, \dots, n_d} \log \left(1 + \left(\llbracket \boldsymbol{\sigma}; U_j \rrbracket_{i_1 \dots i_d} - \mathcal{A}_{i_1 \dots i_d} \right)^2 / \delta^2 \right) \\ &\quad - \sum_{j=1}^{d-t} \sum_{i=1}^R \log(\mathbf{1}_{\text{st}(n_j,1)}(\mathbf{u}_{j,i})) - \sum_{j=d-t+1}^d \log(\mathbf{1}_{\text{st}(n_j,R)}(U_j)) \\ &= \arg \min \sum_{i_1=1, \dots, i_d=1}^{n_1, \dots, n_d} \log \left(1 + \left(\llbracket \boldsymbol{\sigma}; U_j \rrbracket_{i_1 \dots i_d} - \mathcal{A}_{i_1 \dots i_d} \right)^2 / \delta^2 \right) \\ &\quad + \sum_{j=1}^{d-t} \sum_{i=1}^R \iota_{\text{st}(n_j,1)}(\mathbf{u}_{j,i}) + \sum_{j=d-t+1}^d \iota_{\text{st}(n_j,R)}(U_j),\end{aligned}$$

where in the last equality, we have defined $\log(0) = -\infty$. Therefore, from the above deduction, to deal with (2.1) in the presence of Cauchy noise (or even other heavy-tailed noise or outliers), we prefer to solve the following optimization model

$$\begin{aligned}\min \quad \Phi_\delta(\mathcal{A} - \llbracket \boldsymbol{\sigma}; U_j \rrbracket) &:= \frac{\delta^2}{2} \sum_{i_1=1, \dots, i_d=1}^{n_1, \dots, n_d} \log \left(1 + \left(\llbracket \boldsymbol{\sigma}; U_j \rrbracket_{i_1 \dots i_d} - \mathcal{A}_{i_1 \dots i_d} \right)^2 / \delta^2 \right) \\ \text{s.t. } \mathbf{u}_{j,i} &\in \text{st}(n_j, 1), \quad 1 \leq j \leq d-t, \quad 1 \leq i \leq R, \\ U_j &\in \text{st}(n_j, R), \quad d-t+1 \leq j \leq d.\end{aligned} \quad (2.5)$$

Comparing (2.5) with (2.2), we see that the difference is that the least squares loss $\ell_2(t) = t^2/2$ is replaced by the statistically motivated loss function

$$\phi_\delta(t) := \frac{\delta^2}{2} \log(1 + t^2/\delta^2). \quad (2.6)$$

$\phi_\delta(\cdot)$ is called the Cauchy loss. In recent years, various research has been focused on Cauchy loss based models; see, e.g., [12, 17, 19, 27, 32, 34, 38, 49].

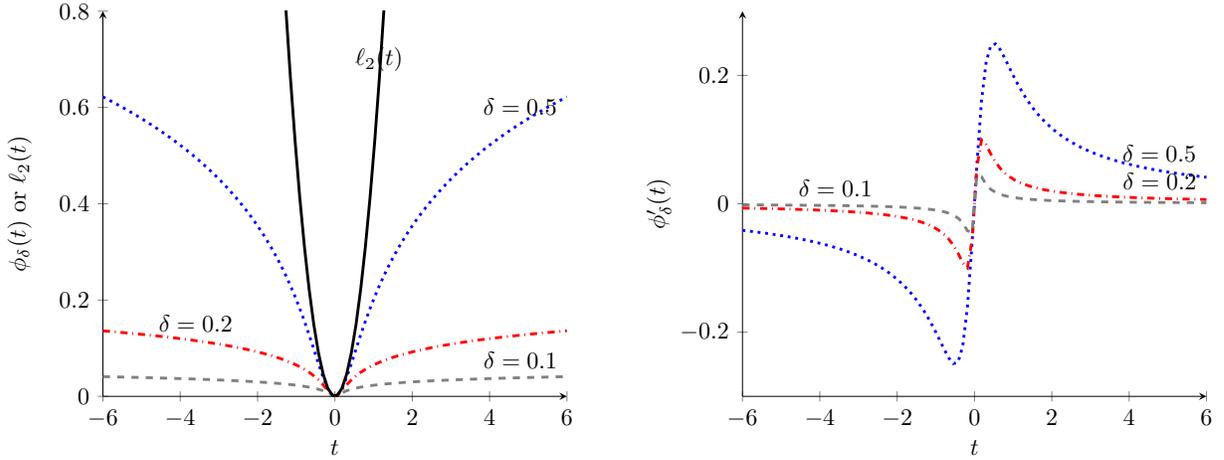


Figure 1: Left: Plots of $\phi_\delta(t) = \frac{\delta^2}{2} \log(1 + t^2/\delta^2)$ with different δ values versus $\ell_2(t) = t^2/2$; Right: Plots of $\phi'_\delta(t) = t/(1 + t^2/\delta^2)$. $\sigma = 0.1$ (the dashed curve), $\sigma = 0.2$ (the dotted-dashed curve), and $\sigma = 0.5$ (the dotted curve); $\ell_2(t)$ (the solid curve).

We discuss some properties of the proposed model (2.5) from the robust statistics viewpoint, which shows (2.5) is not only resistant to Cauchy noise, but may also be resistant to other heavy-tailed noise or outliers. Firstly, we observe that

$$\lim_{|t| \rightarrow +\infty} \phi'_\delta(t) = \lim_{|t| \rightarrow +\infty} \frac{t}{1 + t^2/\delta^2} = 0. \quad (2.7)$$

Such a property is called the redescending property in robust statistics [25], and the minimizer of (2.5) is called a redescending M-estimator. It is known that the redescending M-estimator is robust to heavy-tailed noise and outliers [25]. As a comparison, the derivative of the least squares loss $\ell_2(t) = t^2/2$ is t , whose limit is infinity, which does not have the redescending property. Other loss functions admitting the redescending property include the Welsch loss [13, 14, 21], the Tukey loss [3], the German loss [15], and so on.

Secondly, the parameter δ in (2.6) controls the robustness of the model (2.5). From (2.7), we see that the smaller δ is, the faster $\phi'_\delta(t)$ converges to zero. We plot $\phi'_\delta(t)$ with different δ in the right panel of Fig. 1. On the other hand, taking Taylor expansion of $\phi_\delta(t)$ at 0 yields $\phi_\delta(t) = t^2/2 + o(t^2/\delta^2)$, which shows that $\phi_\delta(t) \approx t^2/2$ as $\delta \rightarrow \infty$. These observations imply that a small δ can enhance the robustness of (2.5). This also reminds us that our model (2.5) is also resistant to Gaussian noise by simply setting a large enough δ . We also plot $\phi_\delta(t)$ with different δ in the left panel of Fig. 1.

Remark 2.1. We discuss several differences between our model (2.5) and some existing robust tensor models. In recent years, robust techniques have been incorporated into tensor decomposition/approximation/recovery/completion/PCA problems, where the L_1 loss function, namely, $\ell_1(t) = |t|$, is frequently employed to deliver robustness. In general, such kind of models can be formulated as [16]

$$\min_{\mathcal{X} \in \mathbb{R}^{n_1 \times \dots \times n_d}} \|\mathbf{L}(\mathcal{X}) - \mathbf{b}\|_1 + \lambda R(\mathcal{X}), \quad (2.8)$$

where \mathbf{L} is a linear operator, and \mathbf{b} has the same size as $\mathbf{L}(\mathcal{X})$; $R(\mathcal{X})$ denotes a certain regularizer that controls the low-rankness of \mathcal{X} , such as the sum of nuclear norms of unfolding matrices of \mathcal{X} [42], and $\lambda > 0$ is the regularization parameter. A special case of (2.8) is the robust tensor PCA, in which \mathbf{L} is the identity operator and \mathbf{b} denotes the observed tensor [16]. It is known that L_1 loss is more suitable for Laplacian noise; on the other hand, one sees that the derivative of $|t|$ does not tend to zero as $|t| \rightarrow +\infty$, meaning that it does not admit the redescending property, while it was pointed out in [33] that the L_1 estimator might behave as bad as the $\ell_2(t)$ estimator in some cases. Comparing with the resulting tensor, (2.8) yields a full tensor of size $n_1 \times \dots \times n_d$, while

ours is compressed into a set of factor matrices, which takes much less storage. Moreover, our orthonormality assumption on some factor matrices is more suitable for certain applications [8–10, 39, 41, 44].

In [1], a robust tensor CP decomposition model has been considered. The differences are that the noise there are required to be sparse, and all the factor matrices are assumed to be columnwisely orthogonal, which are stringent. By using outlier detection techniques, [36] proposed a robust Tucker model. However, the underlying model cannot be clearly formulated as an optimization problem, and the tensor model is different from ours. By using variational inference and Kullback-Leibler divergence, [6] devised a robust algorithm to find CP approximation with orthonormal factors, where the model and the solution method are quite different from ours. In particular, the authors pointed out that their algorithm boils down to the alternating least squares [43] in the absence of outliers. In a recent survey [22], various statistically motivated loss functions are incorporated into tensor CPD, in which the Huber’s loss is considered. As Huber’s loss can be regarded as a smoothed ℓ_1 loss, it does not admit the redescending property as well. The orthonormality is not taken into account in [22]. Note that the idea of employing Cauchy loss has been considered in the authors’ earlier work [49]. Comparing with (2.5), the resulting tensor in [49] is a full tensor and also does not take into account the orthonormality, and the solution method is also different.

The remaining problem is how to solve (2.5) efficiently. For this purpose, several quantitative properties concerning the Cauchy loss for designing and analyzing the solution method are first introduced in the following subsection.

2.1 Quantitative properties concerning $\phi_\delta(\cdot)$

First, we introduce the so-called half-quadratic (HQ) property of $\phi_\delta(\cdot)$, which turns the function into a weighted least squares problem and is crucial for designing the algorithm. Such a property of the Cauchy loss has appeared in the literature; see, e.g., [17, 19], in which the verification is based on the utilization of conjugate functions. While we present a very direct and concise proof. Recall that we have defined $\log(0) = -\infty$.

Lemma 2.1 (Half-quadratic property). *Given $|t| < +\infty$, it holds that*

$$\phi_\delta(t) = \min_{\omega \geq 0} \frac{\omega}{2} t^2 + \frac{\delta^2}{2} \varrho(\omega), \quad (2.9)$$

where $\varrho(\omega) = \omega - \log(\omega) - 1$. Moreover, the minimizer of (2.9) is given by

$$\omega^* = \frac{\delta^2}{\delta^2 + t^2}. \quad (2.10)$$

Proof. First we verify that (2.10) is a minimizer of the right hand-side of (2.9). Denote $g(\omega) := \omega t^2/2 + \delta^2 \varrho(\omega)/2$. As $\varrho(\cdot)$ is convex, it suffices to show that ω^* in (2.10) is a stationary point of $\inf_{\omega \geq 0} g(\omega)$. Since $|t| < +\infty$, we see that the minimizer of $\inf_{\omega \geq 0} g(\omega)$ cannot occur at $\omega = 0$. Thus any stationry point of $\inf_{\omega \geq 0} g(\omega)$ meets

$$g'(\omega) = 0 \Leftrightarrow t^2 + \delta^2 - \frac{\delta^2}{\omega} = 0,$$

and so $\omega = (1 + t^2/\delta^2)^{-1}$, which is exactly (2.10). Inserting this expression into (2.9), we get

$$\begin{aligned} 2g(t) &= \frac{\delta^2}{\delta^2 + t^2} (t^2 + \delta^2) + \delta^2 \log(1 + t^2/\delta^2) - \delta^2 \\ &= \delta^2 \log(1 + t^2/\delta^2), \end{aligned}$$

boiling down to the expression of $\phi_\delta(t)$. The proof is completed. \square

Note that the HQ property has a very clear indication on robustness: Take $t = \mathcal{A}_{i_1 \dots i_d} - \llbracket \boldsymbol{\sigma}; U_j \rrbracket_{i_1 \dots i_d}$ in Lemma 2.1 as the noise; we see that the larger the magnitude of t , the smaller the weight ω it yields, and so the corresponding $\phi_\delta(t)$ is less important in the objective $\Phi_\delta(\cdot)$ in (2.5).

The next two properties are helpful for convergence analysis. Recalling that $\phi'_\delta(t) = \frac{\delta^2 t}{\delta^2 + t^2}$, we have

Proposition 2.2 (Lipschitz gradient). *For any $t_1, t_2 \in \mathbb{R}$ and $\delta > 0$, it holds that*

$$\left| \frac{\delta^2 t_1}{\delta^2 + t_1^2} - \frac{\delta^2 t_2}{\delta^2 + t_2^2} \right| \leq |t_1 - t_2|.$$

Proof. By the mean value theorem, It suffices to show that $|\phi''_\delta(t)| \leq 1$. In fact,

$$|\phi''_\delta(t)| = \left| \frac{\delta^2(\delta^2 - t^2)}{(\delta^2 + t^2)^2} \right| \leq \left| \frac{\delta^2}{\delta^2 + t^2} \right| \leq 1,$$

and the result follows. \square

Proposition 2.3 (Lipschitz-like inequality). *Let $t_1, t_2 \in \mathbb{R}$ be arbitrary, and let $\delta > 0$. Then it holds that*

$$|e| := \left| \delta^2 t_1 \left(\frac{1}{\delta^2 + t_1^2} - \frac{1}{\delta^2 + t_2^2} \right) \right| \leq |t_1 - t_2|.$$

Proof. It is clear that

$$|e| = \left| \sigma^2 t_1 \frac{(t_1 + t_2)(t_1 - t_2)}{(\sigma^2 + t_1^2)(\sigma^2 + t_2^2)} \right| \leq \sigma^2 |t_1| \frac{|t_1| + |t_2|}{(\sigma^2 + t_1^2)(\sigma^2 + t_2^2)} \cdot |t_1 - t_2|.$$

To prove the above relation, it suffices to show the coefficient of $|t_1 - t_2|$ is not greater than 1, i.e.,

$$\varphi(t_1, t_2) := (\sigma^2 + t_1^2)(\sigma^2 + t_2^2) - \sigma^2 |t_1| (|t_1| + |t_2|) \geq 0.$$

In fact,

$$\begin{aligned} \varphi(t_1, t_2) &= \sigma^4 + \sigma^2 t_2^2 + |t_1 t_2| (|t_1 t_2| - \sigma^2) \\ &\geq \sigma^4 + \sigma^2 t_2^2 - \frac{\sigma^4}{4} \geq 0. \end{aligned}$$

Therefore, $|e| \leq |t_1 - t_2|$, as desired. \square

3 HQ-ADMM

By using Lemma 2.1, we equivalently rewrite the objective function $\Phi_\delta(\cdot)$ of (2.5) in what follows. Specifically, since $\Phi_\delta(\cdot)$ is the sum of $\phi_\delta(\cdot)$ functions, taking $t = \mathcal{A}_{i_1 \dots i_d} - \llbracket \boldsymbol{\sigma}; U_j \rrbracket_{i_1 \dots i_d}$ in Lemma 2.1, we have

$$\begin{aligned} &\Phi_\delta(\mathcal{A} - \llbracket \boldsymbol{\sigma}; U_j \rrbracket) \\ &= \frac{1}{2} \min_{\mathcal{W}_{i_1 \dots i_d} \geq 0} \sum_{i_1=1, \dots, i_d}^{n_1, \dots, n_d} \left[\mathcal{W}_{i_1 \dots i_d} (\mathcal{A}_{i_1 \dots i_d} - \llbracket \boldsymbol{\sigma}; U_j \rrbracket_{i_1 \dots i_d})^2 + \delta^2 \varrho(\mathcal{W}_{i_1 \dots i_d}) \right], \end{aligned} \quad (3.11)$$

where we denote $\mathcal{W} = (\mathcal{W}_{i_1 \dots i_d}) \in \mathbb{R}^{n_1 \times \dots \times n_d}$ as a tensor variable. From Lemma 2.1, we see that the optimizer is $\mathcal{W}_{i_1 \dots i_d} = \delta^2 \left(1 + \left(\llbracket \boldsymbol{\sigma}; U_j \rrbracket_{i_1 \dots i_d} - \mathcal{A}_{i_1 \dots i_d} \right)^2 / \delta^2 \right)^{-1}$. As explained in the paragraph below Lemma 2.1, \mathcal{W} can be interpreted as weights to the problem. From the expression of \mathcal{W} , we see that the larger the noise is, the smaller the weight gives to the problem. Such a mechanism helps mitigate heavy-tailed noise or outliers.

In view of (3.11), a straightforward idea to solve (2.5) (with the objective replaced by (3.11)) is to employing an alternating minimization method (AM) by iteratively updating $\boldsymbol{\sigma}, U_j$, and \mathcal{W} . In fact, applying AM to

solve Cauchy loss-based problems have been considered in the literature; see, e.g., [17, 19]. However, for our problem, this would result in that the subproblems related to U_j do not have closed-form solutions. [32] also applied AM to solve Cauchy loss-based problem; however, as their proposed model is unconstrained and the objective function is smooth, AM yields closed-form solutions to each subproblem. [38] incorporated Cauchy loss into models for image processing. However, the problem is convexified by imposing a quadratic term, which results in that the Cauchy loss related subproblem admits a unique solution that can be analytically solved by solving a cubic equation. If the subproblem is nonconvex, then numerical methods have to be applied to solving the Cauchy loss related subproblem, as pointed out in [38], which might result in inefficiency. For other Cauchy loss based image processing problems, [12, 27, 34] proposed to use the conventional alternating direction method of multipliers (ADMM) directly. However, without noticing the HQ property, in the ADMM, solving the Cauchy loss related subproblem also does not admit a closed-form solution. As a result, solving such a subproblem still requires an iterative method. [49] used a linearization technique, which ignored the HQ property.

In view of the above limitations in dealing with Cauchy loss-based problems, in this section, by combining the HQ property and the ADMM framework, we proposed a new method, termed as HQ-ADMM, to solve our model (2.5). The advantage of HQ-ADMM is that all the subproblems involved in the algorithm admit closed-form solutions. In what follows, we derive our method step by step.

