ALTERNATING MAHALANOBIS DISTANCE MINIMIZATION FOR ACCURATE AND WELL-CONDITIONED CP DECOMPOSITION*

NAVJOT SINGH† AND EDGAR SOLOMONIK†

Abstract. Canonical polyadic decomposition (CPD) is prevalent in chemometrics, signal processing, data mining, and many more fields. While many algorithms have been proposed to compute the CPD, alternating least squares (ALS) remains one of the most widely used algorithms for computing the decomposition. Recent works have introduced the notion of eigenvalues and singular values of a tensor and explored applications of eigenvectors and singular vectors in signal processing, data analytics, and various other fields. We introduce a new formulation for deriving singular values and vectors of a tensor by considering the critical points of a function differently from previous works. Computing these critical points in an alternating manner motivates an alternating optimization algorithm which corresponds to the ALS algorithm in the matrix case. However, for tensors with order greater than or equal to 3, it minimizes an objective function which is different from the commonly used least squares loss. Alternating optimization of this new objective leads to simple updates to the factor matrices with the same asymptotic computational cost as ALS. We show that a subsweep of this algorithm can achieve a superlinear convergence rate for exact CPD when the known rank is not larger than the mode lengths of the input tensor. We verify our theoretical arguments experimentally. We then view the algorithm as optimizing a Mahalanobis distance with respect to each factor with the ground metric dependent on the other factors. This perspective allows us to generalize our approach to interpolate between updates corresponding to the ALS and the new algorithm to manage the tradeoff between stability and fitness of the decomposition. Our experimental results show that for approximating synthetic and real-world tensors, this algorithm and its variants converge to a better conditioned decomposition with comparable and sometimes better fitness as compared to the ALS algorithm.

Key words. tensor decomposition, CP decomposition, alternating least squares, eigenvalues, singular values, Mahalanobis distance, condition number

MSC codes. 15A69, 15A72, 65K10, 65Y20, 65Y04, 68W25

DOI. 10.1137/22M1490739

1. Introduction. The canonical polyadic or CANDECOMP/PARAFAC (CP) tensor decomposition [21, 24] is used for analysis and compression of multiparameter datasets and is prevalent in tensor methods for scientific simulation [15, 40, 44, 54]. For an order 3 tensor \mathcal{T} , indexed as t_{ijk} , a rank-R CP decomposition with factor matrices A, B, and C is defined as

$$\mathcal{T} = [\![\boldsymbol{A}, \boldsymbol{B}, \boldsymbol{C}]\!], \quad t_{ijk} = \sum_{r=1}^R a_{ir} b_{jr} c_{kr},$$

where A is indexed as a_{ir} , and similarly for B and C. Determining the CP rank and finding an approximate CP decomposition of a tensor, so as to minimize

(1.1)
$$f(\boldsymbol{A}, \boldsymbol{B}, \boldsymbol{C}) = \frac{1}{2} \left\| \boldsymbol{\mathcal{T}} - [\boldsymbol{A}, \boldsymbol{B}, \boldsymbol{C}] \right\|_{F}^{2},$$

are NP-hard problems [23]. The CP decomposition of a tensor can be computed via various optimization algorithms, such as alternating least squares (ALS) [6, 22, 26, 53]

Funding: This work was supported by the US NSF OAC SSI program grant 1931258.

A2781

^{*}Submitted to the journal's Methods and Algorithms for Scientific Computing section April 15, 2022; accepted for publication (in revised form) June 22, 2023; published electronically November 3, 2023.

 $[\]rm https://doi.org/10.1137/22M1490739$

[†]Department of Computer Science, University of Illinois at Urbana-Champaign, Urbana, IL 61801 USA (navjot2@illinois.edu, solomon2@illinois.edu).

which aims to minimize the objective (1.1) in an alternating manner by considering all except one factor matrix fixed. There have been several attempts to improve the performance of the ALS algorithm [38, 42, 45, 50, 56]. Several methods which aim to minimize (1.1) with respect to all the factor matrices use gradient-based information [1, 46, 48, 57, 60, 63] to update all the factors simultaneously. Another set of methods optimize for all the factors simultaneously by formulating (1.1) as a nonlinear least squares problem by considering the entries of all the factors as variables. In addition to gradient information, these iterative methods use second order information to compute the next step which requires a system solve [46, 62] and can be achieved by an implicit conjugate gradient algorithm [55, 57].

Tensor eigenvalue problems are relevant in the context of solving multilinear systems, simulating quantum systems, exponential data fitting, and many other application areas [49]. However, the study of tensor eigenvalues and tensor singular values is at a relatively nascent stage. [34, 37] provide a definition and introduction to eigenvalues and singular values of a tensor. Computing eigenvalues of a tensor is a hard problem and can be solved via iterative methods for special cases such as computing a subset of eigenvalues of a tensor [32] or computing the real eigenvalues pairs of a real symmetric tensor [16]. The tensor eigenvalue problem is motivated by applications like blind source separation [8] and independent component analysis [25], which also motivate the closely related problem of diagonalizing a tensor. The concept of tensor diagonalization was introduced in [14], where approximate diagonalization of the tensor is considered by minimizing the sum of squares of off-diagonal entries or maximizing the sum of squares of diagonal entries of the tensor. There have been many follow-up works [35, 36, 61, 65] which consider approximate diagonalization of the tensor by invertible and orthogonal transformations.

In this work, we introduce a formulation for computing the singular values and vectors of a tensor by considering a logarithmic penalty function instead of Lagrangian variables [37] and computing the critical points of this function. This formulation, when generalized to computing invariant subspaces of a matrix, leads to diagonalization of the matrix and can be linked to the singular values and vectors of the matrix. When extended to tensors with order greater than or equal to 3, this formulation leads to another notion of diagonalization of the tensor which is different from the one introduced in prior work. The critical points of this new function spectrally diagonalize the tensor; i.e., the transformed equidimensional tensor of mode length R has R elementary eigenvectors with unit eigenvalues. Computing these stationary points alternatively motivates an alternating optimization algorithm for computing the CP decomposition of a tensor.

1.1. Motivation: Eigenvectors via Lagrangian optimization. In the case of low-rank matrix approximation, the Eckart-Young-Mirsky theorem shows that the best low-rank approximation may be obtained from the singular value decomposition. This connection relates low-rank factors to critical points of the bilinear form $f_{\mathbf{M}}(\mathbf{x}, \mathbf{y}) = \mathbf{x}^T \mathbf{M} \mathbf{y}$ with $\|\mathbf{x}\| \neq 0$, $\|\mathbf{y}\| \neq 0$. For tensors of order 3 and higher, tensor singular values have been similarly derived from critical points of multilinear forms. In particular, Lim [37] derives singular vectors and singular values by imposing constraints $\|\mathbf{x}\| = \|\mathbf{y}\| = 1$ and considering the critical points of the Lagrangian function. The same results may be obtained by instead considering a logarithmic interior point barrier function for the constraints, $\|\mathbf{x}\| \neq 0$, $\|\mathbf{y}\| \neq 0$, so

$$f_{M}(x, y) = x^{T} M y - \log(\|x\| \|y\|),$$

 $\nabla f_{M}(x, y) = 0 \Rightarrow M y = x/\|x\|^{2}, M^{T} x = y/\|y\|^{2}.$

Consequently, with $\sigma = 1/(\|\boldsymbol{x}\| \|\boldsymbol{y}\|)$ and $\boldsymbol{u} = \boldsymbol{x}/\|\boldsymbol{x}\|$, $\boldsymbol{v} = \boldsymbol{y}/\|\boldsymbol{y}\|$, we have $\boldsymbol{M}\boldsymbol{v} = \sigma\boldsymbol{u}$ and $\boldsymbol{M}^T\boldsymbol{u} = \sigma\boldsymbol{v}$. The use of a coefficient for the barrier function only affects the scaling of any critical point vectors, \boldsymbol{x} and \boldsymbol{y} . For tensors of order 3 and higher, tensor singular values can be similarly derived from critical points of

$$f_{\mathcal{T}}(\boldsymbol{x}^{(1)}, \dots, \boldsymbol{x}^{(N)}) = \sum_{i_1 \dots i_N} t_{i_1 \dots i_N} x_{i_1}^{(1)} \cdots x_{i_N}^{(N)} - \left(\sum_{i=1}^N \log(\|\boldsymbol{x}^{(i)}\|) \right),$$

$$\nabla f_{\mathcal{T}}(\boldsymbol{x}^{(1)}, \dots, \boldsymbol{x}^{(N)}) = \mathbf{0} \implies \frac{x_{i_j}^{(j)}}{\|\boldsymbol{x}^{(j)}\|^2} = \sum_{i_1 \dots \hat{i}_j \dots i_N} t_{i_1 \dots i_N} x_{i_1}^{(1)} \cdots \hat{x}_{i_j}^{(j)} \cdots x_{i_N}^{(N)},$$

where $i ldots \hat{j} ldots k$ implies j is omitted from the sequence. The use of the 2-norm in the above definitions leads to l^2 singular vectors [37], and with a symmetric tensor and each $\boldsymbol{x}^{(i)} = \boldsymbol{x}^{(j)}$ it yields Z-eigenvectors [34]. Choosing another vector norm in $\{3, \ldots, N\}$ leads to other notions of singular vectors and eigenvectors [34].

In the previous literature, there have been several works that show correspondence between CP decomposition and eigenvectors or singular vectors of a tensor. For latent variable models, under certain assumptions, an orthogonal CP decomposition is also an eigendecomposition of a symmetric tensor [3]. It has been shown that orthogonally decomposable symmetric tensors have orthogonal eigenvectors and that the CP decomposition of these tensors is given by orthogonal eigenvectors [51]. An extension of this result to nonsymmetric tensors with orthogonal singular vectors defining the CP decomposition has also been explored [52].

In this work, motivated by an efficient iterative scheme, we consider an extension of the variational notion of one singular vector tuple to many.

1.2. Convergence results. The new alternating update scheme is highly effective at finding an exact CP decomposition, if one exists. In particular, we show that the method achieves a superlinear rate of local convergence to exact CP decompositions if the CP rank is not greater than any mode length of the input tensor. In section 4, we prove that the convergence order is α per subproblem or α^N per sweep of alternating updates, where α is the unique real root of the polynomial $x^{N-1} - \sum_{i=0}^{N-2} x^i$. For N=3, $\alpha=(1+\sqrt{5})/2$, while for higher N, α increases. A superlinear convergence rate for CP decomposition is also achievable via general optimization algorithms such as Gauss–Newton [55, 57] which does not involve any restrictions on the CP rank. However, the alternating optimization scheme we propose has a much lower per-iteration cost (about the same as ALS, which achieves only linear convergence).

For a given tensor and any choice of rank, the critical points are generally not unique (as in the case of matrices). Theoretical characterization of the conditions under which critical points of (3.3) exist in this scenario remains an open problem, as it requires proving existence of roots of a system of nonlinear equations that have the same number of variables and equations. Note that the problem of proving if the best CP rank approximation exists also requires existence of a solution of system of nonlinear equations. The best CP rank approximation may not exist, which has led to the notion of border rank [31]. Consequently, establishing existence of critical points for our scenario is likely also nontrivial. Therefore, with the assumption that a critical point exists, we show in section 4.1 that the proposed iterative scheme achieves local convergence to it (in this case, at a linear rate).

We have performed numerical experiments to confirm the rate of convergence of the algorithm for computing CP decomposition of different tensors with known rank. The observed rate of convergence values agrees with the theoretical rate of convergence with an error of about 0.2% for order 3 and an error of about 1% for order 4 tensors. The numerical experiments in section 6 also confirm the convergence analysis of the algorithm to stationary points for approximate decomposition of tensors as described in Lemma 4.4.

1.3. Generalizations and experimental evaluation. The proposed algorithm may also be used for approximate CP decomposition but does not minimize the Frobenius norm of the residual directly. Instead, when optimizing for \boldsymbol{A} , the algorithm minimizes

$$\left\| (oldsymbol{I} \otimes oldsymbol{B}^\dagger \otimes oldsymbol{C}^\dagger) \operatorname{vec}(oldsymbol{\mathcal{T}} - \llbracket oldsymbol{A}, oldsymbol{B}, oldsymbol{C}
rbracket)
ight\|_F,$$

where $\text{vec}(\mathcal{T})$ is vectorization of the tensor \mathcal{T} and \otimes is the Kronecker product operation on two matrices. For a fixed residual error in the decomposition, the magnitude of this error metric will generally depend on the conditioning of A, B, and C. Hence, this alternating minimization procedure tends to converge to well-conditioned factors (and well-conditioned CP decompositions [10, 66]).

We generalize this method by considering a Mahalanobis distance between the input and reconstructed tensors. The original motivation for Mahalanobis distance minimization of tensors came from the work on minimization of Wasserstein distance between tensors for nonnegative CP decomposition [2]. This generalization allows us to reformulate each update to a factor of the above-introduced algorithm as a minimizer of a Mahalanobis distance [19] with the ground metric dependent on the other remaining factors. This reformulation helps extend the introduced algorithm to any CP rank by using the same ground metric. Moreover, we are able to interpolate between the introduced algorithm and ALS by interpolating the ground metric. Our experiments in section 6 suggest that the interpolated updates help manage the tradeoff between fitness and conditioning of the decomposition. We measure conditioning of the decomposition by computing the normalized CP decomposition condition number [10]. The condition number can be computed in an efficient manner for decompositions with small CP rank by reducing the size of the matrix for which the smallest singular value needs to be computed. The reduction in size is performed by removing some components of the matrix which contribute to larger singular value components and henceforth reducing the computational cost significantly. We present a proof of this reduction in Appendix A which is specific to CP decomposition and is different from the proof presented for structured block term decompositions [20]. By using this efficient approach, we are able to track the condition number of the decomposition at each iteration of the algorithms. For synthetic as well as real-world tensors, we observe that by utilizing hybrid updates of the introduced algorithm, we can find decompositions with a condition number lower by a factor as large as 10^4 with a change in relative residual of only about 0.01 when compared to ALS.

- 2. Background. We introduce the notation and definitions used in the subsequent sections here along with a brief introduction to the ALS algorithm for computing CP decomposition [12, 21].
- **2.1.** Notation and definitions. We use tensor algebra notation in both elementwise form and specialized form for tensor operations [31]. For vectors, bold

lowercase letters are used, e.g., \boldsymbol{x} . For matrices, bold uppercase letters are used, e.g., \boldsymbol{X} . For tensors, bold calligraphic fonts are used, e.g., $\boldsymbol{\mathcal{X}}$. An order N tensor corresponds to an N-dimensional array with dimensions $s_1 \times \cdots \times s_N$. Elements of vectors, matrices, and tensors are denoted in subscript, e.g., x_i for a vector \boldsymbol{x} , x_{ij} for a matrix \boldsymbol{X} and x_{ijkl} for an order 4 tensor $\boldsymbol{\mathcal{X}}$. The ith column of a matrix \boldsymbol{X} is denoted by \boldsymbol{x}_i . The mode-n matrix product of a tensor $\boldsymbol{\mathcal{X}} \in \mathbb{R}^{s_1 \times \cdots \times s_N}$ with a matrix $\boldsymbol{A} \in \mathbb{R}^{J \times s_n}$ is denoted by $\boldsymbol{\mathcal{X}} \times_n \boldsymbol{A}$, with the result having dimensions $s_1 \times \cdots \times s_{n-1} \times J \times s_{n+1} \times \cdots \times s_N$. Matricization is the process of reshaping a tensor into a matrix. Given a tensor $\boldsymbol{\mathcal{X}}$, the mode-n matricized version is denoted by $\boldsymbol{X}_{(n)} \in \mathbb{R}^{s_n \times K}$, where $K = \prod_{m=1, m \neq n}^N s_m$. We use parenthesized superscripts as labels for different tensors and matrices; e.g., $\boldsymbol{A}^{(1)}$ and $\boldsymbol{A}^{(2)}$ are different matrices.

The Hadamard product of two matrices $\boldsymbol{U}, \boldsymbol{V} \in \mathbb{R}^{I \times J}$ resulting in matrix $\boldsymbol{W} \in \mathbb{R}^{I \times J}$ is denoted by $\boldsymbol{W} = \boldsymbol{U} * \boldsymbol{V}$, where $w_{ij} = u_{ij}v_{ij}$. The inner product of matrices $\boldsymbol{U}, \boldsymbol{V}$ is denoted by $\langle \boldsymbol{U}, \boldsymbol{V} \rangle = \sum_{i,j} u_{ij}v_{ij}$, and similarly for tensors. The outer product of K vectors $\boldsymbol{u}^{(1)}, \dots, \boldsymbol{u}^{(K)}$ of corresponding sizes s_1, \dots, s_K is denoted by $\boldsymbol{\mathcal{X}} = \boldsymbol{u}^{(1)} \circ \dots \circ \boldsymbol{u}^{(K)}$, where $\boldsymbol{\mathcal{X}} \in \mathbb{R}^{s_1 \times \dots \times s_K}$ is an order K tensor.

For matrices $A \in \mathbb{R}^{I \times K} = [a_1, \dots, a_K]$ and $B \in \mathbb{R}^{J \times K} = [b_1, \dots, b_K]$, their Khatri–Rao product results in a matrix of size $(IJ) \times K$ defined by $A \odot B = [a_1 \otimes b_1, \dots, a_K \otimes b_K]$, where $a \otimes b$ denotes the Kronecker product of the two vectors. We define the Mahalanobis norm for a matrix A with ground metric M as $\|A\|_M^2 = \text{vec}(A)^T M \text{vec}(A)$, and similarly for a tensor $\mathcal{T}, \|\mathcal{T}\|_M^2 = \text{vec}(\mathcal{T})^T M \text{vec}(\mathcal{T})$. To ease the notation for N Khatri–Rao products, we use $\bigcirc_{n=1}^N = A^{(N)} \odot \ldots \odot A^{(1)}$, and similarly for Kronecker products of N matrices, $\bigotimes_{n=1}^N A^{(n)} = A^{(N)} \otimes \ldots \otimes A^{(1)}, *_{n=1}^N A^{(n)} = A^{(N)} * \ldots *_{n=1}^N A^{(n)}$. We use $\sigma_{\min}(P)$ to denote the minimum singular value of the matrix P.