Note that (3.11) is quadratic with respect to each U_j , leading to the following formulation

$$\Phi_\delta(\mathcal{A} - \llbracket \boldsymbol{\sigma}; U_j \rrbracket) = \frac{1}{2} \min_{\mathcal{W}_{i_1 \dots i_d} \geq 0} \|\sqrt{\mathcal{W}} \otimes (\mathcal{A} - \llbracket \boldsymbol{\sigma}; U_j \rrbracket)\|_F^2 + \frac{\delta^2}{2} \sum_{i_1=1, \dots, i_d}^{n_1, \dots, n_d} \varrho(\mathcal{W}_{i_1 \dots i_d}),$$

where $\sqrt{\mathcal{W}} = (\sqrt{\mathcal{W}_{i_1 \dots i_d}}) \in \mathbb{R}^{n_1 \times \dots \times n_d}$ and ‘ \otimes ’ denotes the Hadamard product. With this expression at hand, by introducing a slack variable $\mathcal{T} \in \mathbb{R}^{n_1 \times \dots \times n_d}$, we rewrite (2.5) as

$$\begin{aligned} \min_{\boldsymbol{\sigma}, U_j, \mathcal{T}, \mathcal{W}} \Phi_\delta(\mathcal{A} - \mathcal{T}) &= \frac{1}{2} \|\sqrt{\mathcal{W}} \otimes (\mathcal{A} - \mathcal{T})\|_F^2 + \frac{\delta^2}{2} \sum_{i_1=1, \dots, i_d}^{n_1, \dots, n_d} \varrho(\mathcal{W}_{i_1 \dots i_d}) \\ \text{s.t. } \mathcal{T} &= \llbracket \boldsymbol{\sigma}; U_j \rrbracket, \mathcal{W} \geq 0, \\ \mathbf{u}_{j,i}^\top \mathbf{u}_{j,i} &= 1, 1 \leq j \leq d-t, 1 \leq i \leq R, \\ U_j^\top U_j &= I, d-t+1 \leq j \leq d. \end{aligned} \quad (3.12)$$

By introducing a Lagrangian multiplier $\mathcal{Y} \in \mathbb{R}^{n_1 \times \dots \times n_d}$, the augmented Lagrangian function of (3.12) is given by

$$\begin{aligned} L_\tau(\boldsymbol{\sigma}, U_j, \mathcal{T}, \mathcal{Y}, \mathcal{W}) &:= \frac{1}{2} \|\sqrt{\mathcal{W}} \otimes (\mathcal{A} - \mathcal{T})\|_F^2 + \frac{\delta^2}{2} \sum_{i_1=1, \dots, i_d}^{n_1, \dots, n_d} \varrho(\mathcal{W}_{i_1 \dots i_d}) \\ &\quad - \langle \mathcal{Y}, \llbracket \boldsymbol{\sigma}; U_j \rrbracket - \mathcal{T} \rangle + \frac{\tau}{2} \|\llbracket \boldsymbol{\sigma}; U_j \rrbracket - \mathcal{T}\|_F^2, \end{aligned} \quad (3.13)$$

where $\tau > 0$. In what follows, for notational convenience we denote $(\mathcal{Y} + \tau\mathcal{T}) \otimes_{l \neq j}^d \mathbf{u}_{l,i} \in \mathbb{R}^{n_j}$ as the gradient of $\langle \mathcal{Y} + \tau\mathcal{T}, \otimes_{l=1}^d \mathbf{u}_{l,i} \rangle$ with respect to $\mathbf{u}_{j,i}$. Then, the last two terms of (3.13) can be rewritten as

$$\begin{aligned} - \langle \mathcal{Y}, \llbracket \boldsymbol{\sigma}; U_j \rrbracket - \mathcal{T} \rangle + \frac{\tau}{2} \|\llbracket \boldsymbol{\sigma}; U_j \rrbracket - \mathcal{T}\|_F^2 &= \langle \mathcal{Y}, \mathcal{T} \rangle + \frac{\tau}{2} \|\mathcal{T}\|_F^2 - \langle \mathcal{Y} + \tau\mathcal{T}, \llbracket \boldsymbol{\sigma}; U_j \rrbracket \rangle + \frac{\tau}{2} \boldsymbol{\sigma}^\top \boldsymbol{\sigma} \\ &= \langle \mathcal{Y}, \mathcal{T} \rangle + \frac{\tau}{2} \|\mathcal{T}\|_F^2 - \left\langle \mathcal{Y} + \tau\mathcal{T}, \sum_{i=1}^R \sigma_i \otimes_{j=1}^d \mathbf{u}_{j,i} \right\rangle + \frac{\tau}{2} \boldsymbol{\sigma}^\top \boldsymbol{\sigma} \\ &= \langle \mathcal{Y}, \mathcal{T} \rangle + \frac{\tau}{2} \|\mathcal{T}\|_F^2 - \sum_{i=1}^R \sigma_i \left\langle (\mathcal{Y} + \tau\mathcal{T}) \otimes_{l \neq j}^d \mathbf{u}_{l,i}, \mathbf{u}_{j,i} \right\rangle + \frac{\tau}{2} \boldsymbol{\sigma}^\top \boldsymbol{\sigma}, \end{aligned} \quad (3.14)$$

where the first equality is due to Proposition 2.1.

Before presenting the algorithm, we first derive the stationary point system. To this end, we further define Lagrangian multipliers $\eta_{j,i} \in \mathbb{R}$, $1 \leq j \leq d-t$, $1 \leq i \leq R$ attached to the constraints $\mathbf{u}_{j,i}^\top \mathbf{u}_{j,i} = 1$, and $\Lambda_j \in \mathbb{R}^{R \times R}$, $d-t+1 \leq j \leq d$ attached to $U_j^\top U_j = I$, where Λ_j 's are symmetric matrices. Denote

$$\begin{aligned} \hat{L}_\tau(\boldsymbol{\sigma}, U_j, \mathcal{T}, \mathcal{Y}, \mathcal{W}) &:= L_\tau(\boldsymbol{\sigma}, U_1, \dots, U_d, \mathcal{T}, \mathcal{Y}, \mathcal{W}) \\ &+ \sum_{j,i=1}^{d-t,R} \eta_{j,i} (\mathbf{u}_{j,i}^\top \mathbf{u}_{j,i} - 1) + \sum_{j=d-t+1}^d \langle \Lambda_j, U_j^\top U_j - I \rangle. \end{aligned} \quad (3.15)$$

Thus taking derivative of $\hat{L}(\cdot)$ with respect to each $\mathbf{u}_{j,i}$, $1 \leq j \leq d-t$, $1 \leq i \leq R$ and noticing (3.14) yields

$$\sigma_i(\mathcal{Y} + \tau\mathcal{T}) \bigotimes_{l \neq j}^d \mathbf{u}_{l,i} = \eta_{j,i} \mathbf{u}_{j,i}, \quad 1 \leq j \leq d-t, \quad 1 \leq i \leq R. \quad (3.16)$$

Since $\mathbf{u}_{j,i}$'s are normalized, we get $\eta_{j,i} = \sigma_i \langle \mathcal{Y} + \tau\mathcal{T}, \bigotimes_{l=1}^d \mathbf{u}_{l,i} \rangle$. On the other hand, noticing the representation (3.14), taking derivative of $\hat{L}(\cdot)$ with respect to $\boldsymbol{\sigma}$ gives that $\sigma_i = \langle \mathcal{Y} + \tau\mathcal{T}, \bigotimes_{l=1}^d \mathbf{u}_{l,i} \rangle / \tau$, which together with the expression of $\eta_{j,i}$ gives $\eta_{j,i} = \sigma_i^2 \tau$; therefore, (3.16) is in fact as follow

$$(\mathcal{Y} + \tau\mathcal{T}) \bigotimes_{l \neq j}^d \mathbf{u}_{l,i} = \sigma_i \tau \mathbf{u}_{j,i}, \quad 1 \leq j \leq d-t, \quad 1 \leq i \leq R. \quad (3.17)$$

Next, taking derivative with respect to $\mathbf{u}_{j,i}$, $d-t+1 \leq j \leq d$, $1 \leq i \leq R$ and noticing (3.14) gives

$$\sigma_i(\mathcal{Y} + \tau\mathcal{T}) \bigotimes_{l \neq j}^d \mathbf{u}_{l,i} = \sum_{r=1}^R (\Lambda_j)_{i,r} \mathbf{u}_{j,r}, \quad 1 \leq j \leq d-t, \quad 1 \leq i \leq R. \quad (3.18)$$

Denote $\mathcal{E} \in \mathbb{R}^{n_1 \times \dots \times n_d}$ as the all-one tensor; taking derivative with respect to \mathcal{T} and rearranging terms yields

$$\begin{aligned} \mathcal{W} \otimes (\mathcal{T} - \mathcal{A}) + \mathcal{Y} - \tau([\boldsymbol{\sigma}; U_j] - \mathcal{T}) &= 0 \\ \Leftrightarrow (\mathcal{W} + \tau\mathcal{E}) \otimes \mathcal{T} &= \mathcal{W} \otimes \mathcal{A} - \mathcal{Y} + \tau[\boldsymbol{\sigma}; U_j]. \end{aligned} \quad (3.19)$$

As a result, taking (3.17), (3.18), (3.19) and Lemma 2.1 into account, any stationary point $\{\boldsymbol{\sigma}, U_j, \mathcal{T}, \mathcal{Y}, \mathcal{W}\}$ satisfies the following system

$$\begin{cases} (\mathcal{Y} + \tau\mathcal{T}) \bigotimes_{l \neq j}^d \mathbf{u}_{l,i} = \sigma_i \tau \mathbf{u}_{j,i}, & 1 \leq j \leq d-t, \quad 1 \leq i \leq R, \\ \mathbf{u}_{j,i}^\top \mathbf{u}_{j,i} = 1, & 1 \leq j \leq d-t, \quad 1 \leq i \leq R, \\ \sigma_i(\mathcal{Y} + \tau\mathcal{T}) \bigotimes_{l \neq j}^d \mathbf{u}_{l,i} = \sum_{r=1}^R (\Lambda_j)_{i,r} \mathbf{u}_{j,r}, & 1 \leq j \leq d-t, \quad 1 \leq i \leq R, \\ U_j^\top U_j = I, & d-t+1 \leq j \leq d, \\ (\mathcal{W} + \tau\mathcal{E}) \otimes \mathcal{T} = \mathcal{W} \otimes \mathcal{A} - \mathcal{Y} + \tau[\boldsymbol{\sigma}; U_j], \\ [\boldsymbol{\sigma}; U_j] = \mathcal{T}, \\ \mathcal{W}_{i_1 \dots i_d} = \delta^2 \left(\delta^2 + (\mathcal{T}_{i_1 \dots i_d} - \mathcal{A}_{i_1 \dots i_d})^2 \right)^{-1}. \end{cases} \quad (3.20)$$

HQ-ADMM framework Combining the HQ property and the ADMM, our HQ-ADMM computes the following subproblems at each iterate

$$\begin{cases} U_j^{k+1} \in \arg \min_{\|\mathbf{u}_{j,i}\|=1, 1 \leq i \leq R} L_\tau(\boldsymbol{\sigma}^k, U_1^{k+1}, \dots, U_{j-1}^{k+1}, U_j, U_{j+1}^k, \dots, U_d^k, \mathcal{T}^k, \mathcal{Y}^k, \mathcal{W}^k), & 1 \leq j \leq d-t, \\ U_j^{k+1} \in \arg \min_{U_j^\top U_j = I} L_\tau(\boldsymbol{\sigma}^k, U_1^{k+1}, \dots, U_{j-1}^{k+1}, U_j, U_{j+1}^k, \dots, U_d^k, \mathcal{T}^k, \mathcal{Y}^k, \mathcal{W}^k), & d-t+1 \leq j \leq d, \\ \mathcal{T}^{k+1} = \arg \min_{\mathcal{T}} L_\tau(\boldsymbol{\sigma}^k, U_j^{k+1}, \mathcal{T}, \mathcal{Y}^k, \mathcal{W}^k), \\ \mathcal{Y}^{k+1} = \mathcal{Y}^k - \tau([\boldsymbol{\sigma}^k; U_j^{k+1}] - \mathcal{T}^{k+1}), \\ \boldsymbol{\sigma}^{k+1} = \arg \min_{\boldsymbol{\sigma}} L_\tau(\boldsymbol{\sigma}, U_j^{k+1}, \mathcal{T}^{k+1}, \mathcal{Y}^{k+1}, \mathcal{W}^k), \\ \mathcal{W}^{k+1} = \arg \min_{\mathcal{W}} L_\tau(\boldsymbol{\omega}^{k+1}, U_j^{k+1}, \mathcal{T}^{k+1}, \mathcal{Y}^{k+1}, \mathcal{W}). \end{cases}$$

Comparing with the standard ADMM framework, HQ-ADMM involves an additional subproblem to update the weights \mathcal{W} . In what follows, we present how to solve each subproblem.

U_j -subproblems For notational convenience, let

$$\mathbf{v}_{j,i}^{k+1} := (\mathcal{Y}^k + \tau\mathcal{T}^k) \mathbf{u}_{1,i}^{k+1} \otimes \dots \otimes \mathbf{u}_{j-1,i}^{k+1} \otimes \mathbf{u}_{j+1,i}^k \otimes \dots \otimes \mathbf{u}_{d,i}^k$$

represent the gradient of $\langle \mathcal{Y}^k + \tau \mathcal{T}^k, \bigotimes_{l=1}^d \mathbf{u}_{l,i} \rangle$ with respect to $\mathbf{u}_{j,i}$ at the point $(\mathbf{u}_{1,i}^{k+1}, \dots, \mathbf{u}_{j-1,i}^{k+1}, \mathbf{u}_{j,i}^k, \dots, \mathbf{u}_{d,i}^k)$. Denote $V_j^{k+1} := [\mathbf{v}_{j,1}^{k+1}, \dots, \mathbf{v}_{j,R}^{k+1}] \in \mathbb{R}^{n_j \times R}$.

When $1 \leq j \leq d-t$, from the definition of $L_\tau(\cdot)$, $\mathbf{v}_{j,i}$, and noticing the expression (3.14), we have that each column of U_j can be updated as follows

$$\mathbf{u}_{j,i}^{k+1} = \arg \min_{\|\mathbf{u}_{j,i}\|=1} -\sigma_i^k \langle \mathbf{v}_{j,i}^{k+1}, \mathbf{u}_{j,i} \rangle \Leftrightarrow \mathbf{u}_{j,i}^{k+1} = \mathbf{v}_{j,i}^{k+1} / \|\mathbf{v}_{j,i}^{k+1}\|, \quad 1 \leq i \leq R.$$

However, for the convenience of convergence analysis we compute the following instead

$$\mathbf{u}_{j,i}^{k+1} = \tilde{\mathbf{v}}_{j,i}^{k+1} / \|\tilde{\mathbf{v}}_{j,i}^{k+1}\|, \quad \text{where } \tilde{\mathbf{v}}_{j,i}^{k+1} = \sigma_i^k \mathbf{v}_{j,i}^{k+1} + \alpha \mathbf{u}_{j,i}^k, \quad 1 \leq i \leq R; \quad (3.21)$$

here $\alpha > 0$ is an arbitrary constant. Note that $\mathbf{u}_{j,1}^{k+1}, \dots, \mathbf{u}_{j,R}^{k+1}$ can be updated simultaneously.

When $d-t+1 \leq j \leq d$, from the definition of L_τ , $\mathbf{v}_{j,i}^{k+1}$, V_j^{k+1} and recalling (3.14), it follows

$$U_j^{k+1} = \arg \min_{U_j^\top U_j = I} - \sum_{i=1}^R \langle \sigma_i \mathbf{v}_{j,i}^{k+1}, \mathbf{u}_{j,i} \rangle = \arg \max_{U_j^\top U_j = I} \langle V_j^{k+1} \cdot \text{diag}(\boldsymbol{\sigma}^k), U_j \rangle,$$

where $\text{diag}(\boldsymbol{\sigma}) = \text{diag}[\sigma_1, \dots, \sigma_R] \in \mathbb{R}^{R \times R}$ is a diagonal matrix. Similar to (3.21), we in fact compute the following problem instead

$$U_j^{k+1} = \arg \max_{U_j^\top U_j = I} \langle \tilde{V}_j^{k+1}, U_j \rangle, \quad \text{where } \tilde{V}_j^{k+1} = V_j^{k+1} \cdot \text{diag}(\boldsymbol{\sigma}^k) + \alpha U_j^k. \quad (3.22)$$

The above problem is to compute the polar decomposition of \tilde{V}_j^{k+1} , which admits a closed-form solution. Specifically, assume $\tilde{V}_j^{k+1} = P \Xi Q^\top$ is the SVD of \tilde{V}_j^{k+1} , where $P \in \mathbb{R}^{n_j \times R}$, $\Lambda, Q \in \mathbb{R}^{R \times R}$, $P^\top P = I$, $Q^\top Q = Q Q^\top = I$, $\Xi = \text{diag}(\lambda_1, \dots, \lambda_R)$ with λ_i being the singular value of \tilde{V}_j^{k+1} . Then $U_j^{k+1} = P Q^\top$. Moreover, letting $H_j^{k+1} := Q \Xi Q^\top$. Then we see that (3.22) gives the following relation

$$\tilde{V}_j^{k+1} = U_j^{k+1} H_j^{k+1}. \quad (3.23)$$

\mathcal{T} -, $\boldsymbol{\sigma}$ - and \mathcal{W} -subproblems From (3.19), we have that

$$\mathcal{T}_{i_1 \dots i_d}^{k+1} = \left(\mathcal{W}_{i_1 \dots i_d}^k \mathcal{A}_{i_1 \dots i_d} - \mathcal{Y}_{i_1 \dots i_d}^k + \tau \llbracket \boldsymbol{\sigma}^k; U_j^{k+1} \rrbracket_{i_1 \dots i_d} \right) / (\mathcal{W}_{i_1 \dots i_d}^k + \tau). \quad (3.24)$$

To compute $\boldsymbol{\sigma}^{k+1}$, from the expression of (3.14) it is easily seen that

$$\sigma_i^{k+1} = (\mathcal{Y}^{k+1} + \tau \mathcal{T}^{k+1}) \bigotimes_{j=1}^d \mathbf{u}_{j,i}^{k+1} / \tau, \quad 1 \leq i \leq R. \quad (3.25)$$

To compute \mathcal{W}^{k+1} , similar to (3.20) we have

$$\mathcal{W}_{i_1 \dots i_d}^{k+1} = \delta^2 \left(\delta^2 + (\mathcal{T}_{i_1 \dots i_d}^{k+1} - \mathcal{A}_{i_1 \dots i_d})^2 \right)^{-1}. \quad (3.26)$$

In summary, the HQ-ADMM is described in Algorithm 1, where each subproblem admits a closed-form solution.