2.2. ALS for CP decomposition. The CP tensor decomposition [21, 24] for an input tensor $\mathcal{X} \in \mathbb{R}^{I_1 \times \cdots \times I_N}$ is denoted by

$$\mathcal{X} pprox \left[\mathbf{A}^{(1)}, \cdots, \mathbf{A}^{(N)} \right], \quad ext{where} \quad \mathbf{A}^{(i)} = \left[\mathbf{a}_1^{(i)}, \cdots, \mathbf{a}_r^{(i)} \right],$$

and serves to approximate a tensor by a sum of R tensor products of vectors,

$$\mathcal{X} pprox \sum_{r=1}^{R} \boldsymbol{a}_r^{(1)} \circ \cdots \circ \boldsymbol{a}_r^{(N)}.$$

It is sometimes useful to normalize the factor matrices so that each column of the factors has a unit 2-norm and the weights are absorbed into a vector $\mathbf{d} \in \mathbb{R}^R$, given as

$$\mathcal{X} \approx \sum_{r=1}^{R} \boldsymbol{a}_r^{(1)} \circ \cdots \circ \boldsymbol{a}_r^{(N)} = \sum_{r=1}^{R} d_r \bar{\boldsymbol{a}}_r^{(1)} \circ \cdots \circ \bar{\boldsymbol{a}}_r^{(N)},$$

where $\bar{\boldsymbol{A}}^{(n)}$ are column normalized for all n and denoted by

$$\boldsymbol{\mathcal{X}} pprox \left[\boldsymbol{D}; \bar{\boldsymbol{A}}^{(1)}, \dots, \bar{\boldsymbol{A}}^{(N)} \right],$$

where D is a diagonal matrix with d on the diagonal. The CP-ALS method aims to minimize the nonlinear least squares problem

(2.1)
$$f(\boldsymbol{A}^{(1)}, \dots, \boldsymbol{A}^{(N)}) := \frac{1}{2} \left| \left| \boldsymbol{\mathcal{X}} - [\![\boldsymbol{A}^{(1)}, \dots, \boldsymbol{A}^{(N)}]\!] \right| \right|_{F}^{2}$$

by alternatively minimizing a sequence of least squares problems for each of the factor matrices $A^{(n)}$. This results in linear least squares problems for each row,

$$A_{\text{new}}^{(n)} P^{(n)T} \cong X_{(n)},$$

where the matrix $\mathbf{P}^{(n)} \in \mathbb{R}^{D_n \times R}$, where $D_n = \prod_{i=1, i \neq n}^N I_i$, is formed by Khatri–Rao products of the other factor matrices,

(2.2)
$$\boldsymbol{P}^{(n)} = \bigodot_{m=1, m \neq n}^{N} \boldsymbol{A}^{(m)}.$$

These linear least squares problems are often solved via the normal equations [31],

$$oldsymbol{A}_{ ext{new}}^{(n)} oldsymbol{\Gamma}^{(n)} \leftarrow oldsymbol{X}_{(n)} oldsymbol{P}^{(n)},$$

where $\Gamma \in \mathbb{R}^{R \times R}$ can be computed via

(2.3)
$$\Gamma^{(n)} = \underset{i=1, i \neq n}{\overset{N}{*}} S^{(i)}$$

with each $\mathbf{S}^{(i)} = \mathbf{A}^{(i)T} \mathbf{A}^{(i)}$. The matricized tensor times Khatri–Rao product or MT-TKRP computation $\mathbf{M}^{(n)} = \mathbf{X}_{(n)} \mathbf{P}^{(n)}$ is the main computational bottleneck of CP-ALS [5]. For a rank-R CP decomposition, this computation has the cost of $O(I^N R)$ if $I_n = I$ for all $n \in \{1, ..., N\}$. There have been various developments to optimize computation of MTTKRP, like the dimension tree algorithm [4, 27, 28, 29, 47, 67] for dense tensors and sparse MTTKRP [13] for sparse tensors.

- 3. Exact CP decomposition. In this section, we provide a motivation for the derivation of the new alternating optimization scheme to compute exact CP decomposition when the rank is known and is not greater than any of the mode lengths of the input tensor. The method is derived by extending the notion of computing singular vectors via Lagrangian optimization by considering more than one vector at a time by computing the fixed points of a more general function. These fixed points lead to a new alternating update scheme which solves a different least squares problem as compared to ALS. This new scheme has a leading order cost which is the same as ALS, depends on the MTTKRP operation, and therefore is efficient in computation.
- 3.1. Tensor spectral diagonalization via Lagrangian optimization. An equidimensional tensor of order N with mode length s is said to be spectrally diagonalized if it has s unit eigenvalues with s corresponding eigenvectors that are columns of identity. We may obtain a spectral diagonalization of the matrix \boldsymbol{A} by considering the critical points of a generalization of $\boldsymbol{x}^T \boldsymbol{A} \boldsymbol{y} = \langle \boldsymbol{A}, \boldsymbol{x} \boldsymbol{y}^T \rangle$ to the matrix case,

$$f(\boldsymbol{X}, \boldsymbol{Y}) = \langle \boldsymbol{A}, \boldsymbol{X} \boldsymbol{Y}^T \rangle \text{ s.t. } \det(\boldsymbol{X}^T \boldsymbol{X}) \neq 0, \det(\boldsymbol{Y}^T \boldsymbol{Y}) \neq 0.$$

Transforming the inequality constraint into a logarithmic barrier function, we obtain

(3.1)
$$\mathcal{L}_f(\boldsymbol{X}, \boldsymbol{Y}) = \langle \boldsymbol{A}, \boldsymbol{X} \boldsymbol{Y}^T \rangle - \frac{1}{2} \left(\log(\det(\boldsymbol{X}^T \boldsymbol{X})) + \log(\det(\boldsymbol{Y}^T \boldsymbol{Y})) \right)$$

$$= \operatorname{tr}(\boldsymbol{X}^T \boldsymbol{A} \boldsymbol{Y}) - \frac{1}{2} \operatorname{tr} \left(\log(\boldsymbol{X}^T \boldsymbol{X}) + \log(\boldsymbol{Y}^T \boldsymbol{Y}) \right).$$

The above equivalence is due to Jacobi's formula, as $\det(\mathbf{X}^T \mathbf{X}) \neq 0$, $\det(\mathbf{Y}^T \mathbf{Y}) \neq 0$. The critical points of \mathcal{L}_f satisfy

$$AYX^T \cong I \text{ and } A^TXY^T \cong I.$$

The above least squares equations arise from the fact that derivative of $\frac{1}{2}\operatorname{tr}\log(\boldsymbol{X}^T\boldsymbol{X})$ with respect to \boldsymbol{X} is $\boldsymbol{X}^{\dagger T}$. These critical points diagonalize \boldsymbol{A} in the sense that $\boldsymbol{X}^T\boldsymbol{A}\boldsymbol{Y}=\boldsymbol{I}$. In the case of a tensor $\boldsymbol{\mathcal{T}}$, a critical point $(\boldsymbol{X}^{(1)},\ldots,\boldsymbol{X}^{(N)})$ of (3.3)

$$f(\boldsymbol{X}^{(1)}, \dots, \boldsymbol{X}^{(N)}) = \langle \boldsymbol{\mathcal{T}}, [\![\boldsymbol{X}^{(1)}, \dots \boldsymbol{X}^{(N)}]\!] \rangle \text{ s.t. } \det(\boldsymbol{X}^{(n)T} \boldsymbol{X}^{(n)}) \neq 0 \,\forall n \in \{1, \dots, N\},$$

$$\mathcal{L}_f(\boldsymbol{X}^{(1)}, \dots, \boldsymbol{X}^{(N)}) = \sum_{r=1}^R \sum_{i_1 \dots i_N} t_{i_1 \dots i_N} x_{i_1 r}^{(1)} \dots x_{i_N r}^{(N)} - \frac{1}{2} \operatorname{tr} \left(\sum_{i=1}^N \log(\boldsymbol{X}^{(i)T} \boldsymbol{X}^{(i)}) \right)$$

gives invariant subspaces of \mathcal{T} in the sense that, for all $i \in \{1, ..., N\}$,

$$orall oldsymbol{v} \in \operatorname{span} \left\{ igotimes_{j
eq i}^{} oldsymbol{x}_1^{(j)}, \ldots, igotimes_{j
eq i}^{} oldsymbol{x}_R^{(j)}
ight\}, \quad oldsymbol{T}_{(i)} oldsymbol{v} \in \operatorname{span} \left\{ oldsymbol{x}_1^{(i)}, \ldots, oldsymbol{x}_R^{(i)}
ight\},$$

where $T_{(i)}$ is the mode-*i* matricization (unfolding) of the tensor \mathcal{T} . Note that the objective function in (3.3) is neither convex nor concave with respect to any of the factors X, Y, unlike the least squares loss function minimized in ALS which is convex with respect to the factor matrices.

The reconstructed tensor $\tilde{\mathcal{T}}$, where $\tilde{\mathcal{T}} = [\![\boldsymbol{Y}^{(1)}, \dots \boldsymbol{Y}^{(N)}]\!]$, $\boldsymbol{Y}^{(n)} = \boldsymbol{X}^{(n)\dagger}^T$ for all $n \in \{1, \dots, N\}$, captures the action of $\boldsymbol{T}_{(i)}$ on an invariant subspace as described above; the application of each matricization may be performed with bounded backward error. In section 5.1.1, we show that the backward error is bounded by $\|\boldsymbol{T}_{(i)}\boldsymbol{z}^{\perp}\|$, where \boldsymbol{z}^{\perp} is the projection of \boldsymbol{z} onto the orthogonal complement of the column span of $\sum_{j=1,j\neq n}^{N} \boldsymbol{X}^{(j)}$. We show that this bound also holds for ALS, and in general for a family of algorithms based on alternating minimization of Mahalanobis distance [19] between the input and reconstructed tensors.

Another observation regarding the critical points of (3.3) is that each matricization of the tensor reconstructed from a critical point, $\mathcal{X} = [\![\boldsymbol{X}^{(1)}, \dots \boldsymbol{X}^{(N)}]\!]$, is a right inverse of the corresponding matricization of the input tensor \mathcal{T} when the CP rank is equal to the mode length, and more generally,

$$\boldsymbol{T}_{(i)}\boldsymbol{X}_{(i)}^T = \boldsymbol{\Pi}^{(i)} \quad \forall i \in \{1,\dots,N\},$$

where each $\Pi^{(i)}$ is a projector onto the column space of $X^{(i)}$. This \mathcal{X} satisfies some but not all of the properties of previously proposed generalizations of the Moore–Penrose inverse to tensors [36, 58].

Further, the critical point gives a transformation that spectrally diagonalizes \mathcal{T} ,

$$z_{j_1...j_N} = \sum_{i_1...i_N} t_{i_1...i_N} x_{i_1j_1}^{(1)} \cdots x_{i_Nj_N}^{(N)},$$

so that \mathbb{Z} has R eigenvectors that are elementary vectors with unit eigenvalues (for any tensor eigenvector/eigenvalue definition, i.e., l^p eigenvector for any choice of p [34]), since

$$z_{j_k j_k \dots j_p j_k \dots j_k} = \delta_{j_k j_p} \quad \forall \, p \neq k \in \{1, \dots, N\}.$$

Beyond these properties, we show that $\boldsymbol{X}^{(1)\dagger^T},\dots,\boldsymbol{X}^{(N)\dagger^T}$ may be used to obtain an exact CP decomposition or an effective low-rank approximate CP decomposition.

3.2. Alternating optimization method. The critical points defined above may be computed efficiently by a method similar to ALS. For a CP decomposition $[\![A,B,C]\!]$ of tensor \mathcal{T} , ALS solves a set of overdetermined linear equations at each step to minimize the Frobenius norm error relative to one factor; e.g., it solves for A in

$$(\boldsymbol{C}\odot\boldsymbol{B})\boldsymbol{A}^T\cong\boldsymbol{T}_{(1)}^T.$$

The update rule for each subproblem may be written as a product of the pseudoinverse of the Khatri–Rao product of two factors and an unfolding of the tensor, i.e.,

$$\boldsymbol{A} = \boldsymbol{T}_{(1)} (\boldsymbol{C} \odot \boldsymbol{B})^{\dagger T}.$$

Some of the major advantages of the ALS algorithm are its guaranteed monotonic decrease in residual, low per-iteration computational cost and its amenability to parallelization. It has been shown that ALS achieves linear local convergence to minima of the CP residual norm under certain conditions [64].

To obtain a critical point in the high order tensor function (3.3), we propose a different alternating update scheme, which finds the solution U to the linear least squares problem,

$$(\boldsymbol{T}_{(1)}(\boldsymbol{W}\odot\boldsymbol{V}))\boldsymbol{U}^T\cong\boldsymbol{I}.$$

With $A^T = U^{\dagger}$, $B^T = V^{\dagger}$, and $C^T = W^{\dagger}$, we observe that the update is similar to that of ALS,

$$\boldsymbol{A} = \boldsymbol{T}_{(1)}(\boldsymbol{C}^{\dagger T} \odot \boldsymbol{B}^{\dagger T}).$$

A stationary point of this alternating update scheme provides a critical point of (3.3).

We first provide a complete description of the alternating update scheme proposed above. To compute the decomposition of a tensor \mathcal{X} , the algorithm maintains a CP decomposition given by

$$\left[\!\left[oldsymbol{A}^{(1)},\ldots,oldsymbol{A}^{(N)}
ight]\!
ight]$$

and updates each $A^{(n)}$ in an alternating manner,

(3.4)
$$\boldsymbol{A}^{(n)} = \boldsymbol{X}_{(n)} \left(\bigodot_{m=1, m \neq n}^{N} \boldsymbol{A}^{(m)\dagger T} \right).$$

This update may be written in elementwise form in terms of the pseudoinverses $U^{(m)} = A^{(m)\dagger}$ as

$$a_{i_n r}^{(n)} = \sum_{i_1 \dots \hat{i}_n \dots i_N} x_{i_1 \dots i_N} \prod_{m=1, m \neq n}^N u_{r i_m}^{(m)}.$$

Algorithm 3.1 details each sweep of such updates. As with the ALS, it is advisable to recalibrate the norms of the columns of each factor before the subsweep corresponding to the *n*th factor so that $\|\boldsymbol{a}_r^{(k)}\| = 1$ for all $k \neq n$ and r. With this recalibration, a valid convergence criterion is to check whether the magnitude of change

Algorithm 3.1 Basic description of the new alternating update scheme.

```
1: Input: Tensor \mathcal{X} \in \mathbb{R}^{I_1 \times \cdots \times I_N}, rank R \leq I_n \forall n \in \{1, \dots, N\}

2: Initialize \{A^{(1)}, \dots, A^{(N)}\} so each A^{(n)} \in \mathbb{R}^{I_n \times R} is random

3: while not converged do

4: for n \in \{1, \dots, N\} do

5: A^{(n)} = X_{(n)}^{(N)} \left( \bigcirc_{m=1, m \neq n}^{N} A^{(m) \dagger T} \right)

6: end for

7: end while

8: return factor matrices \{A^{(1)}, \dots, A^{(N)}\}
```

in the factors exceeds a predefined threshold at each sweep. The method is invariant to the rescaling in exact arithmetic, but this calibration helps reduce the effects of round-off error as in [64]. This calibration is also cost efficient, since the pseudoinverse of only one matrix changes per subweep. This makes the algorithm accessible to all the optimizations involved in computing the MTTKRP in each subsweep of ALS such as the dimension tree algorithm [4, 27, 28, 29, 47, 67].

3.3. Cost analysis. The cost of each sweep of Algorithm 3.1 corresponds to the cost of computing the pseudoinverse of each factor, as well as a set of N MTTKRP operations computed across different modes. For computing a set of MTTKRP operations for a sweep of the algorithm, a dimension tree [4, 28, 47] or multisweep dimension tree [39] may be used as in the ALS algorithm. Dimension trees provide an efficient way to compute the set of MTTKRP in one sweep by storing and reusing partial computations performed in computing MTTKRP in each subsweep of the sweep. A multisweep dimension tree improves the computation by reusing intermediates across sweeps, reducing the cost per sweep even further. The overall per-sweep cost with a multisweep dimension tree [39] is then given by

$$\frac{2N}{N-1}\Biggl(\prod_{n=1}^{N}I_{n}\Biggr)R+O\left(\left(\sum_{n=1}^{N}I_{n}\right)R^{2}\right).$$

The cost of an ALS sweep with a multisweep dimension tree [39] is

$$\frac{2N}{N-1} \left(\prod_{n=1}^{N} I_n \right) R + O(NR^3),$$

which is less expensive, as solving an overdetermined system via normal equations is cheaper than computing the pseudoinverse of a matrix. If \mathcal{X} is sparse, the method can benefit from existing work on efficiently performing MTTKRP with a sparse tensor [13] and therefore has the same leading order cost per sweep as that of ALS for sparse tensors as well.

4. Convergence rate for exact decomposition. In this section, we theoretically analyze the asymptotic rate of local convergence of Algorithm 3.1 for when an exact CP decomposition of rank R less than equal to the length of all modes of the tensor exists. To derive the rate of convergence for an exact decomposition, we relate the distance between the computed factor in each subproblem and the true factor, to the error in the other factor matrices. The following lemma states the error in

computing $\mathbf{A}^{(N)}$ but can be trivially extended to any $\mathbf{A}^{(n)}$ for $n \in \{1, ..., n\}$. We consider the error in the normalized factor and the error in the magnitude of CP decomposition components modulo column scaling.

LEMMA 4.1. Suppose $\mathcal{X} = [\![\boldsymbol{D}; \boldsymbol{A}^{(1)}, \ldots, \boldsymbol{A}^{(N)}]\!]$, where each $\boldsymbol{A}^{(i)} \in \mathbb{R}^{s_i \times R}$ with $s_i \geq R$ is full rank and has normalized columns, i.e., $\|\boldsymbol{a}_j^{(i)}\|_2 = 1$ for all i, j and $\bar{\boldsymbol{A}}^{(n)} = \boldsymbol{A}^{(n)} + \boldsymbol{\Delta}^{(n)}$ also has normalized columns and satisfies $\|\boldsymbol{\Delta}^{(n)}\|_F = \epsilon_n$ for $n = 1, \ldots, N-1$. Then there exists $\epsilon > 0$ such that if $\epsilon_n < \epsilon$ for $n = 1, \ldots, N-1$, then

$$\tilde{\boldsymbol{A}}^{(N)} = \boldsymbol{X}_{(N)} \left(\bar{\boldsymbol{A}}^{(1)\dagger T} \odot \cdots \odot \bar{\boldsymbol{A}}^{(N-1)\dagger T} \right),$$

where $\tilde{\boldsymbol{A}}^{(N)} = \bar{\boldsymbol{A}}^{(N)} \bar{\boldsymbol{D}}$ ensures that $\bar{\boldsymbol{A}}^{(N)}$ are normalized and satisfies

$$\|\bar{\boldsymbol{A}}^{(N)} - \boldsymbol{A}^{(N)}\|_F = O\left(\prod_{n=1}^{N-1} \epsilon_n\right)$$

and $\|\bar{\boldsymbol{D}} - \boldsymbol{D}\|_F = O(\epsilon)$.

Proof. Let $\epsilon < 1$ and $\epsilon < \min_n(\sigma_{\min}(\boldsymbol{A}^{(n)}))$ for each n; therefore $\bar{\boldsymbol{A}}^{(n)}$ is full rank for each $n = 1, \ldots, N-1$. Substituting the decomposition of $\boldsymbol{\mathcal{X}}$ into the computed solution, we obtain

$$\tilde{\boldsymbol{A}}^{(N)} = \boldsymbol{A}^{(N)} \boldsymbol{D} \left((\bar{\boldsymbol{A}}^{(1)\dagger} (\bar{\boldsymbol{A}}^{(1)} - \boldsymbol{\Delta}^{(1)})) * \cdots * (\bar{\boldsymbol{A}}^{(N-1)\dagger} (\bar{\boldsymbol{A}}^{(N-1)} - \boldsymbol{\Delta}^{(N-1)})) \right)^{T}$$

$$= \boldsymbol{A}^{(N)} \boldsymbol{D} \left((\boldsymbol{I} - \bar{\boldsymbol{A}}^{(1)\dagger} \boldsymbol{\Delta}^{(1)}) * \cdots * (\boldsymbol{I} - \bar{\boldsymbol{A}}^{(N-1)\dagger} \boldsymbol{\Delta}^{(N-1)}) \right)^{T}$$

$$= \boldsymbol{A}^{(N)} \boldsymbol{D} \left(\boldsymbol{S} + (-1)^{N-1} \bar{\boldsymbol{A}}^{(1)\dagger} \boldsymbol{\Delta}^{(1)} * \cdots * \bar{\boldsymbol{A}}^{(N-1)\dagger} \boldsymbol{\Delta}^{(N-1)} \right)^{T},$$

where S includes all cross terms of the Hadamard products, which must be diagonal since any such term includes a Hadamard product with an identity matrix. Since

$$\|\boldsymbol{I} - \boldsymbol{S}\|_F = O(\max_n(\epsilon_n)) = O(\epsilon),$$

S is full rank for sufficiently small ϵ . Let

$$\boldsymbol{\Delta} = \left((-1)^{N-1} \bar{\boldsymbol{A}}^{(1)\dagger} \boldsymbol{\Delta}^{(1)} * \cdots * \bar{\boldsymbol{A}}^{(N-1)\dagger} \boldsymbol{\Delta}^{(N-1)} \right)^{T}.$$

Now, the norm calibration diagonal matrix is defined so that $\bar{d}_{ii} = \|\tilde{\boldsymbol{a}}_i^{(N)}\|_2$. Since

$$\tilde{\boldsymbol{A}}^{(N)} = \boldsymbol{A}^{(N)} \boldsymbol{D} (\boldsymbol{S} + \boldsymbol{\Delta}) = \boldsymbol{A}^{(n)} \boldsymbol{D} + \boldsymbol{A}^{(n)} \boldsymbol{D} (\boldsymbol{S} + \boldsymbol{\Delta} - \boldsymbol{I})$$

and $\|S + \Delta - I\|_2 = O(\epsilon)$, we have

$$\bar{d}_{ii} = \|\tilde{\boldsymbol{a}}_{i}^{(N)}\|_{2} \leq \|\boldsymbol{a}_{i}^{(N)}\|_{2} d_{ii} + \|\boldsymbol{A}^{(N)}\boldsymbol{D}\|_{F} \|\boldsymbol{S} + \boldsymbol{\Delta} - \boldsymbol{I}\|_{2}$$
$$= \|\boldsymbol{a}_{i}^{(N)}\|_{2} d_{ii} + O(\epsilon) = d_{ii} + O(\epsilon).$$

Consequently, $\|\bar{\boldsymbol{D}} - \boldsymbol{D}\|_2 = O(\epsilon)$. Further, we can obtain a tighter bound (in terms of $O(\|\boldsymbol{\Delta}\|_F)$ instead of $O(\epsilon)$ by considering $\tilde{\boldsymbol{A}}^{(N)} = \boldsymbol{A}^{(N)} \boldsymbol{D} \boldsymbol{S} (\boldsymbol{I} + \boldsymbol{S}^{-1} \boldsymbol{\Delta})$, so

$$\bar{d}_{ii} = \|\tilde{\boldsymbol{a}}_{i}^{(N)}\|_{2} \leq \|\boldsymbol{a}_{i}^{(N)}\|_{2}d_{ii}s_{ii} + \|\boldsymbol{A}^{(N)}\boldsymbol{D}\boldsymbol{\Delta}\|_{F} = d_{ii}s_{ii} + \|\boldsymbol{D}\|_{2}\|\boldsymbol{\Delta}\|_{F}.$$

This bound allows us to get the desired result for the error in the factor matrices,

$$\| \boldsymbol{A}^{(N)} - \bar{\boldsymbol{A}}^{(N)} \|_{F} = \| \boldsymbol{A}^{(N)} - \tilde{\boldsymbol{A}}^{(N)} \bar{\boldsymbol{D}}^{-1} \|_{F} = \| \boldsymbol{A}^{(N)} - \boldsymbol{A}^{(N)} \boldsymbol{D} \boldsymbol{S} (\boldsymbol{I} + \boldsymbol{S}^{-1} \boldsymbol{\Delta}) \bar{\boldsymbol{D}}^{-1} \|_{F}$$

$$= O \left(\| \boldsymbol{I} - \boldsymbol{D} \boldsymbol{S} (\boldsymbol{I} + \boldsymbol{S}^{-1} \boldsymbol{\Delta}) \bar{\boldsymbol{D}}^{-1} \|_{F} \right).$$

$$(4.1)$$

Since we have that $\|\bar{\boldsymbol{D}} - \boldsymbol{D}\boldsymbol{S}\|_F = \|\boldsymbol{D}\|_2 \|\boldsymbol{\Delta}\|_F$,

$$\begin{split} \left\| \boldsymbol{I} - \boldsymbol{D} \boldsymbol{S} (\boldsymbol{I} + \boldsymbol{S}^{-1} \boldsymbol{\Delta}) \bar{\boldsymbol{D}}^{-1} \right\|_{F} &= \left\| \boldsymbol{I} - \boldsymbol{D} \boldsymbol{S} \bar{\boldsymbol{D}}^{-1} + \boldsymbol{D} \boldsymbol{\Delta} \bar{\boldsymbol{D}}^{-1} \right\|_{F} \\ &= \left\| \bar{\boldsymbol{D}} \bar{\boldsymbol{D}}^{-1} - \boldsymbol{D} \boldsymbol{S} \bar{\boldsymbol{D}}^{-1} + \boldsymbol{D} \boldsymbol{\Delta} \bar{\boldsymbol{D}}^{-1} \right\|_{F} \\ &= O(\|\boldsymbol{\Delta}\|_{F}) + \left\| \boldsymbol{D} \boldsymbol{\Delta} \bar{\boldsymbol{D}}^{-1} \right\|_{F} = O(\|\boldsymbol{\Delta}\|_{F}). \end{split}$$

Since $\|\Delta\|_F = O(\prod_{n=1}^{N-1} \epsilon_n)$, this completes the proof.

The above lemma states that in Algorithm 3.1, the error in the updated CP decomposition factor relative to the true CP decomposition factor is bounded by the product of errors in the previous N-1 factors. Using this error bound, we derive the convergence rate for Algorithm 3.1.

LEMMA 4.2. For any algorithm where the error in the update is of the order of the product of errors in the previous k updates, the rate of convergence is equal to the positive root of the polynomial $\alpha^k - \sum_{i=0}^{k-1} \alpha^i$.

Proof. Let the error at the *n*th iteration be given as e_n . The worst case error at the *n*th iteration then satisfies the following due to Lemma 4.1:

(4.2)
$$e_n = L \prod_{i=1}^k e_{n-i}, \quad \text{where } L \text{ is a constant.}$$

The above recurrence can be solved by assuming that the error satisfies the following asymptotic relation:

$$(4.3) e_n = Ce_{n-1}^{\alpha},$$

where C is some constant and α is the rate of convergence. From (4.2) and (4.3),

$$\begin{split} Ce_{n-1}^{\alpha} &= L \prod_{i=1}^{k} e_{n-i}, \\ \text{by assumption in (4.3), } C^{\frac{1}{\alpha^p}} e_{n-1-p} &= e_{n-1}^{\frac{1}{\alpha^p}} \quad \forall p \in \{1, \dots k\}, \\ \frac{C^{\sum_{i=0}^{k-1} \frac{1}{\alpha^i}}}{L} &= e_{n-1}^{(-\alpha + \sum_{i=0}^{k-1} \frac{1}{\alpha^i})}. \end{split}$$

Since the left-hand side is constant for $n \to \infty$ due to the assumed convergence rate,

$$\alpha^k - \sum_{i=0}^{k-1} \alpha^i = 0.$$

Now, with all the pieces together we can show that for exact CP decomposition, Algorithm 3.1 locally converges at a rate which is given in the following theorem.

THEOREM 4.3. Suppose $\mathcal{X} = [\![\boldsymbol{D}; \boldsymbol{A}^{(1)}, \dots, \boldsymbol{A}^{(N)}]\!]$, where each $\boldsymbol{A}^{(i)} \in \mathbb{R}^{s_i \times R}$ with $s_i \geq R$ is full rank and has normalized columns. Algorithm 3.1 for computing the

exact CP decomposition of \mathcal{X} converges locally with a rate of convergence equal to α^N , where α is the unique real root of the polynomial $x^{N-1} - \sum_{i=0}^{N-2} x^i$.

Proof. Consider the computation of the CP decomposition of the tensor \mathcal{X} with exact rank R and initial guess $\bar{A}^{(i)}$ with normalized columns such that $\bar{A}^{(i)} = A^{(i)} + \Delta^{(i)}$ with $\|\Delta^{(i)}\|_2 < \epsilon$ sufficiently small for $i = 1, \ldots, N-1$ (as described in Lemma 4.1). Let the error in CP decomposition at the nth iteration be given as

$$E_n = \max \left\{ \|\bar{\boldsymbol{D}} - \boldsymbol{D}\|_F, \|\bar{\boldsymbol{A}}^{(1)} - \boldsymbol{A}^{(1)}\|_F, \dots, \|\bar{\boldsymbol{A}}^{(N)} - \boldsymbol{A}^{(N)}\|_F \right\}.$$

Also, let the error in the kth subiteration of the nth iteration be given as $\epsilon_n^{(k)}$. From Lemma 4.1, we know that the error in CP decomposition is bounded by the maximum error in factor matrices. Since the error decreases at each subiteration, there exists n such that $E_n = O(\epsilon_n^{(1)})$.

From Lemma 4.1, we know that the error in a subsweep of the algorithm modulo the column scaling is of the order of the product of errors in previous N-1 subsweeps; therefore the error at the (n+1)th iteration is bounded by the error in the first factor matrix, given by

$$E_{n+1} = O\left(\prod_{k=1}^{N-1} \epsilon_n^{(k)}\right).$$

By using Lemma 4.2, we know that the error in a subiteration is given by the following recurrence, where α is the positive root of $x^{N-1} - \sum_{i=0}^{N-2} x^i$:

$$\epsilon_n^{(k+1)} = O((\epsilon_n^{(k)})^{\alpha}) \quad \forall k = 1, \dots, N-1.$$

Therefore, the error at the (n+1)th iteration of Algorithm 3.1 can be expressed as

$$E_{n+1} = O\left(\prod_{i=1}^{N-1} \epsilon_1^{\alpha^i}\right) = O\left(E_n^{\sum_{i=1}^{N-1} \alpha^i}\right).$$

The fact that α is a root of the polynomial $x^{N-1} - \sum_{i=0}^{N-2} x^i$ implies that $\alpha^{N-1} - \sum_{i=0}^{N-2} \alpha^i = 0$; that is, $\sum_{i=1}^{N-1} \alpha^i = \alpha^N$. Therefore,

$$E_{n+1} = O\left(E_n^{\alpha^N}\right).$$

This completes the proof to show that Algorithm 3.1 locally converges superlinearly for exact CP rank cases. We verify our theoretical results in the section 6.

4.1. Convergence to other stationary points. The result in Theorem 4.3 can be generalized to the case where a tensor \mathcal{X} can be represented as the sum of two tensors, \mathcal{T} and \mathcal{E} , where \mathcal{T} has an underlying CP decomposition structure of rank R and \mathcal{E} has a CP decomposition that is mostly orthogonal to the decomposition of \mathcal{T} . For such an input tensor \mathcal{X} , we show that Algorithm 3.1 with CP rank R locally converges to the underlying CP factors, provided that a stationary point associated with \mathcal{T} exists. We analyze the convergence rate of the algorithm and show that this is a generalization of the previous result, since we converge to a subset of the CP factors with same convergence rate as in Theorem 4.3, if the factors of \mathcal{T} are in the orthogonal complement of the column space of corresponding CP factors of \mathcal{E} . We show that the error in the normalized factor matrix is a summation of two terms which have separate upper bounds. We keep these separate since the upper bounds are in a different parameter space for each term.

LEMMA 4.4. For a given tensor \mathcal{X} , assume there exists a stationary point of Algorithm 3.1 yielding a positive diagonal matrix and factors $(\mathbf{A}^{(1)}, \dots, \mathbf{A}^{(N)})$, where each $\mathbf{A}^{(i)} \in \mathbb{R}^{s_i \times R}$ with $s_i \geq R$ is full rank and has normalized columns; i.e., $\|\mathbf{a}_j^{(i)}\|_2 = 1$ for all i, j, and their pseudoinverse transposes $(\mathbf{U}^{(1)}, \dots, \mathbf{U}^{(N)})$, so $\mathbf{A}^{(n)\dagger} = \mathbf{U}^{(n)T}$. The stationary point conditions imply that for all $n \in \{1, \dots, N\}$, we have

$$m{A}^{(n)}m{D} = m{X}_{(n)} \left(igodots_{m=1,m
eq n}^{N} m{U}^{(m)}
ight).$$

Further, assume that $\mathcal{X} = [\![\boldsymbol{D}; \boldsymbol{A}^{(1)}, \dots, \boldsymbol{A}^{(N)}]\!] + [\![\hat{\boldsymbol{D}}; \hat{\boldsymbol{A}}^{(1)}, \dots, \hat{\boldsymbol{A}}^{(N)}]\!]$ with $\|\hat{\boldsymbol{A}}^{(n)}T\boldsymbol{U}^{(n)}\|_F \leq \epsilon_{\perp}$, and define $k = \|\hat{\boldsymbol{D}}\|_2\|\boldsymbol{D}^{-1}\|_2$. Given approximations $\tilde{\boldsymbol{A}}^{(n)} = \boldsymbol{A}^{(n)} + \boldsymbol{\Delta}_A^{(n)}$ and $\tilde{\boldsymbol{U}}^{(n)} = \tilde{\boldsymbol{A}}^{(n)\dagger T} = \boldsymbol{U}^{(n)} + \boldsymbol{\Delta}_U^{(n)}$ with normalized columns, there exists $\epsilon > 0$ such that if $\|\boldsymbol{\Delta}_A^{(n)}\|_F, \|\boldsymbol{\Delta}_U^{(n)}\|_F \leq \epsilon_n \leq \epsilon$ for all $n \in \{1, \dots, N-1\}$, then

$$\tilde{\boldsymbol{A}}^{(N)} = \boldsymbol{X}_{(N)} \left(\tilde{\boldsymbol{U}}^{(1)} \odot \cdots \odot \tilde{\boldsymbol{U}}^{(N-1)} \right)$$

satisfies $\|\tilde{\boldsymbol{A}}^{(N)}\bar{\boldsymbol{D}}^{-1}-\boldsymbol{A}^{(N)}\|_F=P+Q$, where $P=O(\prod_{i=1}^{N-1}\epsilon_i)$, $Q=O(k\epsilon\epsilon_\perp)$, and $\bar{\boldsymbol{D}}$ normalizes $\tilde{\boldsymbol{A}}^{(N)}$, i.e., $\bar{d}_{ii}=\|\tilde{\boldsymbol{a}}_i^{(N)}\|_2$. Further, $\|\bar{\boldsymbol{D}}-\boldsymbol{D}\|_F=O(\epsilon)$.