Remark 3.1. 1. HQ-ADMM can be applied to a more general form of (2.5). Specifically, consider the data-fitting term given by $\Phi_\delta(\mathbf{L}(\llbracket \boldsymbol{\sigma}; U_j \rrbracket) - \mathbf{b})$, where \mathbf{L} is a linear operator, and \mathbf{b} denote the observed data of the same size as $\mathbf{L}(\llbracket \boldsymbol{\sigma}; U_j \rrbracket)$. When \mathbf{L} represents the identity operator and \mathbf{b} denotes \mathcal{A} , the data-fitting term boils down to the objective of (2.5). When $\Phi_\delta(\mathbf{L}(\llbracket \boldsymbol{\sigma}; U_j \rrbracket) - \mathbf{b}) = \Phi_\delta(\boldsymbol{\Omega} \circledast (\llbracket \boldsymbol{\sigma}; U_j \rrbracket - \mathcal{A}))$, where $\boldsymbol{\Omega} \in \mathbb{R}^{n_1 \times \dots \times n_d}$ is a given 0-1 tensor with $\Omega_{i_1 \dots i_d} = 1$ if $\mathcal{A}_{i_1 \dots i_d}$ being observed while $\Omega_{i_1 \dots i_d} = 0$ if $\mathcal{A}_{i_1 \dots i_d}$ missing, it can be used to deal with robust tensor approximation with incomplete data. When \mathbf{L} is formed by a set of input data tensors, and each entry of \mathbf{b} denotes the output score of the corresponding input data tensor, it is the objective

Require: $U_j^0 = [\mathbf{u}_{j,i}^0, \dots, \mathbf{u}_{j,R}^0]$, with $\|\mathbf{u}_{j,i}\| = 1$, $1 \leq j \leq d-t$, $1 \leq i \leq R$; $(U^0)^\top U_j^0 = I$, $d-t+1 \leq j \leq d$; $\boldsymbol{\sigma}^0$, \mathcal{T}^0 , \mathcal{Y}^0 , \mathcal{W}^0 , $\alpha > 0$, $\tau > 0$, $\delta > 0$.

- 1: **for** $k = 0, 1, \dots$, **do**
- 2: Compute $\mathbf{u}_{j,i}^{k+1}$ via (3.21), $1 \leq j \leq d-t$, $1 \leq i \leq R$
- 3: Compute U_j^{k+1} via (3.22), $d-t+1 \leq j \leq d$
- 4: Compute \mathcal{T}^{k+1} via (3.24),
- 5: Compute $\mathcal{Y}^{k+1} = \mathcal{Y}^k - \tau (\llbracket \boldsymbol{\sigma}^k; U_j^{k+1} \rrbracket - \mathcal{T}^{k+1})$,
- 6: Compute $\boldsymbol{\sigma}^{k+1}$ via (3.25),
- 7: Compute \mathcal{W}^{k+1} via (3.26).
- 8: **end for**

of the (robust) tensor regression problem. To minimize $\Phi_\delta(\mathbf{L}(\llbracket \boldsymbol{\sigma}; U_j \rrbracket) - \mathbf{b})$ over orthonormal constraints, similar to (3.12), one can also formulate the problem as

$$\begin{aligned} \min_{\boldsymbol{\sigma}, U_j, \mathcal{T}, \mathbf{w}} \quad & \Phi_\delta(\mathbf{L}(\mathcal{T}) - \mathbf{b}) = \frac{1}{2} \|\sqrt{\mathbf{w}} \circledast (\mathbf{L}(\mathcal{T}) - \mathbf{b})\|_F^2 + \frac{\delta^2}{2} \sum_{i=1} \varrho(\mathbf{w}_i) \\ \text{s.t.} \quad & \mathcal{T} = \llbracket \boldsymbol{\sigma}; U_j \rrbracket, \quad \mathbf{w} \geq 0, \\ & \mathbf{u}_{j,i}^\top \mathbf{u}_{j,i} = 1, \quad 1 \leq j \leq d-t, 1 \leq i \leq R, \\ & U_j^\top U_j = I, \quad d-t+1 \leq j \leq d, \end{aligned}$$

where \mathbf{w} is the same size as \mathbf{b} defined similar to that in (3.11). The framework of HQ-ADMM then applies as well.

2. The idea of combining HQ property and ADMM framework can also be extended to solve other Cauchy loss based problems such as those studied in [27, 34]. Specifically, for problems of the form

$$\min_{\mathbf{x}} \Phi_\delta(L\mathbf{x} - \mathbf{b}) + R(\mathbf{x}),$$

where L is a matrix, \mathbf{b} is a vector of proper size, one can also convert it to

$$\min_{\mathbf{x}, \mathbf{w}} \|\sqrt{\mathbf{w}} \circledast (L\mathbf{y} - \mathbf{b})\|_F^2 + \sum_{i=1} \varrho(\mathbf{w}_i) + R(\mathbf{x}), \quad \text{s.t. } \mathbf{x} = \mathbf{y},$$

with \mathbf{w} defined similar to that in (3.11); an algorithm in the spirit of HQ-ADMM can be applied to solve it.

3. An alternative way to obtain closed-form solutions in ADMM for solving (2.5) is to use a linearization technique. For example, one can apply a linearized ADMM to solve the original problem (2.5) instead of the equivalent form (3.12), in which one also replace $\llbracket \boldsymbol{\sigma}; U_j \rrbracket$ by \mathcal{T} ; however, to solve the \mathcal{T} -subproblem, i.e., $\min_{\mathcal{T}} \Phi_\delta(\mathcal{T} - \mathcal{A}) + \langle \mathcal{Y}, \mathcal{T} \rangle + \tau/2 \|\mathcal{T} - \llbracket \boldsymbol{\sigma}; U_j \rrbracket\|_F^2$, which does not admit a closed-form solution, one linearizes $\Phi_\delta(\mathcal{T} - \mathcal{A})$ and then imposes a proximal term. The issue is that by doing this, one does not fully explore the structure of the model, which may lead to inefficiency. On the other hand, extra effort has to be paid to find a suitable step-size for this linearized subproblem.

4 Convergence of HQ-ADMM

This section establishes the convergence of HQ-ADMM. We note that to ensure the convergence, the only requirement is that $\tau \geq \sqrt{10}$. Throughout this section, to simplify the notations, we denote

$$\Delta_{U_j}^{k+1,k} := U_j^{k+1} - U_j^k.$$

The definitions of $\Delta_{\mathcal{T}}^{k+1,k}$, $\Delta_{\mathcal{W}}^{k+1,k}$, and $\Delta_{\mathcal{Y}}^{k+1,k}$ are analogous. In addition, we define the following proximal augmented Lagrangian function

$$\tilde{L}_{\tau}(\boldsymbol{\sigma}, U_j, \mathcal{T}, \mathcal{Y}, \mathcal{W}, \mathcal{T}') := L_{\tau}(\boldsymbol{\sigma}, U_j, \mathcal{T}, \mathcal{Y}, \mathcal{W}) + \frac{2}{\tau} \|\mathcal{T} - \mathcal{T}'\|_F^2,$$

which is needed to study the diminishing property of the terms $\|\Delta_{U_j}^{k+1,k}\|_F$ and $\|\Delta_{\mathcal{T}}^{k+1,k}\|_F$. For convenience we also denote

$$\tilde{L}_{\tau}^{k+1,k} := \tilde{L}_{\tau}(\boldsymbol{\sigma}^{k+1}, U_j^{k+1}, \mathcal{T}^{k+1}, \mathcal{Y}^{k+1}, \mathcal{W}^{k+1}, \mathcal{T}^k). \quad (4.27)$$

We present the first main result in the following, showing that the sequence generated by the algorithm is bounded, and every limit point of the sequence generated by HQ-ADMM is a stationary point. The proof is left to Section 4.1.

Theorem 4.1 (Subsequential convergence). *Let $\{\boldsymbol{\sigma}^k, U_j^k, \mathcal{T}^k, \mathcal{Y}^k, \mathcal{W}^k\}$ be generated by Algorithm 1 with $\tau \geq \sqrt{10}$ and $\alpha > 0$. Then*

1. $\{\boldsymbol{\sigma}^k, U_j^k, \mathcal{T}^k, \mathcal{Y}^k, \mathcal{W}^k\}$ is bounded;
2. the sequence $\{\tilde{L}_{\tau}^{k+1,k}\}$ defined in (4.27) is bounded, nonincreasing and convergent;
3. it holds that

$$\sum_{k=1}^{\infty} \left(\sum_{j=1}^d \|\Delta_{U_j}^{k+1,k}\|_F^2 + \|\Delta_{\mathcal{T}}^{k+1,k}\|_F^2 \right) < +\infty, \quad (4.28)$$

and

$$\|\Delta_{\boldsymbol{\sigma}}^{k+1,k}\| \rightarrow 0, \quad \|\Delta_{\mathcal{W}}^{k+1,k}\|_F \rightarrow 0, \quad \|\llbracket \boldsymbol{\sigma}^k; U_j^k \rrbracket - \mathcal{T}^k\|_F \rightarrow 0. \quad (4.29)$$

Moreover, every limit point $\{\boldsymbol{\sigma}^*, U_j^*, \mathcal{T}^*, \mathcal{Y}^*, \mathcal{W}^*\}$ satisfies the optimality condition (3.20). In particular, $\{\boldsymbol{\sigma}^*, U_j^*\}$ is also a stationary point of the original problem (2.5).

Next, based on the Kurdyka-Lojasiewicz (KL) property [4] which is widely used for proving the global convergence of nonconvex algorithms, we can show that the whole sequence converges to a single limit point. The proof is left to Sect. 4.2.

Theorem 4.2 (Global convergence). *Under the setting of Theorem 4.1, the whole sequence of $\{U_j^k, \mathcal{T}^k\}$ converges to a single limit point, i.e.,*

$$\lim_{k \rightarrow \infty} U_j^k = U_j^*, \quad 1 \leq j \leq d, \quad \lim_{k \rightarrow \infty} \mathcal{T}^k = \mathcal{T}^*.$$

4.1 Proof of Theorem 4.1

To prove the convergence of a nonconvex ADMM, a key step is to upper bound the size of the successive difference of the dual variables by that of the primal variables [23, 29, 47]. For HQ-ADMM, the weight \mathcal{W}^k brings barriers in the estimation of the upper bound. Fortunately, this can be overcome by realizing the relations between \mathcal{W}^k , \mathcal{T}^k and \mathcal{T}^{k-1} by using Lemma 4.1. The resulting estimate is given as follows.

Lemma 4.1. *It holds that*

$$\|\Delta_{\mathcal{Y}}^{k+1,k}\|_F \leq \|\Delta_{\mathcal{T}}^{k+1,k}\|_F + \|\Delta_{\mathcal{T}}^{k,k-1}\|_F.$$

Proof. From (3.24), we have

$$\mathcal{W}^k \circledast (\mathcal{T}^{k+1} - \mathcal{A}) + \mathcal{Y}^k - \tau (\llbracket \boldsymbol{\sigma}^k; U_j^{k+1} \rrbracket - \mathcal{T}^{k+1}) = 0,$$

which together with the definition of \mathcal{Y}^{k+1} yields

$$\mathcal{W}^k \otimes (\mathcal{T}^{k+1} - \mathcal{A}) + \mathcal{Y}^{k+1} = 0. \quad (4.30)$$

Therefore we have

$$\begin{aligned} \|\Delta_{\mathcal{Y}}^{k+1,k}\| &= \left\| \mathcal{W}^k \otimes (\mathcal{T}^{k+1} - \mathcal{A}) - \mathcal{W}^{k-1} \otimes (\mathcal{T}^k - \mathcal{A}) \right\|_F \\ &= \left\| \mathcal{W}^k \otimes (\mathcal{T}^{k+1} - \mathcal{A}) - \mathcal{W}^k \otimes (\mathcal{T}^k - \mathcal{A}) + \mathcal{W}^k \otimes (\mathcal{T}^k - \mathcal{A}) - \mathcal{W}^{k-1} \otimes (\mathcal{T}^k - \mathcal{A}) \right\|_F \\ &\leq \left\| \mathcal{W}^k \otimes (\mathcal{T}^{k+1} - \mathcal{T}^k) \right\|_F + \left\| (\mathcal{W}^k - \mathcal{W}^{k-1}) \otimes (\mathcal{T}^k - \mathcal{A}) \right\|_F \end{aligned} \quad (4.31)$$

Now denote $E_1 := \|\mathcal{W}^k \otimes (\mathcal{T}^{k+1} - \mathcal{T}^k)\|_F$ and $E_2 := \|(\mathcal{W}^k - \mathcal{W}^{k-1}) \otimes (\mathcal{T}^k - \mathcal{A})\|_F$. We first consider E_1 . From the definition of \mathcal{W}^k , we easily see that $\mathcal{W}_{i_1 \dots i_d}^k \leq 1$ for each i_1, \dots, i_d . Therefore,

$$E_1 \leq \|\Delta_{\mathcal{T}}^{k+1,k}\|. \quad (4.32)$$

Next we focus on E_2 . To simplify notations we denote $a_{i_1 \dots i_d}^k := \mathcal{T}_{i_1 \dots i_d}^k - \mathcal{A}_{i_1 \dots i_d}$ and

$$e_{i_1 \dots i_d} := \delta^2 a_{i_1 \dots i_d}^k \left(\frac{1}{\delta^2 + (a_{i_1 \dots i_d}^k)^2} - \frac{1}{\delta^2 + (a_{i_1 \dots i_d}^{k-1})^2} \right).$$

Then E_2 can be expressed as

$$\begin{aligned} E_2^2 &= \sum_{i_1=1, \dots, i_d=1}^{n_1, \dots, n_d} (\mathcal{W}_{i_1 \dots i_d}^{k+1} - \mathcal{W}_{i_1 \dots i_d}^k)^2 (\mathcal{T}_{i_1 \dots i_d} - \mathcal{A}_{i_1 \dots i_d})^2 \\ &= \sum_{i_1=1, \dots, i_d=1}^{n_1, \dots, n_d} \delta^4 (a_{i_1 \dots i_d}^k)^2 \left(\frac{1}{\delta^2 + (a_{i_1 \dots i_d}^k)^2} - \frac{1}{\delta^2 + (a_{i_1 \dots i_d}^{k-1})^2} \right)^2 = \sum_{i_1=1, \dots, i_d=1}^{n_1, \dots, n_d} e_{i_1 \dots i_d}^2. \end{aligned}$$

It follows from Proposition 2.3 that

$$|e_{i_1 \dots i_d}| \leq |a_{i_1 \dots i_d}^k - a_{i_1 \dots i_d}^{k-1}|,$$

and so

$$E_2 \leq \|\mathcal{T}^k - \mathcal{A} - (\mathcal{T}^{k-1} - \mathcal{A})\|_F = \|\Delta_{\mathcal{T}}^{k,k-1}\|_F. \quad (4.33)$$

(4.31) combining with (4.32) and (4.33) yields the desired result. \square

With Lemma 4.1, we then establish a sufficiently decreasing inequality with respect to $\{\tilde{L}_{\tau}^{k+1,k}\}$ defined in (4.27).

Lemma 4.2. *Let the parameter τ satisfy $\tau \geq \sqrt{10}$. Then there holds*

$$\tilde{L}_{\tau}^{k,k-1} - \tilde{L}_{\tau}^{k+1,k} \geq \frac{\alpha}{2} \sum_{j=1}^d \left\| \Delta_{U_j}^{k+1,k} \right\|_F^2 + \frac{1}{\tau} \left\| \Delta_{\mathcal{T}}^{k+1,k} \right\|_F^2, \quad \forall k,$$

where $\alpha > 0$ is defined in (3.21) and (3.22).

Proof. We first consider the decrease caused by U_j . When $1 \leq j \leq d - t$, according to the algorithm, the

expression of $L_\tau(\cdot)$, that $\|\mathbf{u}_{j,i}^k\| = 1$ and recalling the definition of $\mathbf{u}_{j,i}^{k+1}$, $\mathbf{v}_{j,i}^{k+1}$ and $\tilde{\mathbf{v}}_{j,i}^{k+1}$, we have

$$\begin{aligned}
& L_\tau(\boldsymbol{\sigma}^k, U_1^{k+1}, \dots, U_{j-1}^{k+1}, U_j^k, \dots, U_j^d, \mathcal{T}^k, \mathcal{Y}^k, \mathcal{W}^k) - \\
& \quad L_\tau(\boldsymbol{\sigma}, U_1^{k+1}, \dots, U_j^{k+1}, U_{j+1}^k, \dots, U_d^k, \mathcal{T}^k, \mathcal{Y}^k, \mathcal{W}^k) \\
= & \sum_{i=1}^R \langle \sigma_i^k \cdot (\mathcal{Y}^k + \tau \mathcal{T}^k) \mathbf{u}_{1,i}^{k+1} \otimes \dots \otimes \mathbf{u}_{j-1,i}^{k+1} \otimes \mathbf{u}_{j+1,i}^k \otimes \dots \otimes \mathbf{u}_{d,i}^k, \mathbf{u}_{j,i}^{k+1} - \mathbf{u}_{j,i}^k \rangle \\
= & \sum_{i=1}^R \langle \sigma_i^k \cdot \mathbf{v}_{j,i}^{k+1}, \mathbf{u}_{j,i}^{k+1} - \mathbf{u}_{j,i}^k \rangle \\
= & \sum_{i=1}^R \langle \sigma_i^k \cdot \mathbf{v}^{k+1} + \alpha \mathbf{u}_{j,i}^k, \mathbf{u}_{j,i}^{k+1} - \mathbf{u}_{j,i}^k \rangle + \frac{\alpha}{2} \|\mathbf{u}_{j,i}^{k+1} - \mathbf{u}_{j,i}^k\|^2 \\
= & \sum_{i=1}^R \left\langle \tilde{\mathbf{v}}_{j,i}^{k+1}, \frac{\tilde{\mathbf{v}}_{j,i}^{k+1}}{\|\tilde{\mathbf{v}}_{j,i}^{k+1}\|} - \mathbf{u}_{j,i}^k \right\rangle + \frac{\alpha}{2} \|\mathbf{u}_{j,i}^{k+1} - \mathbf{u}_{j,i}^k\|^2 \\
\geq & \frac{\alpha}{2} \sum_{i=1}^R \|\mathbf{u}_{j,i}^{k+1} - \mathbf{u}_{j,i}^k\|^2 = \frac{\alpha}{2} \|\Delta_{U_j}^{k+1,k}\|_F^2, \tag{4.34}
\end{aligned}$$

where the fourth equality follows from the definition of $\mathbf{u}_{j,i}^{k+1}$ and $\tilde{\mathbf{v}}_{j,i}^{k+1}$, and the inequality is due to $\|\mathbf{v}\| \geq \langle \mathbf{v}, \mathbf{u} \rangle$ for any vectors \mathbf{u}, \mathbf{v} of the same size with $\|\mathbf{u}\| = 1$.