Proof. We expand the update as follows:

$$\begin{split} \tilde{\boldsymbol{A}}^{(N)} &= \left([\![\boldsymbol{A}^{(1)}, \dots, \boldsymbol{A}^{(N)} \boldsymbol{D}]\!] + [\![\hat{\boldsymbol{A}}^{(1)}, \dots, \hat{\boldsymbol{A}}^{(N)} \hat{\boldsymbol{D}}]\!] \right)_{(N)} \left(\tilde{\boldsymbol{U}}^{(1)} \odot \cdots \odot \tilde{\boldsymbol{U}}^{(N-1)} \right) \\ &= \boldsymbol{A}^{(N)} \boldsymbol{D} \left(\boldsymbol{A}^{(1)T} \tilde{\boldsymbol{U}}^{(1)} * \cdots * \boldsymbol{A}^{(N-1)T} \tilde{\boldsymbol{U}}^{(N-1)} \right) \\ &+ \hat{\boldsymbol{A}}^{(N)} \hat{\boldsymbol{D}} \left(\hat{\boldsymbol{A}}^{(1)T} \tilde{\boldsymbol{U}}^{(1)} * \cdots * \hat{\boldsymbol{A}}^{(N-1)T} \tilde{\boldsymbol{U}}^{(N-1)} \right) \\ &= \boldsymbol{A}^{(N)} \boldsymbol{D} \left(\left(\boldsymbol{I} + \boldsymbol{A}^{(1)T} \boldsymbol{\Delta}_{\boldsymbol{U}}^{(1)} \right) * \cdots * \left(\boldsymbol{I} + \boldsymbol{A}^{(N-1)T} \boldsymbol{\Delta}_{\boldsymbol{U}}^{(N-1)} \right) \right) \\ &+ \hat{\boldsymbol{A}}^{(N)} \hat{\boldsymbol{D}} \left(\left(\hat{\boldsymbol{A}}^{(1)T} \boldsymbol{U}^{(1)} + \hat{\boldsymbol{A}}^{(1)T} \boldsymbol{\Delta}_{\boldsymbol{U}}^{(1)} \right) * \cdots * \left(\hat{\boldsymbol{A}}^{(N-1)T} \boldsymbol{U}^{(N-1)} \right) \\ &+ \hat{\boldsymbol{A}}^{(N-1)T} \boldsymbol{\Delta}_{\boldsymbol{U}}^{(N-1)} \right) \right). \end{split}$$

By the stationary point condition, we have that

$$\hat{A}^{(N)}\hat{D}\left(\hat{A}^{(1)T}U^{(1)}*\cdots*\hat{A}^{(N-1)T}U^{(N-1)}\right)=0.$$

Consequently, the error reduces to the cross terms of the summations, i.e., 1

$$\begin{split} & \tilde{\boldsymbol{A}}^{(N)} - \boldsymbol{A}^{(N)} \boldsymbol{D} \\ & = \boldsymbol{A}^{(N)} \boldsymbol{D} \left[\overset{N}{\underset{n=1}{\times}} \boldsymbol{A}^{(n)T} \boldsymbol{\Delta}_{\boldsymbol{U}}^{(n)} + \boldsymbol{I} * \sum_{k=1}^{N-1} \left(\sum_{\{m_1, \dots, m_k\} \subset \{1, \dots, N\}} \overset{k}{\underset{l=1}{\times}} \boldsymbol{A}^{(m_l)T} \boldsymbol{\Delta}_{\boldsymbol{U}}^{(m_l)} \right) \right] \\ & + \hat{\boldsymbol{A}}^{(N)} \hat{\boldsymbol{D}} \left[\sum_{k=1}^{N-1} \left(\sum_{\substack{\{m_1, \dots, m_k\} \subset \{1, \dots, N\}, \\ \{w_1, \dots, w_{N-k}\} = \{1, \dots, N\} \setminus \{m_1, \dots, m_k\}} \overset{k}{\underset{l=1}{\times}} \hat{\boldsymbol{A}}^{(m_l)T} \boldsymbol{U}^{(m_l)} \overset{N-k}{\underset{l=1}{\times}} \hat{\boldsymbol{A}}^{(w_l)T} \boldsymbol{\Delta}_{\boldsymbol{U}}^{(w_l)} \right) \right]. \end{split}$$

$$A^{(3)}D\left[A^{(1)T}\Delta_{U}^{(1)}*A^{(2)T}\Delta_{U}^{(2)} + I*(A^{(1)T}\Delta_{U}^{(1)} + A^{(2)T}\Delta_{U}^{(2)})\right]$$

$$+ \hat{A}^{(3)}\hat{D}\left[\hat{A}^{(1)T}U^{(1)}*\hat{A}^{(2)T}\Delta_{U}^{(2)} + \hat{A}^{(2)T}U^{(2)}*\hat{A}^{(1)T}\Delta_{U}^{(1)}\right].$$

¹For N=3, the right-hand side of this formula is

First, since each error term has Frobenius norm $O(\|\Delta^{(n)}\|) = O(\epsilon)$, the column norms of $\tilde{\boldsymbol{A}}^{(N)}$ will yield $\bar{\boldsymbol{D}}$ with $\|\bar{\boldsymbol{D}} - \boldsymbol{D}\|_F = O(\epsilon)$. Further, in order to get the error bound on $\tilde{\boldsymbol{A}}^{(N)}$, we consider the diagonal matrix,

$$oldsymbol{S} = oldsymbol{I} + oldsymbol{I} * \sum_{k=1}^{N-1} \left(\sum_{\{m_1,...,m_k\} \subset \{1,...,N\}} \mathop{\sharp}\limits_{l=1}^k oldsymbol{A}^{(m_l)T} oldsymbol{\Delta}_U^{(m_l)}
ight).$$

We define S as above in order to show that the column norms of $\tilde{A}^{(N)}$, i.e., \bar{D} , are close to DS. Using the above definition of S, we have

$$\begin{split} \boldsymbol{A}^{(N)}\boldsymbol{D} + \boldsymbol{A}^{(N)}\boldsymbol{D} \begin{bmatrix} \sum\limits_{\substack{k=1 \\ k=1}}^{N} \boldsymbol{A}^{(n)T}\boldsymbol{\Delta}_{U}^{(n)} + \boldsymbol{I} * \sum\limits_{k=1}^{N-1} \left(\sum\limits_{\{m_{1},...,m_{k}\} \subset \{1,...,N\}} \overset{k}{\underset{l=1}{\overset{k}{\Rightarrow}}} \boldsymbol{A}^{(m_{l})T}\boldsymbol{\Delta}_{U}^{(m_{l})} \right) \end{bmatrix} \\ = \boldsymbol{A}^{(N)}\boldsymbol{D}\boldsymbol{S} + \boldsymbol{A}^{(N)}\boldsymbol{D}\boldsymbol{S} \begin{bmatrix} \boldsymbol{S}^{-1} & \sum\limits_{n=1}^{N} \boldsymbol{A}^{(n)T}\boldsymbol{\Delta}_{U}^{(n)} \end{bmatrix}. \end{split}$$

Therefore, we have

$$\begin{split} & \hat{\pmb{A}}^{(N)} - \pmb{A}^{(N)} \pmb{D} \pmb{S} = \pmb{A}^{(N)} \pmb{D} \pmb{S} \Bigg[\pmb{S}^{-1} \overset{N}{\underset{n=1}{*}} \pmb{A}^{(n)T} \pmb{\Delta}_{U}^{(n)} \Bigg] \\ & + \hat{\pmb{A}}^{(N)} \hat{\pmb{D}} \left[\sum_{k=1}^{N-1} \left(\sum_{\substack{\{m_{1}, \dots, m_{k}\} \subset \{1, \dots, N\}, \\ \{w_{1}, \dots, w_{N-k}\} = \{1, \dots, N\} \setminus \{m_{1}, \dots, m_{k}\}} \overset{k}{\underset{l=1}{*}} \hat{\pmb{A}}^{(m_{l})T} \pmb{U}^{(m_{l})} \overset{N-k}{\underset{l=1}{*}} \hat{\pmb{A}}^{(w_{l})T} \pmb{\Delta}_{U}^{(w_{l})} \right) \right], \end{split}$$

$$\|\tilde{\boldsymbol{A}}^{(N)} - \boldsymbol{A}^{(N)} \boldsymbol{D} \boldsymbol{S}\|_F = \|\hat{\boldsymbol{P}}\|_F + \|\hat{\boldsymbol{Q}}\|_F,$$

where $\|\hat{\boldsymbol{P}}\|_F = O(\prod_{i=1}^{N-1} \epsilon_i)$ and $\|\hat{\boldsymbol{Q}}\|_F = O(\epsilon \epsilon_\perp)$. The first term in the above equations is a Hadamard product of terms $\boldsymbol{A}^{(i)T}\boldsymbol{\Delta}_U^{(i)}$ and therefore leads to an upper bound of $O(\prod_{i=1}^{N-1} \epsilon_i)$. The second term in the above equations is a summation of terms with products $(\hat{\boldsymbol{A}}^{(i)T}\boldsymbol{U}^{(i)})\hat{\boldsymbol{A}}^{(j)T}\boldsymbol{\Delta}_U^{(j)}$ which results in an upper bound of $O(\epsilon \epsilon_\perp)$.

The same bound follows for each column, so the column norms of $\tilde{\boldsymbol{A}}^{(N)}$, $\bar{\boldsymbol{D}}$ satisfy

$$\|\bar{\boldsymbol{D}} - \boldsymbol{D}\boldsymbol{S}\|_F = \|\bar{\boldsymbol{P}}\|_F + \|\bar{\boldsymbol{Q}}\|_F,$$

where $\|\bar{\boldsymbol{P}}\|_F = O(\prod_{i=1}^{N-1} \epsilon_i)$ and $\|\bar{\boldsymbol{Q}}\|_F = O(\epsilon \epsilon_\perp)$.

Now that we have shown that $\tilde{\boldsymbol{A}}^{(N)}$ is close to $\boldsymbol{A}^{(N)}\boldsymbol{DS}$ and that its column norms $(\bar{\boldsymbol{D}})$ are close to \boldsymbol{DS} , we obtain the final error bound after normalization,

$$\begin{split} \tilde{\boldsymbol{A}}^{(N)} \bar{\boldsymbol{D}}^{-1} - \boldsymbol{A}^{(N)} \boldsymbol{D} \boldsymbol{S} \bar{\boldsymbol{D}}^{-1} &= \boldsymbol{A}^{(N)} \boldsymbol{D} \begin{bmatrix} \sum\limits_{k=1}^{N} \boldsymbol{A}^{(n)T} \boldsymbol{\Delta}_{U}^{(n)} \end{bmatrix} \bar{\boldsymbol{D}}^{-1} \\ &+ \hat{\boldsymbol{A}}^{(N)} \hat{\boldsymbol{D}} \begin{bmatrix} \sum\limits_{k=1}^{N-1} \left(\sum\limits_{\{m_{1}, \dots, m_{k}\} \subset \{1, \dots, N\}, \\ \{w_{1}, \dots, w_{N-k}\} = \{1, \dots, N\} \setminus \{m_{1}, \dots, m_{k}\} \end{bmatrix}} \overset{k}{\underset{l=1}{\overset{k}{=}}} \hat{\boldsymbol{A}}^{(m_{l})T} \boldsymbol{U}^{(m_{l})} \overset{N-k}{\underset{l=1}{\overset{k}{=}}} \hat{\boldsymbol{A}}^{(w_{l})T} \boldsymbol{\Delta}_{U}^{(w_{l})} \end{bmatrix} \bar{\boldsymbol{D}}^{-1}. \\ \text{Since } \boldsymbol{D} \boldsymbol{S} \bar{\boldsymbol{D}}^{-1} &= (\bar{\boldsymbol{D}} - \bar{\boldsymbol{P}} - \bar{\boldsymbol{Q}}) \bar{\boldsymbol{D}}^{-1} = \boldsymbol{I} - (\bar{\boldsymbol{P}} + \bar{\boldsymbol{Q}}) \bar{\boldsymbol{D}}^{-1}, \\ \boldsymbol{A}^{(N)} \boldsymbol{D} \boldsymbol{S} \bar{\boldsymbol{D}}^{-1} &= \boldsymbol{A}^{(N)} - \boldsymbol{A}^{(N)} (\bar{\boldsymbol{P}} + \bar{\boldsymbol{Q}}) \bar{\boldsymbol{D}}^{-1}. \end{split}$$

Therefore,

$$\begin{split} \|\tilde{\boldsymbol{A}}^{(N)}\bar{\boldsymbol{D}}^{-1} - \boldsymbol{A}^{(N)}\|_{F} \\ &= \left\| \boldsymbol{A}^{(N)}\boldsymbol{D} \begin{bmatrix} \sum\limits_{k=1}^{N} \boldsymbol{A}^{(n)T}\boldsymbol{\Delta}_{U}^{(n)} \end{bmatrix} \bar{\boldsymbol{D}}^{-1} \right\|_{F} \\ &+ \left\| \hat{\boldsymbol{A}}^{(N)}\hat{\boldsymbol{D}} \begin{bmatrix} \sum\limits_{k=1}^{N-1} \left(\sum\limits_{\substack{\{m_{1}, \dots, m_{k}\} \subset \{1, \dots, N\}, \\ \{w_{1}, \dots, w_{N-k}\} = \{1, \dots, N\}\} \\ \{m_{1}, \dots, m_{k}\}} & \stackrel{k}{\underset{l=1}{\overset{k}{\rightleftharpoons}}} \hat{\boldsymbol{A}}^{(m_{l})T}\boldsymbol{U}^{(m_{l})} & \stackrel{N-k}{\underset{l=1}{\overset{k}{\rightleftharpoons}}} \hat{\boldsymbol{A}}^{(w_{l})T}\boldsymbol{\Delta}_{U}^{(w_{l})} \\ \end{pmatrix} \right] \bar{\boldsymbol{D}}^{-1} \right\|_{F} \\ &+ \| \boldsymbol{A}^{(N)} \left(\bar{\boldsymbol{P}} + \bar{\boldsymbol{Q}} \right) \bar{\boldsymbol{D}}^{-1} \|_{F} \\ &= P + Q, \end{split}$$

where
$$P = O(\prod_{i=1}^{N-1} \epsilon_i)$$
 and $Q = O(k\epsilon \epsilon_{\perp})$.

This bound shows that Algorithm 3.1 achieves local convergence to fixed points for which $\epsilon_{\perp} = O(1/k)$. In contrast to the exact case, our analysis suggests only a linear convergence rate for approximation with the algorithm.

- 5. Approximate CP decomposition. In the above sections, we have provided a motivation for Algorithm 3.1 to compute a CP decomposition of rank R with R being less than or equal to the smallest mode length of the input tensor. We have shown in Theorem 4.3 that this algorithm exhibits a superlinear local convergence rate for exact CP decomposition problems and achieves a desirable approximation for special input tensors as described in Lemma 4.4. We now focus on the case of finding a good CP approximation for an arbitrary input tensor. We show that Algorithm 3.1 can be viewed as performing coupled minimization of the residual error of the decomposition in terms of a Mahalanobis distance metric [19]. Note that this perspective on the algorithm is different from the one introduced in section 3.1; however, it allows us to formulate an alternating minimization algorithm which generalizes Algorithm 3.1 to any CP rank and to interpolate between the updates of ALS and Algorithm 3.1.
- **5.1.** Mahalanobis distance minimization. Each update of Algorithm 3.1 may be viewed as minimizing a residual error with rescaled components. For an order 3 tensor \mathcal{X} in updating the first factor, it minimizes

(5.1)
$$\|(\boldsymbol{\mathcal{X}} - [\![\boldsymbol{A}, \boldsymbol{B}, \boldsymbol{C}]\!])_{(1)} (\boldsymbol{C}^{\dagger T} \otimes \boldsymbol{B}^{\dagger T})\|_F^2.$$

We can recast the above expression in a way such that the objective function is transformed into a non-Euclidean distance metric with respect to the input tensor. We can rewrite the above expression as

$$\operatorname{vec} \big(\boldsymbol{\mathcal{X}} - [\![\boldsymbol{A}, \boldsymbol{B}, \boldsymbol{C}]\!] \big)^T \boldsymbol{M} \operatorname{vec} \big(\boldsymbol{\mathcal{X}} - [\![\boldsymbol{A}, \boldsymbol{B}, \boldsymbol{C}]\!] \big),$$

where $M = (C^{\dagger T} C^{\dagger} \otimes B^{\dagger T} B^{\dagger} \otimes I)$. This may be viewed as minimizing the Mahalanobis distance metric relative to each factor while keeping the distance metric associated with that factor independent (and then updating it thereafter). This interpretation then enables us to extend Algorithm 3.1 for any CP rank and introduce methods that are a hybrid of Algorithm 3.1 and standard ALS.

5.1.1. Alternating Mahalanobis distance minimization. We consider a variant of Mahalanobis distance [18], which computes the distance between vectors

 $x, y \in \mathbb{R}^n$ as $(d(x, y))^2 = (x - y)^T M (x - y)$ for a given symmetric positive definite matrix $M \in \mathbb{R}^{n \times n}$. The matrix M is called the ground metric matrix. The ground metric generalizes the Euclidean distance to Mahalanobis distance by rotation and scaling of the axes along which the distance is computed. While the underlying ground metric may already be known, it may also be learned via various metric learning techniques [7, 33]. In optimal transport applications and other applications which require computation of distance between probability distributions, a Wasserstein distance is considered instead. Wasserstein distance between tensors with a given ground metric has been considered for nonnegative CP decomposition [2]. The ground metric in Wasserstein distance may be learned via similar metric learning techniques as in Mahalanobis distance [17]. Simultaneous optimization for a ground metric and Wasserstein distance between matrices has been used for nonnegative matrix factorization [68].

We consider minimization of the Mahalanobis distance between tensors with a fixed ground metric which may be updated later. In particular, the objective function minimized for an input tensor $\mathcal{X} \in \mathbb{R}^{I_1 \times \cdots \times I_N}$, factors $\mathbf{A}^{(i)} \in \mathbb{R}^{I_i \times R}$, and ground metric $\mathbf{M} \in \mathbb{R}^{I_1 \dots I_N \times I_1 \dots I_N}$ is

(5.2)
$$f(\mathbf{A}^{(1)}, \dots, \mathbf{A}^{(N)}) = \frac{1}{2} \operatorname{vec}(\mathbf{X} - \mathbf{Y})^{T} \mathbf{M} \operatorname{vec}(\mathbf{X} - \mathbf{Y}),$$
where $\mathbf{Y} = [\mathbf{A}^{(1)}, \dots, \mathbf{A}^{(N)}].$

We restrict the ground metric matrix M to be Kronecker structured defined as

$$\boldsymbol{M} = \bigotimes_{k=1}^{N} \left(\boldsymbol{M}^{(k)} \right)^{-1},$$

where each $(M^{(k)})^{-1} \in \mathbb{R}^{I_k \times I_k}$ may be viewed as a ground metric for each mode of the tensor. This restriction allows us to exploit the computational benefits of the structure and enables us to formulate an efficient alternating minimization algorithm. We consider the objective in (5.2) for a general (fixed) ground metric for alternating optimization, which also allows us to formulate different algorithms for CP decomposition by changing the ground metric. We derive an update for the alternating minimization with respect to the nth factor matrix, given by

(5.3)
$$\mathbf{A}^{(n)} = \min_{\mathbf{A}^{(n)}} \frac{1}{2} \| \mathbf{X}_{(n)} - \mathbf{A}^{(n)} \mathbf{P}^{(n)T} \|_{\mathbf{M}}^{2},$$
where $\mathbf{P}^{(n)} = \bigodot_{m=1, m \neq n}^{N} \mathbf{A}^{(m)}.$

For succinct writing, let $M_{(n)} = \bigotimes_{k=1, k\neq n}^{N} (M^{(k)})^{-1}$. Since the objective function is quadratic in $A^{(n)}$, a minimizer of (5.3) can be found by obtaining obtaining a gradient $G^{(n)}$ with respect to the *n*th factor matrix and setting it to **0**. The gradient is

$$\boldsymbol{G}^{(n)} = \left(\boldsymbol{M}^{(n)}\right)^{-1} \boldsymbol{A}^{(n)} \boldsymbol{P}^{(n)T} \boldsymbol{M}_{(n)} \boldsymbol{P}^{(n)} - \left(\boldsymbol{M}^{(n)}\right)^{-1} \boldsymbol{X}_{(n)} \boldsymbol{M}_{(n)} \boldsymbol{P}^{(n)}.$$

Setting the gradient above to be **0** and equating $M^{(n)}(M^{(n)})^{-1} = I$, we get an update for the *n*th factor given as the solution of the following system:

(5.4)
$$\boldsymbol{A}^{(n)} \left(\boldsymbol{P}^{(n)T} \boldsymbol{M}_{(n)} \boldsymbol{P}^{(n)} \right) = \boldsymbol{X}_{(n)} \boldsymbol{M}_{(n)} \boldsymbol{P}^{(n)}.$$

Using the properties of Khatri–Rao products and Kronecker products, the update for the nth factor matrix reduces to the system of equations

(5.5)
$$\boldsymbol{A}^{(n)}\boldsymbol{Z}^{(n)} = \boldsymbol{X}_{(n)}\boldsymbol{L}^{(n)},$$
 where $\boldsymbol{L}^{(n)} = \bigodot_{k=1,k\neq n}^{N} \left(\boldsymbol{M}^{(k)}\right)^{-1}\boldsymbol{A}^{(k)}$ and $\boldsymbol{Z}^{(n)} = \operatornamewithlimits{*}_{k=1,k\neq n}^{N} \boldsymbol{A}^{(k)T} \left(\boldsymbol{M}^{(k)}\right)^{-1}\boldsymbol{A}^{(k)}.$

The above update leads to the ALS algorithm if $M^{(k)} = I$ for all k. We can retrieve Algorithm 3.1 by defining

(5.6)
$$M^{(k)} = A^{(k)}A^{(k)T} + (I - A^{(k)}A^{(k)\dagger}) \quad \forall k \in \{1, \dots, N\}.$$

The matrix $I - A^{(k)}A^{(k)\dagger}$ is inconsequential when applied to the factor matrices as it is an orthogonal projector that projects onto the null space of the kth factor matrix. It is included to ensure that each $M^{(k)}$ is symmetric positive definite. Since Algorithm 3.1 can be retrieved from the above update, we refer to Algorithm 3.1 as alternating Mahalanobis distance minimization (AMDM).

Let us assume that the iteration involving the above-derived alternating updates to each factor as in (5.4) converges to a critical point. We can then bound the backward error in application of each matricization of the reconstructed tensor $\mathcal{Y} = [A^{(1)}, \cdots, A^{(N)}]$, since from (5.4), for each n, we have

$$egin{aligned} oldsymbol{A}^{(n)} \Big(oldsymbol{P}^{(n)T}oldsymbol{M}_{(n)}oldsymbol{P}^{(n)}\Big) - oldsymbol{X}_{(n)}oldsymbol{M}_{(n)}oldsymbol{P}^{(n)} = oldsymbol{0}, \ \Big(oldsymbol{Y}_{(n)} - oldsymbol{X}_{(n)}\Big)oldsymbol{M}_{(n)}oldsymbol{P}^{(n)} = oldsymbol{0}. \end{aligned}$$

Therefore, we have that

$$\|m{Y}_{(n)}m{z} - m{X}_{(n)}m{z}\| = \|m{X}_{(n)}m{z}^{\perp}\|,$$

where \mathbf{z}^{\perp} is the projection of $\mathbf{z} \in \mathbb{R}^{\prod_{j=1 j \neq n} I_j}$ onto the orthogonal complement of the column span of $\mathbf{M}_{(n)}\mathbf{P}^{(n)}$ or $\bigodot_{j=1,j \neq n}^{N}\mathbf{M}^{(j)}\mathbf{A}^{(j)}$. For ALS, AMDM, and the hybrid methods (introduced in section 5.3) that interpolate between them, it is sufficent to consider the projection onto the orthogonal complement of the column span of $\bigodot_{j=1,j \neq n}^{N}\mathbf{A}^{(j)}$. This is because for each j, the ground metric matrices are chosen such that the column span of $\mathbf{A}^{(j)}$ is an invariant subspace of $\mathbf{M}^{(j)}$.

5.1.2. Comparison of AMDM and ALS for approximate rank-2 CP decomposition. We use the formulation introduced in the previous subsection to generalize AMDM to the case when the CP rank R is greater than the mode lengths and to derive hybrid methods that interpolate between AMDM and ALS. The residual transformation tends to equalize the weight of contribution to the objective function attributed to components of the error associated with different rank-1 parts of the CP decomposition, $[a_i, b_i, c_i]$, without increasing the collinearity of the columns of the factors. We provide an example as an intuition for this assertion.

Consider a tensor $\mathcal{X} = \lambda_1 \mathcal{X}_1 + \lambda_2 \mathcal{X}_2 + \mathcal{N}$, where \mathcal{X}_1 and \mathcal{X}_2 are normalized rank-1 tensors and \mathcal{N} is noise of small magnitude. Assume that $\lambda_1 \gg \lambda_2$ and the rank-1 tensors are highly correlated, i.e., the factors have collinear columns. Let the

current CP decomposition approximation be $\mathbf{\mathcal{Y}} = [\![\bar{\lambda}; \mathbf{A}, \mathbf{B}, \mathbf{C}]\!] = \bar{\lambda}_1 \mathbf{\mathcal{Y}}_1 + \bar{\lambda}_2 \mathbf{\mathcal{Y}}_2$. The least squares objective minimizes

$$\begin{aligned} \|\boldsymbol{\mathcal{X}} - \boldsymbol{\mathcal{Y}}\|_F^2 &= \operatorname{vec}(\boldsymbol{\mathcal{X}} - \boldsymbol{\mathcal{Y}})^T \operatorname{vec}(\boldsymbol{\mathcal{X}} - \boldsymbol{\mathcal{Y}}) \\ &= \operatorname{vec}(\boldsymbol{\mathcal{E}}_1 + \boldsymbol{\mathcal{E}}_2 + \boldsymbol{\mathcal{N}})^T \operatorname{vec}(\boldsymbol{\mathcal{E}}_1 + \boldsymbol{\mathcal{E}}_2 + \boldsymbol{\mathcal{N}}) \\ &= \|\boldsymbol{\mathcal{E}}_1\|_F^2 + \|\boldsymbol{\mathcal{E}}_2\|_F^2 + 2\operatorname{vec}(\boldsymbol{\mathcal{E}}_1)^T \operatorname{vec}(\boldsymbol{\mathcal{E}}_2) + 2\operatorname{vec}(\boldsymbol{\mathcal{N}})^T (\boldsymbol{\mathcal{E}}_1 + \boldsymbol{\mathcal{E}}_2 + \boldsymbol{\mathcal{N}}), \end{aligned}$$

where $\mathcal{E}_1 = \lambda_1 \bar{\mathcal{X}}_1 - \bar{\lambda}_1 \mathcal{Y}_1$ and $\mathcal{E}_2 = \lambda_2 \mathcal{X}_2 - \bar{\lambda}_2 \mathcal{Y}_2$. The ALS algorithm may reduce $\|\mathcal{E}_1\|_F$ and the component of \mathcal{E}_2 in the direction of \mathcal{E}_1 , since it leads to reduction of the terms with larger contribution in the error. This causes an increase in collinearity of the approximated factors and a more ill-conditioned decomposition. On the other hand, the objective in (5.1) can be expressed as

$$\begin{aligned} \|(\boldsymbol{\mathcal{X}} - \boldsymbol{\mathcal{Y}}) \times_1 \boldsymbol{I} \times_2 \boldsymbol{B}^{\dagger T} \times_3 \boldsymbol{C}^{\dagger T}\|_F^2 &= \operatorname{vec}(\boldsymbol{\mathcal{X}} - \boldsymbol{\mathcal{Y}})^T \boldsymbol{M} \operatorname{vec}(\boldsymbol{\mathcal{X}} - \boldsymbol{\mathcal{Y}}) \\ &= \|\boldsymbol{\mathcal{E}}_1\|_{\boldsymbol{M}}^2 + \|\boldsymbol{\mathcal{E}}_2\|_{\boldsymbol{M}}^2 + 2\operatorname{vec}(\boldsymbol{\mathcal{E}}_1)^T \boldsymbol{M} \operatorname{vec}(\boldsymbol{\mathcal{E}}_2) \\ &+ 2\operatorname{vec}(\boldsymbol{\mathcal{N}})^T \boldsymbol{M} (\boldsymbol{\mathcal{E}}_1 + \boldsymbol{\mathcal{E}}_2 + \boldsymbol{\mathcal{N}}), \end{aligned}$$

where $M = C^{\dagger T}C^{\dagger} \otimes B^{\dagger T}B^{\dagger} \otimes I$. The matrix M rescales the components of the error according to the inverse of square of singular values of the factors, since

$$\begin{aligned} \|\boldsymbol{\mathcal{E}}_1\|_{\boldsymbol{M}}^2 &= \operatorname{vec}(\boldsymbol{\mathcal{E}}_1)^T \boldsymbol{M} \operatorname{vec}(\boldsymbol{\mathcal{E}}_1) \\ &= \operatorname{vec}(\bar{\boldsymbol{\mathcal{E}}}_1)^T (\boldsymbol{\Sigma}_{\boldsymbol{C}}^{-2} \otimes \boldsymbol{\Sigma}_{\boldsymbol{B}}^{-2} \otimes \boldsymbol{I}) \operatorname{vec}(\bar{\boldsymbol{\mathcal{E}}}_1), \end{aligned}$$

where $\mathcal{E}_1 = \mathcal{E}_1 \times_1 I \times_2 U_B \times_3 U_C$ with U_B and U_C being the left singular vectors of B and C, respectively. Thus, the error is rotated by the left singular vectors of B and C and then rescaled by the square of the inverse of the singular values of B and C, i.e., Σ_B^{-2} , Σ_C^{-2} . Therefore, if the approximated factors are collinear, the contribution in the direction of the singular vectors of B and C with largest singular value is weighed proportionally less, and similarly the contribution of the error in the direction of those with smaller singular value is weighed more. This reduces the imbalance in error and leads to a better conditioned decomposition.

5.2. Generalizing AMDM to any CP rank. The AMDM algorithm as described in Algorithm 3.1 imposes a constraint that CP rank should be less than or equal to the smallest mode length of the tensor. As mentioned above, the update in (5.5) with the ground metric as defined in (5.6) leads to the AMDM algorithm with the rank constraint. We now describe how this definition of ground metric leads to an AMDM update without any restrictions imposed on the rank. Note that for each $M^{(k)}$ as defined in (5.6),

$$(\boldsymbol{M}^{(k)})^{-1} = \boldsymbol{A}^{(k)\dagger T} \boldsymbol{A}^{(k)\dagger} + (\boldsymbol{I} - \boldsymbol{A}^{(k)} \boldsymbol{A}^{(k)\dagger}).$$

The above equivalence can be verified simply by looking at the singular value decomposition of $M^{(k)}$ and $A^{(k)}$. The first part simply inverts the nonunit singular values of $M^{(k)}$, and the second part is a projector and orthogonal to the first part and is therefore kept as is. The linear system for updating the nth factor matrix as in (5.5) can then be simplified to get

$$\begin{split} \boldsymbol{A}^{(n)}\boldsymbol{Z}^{(n)} &= \boldsymbol{X}_{(n)}\boldsymbol{L}^{(n)},\\ \text{where } \boldsymbol{L}^{(n)} &= \bigodot_{k=1,k\neq n}^{N} \left(\boldsymbol{M}^{(k)}\right)^{-1}\boldsymbol{A}^{(k)} = \bigodot_{k=1,k\neq n}^{N} \boldsymbol{A}^{(k)\dagger T}\\ \text{and } \boldsymbol{Z}^{(n)} &= \operatornamewithlimits{*}_{k=1,k\neq n}^{N} \boldsymbol{A}^{(k)T} \big(\boldsymbol{M}^{(k)}\big)^{-1}\boldsymbol{A}^{(k)} = \operatornamewithlimits{*}_{k=1,k\neq n}^{N} \boldsymbol{A}^{(k)\dagger}\boldsymbol{A}^{(k)}. \end{split}$$

Note that $(\boldsymbol{M}^{(k)})^{-1}$ is not computed explicitly. The above update is equivalent to Algorithm 3.1 when the CP rank $R \leq I_m$ for all $m \in \{1, \dots, N\}$, since in that case $\boldsymbol{A}^{(k)\dagger}\boldsymbol{A}^{(k)}=\boldsymbol{I}$ for all k, and, consequently, $\boldsymbol{Z}^{(n)}=\boldsymbol{I}$ for all n. For the case when the CP rank R is larger than the mode lengths, $\boldsymbol{A}^{(k)\dagger}\boldsymbol{A}^{(k)}$ is a projector onto the row space of $\boldsymbol{A}^{(k)}$, and therefore $\boldsymbol{Z}^{(n)}$ is symmetric positive semidefinite for each n. Since the system is semidefinite, a pivoted Cholesky decomposition followed by a triangular solve must be used for the solution to exist. Alternatively, as in ALS, a regularization term may be introduced to make the system positive definite. The cost of forming this system, $\boldsymbol{Z}^{(n)}$, is of the same leading order as ALS, i.e., $O(IR^2)$ per subsweep. It requires obtaining a pseudoinverse of the factor in the (n-1)th iteration, matrix multiplication, and Hadamard products of the factors and their previously obtained pseudoinverses. The system solve amounts to a computational cost of $O(R^3)$.

5.3. Interpolating between AMDM and ALS. Performing AMDM with an identity ground metric is equivalent to performing ALS. To explore methods that interpolate between AMDM and ALS, we can interpolate the ground metric between the identity matrix and the one associated with AMDM given in (5.6). We can find such ground metrics by decomposing each factor matrix into two low rank matrices such that $\mathbf{A}^{(k)} = \mathbf{A}_1^{(k)} + \mathbf{A}_2^{(k)}$, where $\mathbf{A}_1^{(k)}$ is the best rank-t approximation of $\mathbf{A}^{(k)}$. In other words, $\mathbf{A}_1^{(k)}$ contains the largest t singular values and the corresponding singular vectors of $\mathbf{A}^{(k)}$, and $\mathbf{A}^{(2)}$ contains the rest. Consequently, $\mathbf{A}_1^{(k)}$ and $\mathbf{A}_2^{(k)}$ are orthogonal to each other. The ground metric for each mode can then be defined using only the first part as

(5.7)
$$\mathbf{M}^{(k)} = \mathbf{A}_1^{(k)} \mathbf{A}_1^{(k)T} + \left(\mathbf{I} - \mathbf{A}_1^{(k)} \mathbf{A}_1^{(k)\dagger} \right) \quad \forall k \in \{1, \dots, N\}.$$

Defining a ground metric based on only the first part of the singular value decomposition of the factors leads to hybrid methods, since if the first part is all of the singular value decomposition, then we get back the ground metric in AMDM, and if it is none of the same, then we get back the identity matrix by convention as the orthogonal complement of none is everything.