The decrease of U_j when $d - t + 1 \leq j \leq d$ is similar. From the definition of V_j^{k+1} , It holds that

$$\begin{aligned}
& L_\tau(\boldsymbol{\sigma}^k, U_1^{k+1}, \dots, U_{j-1}^{k+1}, U_j^k, \dots, U_d^k, \mathcal{T}^k, \mathcal{Y}^k, \mathcal{W}^k) - \\
& \quad L_\tau(\boldsymbol{\sigma}^k, U_1^{k+1}, \dots, U_j^{k+1}, U_{j+1}^k, \dots, U_d^k, \mathcal{T}^k, \mathcal{Y}^k, \mathcal{W}^k) \\
= & \sum_{i=1}^R \langle \sigma_i^k \cdot (\mathcal{Y}^k + \tau \mathcal{T}^k) \mathbf{u}_{1,i}^{k+1} \otimes \dots \otimes \mathbf{u}_{j-1,i}^{k+1} \otimes \mathbf{u}_{j+1,i}^k \otimes \dots \otimes \mathbf{u}_{d,i}^k, \mathbf{u}_{j,i}^{k+1} - \mathbf{u}_{j,i}^k \rangle \\
= & \langle V_j^{k+1} \cdot \text{diag}(\boldsymbol{\sigma}^k), U_j^{k+1} - U_j^k \rangle \\
= & \langle V_j^{k+1} \cdot \text{diag}(\boldsymbol{\sigma}^k) + \alpha U_j^k, U_j^{k+1} - U_j^k \rangle + \frac{\alpha}{2} \|U_j^{k+1} - U_j^k\|_F^2 \\
\geq & \frac{\alpha}{2} \|\Delta_{U_j}^{k+1,k}\|_F^2, \tag{4.35}
\end{aligned}$$

where the inequality follows from the definition of U_j^{k+1} in (3.22).

To show the decrease of \mathcal{T} , note that $L_\tau(\cdot)$ is strongly convex with respect to \mathcal{T} , which we can easily deduce that

$$L_\tau(\boldsymbol{\sigma}^k, U_j^{k+1}, \mathcal{T}^k, \mathcal{Y}^k, \mathcal{W}^k) - L_\tau(\boldsymbol{\sigma}^k, U_j^{k+1}, \mathcal{T}^{k+1}, \mathcal{Y}^k, \mathcal{W}^k) \geq \frac{\tau}{2} \|\Delta_{\mathcal{T}}^{k+1,k}\|_F^2. \tag{4.36}$$

Next, it follows from the definition of \mathcal{Y}^{k+1} and Lemma 4.1 that

$$\begin{aligned}
& L_\tau(\boldsymbol{\sigma}^k, U_j^{k+1}, \mathcal{T}^{k+1}, \mathcal{Y}^k, \mathcal{W}^k) - L_\tau(\boldsymbol{\sigma}^k, U_j^{k+1}, \mathcal{T}^{k+1}, \mathcal{Y}^{k+1}, \mathcal{W}^k) \\
= & \langle \mathcal{Y}^{k+1} - \mathcal{Y}^k, \llbracket \boldsymbol{\sigma}^k; U_j^{k+1} \rrbracket - \mathcal{T}^{k+1} \rangle \\
= & -\frac{1}{\tau} \|\Delta_{\mathcal{Y}}^{k+1,k}\|_F^2 \\
\geq & -\frac{2}{\tau} \left(\|\Delta_{\mathcal{T}}^{k+1,k}\|_F^2 + \|\Delta_{\mathcal{T}}^{k,k-1}\|_F^2 \right). \tag{4.37}
\end{aligned}$$

Finally, it follows from the definition of $\boldsymbol{\sigma}^{k+1}$ and \mathcal{W}^{k+1} that

$$L_\tau(\boldsymbol{\sigma}^k, U_j^{k+1}, \mathcal{T}^{k+1}, \mathcal{Y}^{k+1}, \mathcal{W}^k) - L_\tau(\boldsymbol{\sigma}^{k+1}, U_j^{k+1}, \mathcal{T}^{k+1}, \mathcal{Y}^{k+1}, \mathcal{W}^{k+1}) \geq 0. \tag{4.38}$$

As a result, summing up (4.34)–(4.38) yields

$$\begin{aligned}
& L_\tau(\boldsymbol{\sigma}^k, U_j^k, \mathcal{T}^k, \mathcal{Y}^k, \mathcal{W}^k) - L_\tau(\boldsymbol{\sigma}^{k+1}, U_j^{k+1}, \mathcal{T}^{k+1}, \mathcal{Y}^{k+1}, \mathcal{W}^{k+1}) \\
& \geq \frac{\alpha}{2} \sum_{j=1}^d \left\| \Delta_{U_j}^{k+1, k} \right\|_F^2 + \left(\frac{\tau}{2} - \frac{2}{\tau} \right) \left\| \Delta_{\mathcal{T}}^{k+1, k} \right\|_F^2 - \frac{2}{\tau} \left\| \Delta_{\mathcal{T}}^{k, k-1} \right\|_F^2 \\
& \geq \frac{\alpha}{2} \sum_{j=1}^d \left\| \Delta_{U_j}^{k+1, k} \right\|_F^2 + \left(\frac{2}{\tau} + \frac{1}{\tau} \right) \left\| \Delta_{\mathcal{T}}^{k+1, k} \right\|_F^2 - \frac{2}{\tau} \left\| \Delta_{\mathcal{T}}^{k, k-1} \right\|_F^2,
\end{aligned} \tag{4.39}$$

where the last inequality follows from the range of τ . Rearranging the terms of (4.39) gives the desired results. This completes the proof. \square

We then show that $\tilde{L}_\tau^{k, k-1}$ defined in Lemma 4.2 is lower bounded and the sequence $\{\boldsymbol{\sigma}^k, U_j^k, \mathcal{T}^k, \mathcal{Y}^k, \mathcal{W}^k\}$ is bounded as well.

Theorem 4.3. *Under the setting of Lemma 4.2, $\{\tilde{L}_\tau^{k, k-1}\}$ is bounded. The sequence $\{\boldsymbol{\sigma}^k, U_j^k, \mathcal{T}^k, \mathcal{Y}^k, \mathcal{W}^k\}$ generated by Algorithm 1 is bounded as well.*

Proof. Denote $Q^k(\mathcal{T}) := \frac{1}{2} \left\| \sqrt{\mathcal{W}^k} \otimes (\mathcal{T} - \mathcal{A}) \right\|_F^2$; thus we have $\nabla Q^k(\mathcal{T}) = \mathcal{W}^k \otimes (\mathcal{T} - \mathcal{A})$, and it then follows from the quadraticity of $Q^k(\cdot)$ and $\mathcal{Y}^k = -\mathcal{W}^{k-1} \otimes (\mathcal{T}^k - \mathcal{A})$ from (4.30) that

$$\begin{aligned}
& Q^{k-1}(\mathcal{T}^k) - Q^{k-1}(\llbracket \boldsymbol{\sigma}^k; U_j^k \rrbracket) - \langle \mathcal{Y}^k, \llbracket \boldsymbol{\sigma}^k; U_j^k \rrbracket - \mathcal{T}^k \rangle \\
& = \langle \mathcal{W}^{k-1} \otimes (\llbracket \boldsymbol{\sigma}^k; U_j^k \rrbracket - \mathcal{A}), \mathcal{T}^k - \llbracket \boldsymbol{\sigma}^k; U_j^k \rrbracket \rangle \\
& \quad + \frac{1}{2} \left\| \sqrt{\mathcal{W}^{k-1}} \otimes (\llbracket \boldsymbol{\sigma}^k; U_j^k \rrbracket - \mathcal{T}^k) \right\|_F^2 - \langle \mathcal{Y}^k, \llbracket \boldsymbol{\sigma}^k; U_j^k \rrbracket - \mathcal{T}^k \rangle \\
& = \frac{1}{2} \left\| \sqrt{\mathcal{W}^{k-1}} \otimes (\llbracket \boldsymbol{\sigma}^k; U_j^k \rrbracket - \mathcal{T}^k) \right\|_F^2 \\
& \quad + \langle \mathcal{W}^{k-1} \otimes (\llbracket \boldsymbol{\sigma}^k; U_j^k \rrbracket - \mathcal{A}) - \mathcal{W}^{k-1} \otimes (\mathcal{T}^k - \mathcal{A}), \mathcal{T}^k - \llbracket \boldsymbol{\sigma}^k; U_j^k \rrbracket \rangle \\
& = -\frac{1}{2} \left\| \sqrt{\mathcal{W}^{k-1}} \otimes (\llbracket \boldsymbol{\sigma}^k; U_j^k \rrbracket - \mathcal{T}^k) \right\|_F^2 \geq -\frac{1}{2} \left\| \llbracket \boldsymbol{\sigma}^k; U_j^k \rrbracket - \mathcal{T}^k \right\|_F^2,
\end{aligned} \tag{4.40}$$

where the last inequality uses the fact that $0 < \mathcal{W}_{i_1 \dots i_d}^{k-1} \leq 1$.

Based on (4.40), it follows from the proof of Lemma 4.2 that for any $k \geq 2$,

$$\begin{aligned}
& \tilde{L}_\tau^{k-1, k-2} = \tilde{L}_\tau(\boldsymbol{\sigma}^{k-1}, U_j^{k-1}, \mathcal{T}^{k-1}, \mathcal{Y}^{k-1}, \mathcal{W}^{k-1}, \mathcal{T}^{k-2}) \geq \tilde{L}_\tau(\boldsymbol{\sigma}^k, U_j^k, \mathcal{T}^k, \mathcal{Y}^k, \mathcal{W}^{k-1}, \mathcal{T}^{k-1}) \\
& = Q^{k-1}(\mathcal{T}^k) + \frac{\delta^2}{2} \sum_{i_1=1, \dots, i_d=1}^{n_1, \dots, n_d} \varrho(\mathcal{W}_{i_1 \dots i_d}^{k-1}) - \langle \mathcal{Y}^k, \llbracket \boldsymbol{\sigma}^k; U_j^k \rrbracket - \mathcal{T}^k \rangle + \frac{\tau}{2} \left\| \llbracket \boldsymbol{\sigma}^k; U_j^k \rrbracket - \mathcal{T}^k \right\|_F^2 + \frac{2}{\tau} \left\| \Delta_{\mathcal{T}}^{k, k-1} \right\|_F^2 \\
& \geq Q^{k-1}(\llbracket \boldsymbol{\sigma}^k; U_j^k \rrbracket) + \frac{\tau-1}{2} \left\| \llbracket \boldsymbol{\sigma}^k; U_j^k \rrbracket - \mathcal{T}^k \right\|_F^2 + \frac{\delta^2}{2} \sum_{i_1=1, \dots, i_d=1}^{n_1, \dots, n_d} \varrho(\mathcal{W}_{i_1 \dots i_d}^{k-1}) + \frac{2}{\tau} \left\| \Delta_{\mathcal{T}}^{k, k-1} \right\|_F^2 \\
& > -\infty,
\end{aligned} \tag{4.41}$$

where the first inequality follows from (4.40) and the last one is due to the range of τ and $\varrho(\cdot) \geq 0$. Thus $\{\tilde{L}_\tau^{k, k-1}\}$ is a lower bounded sequence. This together with Lemma 4.2 shows that $\{\tilde{L}_\tau^{k, k-1}\}$ is bounded.

We then show the boundedness of $\{\boldsymbol{\sigma}^k, U_j^k, \mathcal{T}^k, \mathcal{Y}^k, \mathcal{W}^k\}$. The boundedness of $\{U_j^k\}$ and $\{\mathcal{W}^k\}$ is obvious. Next, denote $g(\boldsymbol{\sigma}^k)$ as the formulation in line 3 of (4.41) with respect to $\boldsymbol{\sigma}^k$. Since by Proposition 2.1, namely, the orthonormality of $\bigotimes_{j=1}^d \mathbf{u}_{j,i}^k$,

$$\left\| \llbracket \boldsymbol{\sigma}^k; U_j^k \rrbracket - \mathcal{T}^k \right\|_F = \left\| \boldsymbol{\sigma}^k \right\|^2 - 2 \langle \llbracket \boldsymbol{\sigma}^k; U_j^k \rrbracket, \mathcal{T}^k \rangle + \left\| \mathcal{T}^k \right\|_F^2,$$

while $Q^{k-1}(\llbracket \boldsymbol{\sigma}^k; U_j^k \rrbracket)$ is convex with respect to $\boldsymbol{\sigma}^k$, we see that $g(\boldsymbol{\sigma}^k)$ is strongly convex with respect to $\boldsymbol{\sigma}^k$. This together with the boundedness of $\{\tilde{L}_\tau^{k, k-1}\}$ and (4.41) gives the boundedness of $\{\boldsymbol{\sigma}^k\}$. Quite similarly we

have that $\{\mathcal{T}^k\}$ is bounded. Finally, the boundedness of $\{\mathcal{Y}^k\}$ follows from the expression of the \mathcal{T} -subproblem (3.24). As a result, the sequence $\{\boldsymbol{\sigma}^k, U_j^k, \mathcal{T}^k, \mathcal{Y}^k, \mathcal{W}^k\}$ is bounded. This completes the proof. \square

Proof of Theorem 4.1. Lemma 4.2 in connection with Theorem 4.3 yields points 1, 2, and (4.28); (4.28) together with Lemma 4.1 and the definition of \mathcal{Y}^{k+1} , $\boldsymbol{\sigma}^{k+1}$ and \mathcal{W}^{k+1} gives (4.29). On the other hand, since the sequence is bounded, limit points exist. Assume that $\{\boldsymbol{\sigma}^*, U_j^*, \mathcal{T}^*, \mathcal{Y}^*, \mathcal{W}^*\}$ is a limit point with

$$\lim_{l \rightarrow \infty} \{\boldsymbol{\sigma}^{kl}, U_j^{kl}, \mathcal{T}^{kl}, \mathcal{Y}^{kl}, \mathcal{W}^{kl}\} = \{\boldsymbol{\sigma}^*, U_j^*, \mathcal{T}^*, \mathcal{Y}^*, \mathcal{W}^*\}.$$

(4.28), (4.29) then implies that

$$\lim_{l \rightarrow \infty} \{\boldsymbol{\sigma}^{kl+1}, U_j^{kl+1}, \mathcal{T}^{kl+1}, \mathcal{Y}^{kl+1}, \mathcal{W}^{kl+1}\} = \{\boldsymbol{\sigma}^*, U_j^*, \mathcal{T}^*, \mathcal{Y}^*, \mathcal{W}^*\}.$$

Therefore, taking the limit into l with respect to the $\mathbf{u}_{j,i}$ -subproblem (3.21) yields

$$\mathbf{v}_{j,i}^* \boldsymbol{\sigma}_i^* + \alpha \mathbf{u}_{j,i}^* = \|\tilde{\mathbf{v}}_{j,i}^*\| \mathbf{u}_{j,i}^*, \quad 1 \leq j \leq d-t, \quad 1 \leq i \leq R. \quad (4.42)$$

Multiplying both sides by $\mathbf{u}_{j,i}^*$ gives

$$\|\tilde{\mathbf{v}}_{j,i}^*\| = \alpha + \sigma_i^* \langle \mathbf{v}_{j,i}^*, \mathbf{u}_{j,i}^* \rangle = \alpha + \sigma_i^* \left\langle \mathcal{Y}^* + \tau \mathcal{T}^*, \bigotimes_{j=1}^d \mathbf{u}_{j,i}^* \right\rangle = \alpha + \tau (\sigma_i^*)^2, \quad (4.43)$$

where the second equality follows from the definition of $\mathbf{v}_{j,i}$ and the last one is given by passing the limit into the expression of σ_i^{kl+1} (3.25). Thus (4.42) together with (4.43) gives

$$(\mathcal{Y}^* + \tau \mathcal{T}^*) \bigotimes_{l \neq j}^d \mathbf{u}_{l,i}^* = \sigma_i^* \tau \mathbf{u}_{j,i}^*, \quad (4.44)$$

i.e., the first equation of the stationary point system (3.20).

Taking the limit into l with respect to the U_j -subproblem (3.22) and noticing the expression (3.23), we get

$$V_j^* \text{diag}(\boldsymbol{\sigma}^*) + \alpha U_j^* = U_j^* H_j^*,$$

where H_j^* is a symmetric matrix. Writing it columnwisely, we obtain

$$\sigma_i^* (\mathcal{Y}^* + \tau \mathcal{T}^*) \bigotimes_{l \neq j}^d \mathbf{u}_{l,i}^* = \sum_{i=1}^R (H_j^*)_{i,r} \mathbf{u}_{j,r}^* - \alpha \mathbf{u}_{j,i}^*, \quad d-t+1 \leq j \leq d, \quad 1 \leq i \leq R.$$

Denoting $\Lambda_j^* := H_j^* - \alpha I$, the above is exactly the third equality of (3.20). On the other hand, passing the limit into the expression of \mathcal{T}^k (3.24) and \mathcal{W}^k (3.26) respectively gives the \mathcal{T}^* - and \mathcal{W}^* - formulas in (3.20). Finally, the first expression of (4.29) yields $\mathcal{T}^* = \llbracket \boldsymbol{\sigma}^*; U_j^* \rrbracket$. Taking the above pieces together, we have that $\{\boldsymbol{\sigma}^*, U_j^*, \mathcal{T}^*, \mathcal{Y}^*, \mathcal{W}^*\}$ satisfies the stationary point system (3.20).