Note that for each k,

$$\left(oldsymbol{M}^{(k)}
ight)^{-1} = oldsymbol{A}_1^{(k)\dagger T} oldsymbol{A}_1^{(k)\dagger} + \left(oldsymbol{I} - oldsymbol{A}_1^{(k)} oldsymbol{A}_1^{(k)\dagger}
ight).$$

The update for the *n*th factor matrix also becomes a combination of AMDM and ALS where the first part of the singular value decomposition of factors is treated as in AMDM and the second one as in ALS. More precisely, the system of equations solved for updating the *n*th factor matrix is

$$\mathbf{A}^{(n)} \mathbf{Z}^{(n)} = \mathbf{X}_{(n)} \mathbf{L}^{(n)},$$
where $\mathbf{L}^{(n)} = \bigodot_{k=1, k \neq n}^{N} \left(\mathbf{M}^{(k)} \right)^{-1} \mathbf{A}^{(k)} = \bigodot_{k=1, k \neq n}^{N} \left(\mathbf{A}_{1}^{(k) \dagger T} + \mathbf{A}_{2}^{(k)} \right)$
(5.8) and $\mathbf{Z}^{(n)} = \underset{k=1, k \neq n}{\overset{N}{*}} \mathbf{A}^{(k)T} \left(\mathbf{M}^{(k)} \right)^{-1} \mathbf{A}^{(k)} = \underset{k=1, k \neq n}{\overset{N}{*}} \left(\mathbf{A}_{1}^{(k) \dagger} \mathbf{A}_{1}^{(k)} + \mathbf{A}_{2}^{(k)T} \mathbf{A}_{2}^{(k)} \right).$

As mentioned above, in the above update, $\boldsymbol{L}^{(n)}$ and $\boldsymbol{Z}^{(n)}$ can be split into two orthogonal parts as shown. The first part which involves $\boldsymbol{A}_1^{(k)}$ can be characterized as the AMDM part, and the one which involves $\boldsymbol{A}_2^{(k)}$ can be characterized as the ALS part. Note that computing these interpolating updates does not involve extra computational cost asymptotically when compared to ALS. Using these interpolating updates, we describe a hybrid algorithm in the subsection below. This algorithm is shown to perform the best in terms of residual and conditioning in section 6.

5.4. Hybrid algorithm. In this section, we leverage the interpolating updates to describe the hybrid algorithm which interpolates between AMDM and ALS. This algorithm starts by computing AMDM updates and transitions into ALS by implicitly reducing components used in the ground metric or reducing the size of $A_1^{(k)}$ in (5.7) until the ground metric is I. The transition in the hybrid algorithm is chosen to be from AMDM to ALS and not the other way around as ALS updates perturb the condition number less and make more progress when the updates are well-conditioned, and AMDM keeps the condition number bounded but may compromise on the Frobenius norm of the residual. The number of components of $A^{(k)}$ in $A_1^{(k)}$ is controlled by a hyperparameter, t. Note that t controls if the updates are more likely to behave as ALS or AMDM; for example, if $A_1^{(k)}$ has a larger column space, then the updates are more likely to behave like AMDM. This hyperparameter may depend on the input tensor and conditioning of the decomposition and therefore needs to be tuned according to the problem. The details of the algorithm are described below and summarized in Algorithm 5.1.

The hybrid algorithm starts by normalizing the columns of all the factors and absorbing norms in the first factor as described in section 3 and then computing a reduced singular value decomposition of all the factors which costs $O(\sum_{n=1}^{N} I_n R \min(I_n, R))$. At the *n*th subiteration of the algorithm, the right- and left-hand sides of the system, $\boldsymbol{X}_{(n)} \boldsymbol{L}^{(n)}$ and $\boldsymbol{Z}^{(n)}$, are computed using (5.8). Let $\boldsymbol{U}_1^{(m)}, \boldsymbol{V}_1^{(m)}$ be the left and right singular vectors of $\boldsymbol{A}_1^{(m)}$, respectively, $\boldsymbol{U}_2^{(m)}, \boldsymbol{V}_2^{(m)}$ be the left and right singular vectors of $\boldsymbol{A}_2^{(m)}$, respectively, and $\boldsymbol{S}_1^{(m)}, \boldsymbol{S}_2^{(m)}$ be the singular values of $A_1^{(m)}$ and $A_2^{(m)}$, respectively. The computation of $Z^{(n)}$ requires $O(R^2 \min(I_n, R))$ operations. This computation can be made efficient by first storing the singular value decomposition of each of the factors and then computing $\mathbf{Z}_m = \mathbf{V}_1^{(m)} \mathbf{V}_1^{(m)T} + \mathbf{V}_2^{(m)} (\mathbf{S}_2^{(m)})^2 \mathbf{V}_2^{(m)T}$ for each m. Then, $\mathbf{Z}^{(n)}$ can be computed by taking Hadamard products of Z_m as shown in Algorithm 5.1. A further reduction in cost can be obtained by updating \boldsymbol{Z}_m at each subiteration by only computing the singular value decomposition of the previous factor. This amortization is similar to the one used in ALS while computing the left-hand sides, leading to a reduction in cost by a factor of N. The computation of $X_{(n)}L^{(n)}$ can be similarly performed by using $L_m = U_1^{(m)}(S_1^{(m)})^{-1}V_1^{(m)T} + U_2^{(m)}S_2^{(m)}V_2^{(m)T}$. Then $X_{(n)}L^{(n)} = X_{(n)} \bigcirc_{m=1, m\neq n}^{N} L_m$ which is an MTTKRP computation and can be performed efficiently and amortized as in ALC. The performed efficiently and amortized as in ALS. The symmetric semidefinite system solve requires $O(R^3)$. Computing the right-hand side, i.e., performing MTTKRP, is the most computationally expensive operation with a cost of $O(\prod_{n=1}^{N} I_n R)$ for each subsweep. In addition to this, the factor matrix obtained after the solve is normalized, and a reduced singular value decomposition is obtained to update the singular value decomposition which costs $O(I_nR\min(I_n,R))$. Therefore, the asymptotic computational cost of the algorithm is the same as ALS, being $O(\prod_{n=1}^{N} I_n R)$.

We summarize the above-described hybrid algorithm in Algorithm 5.1.

Algorithm 5.1 General AMDM: AMDM with singular value thresholding.

```
1: Input: Tensor \mathcal{X} \in \mathbb{R}^{I_1 \times \cdots \times I_N}, threshold t, rank R
 2: Initialize \{\boldsymbol{A}^{(1)}, \dots, \boldsymbol{A}^{(N)}\} so each \boldsymbol{A}^{(n)} \in \mathbb{R}^{I_n \times R} is random
 3: for n \in \{2, ..., N\} do
          \boldsymbol{A}^{(n)} = \text{normalize}(\boldsymbol{A}^{(n)})
  4:
          U^{(n)} = \min(I_n, R) left singular vectors of A^{(n)}
  5:
          V^{(n)} = \min(I_n, R) right singular vectors of A^{(n)}
  6:
          s^{(n)} = \min(I_n, R) singular values of A^{(n)}
 8: end for
 9: while Convergence do
10:
          for n \in \{1, ..., N\} do
11:
               for m \in \{1, ..., N\}, m \neq n do
                  \mathbf{s}_{\mathrm{ps}}^{(m)} = \text{first } t \text{ values of } \mathbf{s}^{(m)} \text{ inverted and others as it is}
12:
                  egin{aligned} oldsymbol{L}_m &= oldsymbol{U}^{(m)} \mathbf{diag}(oldsymbol{s}_{	ext{ps}}^{(m)}) oldsymbol{V}^{(m)T} \ oldsymbol{Z}_m &= oldsymbol{V}^{(m)} \mathbf{diag}(oldsymbol{s}_{	ext{ps}}^{(m)} * oldsymbol{s}^{(m)}) oldsymbol{V}^{(m)T} \end{aligned}
13:
14:
15:
               end for
              Compute \boldsymbol{Z}^{(n)} = *_{m=1,m\neq n}^{N} \boldsymbol{Z}_{m}

Compute \boldsymbol{X}_{(n)} \boldsymbol{L}^{(n)} where \boldsymbol{L}^{(n)} = \bigodot_{m=1,m\neq n}^{N} \boldsymbol{L}_{m}
16:
17:
              Solve for A^{(n)} in A^{(n)}Z^{(n)} = X_{(n)}L^{(n)} as in (5.8)
18:
              if n is N: Check Convergence, if converged:
19:
               \mathbf{A}^{(n)} = \text{normalize}(\mathbf{A}^{(n)})
20:
               Update U^{(n)} = \min(I_n, R) left singular vectors of A^{(n)}
21:
               Update V^{(n)} = \min(I_n, R) right singular vectors of A^{(n)}
22:
               Update s^{(n)} = \min(I_n, R) singular values of A^{(n)}
23:
24:
          end for
25: end while
26: return factor matrices \{A^{(1)}, \dots, A^{(N)}\}
```

6. Numerical experiments. We perform numerical experiments to demonstrate the convergence behavior of the AMDM algorithm compared to ALS for various tensors which include synthetic examples and tensors arising in different applications. Both algorithms are implemented in Python and are publicly available.² The ALS implementation used in this work has been used in several previous works [38, 39, 55]. We use absolute residual and fitness of the decomposition in the Frobenius norm as metrics to measure the closeness of the decomposition to the input tensor. For an input tensor \mathcal{X} , these are given as

$$r = \|\mathcal{X} - \mathcal{Y}\|_F$$
 and $f = 1 - \frac{\|\mathcal{X} - \mathcal{Y}\|_F}{\|\mathcal{X}\|_F}$,

respectively, where $\mathbf{\mathcal{Y}} = [\![\mathbf{A}^{(1)}, \dots, \mathbf{A}^{(N)}]\!]$ is the approximated tensor. For measuring the stability of the decomposition or the degree of overlap of CP components, we use the normalized CP decomposition condition number [10] to measure the sensitivity or degree of overlap of rank-1 components of the decomposition.

²https://github.com/cyclops-community/tensor_decomposition.

The normalized CP condition number is given by the reciprocal of the smallest singular value of Terracini's matrix associated with the CP decomposition. For an equidimensional tensor of order N with mode length s and CP rank R, the size of Terracini's matrix is $s^N \times (N(s-1)+1)R$. For CP rank lower than mode lengths of the tensor, this matrix can be compressed to $R^N \times (N(R-1)+1)R$, and the CP decomposition condition number can be efficiently computed with a cost of $O(R^{N+4})$. The details of computation of the condition number are in Appendix A. A more general result regarding computing the condition number of joint decompositions efficiently is proved in [20]. Our experiments consider two types of synthetic tensors.

Tensor made by random matrices (random tensor). We create these tensors based on known uniformly distributed randomly generated factor matrices $\mathbf{A}^{(n)} \in (0,1)^{s \times R}$, $\mathbf{\mathcal{X}} = [\![\mathbf{A}^{(1)}, \dots, \mathbf{A}^{(N)}]\!]$.

Tensor made by collinear random matrices (collinearity tensor). We use a similar approach as used in [1] to generate factor matrices with a fixed value of collinearity, say C. That means that these tensors are created with randomly generated factors $\mathbf{A}^{(n)} \in \mathbb{R}^{s \times R}$ with the following property:

$$\mathbf{a}_r^{(n)T}\mathbf{a}_z^{(n)} = C$$

s.t. $\|\mathbf{a}_r^{(n)}\| = 1 \ \forall r \neq z \in \{1, \dots, R\}.$

We then set $d_{ii} = i$ for all $i \in \{1, ..., R\}$ to create a tensor, $\mathcal{X} = [D; \mathbf{A}^{(1)}, ..., \mathbf{A}^{(N)}]$. We consider four tensors from various real-world applications.

Sleep-EDF tensor: This dataset has been used to identify sleeping patterns. It comprises electroencephalogram (EEG) and electromyography (EMG) data in the nonrapid eye movement stage of sleep [30].

MGH tensor: This dataset consists of data from Massachusetts General Hospital. It includes combinations of EEG data, respiratory signals, and EMG signals. This dataset was used to analyze sleep using deep neural networks [9].

SCF tensor. We consider the density fitting tensor (Cholesky factor of the twoelectron integral tensor) arising in quantum chemistry. This tensor has been used previously in [55] to compare the efficacy of the Gauss–Newton and ALS algorithms. We leverage the PySCF library [59] to generate the three-dimensional compressed density fitting tensor, representing the compressed restricted Hartree–Fock wave function of water molecule chain systems with a STO-3G basis set. The number of molecules in the system is set to three for this experiment.

Amino acid tensor. This dataset consists of five simple laboratory-made samples. Each sample contains different amounts of tyrosine, tryptophan, and phenylalanine dissolved in phosphate-buffered water. The samples were measured by fluorescence [11].

The experiments are divided broadly into two categories:

Exact decomposition. We create synthetic tensors with known CP rank R and compare the convergence behavior of the AMDM algorithm with the ALS algorithm for exact CP decomposition.

Approximate decomposition. We create synthetic tensors with known CP rank and special structure such as with added noise or as described in Lemma 4.4. We then approximate these tensors with CP rank R which is lower than the underlying decomposition rank. We also consider real-world tensors from different applications with unknown CP rank. For synthetic tensors, we create 10 different tensors and run each with a random initialization for ALS and different variants of AMDM and plot the mean of fitness and condition number at each iteration with confidence intervals

as the minimum and maximum values. The random initialization are factor matrices with entries sampled from a uniform distribution of real numbers between 0 and 1. For real-world tensors, we plot the mean of 5 initializations for larger tensors, Sleep-EDF, MGH, and SCF tensors and the mean of 10 iterations for amino acid tensors with confidence intervals as minimum and maximum values.

6.1. Exact CP decomposition. We compare ALS and Algorithm 3.1 for computing exact CP decomposition of synthetic tensors in Figures 1 and 2 and verify our theoretical results. We create *collinearity* and *random* tensors of specified CP rank to analyze the convergence of these algorithms. We claim that the algorithm has converged for a CP decomposition with exact rank if the absolute residual is below 10^{-9} .

In Figure 1, we create equidimensional synthetic tensors of order N with each mode length s to follow s=100 for order 3, s=53 for order 4, and s=24 for order 5 with fixed CP rank R=20. We initialize the factors with uniformly distributed random matrices for both algorithms and plot the absolute residual for each iteration. For the *collinearity* tensors, collinearity value, i.e., C, is set to be 0.9. We can observe superlinear convergence of Algorithm 3.1 for both cases, while ALS appears to be slowed down by the "swamp" phenomenon for the *collinearity* tensors.

In Figure 2, we create a random equidimensional tensor with mode length s=100 and CP rank R=200. We can observe that Algorithm 5.1 converges to an approximate tensor which is $\epsilon=10^{-10}$ close to the exact tensor while taking only a few iterations to do so. Since R>s, there is a possibility that the synthetic tensor might have a border rank lower than 200, and we may be finding the corresponding factor matrices.

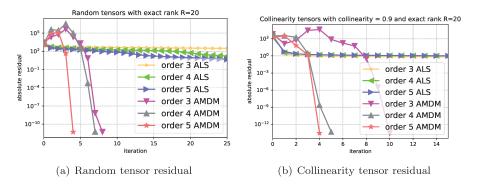


Fig. 1. Superlinear convergence of the AMDM algorithm for exact CP decomposition.

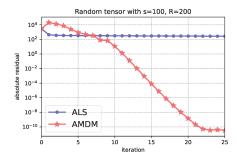


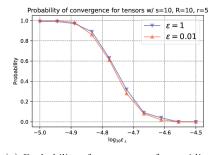
Fig. 2. Linear rate of convergence of AMDM for large rank for exact CP decomposition.

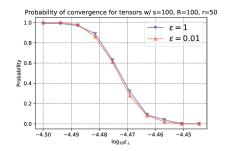
Therefore we cannot guarantee that we find the exact input tensor [43]. However, we do observe a linear convergence rate to obtain this approximation. ALS makes slow progress for this case. A reason for that might again be related to the collinearity of the factors, since when R > s, collinearity is high and ALS is more likely to experience the "swamp" phenomenon [41].

6.2. Approximate CP decomposition. We plot the probability of convergence to the desired decomposition for synthetic tensors as described in Lemma 4.4 with respect to the ϵ_{\perp} to verify our theoretical results in Figure 3. We construct these tensors by constructing the first R/2 columns of the factors with random matrices and then constructing the other half by projecting it onto the orthogonal complement of the column space of the first half and adding Gaussian noise of amplitude ϵ_{\perp} to the same. We construct 100 such tensors, and for each tensor we consider 5 initial guesses which are ϵ away from the desired decomposition as described in 4.4. We plot the probability of convergence over these 100 tensors by considering if at least 1 initial guess is within 10^{-9} of the desired factors. We observe that the probability of convergence to the desired decomposition is 1 when the ϵ_{\perp} is small, irrespective of the size and ϵ , thereby verifying that the first half of CP decomposition is a stationary point. The probability decreases as we increase ϵ_{\perp} ; i.e., the decomposition converges to different stationary points for these tensors.

We compare ALS and different variants of AMDM for computing approximate CP decomposition of synthetic tensors in Figure 4 and application tensors that admit approximations with a low CP rank in Figure 5, Figure 6, Figure 7, and Figure 8. We plot the fitness and the condition number of the CP decomposition to compare ALS and variants of AMDM. The integer associated with AMDM t = # corresponds to the number of singular values inverted for each factor or best rank-t approximation A_1 in (5.7). The hybrid algorithm starts by using threshold t = R in Algorithm 5.1; i.e., it starts with Algorithm 3.1 and gradually decreases the threshold to 0 to recover the ALS algorithm.

In Figure 4, we compute a rank-10 CP decomposition of the collinearity tensor with collinearity C = 0.9 and exact CP rank R = 10 with added Gaussian noise tensor. Each entry of the noise tensor is distributed normally with mean $\mu = 0$ and standard deviation $\sigma = 0.001$. We create 10 such tensors randomly and, for each tensor, initialize the algorithms with a random initial guess. We observe that

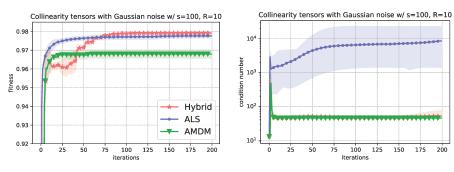




(a) Probability of convergence for equidimen-CP rank= 10 and approximate rank= 5

(b) Probability of convergence for equidimensional tensors with mode length= 10, exact sional tensors with mode length= 100, exact CP rank= 100 and approximate rank= 50

Fig. 3. Probability of convergence for 100 tensors with 5 initial guesses in the setting as described in Lemma 4.4



- (a) Collinearity tensor with Gaussian noise (b) Collinearity tensor with Gaussian noise fitness
 - condition number

Fig. 4. Collinearity tensor of size $100 \times 100 \times 100$ with exact CP rank R = 10 with added Gaussian noise with each entry distributed with mean $\mu = 0$ and standard deviation $\sigma = 0.001$.