Next, we show that $\{\boldsymbol{\sigma}^*, U_j^*\}$ is also a stationary point of problem (2.5). We define its Lagrangian function as $L_{\Phi} := \Phi_{\delta}(\boldsymbol{\sigma}, U_j) - \sum_{j,i=1}^{d-t,R} \eta_{j,i} (\mathbf{u}_{j,i}^{\top} \mathbf{u}_{j,i} - 1) - \sum_{j=d-t+1}^d \langle \Lambda_j, U_j^{\top} U_j - I \rangle$, similar to that in (3.15). Taking derivative yields

$$\begin{cases} \partial_{\mathbf{u}_{j,i}} \Phi_{\delta}(\boldsymbol{\sigma}; U_j) = \eta_{j,i} \mathbf{u}_{j,i} \Leftrightarrow \mathcal{W} \otimes (\llbracket \boldsymbol{\sigma}; U_j \rrbracket - \mathcal{A}) \cdot \sigma_i \bigotimes_{l \neq j} \mathbf{u}_{l,i} = \eta_{j,i} \mathbf{u}_{j,i}, & 1 \leq j \leq d-t, 1 \leq i \leq R, \\ \partial_{\mathbf{u}_{j,i}} \Phi_{\delta}(\boldsymbol{\sigma}, U_j) = \sum_{r=1}^R (\Lambda_j)_{i,r} \mathbf{u}_{j,r} \Leftrightarrow \mathcal{W} \otimes (\llbracket \boldsymbol{\sigma}; U_j \rrbracket - \mathcal{A}) \cdot \sigma_i \bigotimes_{l \neq j} \mathbf{u}_{l,i} = \sum_{r=1}^R (\Lambda_j)_{i,r} \mathbf{u}_{j,r}, & d-t+1 \leq j \leq d, 1 \leq i \leq R, \\ \partial_{\boldsymbol{\sigma}} \Phi_{\delta}(\boldsymbol{\sigma}, U_j) = 0 \Leftrightarrow \langle \mathcal{W} \otimes (\llbracket \boldsymbol{\sigma}; U_j \rrbracket - \mathcal{A}), \bigotimes_{j=1}^d \mathbf{u}_{j,i} \rangle = 0, & 1 \leq i \leq R, \end{cases} \quad (4.45)$$

where $\mathcal{W}_{i_1 \dots i_d} = \delta^2 \left(1 + \left(\llbracket \boldsymbol{\sigma}; U_j \rrbracket_{i_1 \dots i_d} - \mathcal{A}_{i_1 \dots i_d} \right)^2 / \delta^2 \right)^{-1}$; multiplying $\mathbf{u}_{j,i}$ in both sides of the first equality above, and noticing the last equality, we get $\eta_{j,i} = 0$.

Since $\mathcal{T}^* = \llbracket \boldsymbol{\sigma}^*; U_j^* \rrbracket$, the \mathcal{T} -subproblem (3.24) also gives $\mathcal{Y}^* = \mathcal{W}^* \otimes (\llbracket \boldsymbol{\sigma}^*; U_j^* \rrbracket - \mathcal{A})$. This together with (4.44) and that $\mathcal{T}^* \bigotimes_{l \neq j}^d \mathbf{u}_{j,i}^* = \llbracket \boldsymbol{\sigma}^*; U_j^* \rrbracket \bigotimes_{l \neq j}^d \mathbf{u}_{j,i}^* = \sigma_i^* \mathbf{u}_{j,i}^*$ gives $\mathcal{W}^* \otimes (\llbracket \boldsymbol{\sigma}^*; U_j^* \rrbracket - \mathcal{A}) \bigotimes_{l \neq j}^d \mathbf{u}_{j,i}^* = 0$, i.e., the first equality of (4.45) by noticing $\eta_{j,i} = 0$. In a similar vein, we get that

$$\sigma_i^* \mathcal{W}^* \otimes (\llbracket \boldsymbol{\sigma}^*; U_j^* \rrbracket - \mathcal{A}) = \sum_{i=1}^R (H_j^*)_{i,r} \mathbf{u}_{j,r}^* - (\alpha + \tau \sigma_i^*) \mathbf{u}_{j,i}^*.$$

Taking $\Lambda_j := H_j^* - (\alpha + \tau \sigma_i^*) I$ gives the second relation of (4.45). The last equality follows directly from $\mathcal{W}^* \otimes (\llbracket \boldsymbol{\sigma}^*; U_j^* \rrbracket - \mathcal{A}) \bigotimes_{l \neq j}^d \mathbf{u}_{j,i}^* = 0$. The proof has been completed. \square

4.2 Proof of Theorem 4.2

To prove Theorem 4.2, we first recall some definitions from nonsmooth analysis. Denote $\text{dom} f := \{\mathbf{x} \in \mathbb{R}^n \mid f(\mathbf{x}) < +\infty\}$.

Definition 4.1 (c.f. [2]). *For $\mathbf{x} \in \text{dom} f$, the Fréchet subdifferential, denoted as $\hat{\partial}f(\mathbf{x})$, is the set of vectors $\mathbf{z} \in \mathbb{R}^n$ satisfying*

$$\liminf_{\substack{\mathbf{y} \neq \mathbf{x} \\ \mathbf{y} \rightarrow \mathbf{x}}} \frac{f(\mathbf{y}) - f(\mathbf{x}) - \langle \mathbf{z}, \mathbf{y} - \mathbf{x} \rangle}{\|\mathbf{x} - \mathbf{y}\|} \geq 0. \quad (4.46)$$

The subdifferential of f at $\mathbf{x} \in \text{dom} f$, written ∂f , is defined as

$$\partial f(\mathbf{x}) := \left\{ \mathbf{z} \in \mathbb{R}^n : \exists \mathbf{x}^k \rightarrow \mathbf{x}, f(\mathbf{x}^k) \rightarrow f(\mathbf{x}), \mathbf{z}^k \in \hat{\partial}f(\mathbf{x}^k) \rightarrow \mathbf{z} \right\}.$$

It is known that $\hat{\partial}f(\mathbf{x}) \subset \partial f(\mathbf{x})$ for each $\mathbf{x} \in \mathbb{R}^n$ [4]. An extended-real-valued function is a function $f : \mathbb{R}^n \rightarrow [-\infty, \infty]$, which is proper if $f(\mathbf{x}) > -\infty$ for all \mathbf{x} and $f(\mathbf{x}) < \infty$ for at least one \mathbf{x} . It is called closed if it is lower semi-continuous (l.s.c. for short). The global convergence relies on the the Kurdyka-Łojasiewicz (KL) property given as follows:

Definition 4.2 (KL property and KL function, c.f. [2,4]). *A proper function f is said to have the KL property at $\bar{\mathbf{x}} \in \text{dom} \partial f := \{\mathbf{x} \in \mathbb{R}^n \mid \partial f(\mathbf{x}) \neq \emptyset\}$, if there exist $\bar{\epsilon} \in (0, \infty]$, a neighborhood \mathcal{N} of $\bar{\mathbf{x}}$, and a continuous and concave function $\psi : [0, \bar{\epsilon}) \rightarrow \mathbb{R}_+$ which is continuously differentiable on $(0, \bar{\epsilon})$ with positive derivatives and $\psi(0) = 0$, such that for all $\mathbf{x} \in \mathcal{N}$ satisfying $f(\bar{\mathbf{x}}) < f(\mathbf{x}) < f(\bar{\mathbf{x}}) + \bar{\epsilon}$, it holds that*

$$\psi'(f(\mathbf{x}) - f(\bar{\mathbf{x}})) \text{dist}(0, \partial f(\mathbf{x})) \geq 1,$$

where $\text{dist}(0, \partial f(\mathbf{x}))$ means the distance from the original point to the set $\partial f(\mathbf{x})$. If a proper and l.s.c. function f satisfies the KL property at each point of $\text{dom} \partial f$, then f is called a KL function.

We then simplify $\tilde{L}_\tau(\cdot)$ by eliminating the variables \mathcal{W} and $\boldsymbol{\sigma}$. First, from the definition of \mathcal{W}^{k+1} and Lemma 2.1, we have that

$$\left\| \sqrt{\mathcal{W}^{k+1}} \otimes (\mathcal{T}^{k+1} - \mathcal{A}) \right\|_F^2 + \delta^2 \sum_{i_1=1, \dots, i_d=1}^{n_1, \dots, n_d} \varrho(\mathcal{W}_{i_1 \dots i_d}^{k+1}) = \Phi_\delta(\mathcal{T}^{k+1} - \mathcal{A}),$$

where $\Phi_\delta(\cdot)$ is defined in (2.5). This eliminate the \mathcal{W} from $\tilde{L}_\tau(\cdot)$. On the other hand, it follows from the definition of $\boldsymbol{\sigma}^{k+1}$ (3.25) that

$$\begin{aligned} & - \langle \mathcal{Y}^{k+1}, [\boldsymbol{\sigma}^{k+1}; U_j^{k+1}] - \mathcal{T}^{k+1} \rangle + \frac{\tau}{2} \left\| [\boldsymbol{\sigma}^{k+1}; U_j^{k+1}] - \mathcal{T}^{k+1} \right\|_F^2 \\ = & \langle \mathcal{Y}^{k+1}, \mathcal{T}^{k+1} \rangle + \frac{\tau}{2} \left\| \mathcal{T}^{k+1} \right\|_F^2 - \frac{1}{2\tau} \sum_{i=1}^R \left((\mathcal{Y}^{k+1} + \tau \mathcal{T}^{k+1}) \bigotimes_{j=1}^d \mathbf{u}_{j,i}^{k+1} \right)^2. \end{aligned}$$

Thus σ is also eliminated. In what follows, whenever necessary, σ_i^k still represents the expression $(\mathcal{Y}^k + \tau \mathcal{T}^k) \bigotimes_{j=1}^d \mathbf{u}_{j,i}^k / \tau$, but we only treat it as a representation instead of a variable.

Then $\tilde{L}_\tau(\boldsymbol{\sigma}^{k+1}, U_j^{k+1}, \mathcal{T}^{k+1}, \mathcal{Y}^{k+1}, \mathcal{W}^{k+1} \mathcal{T}^k)$ can be equivalently written as

$$\begin{aligned} & \tilde{L}_\tau(U_j^{k+1}, \mathcal{T}^{k+1}, \mathcal{Y}^{k+1}, \mathcal{T}^k) \\ = & \frac{1}{2} \Phi_\delta(\mathcal{T}^{k+1} - \mathcal{A}) + \langle \mathcal{Y}^{k+1}, \mathcal{T}^{k+1} \rangle + \frac{\tau}{2} \left\| \mathcal{T}^{k+1} \right\|_F^2 - \frac{1}{2\tau} \sum_{i=1}^R \left((\mathcal{Y}^{k+1} + \tau \mathcal{T}^{k+1}) \bigotimes_{j=1}^d \mathbf{u}_{j,i}^{k+1} \right)^2 + \frac{2}{\tau} \left\| \Delta_{\mathcal{T}^{k+1}, k} \right\|_F^2. \end{aligned}$$

In addition, we denote

$$\tilde{L}_{\tau, \alpha}(U_j, \mathcal{T}, \mathcal{Y}, \mathcal{T}') := \tilde{L}_\tau(U_j, \mathcal{T}, \mathcal{Y}, \mathcal{T}') - \frac{\alpha}{2} \sum_{j=1}^d \left\| U_j \right\|_F^2 + \sum_{j=1, i=1}^{d-t, R} \iota_{\text{st}(n_j, 1)}(\mathbf{u}_{j,i}) + \sum_{j=d-t+1}^d \iota_{\text{st}(n_j, R)}(U_j).$$

We can see that under the constraints of the optimization problem (2.5), $\tilde{L}_{\tau,\alpha}(\cdot) = \tilde{L}_\tau(\cdot) + c$ where c is a constant. This together with Theorem 4.1 shows that $\{\tilde{L}_{\tau,\alpha}(U_j^{k+1}, \mathcal{T}^{k+1}, \mathcal{Y}^{k+1}, \mathcal{T}^k), \}$ is also a bounded and nonincreasing sequence. In addition, we have that $\tilde{L}_{\tau,\alpha}(\cdot)$ is a KL function.

Proposition 4.1. $\tilde{L}_{\tau,\alpha}(U_j, \mathcal{T}, \mathcal{Y}, \mathcal{T}')$ defined above is a proper, l.s.c., and KL function.

Proof. It is clear that $\tilde{L}_{\tau,\alpha}(\cdot)$ is proper and l.s.c.. Next, since the constrained sets in (2.5) are all Stiefel manifolds, items 2 and 6 of [4, Example 2] tell us that they are semi-algebraic sets, and their indicator functions are semi-algebraic functions. Therefore, the indicator functions are KL functions [4, Theorem 3]. On the other hand, the remaining part of $\tilde{L}_{\tau,\alpha}$ (besides the indicator functions) is an analytic function and hence it is KL [4]. As a result, $\tilde{L}_{\tau,\alpha}(U_j, \mathcal{T}, \mathcal{Y}, \mathcal{T}')$ is a KL function. \square

In the sequel, we mainly rely on $\tilde{L}_{\tau,\alpha}(\cdot)$ to prove the global convergence. For convenience, we denote

$$\tilde{L}_{\tau,\alpha}^{k+1,k} := \tilde{L}_{\tau,\alpha}(U_j^{k+1}, \mathcal{T}^{k+1}, \mathcal{Y}^{k+1}, \mathcal{T}^k), \text{ and } \partial\tilde{L}_{\tau,\alpha}^{k+1,k} := \partial\tilde{L}_{\tau,\alpha}(U_j^{k+1}, \mathcal{T}^{k+1}, \mathcal{Y}^{k+1}, \mathcal{T}^k);$$

denote $\Delta_{U_j, \mathcal{T}}^{k+1,k} := (U_j^{k+1}, \mathcal{T}^{k+1}) - (U_j^k, \mathcal{T}^k)$, and

$$\|\Delta_{U_j, \mathcal{T}}^{k+1,k}\|_F := \sqrt{\sum_{j=1}^d \left(\|\Delta_{U_j}^{k+1,k}\|_F^2 + \|\Delta_{\mathcal{T}}^{k+1,k}\|_F^2 \right)}.$$

Lemma 4.3. *There exists a large enough constant $c_0 > 0$, such that*

$$\text{dist}(\mathbf{0}, \partial\tilde{L}_{\tau,\alpha}^{k+1,k}) \leq c_0 \left(\|\Delta_{U_j, \mathcal{T}}^{k+1,k}\|_F + \|\Delta_{U_j, \mathcal{T}}^{k,k-1}\|_F \right). \quad (4.47)$$

Proof. We first consider $\partial_{\mathbf{u}_{j,i}} \tilde{L}_{\tau,\alpha}^{k+1,k}$, $1 \leq j \leq d-t$, $1 \leq i \leq R$, and $\partial_{U_j} \tilde{L}_{\tau,\alpha}^{k+1,k}$, $d-t+1 \leq j \leq d$, respectively. In what follows, we denote

$$\bar{\mathbf{v}}_{j,i}^{k+1} := \sigma_i^{k+1} (\mathcal{Y}^{k+1} + \tau\mathcal{T}^{k+1}) \bigotimes_{l \neq j} \mathbf{u}_{l,i}^{k+1} + \alpha \mathbf{u}_{j,i}^{k+1}, \text{ and } \bar{\mathbf{V}}_j^{k+1} := [\bar{\mathbf{v}}_{j,1}^{k+1}, \dots, \bar{\mathbf{v}}_{j,R}^{k+1}].$$

We also recall $\mathbf{v}_{j,i}^{k+1} := (\mathcal{Y}^k + \tau\mathcal{T}^k) \mathbf{u}_{1,i}^{k+1} \otimes \dots \otimes \mathbf{u}_{j-1,i}^{k+1} \otimes \mathbf{u}_{j+1,i}^k \otimes \dots \otimes \mathbf{u}_{d,i}^k$ and $\tilde{\mathbf{v}}_{j,i}^{k+1} = \sigma_i^k \mathbf{v}_{j,i}^{k+1} + \alpha \mathbf{u}_{j,i}^k$ for later use. In addition, denote $\tilde{\mathbf{V}}_j^{k+1} := [\tilde{\mathbf{v}}_{j,1}^{k+1}, \dots, \tilde{\mathbf{v}}_{j,R}^{k+1}]$.

For $1 \leq j \leq d-t$, one has

$$\begin{aligned} \partial_{\mathbf{u}_{j,i}} \tilde{L}_{\tau,\alpha}^{k+1,k} &= -\sigma_i^{k+1} (\mathcal{Y}^{k+1} + \tau\mathcal{T}^{k+1}) \bigotimes_{l \neq j} \mathbf{u}_{l,i}^{k+1} - \alpha \mathbf{u}_{j,i}^{k+1} + \partial_{\ell_{\text{st}(n_j,1)}}(\mathbf{u}_{j,i}^{k+1}) \\ &= -\bar{\mathbf{v}}_{j,i}^{k+1} + \partial_{\ell_{\text{st}(n_j,1)}}(\mathbf{u}_{j,i}^{k+1}). \end{aligned} \quad (4.48)$$

we then wish to show that

$$\tilde{\mathbf{v}}_{j,i}^{k+1} \in \hat{\partial}_{\ell_{\text{st}(n_j,1)}}(\mathbf{u}_{j,i}^{k+1}) \subset \partial_{\ell_{\text{st}(n_j,1)}}(\mathbf{u}_{j,i}^{k+1}). \quad (4.49)$$

The proof is similar to that of [48, Lemma 6.1]. First, from the definition of $\ell_{\text{st}(n_j,1)}(\cdot)$ and $\hat{\partial}_{\ell_{\text{st}(n_j,1)}}(\cdot)$ in (4.46), it is not hard to see that if $\mathbf{y} \notin \text{st}(n_j, 1)$, then (4.46) clearly holds when $\mathbf{z} = \tilde{\mathbf{v}}_{j,i}^{k+1}$; otherwise if $\mathbf{y} \in \text{st}(n_j, 1)$, i.e., $\|\mathbf{y}\| = 1$, then from the definition of $\mathbf{u}_{j,i}^{k+1}$, we see that

$$\mathbf{u}_{j,i}^{k+1} = \arg \max_{\|\mathbf{y}\|=1} \langle \mathbf{y}, \tilde{\mathbf{v}}_{j,i}^{k+1} \rangle \Leftrightarrow \langle \tilde{\mathbf{v}}_{j,i}^{k+1}, \mathbf{u}_{j,i}^{k+1} - \mathbf{y} \rangle \geq 0, \quad \forall \|\mathbf{y}\| = 1,$$

which together with $\ell_{\text{st}(n_j,1)}(\mathbf{y}) = 0$ and $\ell_{\text{st}(n_j,1)}(\mathbf{u}_{j,i}^{k+1}) = 0$ gives

$$\liminf_{\mathbf{y} \neq \mathbf{u}_{j,i}^{k+1}, \mathbf{y} \rightarrow \mathbf{u}_{j,i}^{k+1}} \frac{\ell_{\text{st}(n_j,1)}(\mathbf{y}) - \ell_{\text{st}(n_j,1)}(\mathbf{u}_{j,i}^{k+1}) - \langle \tilde{\mathbf{v}}_{j,i}^{k+1}, \mathbf{y} - \mathbf{u}_{j,i}^{k+1} \rangle}{\|\mathbf{y} - \mathbf{u}_{j,i}^{k+1}\|} \geq 0.$$

As a result, (4.49) is true, which together with (4.48) shows that

$$\tilde{\mathbf{v}}_{j,i}^{k+1} - \bar{\mathbf{v}}_{j,i}^{k+1} \in \partial_{\mathbf{u}_{j,i}} \tilde{L}_{\tau,\alpha}^{k+1,k}, \quad 1 \leq j \leq d-t, \quad 1 \leq i \leq R.$$

Let $\mathbf{0}$ denote the original. Then by using the triangle inequality and the boundness of $\{\boldsymbol{\sigma}^k, U^k, \mathcal{T}^k, \mathcal{Y}^k\}$, and noticing the definition of $\Delta_{U_j, \mathcal{T}}^{k+1, k}$, there must exist large enough constants $c_1, c_2 > 0$ only depending on τ, α , and the size of $\{\boldsymbol{\sigma}^k, U^k, \mathcal{T}^k, \mathcal{Y}^k\}$, such that

$$\begin{aligned}
& \text{dist}(\mathbf{0}, \partial_{\mathbf{u}_{j,i}} \tilde{L}_{\tau, \alpha}^{k+1, k}) \\
& \leq \left\| \tilde{\mathbf{v}}_{j,i}^{k+1} - \bar{\mathbf{v}}_{j,i}^{k+1} \right\| \\
& \leq c_1 \left(\sum_{j=1}^d \left\| \Delta_{U_j}^{k+1, k} \right\|_F + \left\| \Delta_{\mathcal{T}}^{k+1, k} \right\|_F + \left\| \Delta_{\mathcal{Y}}^{k+1, k} \right\|_F \right) \\
& \leq c_1 \left(\sum_{j=1}^d \left\| \Delta_{U_j}^{k+1, k} \right\|_F + 2 \left\| \Delta_{\mathcal{T}}^{k+1, k} \right\|_F + \left\| \Delta_{\mathcal{T}}^{k, k-1} \right\|_F \right) \\
& \leq c_2 \left(\left\| \Delta_{U_j, \mathcal{T}}^{k+1, k} \right\|_F + \left\| \Delta_{U_j, \mathcal{T}}^{k, k-1} \right\|_F \right), \quad 1 \leq j \leq d-t.
\end{aligned} \tag{4.50}$$

On the other hand, for $d-t+1 \leq j \leq d$, by noticing the definition of \bar{V}_j^{k+1} , we have

$$\partial_{U_j} \tilde{L}_{\tau, \alpha}^{k+1, k} = -\bar{V}_j^{k+1} + \partial_{\text{st}(n_j, R)}(U_j^{k+1}).$$

From the definition of U_j^{k+1} in (3.22) and similar to the above argument, we can show that $\tilde{V}_j^{k+1} \in \partial_{\text{st}(n_j, R)}(U_j^{k+1})$. Thus

$$\tilde{V}_j^{k+1} - \bar{V}_j^{k+1} \in \partial_{U_j} \tilde{L}_{\tau, \alpha}^{k+1, k}, \quad d-t+1 \leq j \leq d.$$

Similar to (4.50), there exists a large enough constant $c_3 > 0$ such that

$$\text{dist}(\mathbf{0}, \partial_{\mathbf{u}_{j,i}} \tilde{L}_{\tau, \alpha}^{k+1, k}) \leq c_3 \left(\left\| \Delta_{U_j, \mathcal{T}}^{k+1, k} \right\|_F + \left\| \Delta_{U_j, \mathcal{T}}^{k, k-1} \right\|_F \right), \quad d-t+1 \leq j \leq d. \tag{4.51}$$

We then consider

$$\nabla_{\mathcal{T}} \tilde{L}_{\tau, \alpha}^{k+1, k} = \mathcal{W}^{k+1} \circledast (\mathcal{T}^{k+1} - \mathcal{A}) + \mathcal{Y}^{k+1} - \tau \left(\llbracket \boldsymbol{\sigma}^{k+1}; U_j^{k+1} \rrbracket - \mathcal{T}^{k+1} \right) + \frac{4}{\tau} (\mathcal{T}^{k+1} - \mathcal{T}^k).$$

Note that \mathcal{W}^{k+1} and $\boldsymbol{\sigma}^{k+1}$ above are only representations instead of variables, which represent (3.26) and (3.25). From the expression of \mathcal{Y}^{k+1} in (4.30), we have

$$\begin{aligned}
\left\| \mathcal{W}^{k+1} \circledast (\mathcal{T}^{k+1} - \mathcal{A}) + \mathcal{Y}^{k+1} \right\|_F &= \left\| (\mathcal{W}^{k+1} - \mathcal{W}^k) \circledast (\mathcal{T}^{k+1} - \mathcal{A}) \right\|_F \\
&\leq \left\| \Delta_{\mathcal{T}}^{k+1, k} \right\|_F,
\end{aligned}$$

where the inequality follows from Proposition 2.3. On the other side,

$$\begin{aligned}
\tau \left\| \llbracket \boldsymbol{\sigma}^{k+1}; U_j^{k+1} \rrbracket - \mathcal{T}^{k+1} \right\|_F &= \tau \left\| \llbracket \boldsymbol{\sigma}^{k+1}; U_j^{k+1} \rrbracket - \llbracket \boldsymbol{\sigma}^k; U_j^{k+1} \rrbracket + \llbracket \boldsymbol{\sigma}^k; U_j^{k+1} \rrbracket - \mathcal{T}^{k+1} \right\|_F \\
&\leq \tau \left\| \llbracket \boldsymbol{\sigma}^{k+1}; U_j^{k+1} \rrbracket - \llbracket \boldsymbol{\sigma}^k; U_j^{k+1} \rrbracket \right\|_F + \left\| \Delta_{\mathcal{Y}}^{k+1, k} \right\|_F \\
&\leq c_4 \left(\left\| \Delta_{U_j, \mathcal{T}}^{k+1, k} \right\|_F + \left\| \Delta_{U_j, \mathcal{T}}^{k, k-1} \right\|_F \right),
\end{aligned} \tag{4.52}$$

where $c_4 > 0$ is large enough. Combining the above pieces shows that there exists a large enough constant $c_5 > 0$ such that

$$\left\| \nabla_{\mathcal{T}} \tilde{L}_{\tau, \alpha}^{k+1, k} \right\|_F \leq c_5 \left(\left\| \Delta_{U_j, \mathcal{T}}^{k+1, k} \right\|_F + \left\| \Delta_{U_j, \mathcal{T}}^{k, k-1} \right\|_F \right). \tag{4.53}$$

Next, it follows from (4.52) that

$$\left\| \nabla_{\mathcal{Y}} \tilde{L}_{\tau, \alpha}^{k+1, k} \right\|_F = \left\| \llbracket \boldsymbol{\sigma}^{k+1}; U_j^{k+1} \rrbracket - \mathcal{T}^{k+1} \right\|_F \leq \frac{c_4}{\tau} \left(\left\| \Delta_{U_j, \mathcal{T}}^{k+1, k} \right\|_F + \left\| \Delta_{U_j, \mathcal{T}}^{k, k-1} \right\|_F \right). \tag{4.54}$$

Finally,

$$\left\| \nabla_{\mathcal{T}'} \tilde{L}_{\tau, \alpha}^{k+1, k} \right\|_F = \frac{4}{\tau} \left\| \Delta_{\mathcal{T}}^{k+1, k} \right\|_F. \tag{4.55}$$

Combining (4.50), (4.51), (4.53), (4.54), (4.55), we get that there exists a large enough constant $c_0 > 0$ independent of k , such that

$$\text{dist}(\mathbf{0}, \partial \tilde{L}_{\tau, \alpha}^{k+1, k}) \leq c_0 \left(\left\| \Delta_{U_j, \mathcal{T}}^{k+1, k} \right\|_F + \left\| \Delta_{U_j, \mathcal{T}}^{k, k-1} \right\|_F \right),$$

as desired. \square

Now we can present the proof concerning global convergence.

Proof of Theorem 4.2. We have mentioned that $\{\tilde{L}_{\tau, \alpha}^{k+1, k}\}$ inherits the properties of $\{\tilde{L}_{\tau}^{k+1, k}\}$, i.e., it is bounded, nonincreasing and convergent. We denote its limit as $\tilde{L}_{\tau, \alpha}^* = \lim_{k \rightarrow \infty} \tilde{L}_{\tau, \alpha}^{k+1, k} = \tilde{L}_{\tau, \alpha}(U_j^*, \mathcal{T}^*, \mathcal{Y}^*, \mathcal{T}^*)$ where $\{U_j^*, \mathcal{T}^*, \mathcal{Y}^*, \mathcal{T}^*\}$ is a limit point. According to Definition 4.2 and Proposition 4.1, there exist an $\epsilon_0 > 0$, a neighborhood of $\{U_j^*, \mathcal{T}^*, \mathcal{Y}^*, \mathcal{T}^*\}$, and a continuous and concave function $\psi(\cdot) : [0, \epsilon_0) \rightarrow \mathbb{R}_+$ such that for all $\{U_j, \mathcal{T}, \mathcal{Y}, \mathcal{T}'\} \in \mathcal{N}$ satisfying $\tilde{L}_{\tau, \alpha}^* < \tilde{L}_{\tau, \alpha}(U_j, \mathcal{T}, \mathcal{Y}, \mathcal{T}') < \tilde{L}_{\tau, \alpha}^* + \epsilon_0$, there holds

$$\psi'(\tilde{L}_{\tau, \alpha}(U_j, \mathcal{T}, \mathcal{Y}, \mathcal{T}') - \tilde{L}_{\tau, \alpha}^*) \text{dist}(0, \partial \tilde{L}_{\tau, \alpha}(U_j, \mathcal{T}, \mathcal{Y}, \mathcal{T}') \geq 1. \quad (4.56)$$

Let $\epsilon_1 > 0$ be such that

$$\mathbb{B}_{\epsilon_1} := \{(U_j, \mathcal{T}, \mathcal{Y}, \mathcal{T}') \mid \|U_j - U_j^*\|_F < \epsilon_1, 1 \leq j \leq d, \|\mathcal{T} - \mathcal{T}^*\|_F < \epsilon_1, \|\mathcal{Y} - \mathcal{Y}^*\|_F < 2\epsilon_1, \|\mathcal{T}' - \mathcal{T}^*\|_F < 2\epsilon_1\} \subset \mathcal{N},$$

and let $\mathbb{B}_{\epsilon_1}^{U_j, \mathcal{T}} := \{(U_j, \mathcal{T}) \mid \|U_j - U_j^*\|_F < \epsilon_1, 1 \leq j \leq d, \|\mathcal{T} - \mathcal{T}^*\|_F < \epsilon_1\}$. From the stationary point system (3.20) and the expression of \mathcal{Y}^{k+1} in (4.30), we have

$$\begin{aligned} \|\mathcal{Y}^k - \mathcal{Y}^*\|_F &= \|\mathcal{W}^{k-1} \otimes (\mathcal{T}^k - \mathcal{A}) - \mathcal{W}^* \otimes (\mathcal{T}^* - \mathcal{A})\|_F \\ &\leq \|\mathcal{W}^{k-1} \otimes (\mathcal{T}^k - \mathcal{A}) - \mathcal{W}^k \otimes (\mathcal{T}^k - \mathcal{A})\|_F + \|\mathcal{W}^k \otimes (\mathcal{T}^k - \mathcal{A}) - \mathcal{W}^* \otimes (\mathcal{T}^* - \mathcal{A})\|_F \\ &= \left\| \Delta_{\mathcal{T}}^{k, k-1} \right\|_F + \left\| \Delta_{\mathcal{T}}^{k, *} \right\|_F \end{aligned} \quad (4.57)$$

where the last inequality follows from Propositions 2.3 and 2.2. On the other hand,

$$\|\mathcal{T}^{k-1} - \mathcal{T}^*\|_F \leq \left\| \Delta_{\mathcal{T}}^{k, k-1} \right\|_F + \left\| \Delta_{\mathcal{T}}^{k, *} \right\|_F. \quad (4.58)$$

As Theorem 4.1 shows that there exists $k_0 > 0$ such that for $k \geq k_0$, $\left\| \Delta_{\mathcal{T}}^{k, k-1} \right\|_F < \epsilon_1$, (4.57) and (4.58) tells us that if $k \geq k_0$ and $(U_j^k, \mathcal{T}^k) \in \mathbb{B}_{\epsilon_1}^{U_j, \mathcal{T}}$, then $\{U_j^k, \mathcal{T}^k, \mathcal{Y}^k, \mathcal{T}^{k-1}\} \in \mathbb{B}_{\epsilon_1} \subset \mathcal{N}$. Such k_0 must exist as $\{U_j^*, \mathcal{T}^*, \mathcal{Y}^*, \mathcal{T}^*\}$ is a limit point. In addition, denote $c_1 := \min\{\alpha/2, 1/\tau\}$; then there exists $k_1 \geq k_0$ such that $(U_j^{k_1}, \mathcal{T}^{k_1}) \in \mathbb{B}_{\epsilon_1/2}^{U_j, \mathcal{T}}$ and

$$\begin{aligned} \frac{c_0}{2\sqrt{c_1 c_2}} \left\| \Delta_{U_j, \mathcal{T}}^{k_1, k_1-1} \right\|_F &< \frac{\epsilon_1}{16}, \quad \frac{c_0}{2\sqrt{c_1 c_2}} \left\| \Delta_{U_j, \mathcal{T}}^{k_1-1, k_1-2} \right\|_F < \frac{\epsilon_1}{16}, \quad \frac{c_2}{2\sqrt{c_1}} \psi(\tilde{L}_{\tau, \alpha}^{k_1, k_1-1} - L_{\tau, \alpha}^*) < \frac{\epsilon_1}{4}, \\ L_{\tau, \alpha}^* &< \tilde{L}_{\tau, \alpha}^{k_1, k_1-1} < L_{\tau, \alpha}^* + \epsilon_0, \end{aligned} \quad (4.59)$$

where c_0 is the constant appeared in Lemma 4.3, and c_2 is a constant such that $c_2 > 16c_0/\sqrt{c_1}$.

In what follows, we use induction method to show that $(U_j^k, \mathcal{T}^k) \in \mathbb{B}_{\epsilon_1}^{U_j, \mathcal{T}}$ for all $k > k_1$. Since $\psi(\cdot)$ in Definition 4.2 is concave, it holds that for any k ,

$$\psi'(\tilde{L}_{\tau, \alpha}^{k, k-1} - L_{\tau, \alpha}^*) \left((\tilde{L}_{\tau, \alpha}^{k, k-1} - \tilde{L}_{\tau, \alpha}^*) - (\tilde{L}_{\tau, \alpha}^{k+1, k} - \tilde{L}_{\tau, \alpha}^*) \right) \leq \psi(\tilde{L}_{\tau, \alpha}^{k, k-1} - \tilde{L}_{\tau, \alpha}^*) - \psi(\tilde{L}_{\tau, \alpha}^{k+1, k} - \tilde{L}_{\tau, \alpha}^*); \quad (4.60)$$

on the other side, from the previous paragraph we see that $(U_j^{k_1}, \mathcal{T}^{k_1}) \in \mathbb{B}_{\epsilon_1/2}^{U_j, \mathcal{T}}$, $\{U_j^{k_1}, \mathcal{T}^{k_1}, \mathcal{Y}^{k_1}, \mathcal{T}^{k_1-1}\} \in \mathbb{B}_{\epsilon_1} \subset \mathcal{N}$, and so (4.56) holds at $\{U_j^{k_1}, \mathcal{T}^{k_1}, \mathcal{Y}^{k_1}, \mathcal{T}^{k_1-1}\}$. Recall $c_1 = \min\{\alpha/2, 1/\tau\}$. From Lemma 4.2 and the relation between \tilde{L}_{τ} and $\tilde{L}_{\tau, \alpha}$, we obtain

$$\begin{aligned} c_1 \left\| \Delta_{U_j, \mathcal{T}}^{k_1+1, k_1} \right\|_F^2 &\leq \tilde{L}_{\tau, \alpha}^{k_1, k_1-1} - \tilde{L}_{\tau, \alpha}^{k_1+1, k_1} \\ &\leq \frac{\psi(\tilde{L}_{\tau, \alpha}^{k_1, k_1-1} - \tilde{L}_{\tau, \alpha}^*) - \psi(\tilde{L}_{\tau, \alpha}^{k_1+1, k_1} - \tilde{L}_{\tau, \alpha}^*)}{\psi'(\tilde{L}_{\tau, \alpha}^{k_1, k_1-1} - \tilde{L}_{\tau, \alpha}^*)} \\ &\leq c_2 \left(\psi(\tilde{L}_{\tau, \alpha}^{k_1, k_1-1} - \tilde{L}_{\tau, \alpha}^*) - \psi(\tilde{L}_{\tau, \alpha}^{k_1+1, k_1} - \tilde{L}_{\tau, \alpha}^*) \right) \cdot c_2^{-1} \text{dist}(0, \partial \tilde{L}_{\tau, \alpha}^{k_1, k_1-1}), \end{aligned}$$

where the second inequality is due to (4.60) while the last one comes from (4.56). Using $\sqrt{ab} \leq \frac{a+b}{2}$ for $a \geq 0, b \geq 0$, invoking (4.47) and noticing the range in (4.59), we obtain

$$\begin{aligned} \sqrt{c_1} \left\| \Delta_{U_j, \mathcal{T}}^{k_1+1, k} \right\|_F &\leq \frac{c_2}{2} \left(\psi(\tilde{L}_{\tau, \alpha}^{k_1, k_1-1} - \tilde{L}_{\tau, \alpha}^*) - \psi(\tilde{L}_{\tau, \alpha}^{k_1+1, k_1} - \tilde{L}_{\tau, \alpha}^*) \right) \\ &\quad + \frac{c_0}{2c_2} \left(\left\| \Delta_{U_j, \mathcal{T}}^{k_1, k_1-1} \right\|_F + \left\| \Delta_{U_j, \mathcal{T}}^{k_1-1, k_1-2} \right\|_F \right) \\ &< \frac{\sqrt{c_1} \epsilon_1}{4} + \frac{\sqrt{c_1} \epsilon_1}{8} < \frac{\sqrt{c_1} \epsilon_1}{2}, \end{aligned}$$

and so

$$\left\| \Delta_{U_j, \mathcal{T}}^{k_1+1, *} \right\|_F \leq \left\| \Delta_{U_j, \mathcal{T}}^{k_1+1, k_1} \right\|_F + \left\| \Delta_{U_j, \mathcal{T}}^{k_1, *} \right\|_F < \frac{\epsilon_1}{2} + \frac{\epsilon_1}{2} = \epsilon_1,$$

namely, $(U_j^{k_1+1}, \mathcal{T}^{k_1+1}) \in \mathbb{B}_{\epsilon_1}^{U_j, \mathcal{T}}$.