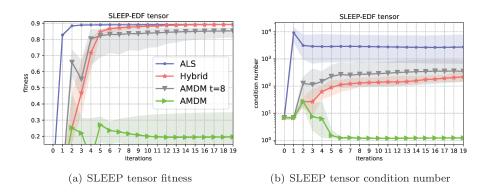


Fig. 5. Sleep-EDF tensor of dimensions $2048 \times 14 \times 129 \times 86$ approximated with CP rank R = 10, where t is the singular value threshold in Algorithm 5.1.

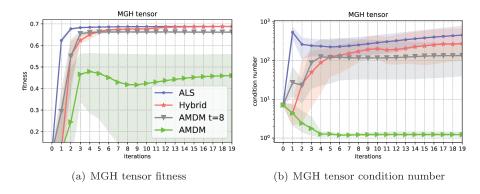


Fig. 6. MGH tensor of dimensions $2048 \times 12 \times 257 \times 43$ approximated with CP rank R = 10, where t is the singular value threshold in Algorithm 5.1.

the mean fitness of the AMDM algorithm is comparable to ALS with tight max-min confidence intervals and a low condition number for all the tensors. On the other hand, ALS maintains a higher fitness with $15-1500\times$ the mean condition number. Note

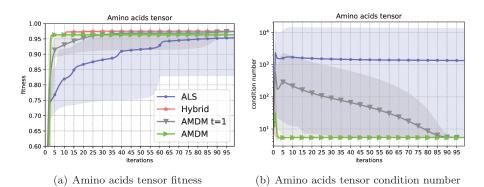


Fig. 7. Amino acid tensor of dimensions $5 \times 61 \times 201$ approximated with CP rank R=3, where t is the singular value threshold in Algorithm 5.1.

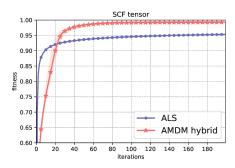


FIG. 8. SCF tensor of size $339 \times 21 \times 21$ approximated with CP rank R = 200, where t is the $\arg\min_i \sigma_{max}/\sigma_i < 100$ threshold in Algorithm 5.1.

that in this approximate CP decomposition setting, the convergence of the AMDM algorithm behaves similarly to ALS; i.e., the rate is linear, as shown in Lemma 4.4. We can clearly see the advantage of the hybrid algorithm where we start with AMDM and decrease the number of singular values inverted by 1 after every 10 iterations. The hybrid algorithm has a mean fitness that is higher than ALS while maintaining a condition number that is very close to AMDM. This suggests that for these types of tensors a more stable decomposition with desirable accuracy can be found via the hybrid algorithm.

In Figure 5 and Figure 6, we compute the CP decomposition of the Sleep-EDF tensor and MGH tensor with CP rank R=10. We consider several variants of the hybrid Algorithm 5.1. For the hybrid algorithm, one less singular value is inverted after every iteration. We clearly see a pattern in both tensors that if lesser singular values are inverted, then the fitness is higher and the CP decomposition condition number is larger. We also see that the hybrid algorithm is able to achieve a fitness as high as ALS while maintaining a lower condition number. For the Sleep-EDF tensor we see this effect being more pronounced as we get a mean condition number that is around an order of magnitude smaller with small confidence intervals while the fitness is the same. For the MGH tensor, we have that the mean condition number of the hybrid algorithm is not very different from the condition number in ALS, and after 50 iterations this difference is also lesser.

In Figure 7, we compute the CP decomposition of the amino acid tensor with rank R=3. We use 10 random initializations for this tensor because of its small size. This tensor is different from Sleep-EDF and MGH as we have the a higher mean fitness for AMDM variants as compared to ALS. The large confidence interval of ALS suggests that it is sensitive to initializations, whereas AMDM and the hybrid algorithm where we decrease the number of singular values inverted by 1 after every 10 iterations have a tighter max-min confidence interval suggesting that they are less susceptible to initialization and maintain a higher fitness with a several orders of magnitude lower condition number. Also, note that all the AMDM variants converge to a similar fitness (the confidence intervals also collapse) with a similar condition number which is close to unity for this case.

In Figure 8, we compute CP decomposition of the SCF tensor with rank R=200 (exceeding 2 of the 3 tensor dimensions). We use a relative tolerance criterion for computing $t= \underset{\sigma_i}{\operatorname{argmin}}_i(\frac{\sigma_{\max}}{\sigma_i} < 100)$ in the hybrid algorithm; i.e., singular values σ_i are inverted only if $\frac{\sigma_{\max}}{\sigma_i} < 100$, where σ_{\max} is the maximum singular value. The number of initialization used here is 5, and the hybrid algorithm outperforms ALS in terms of fitness, as the mean fitness reaches 0.992 fitness in 150 iterations, whereas ALS reaches 0.95 in 200 iterations; the max-min confidence bounds suggests that the hybrid algorithm is a clear winner here.

7. Conclusion. In this work, we have proposed an alternating optimization algorithm, AMDM, to compute a CP decomposition of the tensor. This algorithm achieves superlinear local convergence for exact CP rank problems when the CP rank is smaller than or equal to all the mode lengths of the tensor with the same asymptotic computational cost as that of ALS. For approximating a tensor via CP decomposition, we theoretically show that the algorithm locally converges to the stationary points of (3.3) for tensors with special CP structure. Although the existence of these stationary points for any tensor is an open problem, we empirically confirm that the AMDM algorithm converges to these stationary points for various tensors. Viewing the algorithm as minimizing a Mahalanobis distance helps in generalization of the method for CP rank larger than the mode lengths and interpolation between the AMDM and ALS algorithms. We also formulate an efficient way to compute the CP decomposition condition number to track the condition of the decomposition throughout the algorithm. Our numerical experiments confirm that interpolation of algorithms between AMDM and ALS leads to a better conditioned decomposition without significant difference in fitness as compared to ALS for synthetic tensors as well as most of the tested real-world tensors. We provide an intuitive reasoning of this phenomenon and leave the detailed analysis as a future direction of research.

Appendix A. Computing the condition number of a CP decomposition. It has been shown that the CP decomposition condition number is the reciprocal of the smallest singular value of a matrix called Terracini's matrix. This matrix consists of the orthogonal basis for the tangent space of each of the rank-1 components of the reconstructed tensor. We will refer to the normalized condition number as the condition number of CP decomposition, and we refer the reader to [10, 66] for details about how the notion of condition number of a CP decomposition is defined and derived. Consider an equidimensional order 3 real tensor \mathcal{X} with mode length s. Let the CP decomposition approximation of rank s be given by r0; then Terracini's matrix r0 is r1 is r2, where for all r3, where for all r3, where for all r4, is r5.

$$oldsymbol{V}_i = [oldsymbol{a}_i \otimes oldsymbol{b}_i \otimes oldsymbol{c}_i \quad oldsymbol{Q}_{oldsymbol{a}_i}^ot \otimes oldsymbol{c}_i \quad oldsymbol{a}_i \otimes oldsymbol{b}_i \otimes oldsymbol{C}_i^ot],$$

and $Q_{a_i}^{\perp} \in \mathbb{R}^{s \times (s-1)}$ is an orthogonal basis of the orthogonal complement of a_i , and $Q_{b_i}^{\perp}$, and $Q_{c_i}^{\perp}$ are defined similarly. Consequently, Terrracini's matrix is of size $s^3 \times R(3s-2)$, and the computational cost of computing the smallest singular value via a Krylov subspace method is $O(s^5R^2)$. For an order N tensor, this cost is $O(N^2s^{N+2}R^2)$ and therefore expensive to compute for decompositions with moderately large mode lengths.

The cost of computing the condition number can be decreased significantly for when the rank of the CP decomposition is less than all the mode lengths of the input tensor, i.e., if R < s. Assume that $R \le s$; then since the condition number is invariant to orthogonal transformations [10], for a CP decomposition of an order 3 tensor,

$$\kappa(\llbracket \boldsymbol{D}; \boldsymbol{A}, \boldsymbol{B}, \boldsymbol{C} \rrbracket) = \kappa(\llbracket \boldsymbol{D}; \boldsymbol{Q}_{\boldsymbol{A}}^T \boldsymbol{A}, \boldsymbol{Q}_{\boldsymbol{B}}^T \boldsymbol{B}, \boldsymbol{Q}_{\boldsymbol{C}}^T \boldsymbol{C} \rrbracket),$$

where $Q_A \in \mathbb{R}^{s \times s} = [Q_A^{(1)}Q_A^{(2)}]$ and the columns of $Q_A^{(1)} \in \mathbb{R}^{s \times R}$ are an orthogonal basis of the column space of A, while the columns of $Q_A^{(2)} \in \mathbb{R}^{s \times (s-R)}$ are an orthogonal basis for the orthogonal complement of the column space of A. We define $Q_B = [Q_B^{(1)}Q_B^{(2)}]$ and $Q_C = [Q_C^{(1)}Q_C^{(2)}]$ similarly. The transformed Terracini's matrix $U = [U_1 \dots U_R]$, where for all $i \in \{1, \dots, R\}$,

$$egin{aligned} oldsymbol{U}_i = egin{bmatrix} oldsymbol{Q}_{oldsymbol{A}}^T oldsymbol{a}_i \otimes oldsymbol{Q}_{oldsymbol{B}}^T oldsymbol{b}_i \otimes oldsymbol{Q}_{oldsymbol{C}}^T oldsymbol{c}_i \ oldsymbol{U}_i^{(1)} & oldsymbol{Q}_{oldsymbol{a}_i}^{\perp} \otimes oldsymbol{Q}_{oldsymbol{C}}^T oldsymbol{c}_i \ oldsymbol{Q}_{oldsymbol{A}}^T oldsymbol{a}_i \otimes oldsymbol{Q}_{oldsymbol{C}}^T oldsymbol{c}_i \ oldsymbol{U}_i^{(2)} \ oldsymbol{U}_i^{(2)} \ oldsymbol{Q}_{oldsymbol{A}}^T oldsymbol{a}_i \otimes oldsymbol{Q}_{oldsymbol{C}}^T oldsymbol{c}_i \ oldsymbol{U}_i^{(2)} \ oldsymbol{U}_i^{(2)} \ oldsymbol{U}_i^{(3)} \ oldsymbol{U}_i^{(3)} \ oldsymbol{U}_i^{(3)} \ oldsymbol{Q}_{oldsymbol{C}}^T oldsymbol{c}_i \ oldsymbol{Q}_{oldsymbol{A}}^T oldsymbol{a}_i \otimes oldsymbol{Q}_{oldsymbol{C}}^T oldsymbol{c}_i \ oldsymbol{Q}_i^T oldsymbol{a}_i \otimes oldsymbol{Q}_{oldsymbol{C}}^T oldsymbol{c}_i \ oldsymbol{Q}_i^T oldsymbol{c}_i \otimes oldsymbol{Q}_{oldsymbol{C}}^T oldsymbol{c}_i \ oldsymbol{Q}_i^T oldsymbol{c}_i \otimes oldsymbol{Q}_{oldsymbol{C}}^T oldsymbol{c}_i \ oldsymbol{Q}_i^T oldsymbol{c}_i \otimes oldsymbol{Q}_i^T oldsymbol{c}_i \ oldsymbol{c}_i \otimes oldsymbol{Q}_i^T oldsymbol{c}_i \otimes oldsymbol{Q}_i^T oldsymbol{c}_i \ oldsymbol{C}_i \otimes oldsymbol{C}_i \otimes oldsymbol{C}_i \ oldsymbol{C}_i \otimes old$$

where $\bar{\boldsymbol{Q}}_{\boldsymbol{a}_i}^{\perp} = \boldsymbol{Q}_{\boldsymbol{A}}^T \boldsymbol{Q}_{\boldsymbol{a}_i}^{\perp} \in \mathbb{R}^{s \times (s-1)}$ is an orthogonal basis of the orthogonal complement of $\boldsymbol{Q}_{\boldsymbol{A}}^T \boldsymbol{a}_i$ and $\bar{\boldsymbol{Q}}_{\boldsymbol{b}_i}^{\perp}$, and $\bar{\boldsymbol{Q}}_{\boldsymbol{c}_i}^{\perp}$ are defined similarly. Note that $\boldsymbol{U}_i^{(j)T} \boldsymbol{U}_i^{(k)} = \boldsymbol{0}$ for $j \neq k$, since $\bar{\boldsymbol{Q}}_{\boldsymbol{a}_i}^{\perp T} \boldsymbol{Q}_{\boldsymbol{A}}^T \boldsymbol{a}_i = 0$. Consequently,

$$\sigma_{\min}(\boldsymbol{U}) = \min_{j \in \{2,3,4\}} \sigma_{\min} \bigg(\begin{bmatrix} \boldsymbol{U}_1^{(1)} & \boldsymbol{U}_1^{(j)} & \dots & \boldsymbol{U}_R^{(1)} & \boldsymbol{U}_R^{(j)} \end{bmatrix} \bigg).$$

We analyze $[\boldsymbol{U}_i^{(1)} \, \boldsymbol{U}_i^{(j)}]$ instead of $\boldsymbol{U}_i^{(j)}$ separately in order to have Kronecker products with orthogonal (square) matrices instead of nonsquare matrices with orthonormal columns, thereby simplifying analysis. After the above transformation, we can obtain a reduced form of smaller dimensions each of the four matrices to compute the condition number more efficiently. Note that $\boldsymbol{Q}_A^T \boldsymbol{a}_i = [\boldsymbol{Q}_A^{(1)T} \boldsymbol{a}_i]$, and similarly for $\boldsymbol{Q}_B^T \boldsymbol{b}_i$ and $\boldsymbol{Q}_C^T \boldsymbol{c}_i$. Further, we can choose the columns \boldsymbol{Q}_a^{\perp} so that $\bar{\boldsymbol{Q}}_{a_i}^{\perp} = \boldsymbol{Q}_A^T \boldsymbol{Q}_{a_i}^{\perp} = [\boldsymbol{Q}_{a_i}^{(1)T} \boldsymbol{a}_{a_i}]$, where $\boldsymbol{Q}_{A}^{(1)T} \boldsymbol{a}_i \in \mathbb{R}^{R \times (R-1)}$ is an orthogonal basis of the orthogonal complement of $\boldsymbol{Q}_A^{(1)T} \boldsymbol{a}_i$, and similarly for $\boldsymbol{Q}_{b_i}^{\perp}$ and $\boldsymbol{Q}_{c_i}^{\perp}$. Consequently, for j=2,

$$\begin{split} &\sigma_{\min} \Bigg(\left[\boldsymbol{U}_{1}^{(1)} \quad \boldsymbol{U}_{1}^{(2)} \quad \dots \quad \boldsymbol{U}_{R}^{(1)} \quad \boldsymbol{U}_{R}^{(2)} \right] \Bigg) \\ &= \sigma_{\min} \Bigg(\left[\boldsymbol{Q}_{A}^{T} \boldsymbol{a}_{1} \otimes \boldsymbol{Q}_{B}^{T} \boldsymbol{b}_{1} \otimes \boldsymbol{Q}_{C}^{T} \boldsymbol{c}_{1} \quad \bar{\boldsymbol{Q}}_{\boldsymbol{a}_{1}}^{\perp} \otimes \boldsymbol{Q}_{B}^{T} \boldsymbol{b}_{1} \otimes \boldsymbol{Q}_{C}^{T} \boldsymbol{c}_{1} \dots \right] \Bigg) \\ &= \sigma_{\min} \Bigg(\left[\boldsymbol{Q}_{A}^{(1)T} \boldsymbol{a}_{1} \quad \boldsymbol{Q}_{\boldsymbol{Q}_{A}^{(1)T} \boldsymbol{a}_{1}}^{\perp} \quad \boldsymbol{0} \right] \otimes \left[\boldsymbol{Q}_{B}^{(1)T} \boldsymbol{b}_{1} \right] \otimes \left[\boldsymbol{Q}_{C}^{(1)T} \boldsymbol{c}_{1} \right] \quad \dots \Bigg) \\ &= \min \Bigg\{ \sigma_{\min} \Bigg(\left[\boldsymbol{Q}_{\boldsymbol{Q}_{A}^{(1)T} \boldsymbol{a}_{1}} \otimes \boldsymbol{Q}_{B}^{(1)T} \boldsymbol{b}_{1} \otimes \boldsymbol{Q}_{C}^{(1)T} \boldsymbol{c}_{1} \dots \right] \Bigg), \sigma_{\min} \Bigg(\left[\boldsymbol{Q}_{B}^{(1)T} \boldsymbol{b}_{1} \otimes \boldsymbol{Q}_{C}^{(1)T} \boldsymbol{c}_{1} \dots \right] \Bigg) \Bigg\} \\ &= \sigma_{\min} \Bigg(\left[\boldsymbol{Q}_{\boldsymbol{Q}_{A}^{(1)T} \boldsymbol{a}_{1}} \otimes \boldsymbol{Q}_{B}^{(1)T} \boldsymbol{b}_{1} \otimes \boldsymbol{Q}_{C}^{(1)T} \boldsymbol{c}_{1} \quad \dots \quad \boldsymbol{Q}_{\boldsymbol{Q}_{A}^{(1)T} \boldsymbol{a}_{R}} \otimes \boldsymbol{Q}_{B}^{(1)T} \boldsymbol{b}_{R} \otimes \boldsymbol{Q}_{C}^{(1)T} \boldsymbol{c}_{R} \right] \Bigg) \\ &= \sigma_{\min} \Bigg(\left[\boldsymbol{Q}_{A}^{(1)T} \boldsymbol{a}_{1} \otimes \boldsymbol{Q}_{B}^{(1)T} \boldsymbol{b}_{1} \otimes \boldsymbol{Q}_{C}^{(1)T} \boldsymbol{c}_{1} \quad \boldsymbol{Q}_{\boldsymbol{Q}_{A}^{(1)T} \boldsymbol{a}_{1}} \otimes \boldsymbol{Q}_{B}^{(1)T} \boldsymbol{b}_{1} \otimes \boldsymbol{Q}_{C}^{(1)T} \boldsymbol{c}_{1} \dots \right] \Bigg) \\ &= \sigma_{\min} \Bigg(\left[\bar{\boldsymbol{U}}_{1}^{(1)} \quad \bar{\boldsymbol{U}}_{1}^{(2)} \quad \dots \quad \bar{\boldsymbol{U}}_{R}^{(1)} \quad \bar{\boldsymbol{U}}_{R}^{(2)} \right] \Bigg), \end{split}$$

where $Q_{Q_A^{(1)_T}a_i} = [Q_A^{(1)_T}a_iQ_{Q_A^{(1)_T}a_i}^{\perp}] \in \mathbb{R}^{R \times R}$ is an orthogonal (square) matrix. A similar argument can also be shown for j=3,4 to show that

$$\sigma_{\min}(oldsymbol{U}) = \min_{j \in \{2,3,4\}} \sigma_{\min} \left(egin{array}{ccc} ar{oldsymbol{U}}_1^{(1)} & ar{oldsymbol{U}}_1^{(j)} & \dots & ar{oldsymbol{U}}_R^{(1)} & ar{oldsymbol{U}}_R^{(j)} \end{bmatrix}
ight) = \sigma_{\min}(ar{oldsymbol{U}}),$$

where

$$ar{m{U}}_i = egin{bmatrix} ar{m{u}}_i \otimes ar{m{b}}_i \otimes ar{m{c}}_i & m{Q}_{m{ar{a}}_i}^ot \otimes ar{m{b}}_i \otimes ar{m{c}}_i & ar{m{a}}_i \otimes m{Q}_{ar{ar{b}}_i}^ot \otimes ar{m{c}}_i & ar{m{a}}_i \otimes m{Q}_{ar{ar{c}}_i}^ot \otimes ar{m{c}}_i & ar{m{d}}_i \otimes ar{m{b}}_i \otimes m{Q}_{ar{ar{c}}_i}^ot \end{bmatrix},$$

 $\bar{\boldsymbol{a}}_i = \boldsymbol{Q}_{\boldsymbol{A}}^{(1)T} \boldsymbol{a}_i, \bar{\boldsymbol{b}}_i = \boldsymbol{Q}_{\boldsymbol{B}}^{(1)T} \boldsymbol{b}_i$, and $\bar{\boldsymbol{c}}_i = \boldsymbol{Q}_{\boldsymbol{C}}^{(1)T} \boldsymbol{c}_i$. The size of the above reduced Terracini's matrix is $R^3 \times R(3R-2)$; hence a direct computation of the singular value decomposition can be used to compute the condition number with a cost of $O(R^7)$ for R < s.

All the above arguments can be easily generalized to an order N nonequidimensional tensor. Therefore, we showed that the condition number of CP decomposition is invariant to the following transformation:

$$\kappa(\llbracket \boldsymbol{D}; \boldsymbol{A}_1 \dots, \boldsymbol{A}_N \rrbracket) = \kappa\left(\llbracket \boldsymbol{D}; \boldsymbol{Q}_{\boldsymbol{A}_1}^{(1)T} \boldsymbol{A}_1, \dots, \boldsymbol{Q}_{\boldsymbol{A}_N}^{(N)T} \boldsymbol{A}_N \rrbracket\right),$$

where for all $i \in \{1, ..., N\}$, the columns of $Q_{A_i}^{(1)}$ are an orthonormal basis of the column space of A_i .

Acknowledgments. The authors would like to thank Ardavan (Ari) Afshar for detailed discussions about his work on minimizing Wasserstein distance between tensors, from which this work is derived. The authors would also like to thank Jimeng Sun, Cheng Qian, and Chaoqi Yang for fruitful discussions and for providing datasets which motivated this work. We would also like to thank Nick Dewaele and Nick Vannieuwenhoven for their useful comments and feedback on the manuscript. We would also like to thank the anonymous reviewers for their valuable feedback and suggestions which improved the manuscript by a substantial amount.

REFERENCES

- [1] E. ACAR, D. M. DUNLAVY, AND T. G. KOLDA, A scalable optimization approach for fitting canonical tensor decompositions, J. Chemometrics, 25 (2011), pp. 67–86.
- [2] A. Afshar, K. Yin, S. Yan, C. Qian, J. C. Ho, H. Park, and J. Sun, Swift: Scalable Wasserstein factorization for sparse nonnegative tensors, in Proceedings of the AAAI Conference, 2021.
- [3] A. Anandkumar, R. Ge, D. J. Hsu, S. M. Kakade, and M. Telgarsky, *Tensor decompositions for learning latent variable models*, J. Mach. Learn. Res., 15 (2014), pp. 2773–2832.
- [4] G. BALLARD, K. HAYASHI, AND R. KANNAN, Parallel Nonnegative CP Decomposition of Dense Tensors, preprint, arXiv:1806.07985, 2018.
- [5] G. BALLARD, N. KNIGHT, AND K. ROUSE, Communication lower bounds for matricized tensor times Khatri-Rao product, in Proceedings of the International Parallel and Distributed Processing Symposium (IPDPS), IEEE, 2018, pp. 557–567.
- [6] C. Battaglino, G. Ballard, and T. G. Kolda, A practical randomized CP tensor decomposition, SIAM J. Matrix Anal. Appl., 39 (2018), pp. 876–901.
- [7] A. Bellet, A. Habrard, and M. Sebban, A Survey on Metric Learning for Feature Vectors and Structured Data, preprint, arXiv:1306.6709, 2013.
- [8] A. BELOUCHRANI, K. ABED-MERAIM, J.-F. CARDOSO, AND E. MOULINES, A blind source separation technique using second-order statistics, IEEE Trans. Signal Process., 45 (1997), pd. 434–444.
- [9] S. BISWAL, H. SUN, B. GOPARAJU, M. B. WESTOVER, J. SUN, AND M. T. BIANCHI, Expertlevel sleep scoring with deep neural networks, J. Amer. Med. Inform. Assoc., 25 (2018), pp. 1643–1650.
- [10] P. Breiding and N. Vannieuwenhoven, The condition number of joint decompositions, SIAM J. Matrix Anal. Appl., 39 (2018), pp. 287–309.
- [11] R. Bro, PARAFAC tutorial and applications, Chemometrics Intell. Lab. Syst., 38 (1997), pp. 149–171.
- [12] J. D. CARROLL AND J.-J. CHANG, Analysis of individual differences in multidimensional scaling via an n-way generalization of Eckart-Young decomposition, Psychometrika, 35 (1970), pp. 283-319.
- [13] J. CHOI, X. LIU, S. SMITH, AND T. SIMON, Blocking optimization techniques for sparse tensor computation, in Proceedings of the International Parallel and Distributed Processing Symposium (IPDPS), IEEE, 2018, pp. 568–577.
- [14] P. COMON, Tensor diagonalization, a useful tool in signal processing, IFAC Proc., 27 (1994), pp. 77–82.
- [15] F. Cong, Q.-H. Lin, L.-D. Kuang, X.-F. Gong, P. Astikainen, and T. Ristaniemi, Tensor decomposition of EEG signals: A brief review, J. Neurosci. Methods, 248 (2015), pp. 59–69.
- [16] C.-F. Cui, Y.-H. Dai, and J. Nie, All real eigenvalues of symmetric tensors, SIAM J. Matrix Anal. Appl., 35 (2014), pp. 1582–1601.
- [17] M. CUTURI AND D. AVIS, Ground metric learning, J. Mach. Learn. Res., 15 (2014), pp. 533–564.
- [18] R. DE MAESSCHALCK, D. JOUAN-RIMBAUD, AND D. L. MASSART, The Mahalanobis distance, Chemometrics Intell. Lab. Syst., 50 (2000), pp. 1–18.
- [19] P. C. MAHALANOBIS, On the generalised distance in statistics, Proc. Natl. Instit. Sci. India, 2 (1936), pp. 49–55.
- [20] N. Dewaele, P. Breiding, and N. Vannieuwenhoven, The Condition Number of Many Tensor Decompositions Is Invariant Under Tucker Compression, manuscript, 2021.
- [21] R. A. HARSHMAN, Foundations of the PARAFAC Procedure: Models and Conditions for an Explanatory Multimodal Factor Analysis, University of California at Los Angeles, Los Angeles, CA, 1970.
- [22] K. HAYASHI, G. BALLARD, J. JIANG, AND M. TOBIA, Shared Memory Parallelization of MT-TKRP for Dense Tensors, preprint, arXiv:1708.08976, 2017.
- [23] C. J. HILLAR AND L.-H. LIM, Most tensor problems are NP-hard, J. ACM, 60 (2013), pp. 45:1–45:39.
- [24] F. L. HITCHCOCK, The expression of a tensor or a polyadic as a sum of products, Stud. Appl. Math., 6 (1927), pp. 164–189.
- [25] A. HYVÄRINEN, Survey on Independent Component Analysis, manuscript, 1999.
- [26] L. KARLSSON, D. KRESSNER, AND A. USCHMAJEW, Parallel algorithms for tensor completion in the CP format, Parallel Comput., 57 (2016), pp. 222–234.
- [27] O. KAYA, High Performance Parallel Algorithms for Tensor Decompositions, Ph.D. thesis, Université de Lyon, 2017.

- [28] O. KAYA AND Y. ROBERT, Computing dense tensor decompositions with optimal dimension trees, Algorithmica, 81 (2019), pp. 2092–2121.
- [29] O. KAYA AND B. UÇAR, Parallel CP Decomposition of Sparse Tensors Using Dimension Trees, Ph.D. thesis, Inria-Research Centre Grenoble-Rhône-Alpes, 2016.
- [30] B. Kemp, A. H. Zwinderman, B. Tuk, H. A. Kamphuisen, and J. J. Oberye, Analysis of a sleep-dependent neuronal feedback loop: The slow-wave microcontinuity of the EEG, IEEE Trans. Biomed. Eng., 47 (2000), pp. 1185–1194.
- [31] T. G. KOLDA AND B. W. BADER, Tensor decompositions and applications, SIAM Rev., 51 (2009), pp. 455-500.
- [32] T. G. KOLDA AND J. R. MAYO, Shifted power method for computing tensor eigenpairs, SIAM J. Matrix Anal. Appl., 32 (2011), pp. 1095–1124.
- [33] B. Kulis, Metric learning: A survey, Found. Trends Mach. Learn., 5 (2013), pp. 287–364.
- [34] G. Li, L. Qi, And G. Yu, The Z-eigenvalues of a symmetric tensor and its application to spectral hypergraph theory, Numer. Linear Algebra Appl., 20 (2013), pp. 1001–1029.
- [35] J. LI, K. USEVICH, AND P. COMON, Globally convergent Jacobi-type algorithms for simultaneous orthogonal symmetric tensor diagonalization, SIAM J. Matrix Anal. Appl., 39 (2018), pp. 1–22.
- [36] M. LIANG AND B. ZHENG, Further results on Moore-Penrose inverses of tensors with application to tensor nearness problems, Comput. Math. Appl., 77 (2019), pp. 1282–1293.
- [37] L.-H. LIM, Singular values and eigenvalues of tensors: A variational approach, in Proceedings of the 1st IEEE International Workshop on Computational Advances in Multi-Sensor Adaptive Processing, IEEE, 2005, pp. 129–132.
- [38] L. MA AND E. SOLOMONIK, Accelerating Alternating Least Squares for Tensor Decomposition by Pairwise Perturbation, preprint, arXiv:1811.10573, 2018.
- [39] L. MA AND E. SOLOMONIK, Efficient parallel CP decomposition with pairwise perturbation and multi-sweep dimension tree, in Proceedings of the International Parallel and Distributed Processing Symposium (IPDPS), IEEE, 2021, pp. 412–421.
- [40] K. MARUHASHI, F. GUO, AND C. FALOUTSOS, Multiaspectforensics: Pattern mining on large-scale heterogeneous networks with tensor analysis, in Proceedings of the International Conference on Advances in Social Networks Analysis and Mining, IEEE, 2011, pp. 203–210.
- [41] B. C. MITCHELL AND D. S. BURDICK, Slowly converging PARAFAC sequences: Swamps and two-factor degeneracies, J. Chemometrics, 8 (1994), pp. 155–168.
- [42] D. MITCHELL, N. YE, AND H. DE STERCK, Nesterov Acceleration of Alternating Least Squares for Canonical Tensor Decomposition, preprint, arXiv:1810.05846, 2018.
- [43] M. J. Mohlenkamp, Musings on multilinear fitting, Linear Algebra Appl., 438 (2013), pp. 834–852.
- [44] K. R. MURPHY, C. A. STEDMON, D. GRAEBER, AND R. BRO, Fluorescence spectroscopy and multi-way techniques. PARAFAC, Anal. Methods, 5 (2013), pp. 6557–6566.
- [45] D. NION AND L. DE LATHAUWER, An enhanced line search scheme for complex-valued tensor decompositions. Application in DS-CDMA, Signal Process., 88 (2008), pp. 749-755.
- [46] P. PAATERO, A weighted non-negative least squares algorithm for three-way PARAFAC factor analysis, Chemometrics Intell. Lab. Syst., 38 (1997), pp. 223–242.
- [47] A.-H. Phan, P. Tichavský, and A. Cichocki, Fast alternating LS algorithms for high order CANDECOMP/PARAFAC tensor factorizations, IEEE Trans. Signal Process., 61 (2013), pp. 4834–4846.
- [48] A.-H. Phan, P. Tichavský, and A. Cichocki, Low complexity damped Gauss-Newton algorithms for CANDECOMP/PARAFAC, SIAM J. Matrix Anal. Appl., 34 (2013), pp. 126–147.
- [49] L. QI, H. CHEN, AND Y. CHEN, Tensor Eigenvalues and their Applications, Springer, Berlin, 2018.
- [50] M. RAJIH, P. COMON, AND R. A. HARSHMAN, Enhanced line search: A novel method to accelerate PARAFAC, SIAM J. Matrix Anal. Appl., 30 (2008), pp. 1128-1147.
- [51] E. ROBEVA, Orthogonal decomposition of symmetric tensors, SIAM J. Matrix Anal. Appl., 37 (2016), pp. 86–102.
- [52] E. ROBEVA AND A. SEIGAL, Singular vectors of orthogonally decomposable tensors, Linear Multilinear Algebra, 65 (2017), pp. 2457–2471.
- [53] M. D. SCHATZ, T. M. LOW, R. A. VAN DE GEIJN, AND T. G. KOLDA, Exploiting symmetry in tensors for high performance: Multiplication with symmetric tensors, SIAM J. Sci. Comput., 36 (2014), pp. C453–C479.

- [54] 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, pp. 3551–3582.
- [55] N. Singh, L. Ma, H. Yang, and E. Solomonik, Comparison of accuracy and scalability of Gauss-Newton and alternating least squares for CANDECOMC/PARAFAC decomposition, SIAM J. Sci. Comput., 43 (2021), pp. C290-C311.
- [56] L. SORBER, I. DOMANOV, M. VAN BAREL, AND L. LATHAUWER, Exact line and plane search for tensor optimization, Comput. Optim. Appl., 63 (2016), pp. 121–142.
- [57] L. SORBER, M. VAN BAREL, AND L. DE LATHAUWER, Optimization-based algorithms for tensor decompositions: Canonical polyadic decomposition, decomposition in rank-(l_r,l_r,1) terms, and a new generalization, SIAM J. Optim., 23 (2013), pp. 695–720.
- [58] L. Sun, B. Zheng, C. Bu, and Y. Wei, Moore-Penrose inverse of tensors via Einstein product, Linear Multilinear Algebra, 64 (2016), pp. 686–698.
- [59] Q. Sun, T. C. Berkelbach, N. S. Blunt, G. H. Booth, S. Guo, Z. Li, J. Liu, J. D. Mc-Clain, E. R. Sayfutyarova, S. Sharma, et al., PySCF: The Python-based simulations of chemistry framework, Wiley Interdiscip. Rev. Comput. Mol. Sci., 8 (2018), e1340.
- [60] P. Tichavský, A. H. Phan, and A. Cichocki, A further improvement of a fast damped Gauss-Newton algorithm for CANDECOMP-PARAFAC tensor decomposition, in Proceedings of the International Conference on Acoustics, Speech and Signal Processing, IEEE, 2013, pp. 5964–5968.
- [61] P. TICHAVSKÝ, A. H. PHAN, AND A. CICHOCKI, Non-orthogonal tensor diagonalization, Signal Process., 138 (2017), pp. 313–320.
- [62] G. Tomasi and R. Bro, PARAFAC and missing values, Chemometrics Intell. Lab. Syst., 75 (2005), pp. 163–180.
- [63] G. TOMASI AND R. BRO, A comparison of algorithms for fitting the PARAFAC model, Comput. Statist. Data Anal., 50 (2006), pp. 1700–1734.
- [64] A. USCHMAJEW, Local convergence of the alternating least squares algorithm for canonical tensor approximation, SIAM J. Matrix Anal. Appl., 33 (2012), pp. 639–652.
- [65] K. USEVICH, J. LI, AND P. COMON, Approximate matrix and tensor diagonalization by unitary transformations: Convergence of Jacobi-type algorithms, SIAM J. Optim., 30 (2020), pp. 2998–3028.
- [66] N. VANNIEUWENHOVEN, Condition numbers for the tensor rank decomposition, Linear Algebra Appl., 535 (2017), pp. 35–86.
- [67] N. VANNIEUWENHOVEN, K. MEERBERGEN, AND R. VANDEBRIL, Computing the gradient in optimization algorithms for the CP decomposition in constant memory through tensor blocking, SIAM J. Sci. Comput., 37 (2015), pp. C415-C438.
- [68] G. Zen, E. Ricci, and N. Sebe, Simultaneous ground metric learning and matrix factorization with earth mover's distance, in Proceedings of the 22nd International Conference on Pattern Recognition, IEEE, 2014, pp. 3690–3695.