Now assume that $(U_j^k, \mathcal{T}^k) \in \mathbb{B}_{\epsilon_1}^{U_j, \mathcal{T}}$ for $k = k_1, \dots, K$. This implies that (4.56) is true at $\{U_j^k, \mathcal{T}^k, \mathcal{Y}^k, \mathcal{T}^{k-1}\}$, and similarly to the above analysis, we have

$$\sqrt{c_1} \left\| \Delta_{U_j, \mathcal{T}}^{k+1, k} \right\|_F \leq \frac{c_2}{2} \left(\psi(\tilde{L}_{\tau, \alpha}^{k, k-1} - \tilde{L}_{\tau, \alpha}^*) - \psi(\tilde{L}_{\tau, \alpha}^{k+1, k} - \tilde{L}_{\tau, \alpha}^*) \right) + \frac{c_0}{2c_2} \left(\left\| \Delta_{U_j, \mathcal{T}}^{k, k-1} \right\|_F + \left\| \Delta_{U_j, \mathcal{T}}^{k-1, k-2} \right\|_F \right), \quad k = k_1, \dots, K. \quad (4.61)$$

We then show that $(U_j^{K+1}, \mathcal{T}^{K+1}) \in \mathbb{B}_{\epsilon_1}^{U_j, \mathcal{T}}$. Summing (4.61) for $k = k_1, \dots, K$ yields

$$\begin{aligned} \sqrt{c_1} \sum_{k=k_1}^K \left\| \Delta_{U_j, \mathcal{T}}^{k+1, k} \right\|_F &\leq \frac{c_2}{2} \left(\psi(\tilde{L}_{\tau, \alpha}^{k_1, k_1-1} - \tilde{L}_{\tau, \alpha}^*) - \psi(\tilde{L}_{\tau, \alpha}^{K+1, K} - \tilde{L}_{\tau, \alpha}^*) \right) + \frac{c_0}{2c_2} \sum_{k=k_1}^K \left(\left\| \Delta_{U_j, \mathcal{T}}^{k, k-1} \right\|_F + \left\| \Delta_{U_j, \mathcal{T}}^{k-1, k-2} \right\|_F \right) \\ &\leq \frac{c_2}{2} \left(\psi(\tilde{L}_{\tau, \alpha}^{k_1, k_1-1} - \tilde{L}_{\tau, \alpha}^*) - \psi(\tilde{L}_{\tau, \alpha}^{K+1, K} - \tilde{L}_{\tau, \alpha}^*) \right) \\ &\quad + \frac{c_0}{c_2} \sum_{k=k_1}^{K-1} \left\| \Delta_{U_j, \mathcal{T}}^{k+1, k} \right\|_F + \frac{2c_0}{c_2} \left\| \Delta_{U_j, \mathcal{T}}^{k_1, k_1-1} \right\|_F + \frac{c_0}{c_2} \left\| \Delta_{U_j, \mathcal{T}}^{k_1-1, k_1-2} \right\|_F. \end{aligned} \quad (4.62)$$

Rearranging the terms, noticing (4.59) and noticing that $\frac{c_2}{c_0} > \frac{\sqrt{c_1}}{16}$, we have

$$\frac{15\sqrt{c_1}}{16} \sum_{k=k_1}^K \left\| \Delta_{U_j, \mathcal{T}}^{k+1, k} \right\|_F \leq \frac{\sqrt{c_1}}{4} \epsilon_1 + \frac{\sqrt{c_1} \epsilon_1}{16} + \frac{\sqrt{c_1} \epsilon_1}{16},$$

and so

$$\begin{aligned} \left\| \Delta_{U_j, \mathcal{T}}^{K+1, *} \right\|_F &\leq \left\| \Delta_{U_j, \mathcal{T}}^{K+1, k_1} \right\|_F + \left\| \Delta_{U_j, \mathcal{T}}^{k_1, *} \right\|_F \\ &< \sum_{k=k_1}^K \left\| \Delta_{U_j, \mathcal{T}}^{k+1, k} \right\|_F + \frac{\epsilon_1}{2} \\ &< \frac{3\epsilon_1}{8} + \frac{\epsilon_1}{2} < \epsilon_1. \end{aligned}$$

Thus induction method implies that $(U_j^k, \mathcal{T}^k) \in \mathbb{B}_{\epsilon_1}^{U_j, \mathcal{T}}$ for all $k \geq k_1$, i.e., $\{U_j^k, \mathcal{T}^k, \mathcal{Y}^k, \mathcal{T}^{k-1}\} \in \mathcal{N}$, $k \geq k_1$. As a result, (4.61) holds for all $k \geq k_1$, so does (4.62). Therefore, letting $K \rightarrow \infty$ in (4.62) yields

$$\sum_{k=1}^{\infty} \left\| \Delta_{U_j, \mathcal{T}}^{k+1, k} \right\|_F < +\infty,$$

which shows that $\{U_j^k, \mathcal{T}^k\}$ is a Cauchy sequence and hence converges. Since (U_j^*, \mathcal{T}^*) in Theorem 4.1 is a limit point, the whole sequence converges to (U_j^*, \mathcal{T}^*) . This completes the proof. \square

5 Numerical Experiments

We evaluate the robustness of model (2.5) solved by HQ-ADMM in this section using synthetic and real data. The least squares based model (2.2) is used as a comparison. (2.2) is solved by the alternating least squares

(ALS) method. All the computations are conducted on an Intel i7-7770 CPU desktop computer with 32 GB of RAM. The supporting software is Matlab R2015b. The Matlab package Tensorlab [45] is employed for tensor operations. The Matlab code of HQ-ADMM is available at https://github.com/yuningyang19/hqadmm_rota.

The stopping criterion for HQ-ADMM is $\left| \left\| \left[\left[\boldsymbol{\sigma}^{k+1}; U_j^{k+1} \right] - \mathcal{A} \right\|_F - \left\| \left[\boldsymbol{\sigma}^k; U_j^k \right] - \mathcal{A} \right\|_F \right| \leq 10^{-6}$ or $k \geq 2000$ for practical reasons. The parameter α in HQ-ADMM is set to 10^{-8} , $\tau \in \{0.7, 1\}$; $\delta = 0.05$.

Synthetic data We consider randomly generated tensors contaminated by different kinds of noises listed in the following

- $\mathcal{A} = \mathcal{A}_0 / \|\mathcal{A}_0\|_F + \beta \cdot \mathcal{N} / \|\mathcal{N}\|_F$, where \mathcal{A}_0 is the ground truth tensor specified later, and \mathcal{N} denotes the Cauchy noise, with scale parameter $\delta = 0.05$. $\beta = 0.5$;
- $\mathcal{A} = \mathcal{A}_0 / \|\mathcal{A}_0\|_F + \mathcal{O}$. Here \mathcal{O} denotes sparse outliers, with sparsity 0.1, i.e., 10% of the entries of \mathcal{A}_0 are contaminated by outliers. Outliers are drawn uniformly from $[0, 10]$;
- $\mathcal{A} = \mathcal{A}_0 / \|\mathcal{A}_0\|_F + \beta \cdot \mathcal{N} / \|\mathcal{N}\|_F$, where \mathcal{N} denotes Gaussian noise, with $\beta = 0.1$.

The ground truth tensor $\mathcal{A}_0 = \sum_{i=1}^R \sigma_i \otimes_{j=1}^d \mathbf{u}_{j,i}$, where U_j are randomly drawn from a uniformly distribution in $[-1, 1]$. U_j , $d - t + 1 \leq j \leq d$, are then made to be columnwisely orthonormal, while the remaining U_j are columnwisely normalized. σ_i are drawn from Gaussian distribution. For convenience, we set $d = 3$ or 4 , $n_1 = \dots = n_d$, and $R = 5$ in all the experiments in this part. The initializers for HQ-ADMM and ALS are randomly generated. The reported results are averaged over 50 instances for each case.

Table 1: Comparison of HQ-ADMM for (2.5) and ALS for (2.2) when the ground truth tensor is contaminated by Cauchy noise.

n	(d, t)	HQ-ADMM for (2.5)			ALS for (2.2)		
		err.	iter.	time	err.	iter.	time
10	(3, 1)	5.57E-02	395	0.16	4.29E-01	149	0.04
20	(3, 1)	4.66E-02	315	0.21	4.20E-01	147	0.05
50	(3, 1)	4.30E-02	45	0.09	4.33E-01	309	0.27
80	(3, 1)	3.05E-02	71	0.77	4.31E-01	190	1.16
90	(3, 1)	3.04E-02	47	0.76	4.29E-01	152	1.28
100	(3, 1)	3.21E-02	86	1.62	4.41E-01	210	1.82
10	(3, 2)	5.25E-02	453	0.19	3.84E-01	33	0.01
20	(3, 2)	2.93E-02	137	0.10	4.12E-01	17	0.01
60	(3, 2)	2.25E-02	200	1.03	4.42E-01	11	0.04
80	(3, 2)	2.20E-02	58	0.60	4.18E-01	11	0.07
90	(3, 2)	2.02E-02	136	2.11	4.33E-01	14	0.11
100	(3, 2)	2.57E-02	96	1.84	4.23E-01	10	0.09
80	(3, 3)	1.39E-02	35	0.34	1.41E+00	2	0.02
100	(3, 3)	2.08E-02	89	1.69	1.41E+00	2	0.03
10	(4, 1)	3.86E-02	64	0.08	4.12E-01	341	0.21
20	(4, 1)	7.98E-02	40	0.16	4.45E-01	613	1.02
30	(4, 1)	7.37E-02	28	0.71	4.25E-01	485	6.55
40	(4, 1)	5.08E-02	25	1.62	4.47E-01	637	16.68
10	(4, 2)	4.98E-02	75	0.09	4.56E-01	299	0.19
20	(4, 2)	1.11E-01	53	0.20	4.73E-01	527	0.94
30	(4, 2)	7.33E-02	36	1.09	4.76E-01	394	6.06
40	(4, 2)	6.85E-02	27	1.75	4.70E-01	705	19.25
10	(4, 3)	9.57E-02	100	0.12	4.83E-01	664	0.41
20	(4, 3)	8.60E-02	69	0.27	5.00E-01	707	1.04
30	(4, 3)	1.29E-01	35	0.98	5.18E-01	645	9.72
40	(4, 3)	1.40E-01	30	1.86	5.41E-01	878	22.68

Table 2: Comparison of HQ-ADMM for (2.5) and ALS for (2.2) when the ground truth tensor is contaminated by outliers.

n	(d, t)	HQ-ADMM for (2.5)			ALS for (2.2)		
		err.	iter.	time	err.	iter.	time
10	(3, 1)	4.54E-01	89	0.04	1.40E+00	150	0.04
20	(3, 1)	5.95E-02	46	0.04	1.41E+00	251	0.09
50	(3, 1)	1.99E-02	31	0.10	1.41E+00	757	0.95
80	(3, 1)	2.21E-02	27	0.55	1.41E+00	1456	12.17
90	(3, 1)	3.52E-02	28	0.70	1.41E+00	1204	11.59
100	(3, 1)	2.82E-02	31	0.91	1.41E+00	1390	15.44
10	(3, 2)	4.32E-01	56	0.03	1.41E+00	120	0.04
20	(3, 2)	6.13E-02	35	0.04	1.41E+00	314	0.15
50	(3, 2)	7.50E-03	25	0.07	1.41E+00	592	0.69
80	(3, 2)	7.40E-03	25	0.42	1.41E+00	820	6.05
90	(3, 2)	6.66E-03	26	0.65	1.41E+00	828	7.80
100	(3, 2)	8.16E-03	27	0.90	1.41E+00	928	11.99
80	(3, 3)	6.08E-03	25	0.42	1.41E+00	2	0.02
100	(3, 3)	6.72E-03	27	0.80	1.41E+00	2	0.04
10	(4, 1)	1.04E-01	76	0.23	1.42E+00	187	0.14
20	(4, 1)	2.91E-02	34	0.28	1.41E+00	439	1.02
30	(4, 1)	4.40E-02	28	1.06	1.41E+00	1173	18.40
40	(4, 1)	6.09E-02	27	2.00	1.41E+00	885	26.09
10	(4, 2)	1.31E-01	67	0.08	1.41E+00	246	0.16
20	(4, 2)	5.23E-02	28	0.13	1.41E+00	729	1.12
30	(4, 2)	6.17E-02	27	0.85	1.41E+00	697	12.68
40	(4, 2)	3.36E-02	29	1.88	1.41E+00	1047	29.12
10	(4, 3)	1.40E-01	64	0.08	1.41E+00	208	0.13
20	(4, 3)	8.14E-02	29	0.12	1.41E+00	622	0.92
30	(4, 3)	8.45E-02	38	1.15	1.41E+00	900	14.85
40	(4, 3)	1.13E-01	30	2.12	1.41E+00	846	24.38

Comparisons of HQ-ADMM for solving (2.5) and ALS for solving (2.2) with Cauchy noise are reported in Table 1, where $\text{err.} = \|\mathcal{A}_0/\|\mathcal{A}_0\|_F - \mathcal{A}^*/\|\mathcal{A}^*\|_F\|_F$, with $\mathcal{A}^* = \llbracket \boldsymbol{\sigma}^*; U_j^* \rrbracket$ the tensor generated by the algorithm. ‘‘iter.’’ denotes the number of iterates, and ‘‘time’’ stands for the CPU time consumed by the algorithm. From the ‘‘err.’’ columns, we see that in all cases, HQ-ADMM performs much better than ALS; in particular, ‘‘err.’’ of HQ-ADMM is smaller than 0.1 in almost all cases, which confirms that the proposed model and algorithm are consistent with Cauchy noise. Considering the efficiency, we see that HQ-ADMM all converges within 500 iterates, and it consumes 1 ~ 2 seconds. Comparing with ALS, when $d = 3$, ALS is more efficient in most cases, while HQ-ADMM outperforms ALS when $d = 4$. Thus HQ-ADMM is efficient.

The cases contaminated by outliers are reported in Table 2, from which we can still observe that HQ-ADMM for solving (2.5) is consistent with outliers, owing to the redescending property of the Cauchy loss. HQ-ADMM outperforms ALS in terms of the iterates and CPU time.

The cases with Gaussian noise are reported in Table 3. It is known that model (2.2) is consistent with Gaussian noise, which can be seen from the table. We also observe that (2.5) is consistent with Gaussian noise from the third column, although the results are slightly worse than (2.2), as reported in the table. However, it is interesting to see that in some cases, namely, $(n, d, t) = (80, 3, 1), (30, 4, 1), (40, 4, 1), (20, 4, 2), (30, 4, 3)$, HQ-ADMM for (2.5) is slightly better than ALS for (2.2). HQ-ADMM still shows its efficiency, and is more stable than ALS, as ALS needs much more iterates when $t = 1$.

Table 3: Comparison of HQ-ADMM for (2.5) and ALS for (2.2) when the ground truth tensor is contaminated by Gaussian noise.

n	(d, t)	HQ-ADMM for (2.5)			ALS for (2.2)		
		err.	iter.	time	err.	iter.	time
10	(3, 1)	4.51E-02	198	0.09	4.09E-02	676	0.18
20	(3, 1)	3.62E-02	53	0.04	2.73E-02	564	0.19
50	(3, 1)	2.24E-02	30	0.08	2.18E-02	550	0.58
80	(3, 1)	2.14E-02	34	0.57	2.72E-02	716	5.78
90	(3, 1)	2.70E-02	33	0.79	2.44E-02	696	6.69
100	(3, 1)	2.79E-02	34	0.98	2.28E-02	712	7.75
10	(3, 2)	3.89E-02	296	0.13	3.48E-02	16	0.01
20	(3, 2)	2.15E-02	65	0.05	1.87E-02	17	0.01
50	(3, 2)	7.99E-03	24	0.07	7.67E-03	14	0.02
80	(3, 2)	4.90E-03	24	0.40	4.82E-03	20	0.15
90	(3, 2)	4.68E-03	25	0.60	4.34E-03	41	0.40
100	(3, 2)	3.85E-03	24	0.72	3.85E-03	7	0.10
10	(4, 1)	1.01E-01	673	0.83	8.62E-02	613	0.42
20	(4, 1)	7.46E-02	67	0.31	6.21E-02	699	1.33
30	(4, 1)	6.22E-02	29	1.05	6.61E-02	692	11.90
40	(4, 1)	8.68E-02	27	1.92	1.11E-01	858	24.49
10	(4, 2)	1.39E-02	45	0.15	1.74E-02	20	0.02
20	(4, 2)	4.75E-03	23	0.20	9.09E-03	17	0.05
30	(4, 2)	5.42E-03	26	0.91	2.71E-03	14	0.25
40	(4, 2)	2.26E-03	26	2.10	1.96E-03	41	1.24
10	(4, 3)	1.29E-02	48	0.17	1.23E-02	10	0.01
20	(4, 3)	4.93E-03	24	0.21	4.73E-03	10	0.04
30	(4, 3)	2.72E-03	25	0.98	2.88E-03	30	0.53
40	(4, 3)	1.95E-03	26	2.15	1.92E-03	21	0.67

Table 4: HQ-ADMM for video surveillance with different R . The last column shows the compressed ratio of the compressed background factors D, U, V to the sum of background frames $B_r, 1 \leq r \leq l$.

R	iter.	time	$\frac{R(1000+144+176)}{1000*144*176}$
10	43	33.86	0.05%
20	31	26.02	0.1%
30	26	21.58	0.16%
40	43	38.13	0.21%
50	31	28.78	0.26%

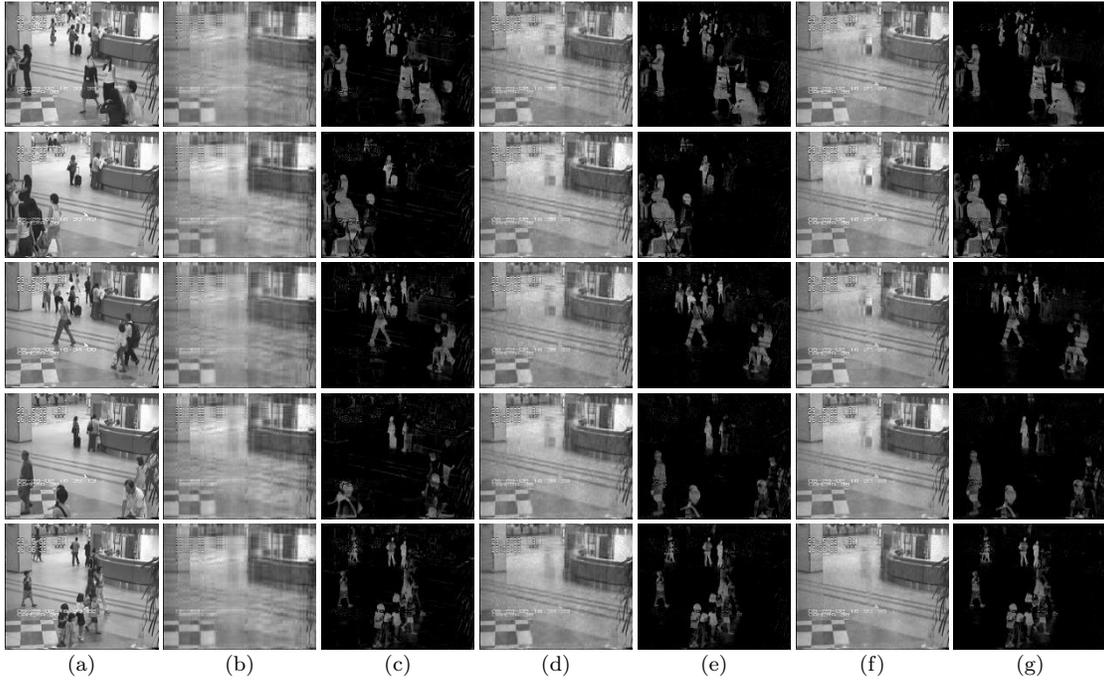


Figure 2: Some extracting frames by HQ-ADMM from the video airport. Column (a): The original frames; Columns (b) and (c): Extracted with $R = 10$; Columns (d) and (e): Extracted with $R = 30$; Columns (f) and (g): Extracted with $R = 50$.

Simultaneous foreground-background extraction and compression Foreground-background extraction finds applications in video surveillance, where the aim is to detect moving objects such as human beings from static background. As the background changes little in the video, it is reasonable to project the background frames to a low dimensional subspace to compress the data. We show how this problem can be fitted into our model (2.5). Assume that a gray video consists of l frames, each of size $m \times n$, resulting into a third-order tensor $\mathcal{A} \in \mathbb{R}^{l \times m \times n}$. Let A_i denotes its i -th frame. Our goal is to decompose it as $A_r = B_r + F_r$, in which B_r and F_r denote the back-/foreground frames, respectively. Under the assumption that B_r 's lie in a low dimensional subspace with commonalities, we write $B_r = UD_rV^\top = \sum_{i=1}^R (D_r)_{ii} \mathbf{u}_i \mathbf{v}_i^\top$, $1 \leq r \leq l$, where $U = [\mathbf{u}_1, \dots, \mathbf{u}_R]$, $V = [\mathbf{v}_1, \dots, \mathbf{v}_R]$ are orthonormal matrices, D_r is diagonal, and R is a parameter. On the other hand, the foreground is often sparse and can be recognized as outliers. Therefore, the Cauchy loss can be employed to control the effect of outliers. Denoting

$$\phi_\delta(A_r - UD_rV^\top) := \sum_{s=1, t=1}^{m, n} \frac{\delta^2}{2} \log \left(1 + ((A_r)_{st} - (UD_rV^\top)_{st})^2 / \delta^2 \right),$$

the problem can be modeled as

$$\min_{U^\top U=I, V^\top V=I} \sum_{r=1}^l \phi_\delta(A_r - UD_rV^\top).$$

If we further denote $D \in \mathbb{R}^{l \times R}$ where the r -th row is exactly the diagonal entries of D_r , the it can be written in the form of (2.5), i.e.,

$$\min_{U^\top U=I, V^\top V=I} \Phi_\delta(\mathcal{A} - \llbracket D, U, V \rrbracket),$$

where σ is absorbed into D .

The tested video ‘‘airport’’ was downloaded from http://perception.i2r.a-star.edu.sg/bk_model/bk_index.html. The video consists of 4583 frames, each of size 144×176 . We use 1000 frames, resulting into a tensor $\mathcal{A} \in \mathbb{R}^{1000 \times 144 \times 176}$. \mathcal{A} is then normalized for conveniently choosing parameters, where we set $\delta = 0.05$, $\tau = 1$, and $\alpha = 10^{-8}$. The parameter R varies in $\{10, 20, 30, 40, 50\}$. The quantitative results

are reported in Table 4, in which we can see that HQ-ADMM stops around $30 \sim 40$ iterates, and consumes around 30 seconds, which demonstrates the efficiency of the algorithm. The last column shows the compressed ratio of the compressed background factors D, U, V to the sum of background frames B_r , $1 \leq r \leq l$, from which we observe that the ratio is very high, resulting into low storage space. Some extracted frames with $R \in \{10, 30, 50\}$ are illustrated in Fig. 2. From the figures, we see that even when $R = 10$, HQ-ADMM can successfully separate the back-/foreground; of course, when $R \geq 30$, the extracted frames are of higher quality, in that the background frames reconstructed from UD_rV^\top are more clear.

6 Conclusions

Heavy-tailed noise and outliers often contaminate real-world data. In the context of tensor canonical polyadic approximation problem with one or more latent factor matrices having orthonormal columns, most existing models rely on the least squares loss, which is not resistant to heavy-tailed noise or outliers. To gain robustness, a Cauchy loss based robust orthogonal tensor approximation model was proposed in this work. To efficiently solve this model, by exploring its half-quadratic property, a new algorithm, termed as HQ-ADMM, was developed under the framework of alternating direction method of multipliers. Its global convergence was then established, thanks to some nice properties of the Cauchy loss. Numerical experiments on synthetic as well as real data demonstrate the efficiency and robustness of the proposed model and algorithm. In future work, it would be interesting to incorporate other robust losses in the orthogonal tensor approximation problem and to apply HQ-ADMM to solve other Cauchy loss based problems, as noted in Remark 3.1.

References

- [1] A. Anandkumar, P. Jain, Y. Shi, and U. N. Niranjan. Tensor vs. matrix methods: Robust tensor decomposition under block sparse perturbations. In *Artificial Intelligence and Statistics*, pages 268–276, 2016. 6
- [2] H. Attouch, J. Bolte, and B. F. Svaiter. Convergence of descent methods for semi-algebraic and tame problems: proximal algorithms, forward–backward splitting, and regularized Gauss–Seidel methods. *Math. Program.*, 137(1-2):91–129, 2013. 17
- [3] Al. Beaton and J. Tukey. The fitting of power series, meaning polynomials, illustrated on band-spectroscopic data. *Technometrics*, 16(2):147–185, 1974. 5
- [4] J. Bolte, S. Sabach, and M. Teboulle. Proximal alternating linearized minimization for nonconvex and nonsmooth problems. *Math. Program.*, 146(1-2):459–494, 2014. 12, 17, 18
- [5] J. Chen and Y. Saad. On the tensor SVD and the optimal low rank orthogonal approximation of tensors. *SIAM J. Matrix Anal. Appl.*, 30(4):1709–1734, 2009. 2, 3
- [6] L. Cheng, Y.-C. Wu, and H. V. Poor. Probabilistic tensor canonical polyadic decomposition with orthogonal factors. *IEEE Trans. Signal Process.*, 65(3):663–676, 2016. 6
- [7] A. Cichocki, D. Mandic, L. De Lathauwer, G. Zhou, Q. Zhao, C. Caiafa, and H. A. Phan. Tensor decompositions for signal processing applications: From two-way to multiway component analysis. *IEEE Signal Process. Mag.*, 32(2):145–163, 2015. 1
- [8] A. L. F. De Almeida, A. Y. Kibangou, S. Miron, and D. C. Araújo. Joint data and connection topology recovery in collaborative wireless sensor networks. In *Proc. of the IEEE Int. Conference on Acoustics, Speech and Signal Processing (ICASSP 2013)*, pages 5303–5307. IEEE, 2013. 2, 6

- [9] L. De Lathauwer. Algebraic methods after prewhitening. In *Handbook of Blind Source Separation*, pages 155–177. Elsevier, 2010. 2, 6
- [10] L. De Lathauwer. A short introduction to tensor-based methods for factor analysis and blind source separation. In *Proc. of the IEEE Int. Symp. on Image and Signal Processing and Analysis (ISPA 2011)*, pages 558–563. IEEE, 2011. 2, 6
- [11] L. De Lathauwer, B. De Moor, and J. Vandewalle. A multilinear singular value decomposition. *SIAM J. Matrix Anal. Appl.*, 21:1253–1278, 2000. 2
- [12] M. Ding, T.-Z. Huang, T.-H. Ma, X.-L. Zhao, and J.-H. Yang. Cauchy noise removal using group-based low-rank prior. *Applied Math. Comput.*, 372:124971, 2020. 4, 8
- [13] Y. Feng, J. Fan, and J. Suykens. A statistical learning approach to modal regression. *J. Mach. Learn. Res.*, 21(2):1–35, 2020. 5
- [14] Y. Feng, X. Huang, L. Shi, Y. Yang, and J. Suykens. Learning with the maximum correntropy criterion induced losses for regression. *J. Mach. Learn. Res.*, 16:993–1034, 2015. 5
- [15] S. Ganan and D. McClure. Bayesian image analysis: An application to single photon emission tomography. *Amer. Statist. Assoc.*, pages 12–18, 1985. 5
- [16] D. Goldfarb and Z. Qin. Robust low-rank tensor recovery: Models and algorithms. *SIAM J. Matrix Anal. Appl.*, 35(1):225–253, 2014. 5
- [17] N. Guan, T. Liu, Y. Zhang, D. Tao, and L. S. Davis. Truncated Cauchy non-negative matrix factorization. *IEEE Trans. Pattern Anal. Mach. Intell.*, 41(1):246–259, 2017. 4, 6, 8
- [18] Y. Guan and D. Chu. Numerical computation for orthogonal low-rank approximation of tensors. *SIAM J. Matrix Anal. Appl.*, 40(3):1047–1065, 2019. 2, 3
- [19] R. He, W.-S. Zheng, and B.-G. Hu. Maximum correntropy criterion for robust face recognition. *IEEE Trans. Pattern Anal. Mach. Intell.*, 33(8):1561–1576, 2010. 4, 6, 8
- [20] C. J. Hillar and L.-H. Lim. Most tensor problems are NP-hard. *J. ACM*, 60(6):45:1–45:39, 2013. 2
- [21] P. Holland and R. Welsch. Robust regression using iteratively reweighted least-squares. *Commun. Stat.-Theory Methods*, 6(9):813–827, 1977. 5
- [22] D. Hong, T. G. Kolda, and J. A. Duersch. Generalized canonical polyadic tensor decomposition. *SIAM Rev.*, 62(1):133–163, 2020. 6
- [23] M. Hong, Z.-Q. Luo, and M. Razaviyayn. Convergence analysis of alternating direction method of multipliers for a family of nonconvex problems. *SIAM J. Optim.*, 26(1):337–364, 2016. 12
- [24] S. Hu and K. Ye. Linear convergence of an alternating polar decomposition method for low rank orthogonal tensor approximations. *arXiv preprint arXiv:1912.04085*, 2019. 2
- [25] P. J. Huber. *Robust statistics*, volume 523. John Wiley & Sons, 2004. 2, 5
- [26] M. Ishteva, P.-A. Absil, and P. Van Dooren. Jacobi algorithm for the best low multilinear rank approximation of symmetric tensors. *SIAM J. Matrix Anal. Appl.*, 34(2):651–672, 2013. 2
- [27] G. Kim, J. Cho, and M. Kang. Cauchy noise removal by weighted nuclear norm minimization. *J. Sci. Comput.*, 83:15, 2020. 4, 8, 11
- [28] T. G. Kolda and B. W. Bader. Tensor decompositions and applications. *SIAM Rev.*, 51:455–500, 2009. 1, 3

- [29] G. Li and T. K. Pong. Global convergence of splitting methods for nonconvex composite optimization. *SIAM J. Optim.*, 25(4):2434–2460, 2015. [12](#)
- [30] J. Li, K. Usevich, and P. Comon. Globally convergent Jacobi-type algorithms for simultaneous orthogonal symmetric tensor diagonalization. *SIAM J. Matrix Anal. Appl.*, 39(1):1–22, 2018. [2](#)
- [31] J. Li and S. Zhang. Polar decomposition based algorithms on the product of stiefel manifolds with applications in tensor approximation. *arXiv preprint arXiv:1912.10390*, 2019. [2](#)
- [32] X. Li, Q. Lu, Y. Dong, and D. Tao. Robust subspace clustering by cauchy loss function. *IEEE Trans. Neural Netw. Learn. Syst.*, 30(7):2067–2078, 2018. [4](#), [8](#)
- [33] R. Maronna, O. Bustos, and V. Yohai. Bias-and efficiency-robustness of general m-estimators for regression with random carriers. In *Smoothing Techniques for Curve Estimation*, pages 91–116. Springer, 1979. [5](#)
- [34] J.-J. Mei, Y. Dong, T.-Z. Huang, and W. Yin. Cauchy noise removal by nonconvex admm with convergence guarantees. *J. Sci. Comput.*, 74(2):743–766, 2018. [4](#), [8](#), [11](#)
- [35] J. Pan and M. K. Ng. Symmetric orthogonal approximation to symmetric tensors with applications to image reconstruction. *Numer. Linear Algebra Appl.*, 25(5):e2180, 2018. [2](#)
- [36] V. Pravdova, F. Estienne, B. Walczak, and D. L. Massart. A robust version of the Tucker3 model. *Chemometr. Intell. Lab. Syst.*, 59(1):75 – 88, 2001. [6](#)
- [37] B. Savas and L.-H. Lim. Quasi-Newton methods on grassmannians and multilinear approximations of tensors. *SIAM J. Sci. Comput.*, 32(6):3352–3393, 2010. [2](#)
- [38] F. Sciacchitano, Y. Dong, and T. Zeng. Variational approach for restoring blurred images with cauchy noise. *SIAM J. Imag. Sci.*, 8(3):1894–1922, 2015. [4](#), [8](#)
- [39] A. Shashua and A. Levin. Linear image coding for regression and classification using the tensor-rank principle. In *CVPR*, volume 1, pages I–I. IEEE, 2001. [1](#), [6](#)
- [40] N. D. Sidiropoulos, L. De Lathauwer, X. Fu, K. Huang, E. E. Papalexakis, and C. Faloutsos. Tensor decomposition for signal processing and machine learning. *IEEE Trans. Signal Process.*, 65(13):3551–3582. [1](#)
- [41] N. D. Sidiropoulos, G. B. Giannakis, and R. Bro. Blind parafac receivers for ds-cdma systems. *IEEE Trans. Signal Process.*, 48(3):810–823, 2000. [2](#), [6](#)
- [42] M. Signoretto, Q. T. Dinh, L. De Lathauwer, and J. A. K. Suykens. Learning with tensors: a framework based on convex optimization and spectral regularization. *Mach. Learn.*, 94(3):303–351, 2014. [5](#)
- [43] M. Sørensen, L. De Lathauwer, P. Comon, S. Icart, and L. Deneire. Canonical polyadic decomposition with a columnwise orthonormal factor matrix. *SIAM J. Matrix Anal. Appl.*, 33(4):1190–1213, 2012. [2](#), [3](#), [6](#)
- [44] M. Sørensen, L. De Lathauwer, and L. Deneire. PARAFAC with orthogonality in one mode and applications in DS-CDMA systems. In *Proc. of the IEEE Int. Conference on Acoustics, Speech and Signal Processing (ICASSP 2010)*, pages 4142–4145. IEEE, 2010. [2](#), [6](#)
- [45] N. Vervliet, O. Debals, L. Sorber, M. Van Barel, and L. De Lathauwer. Tensorlab 3.0, Mar. 2016. Available online. [22](#)
- [46] L. Wang, M. T. Chu, and B. Yu. Orthogonal low rank tensor approximation: Alternating least squares method and its global convergence. *SIAM J. Matrix Anal. and Appl.*, 36(1):1–19, 2015. [2](#), [3](#)

- [47] Y. Wang, W. Yin, and J. Zeng. Global convergence of ADMM in nonconvex nonsmooth optimization. *J. Sci. Comput.*, 78(1):29–63, 2019. [12](#)
- [48] Y. Yang. The epsilon-alternating least squares for orthogonal low-rank tensor approximation and its global convergence. *arXiv preprint arXiv:1911.10921*, 2019. [2](#), [3](#), [18](#)
- [49] Y. Yang, Y. Feng, and J. A. K. Suykens. Robust low-rank tensor recovery with regularized re-descending M-estimator. *IEEE Trans. Neural Netw. Learn. Syst.*, 27(9):1933–1946, 2015. [4](#), [6](#), [8](#)