



A New Alternating Optimization Algorithm for CP Decomposition

Navjot Singh and Edgar Solomonik

 ·  ·  ·  @CS@Illinois

Department of Computer Science
University of Illinois at Urbana-Champaign

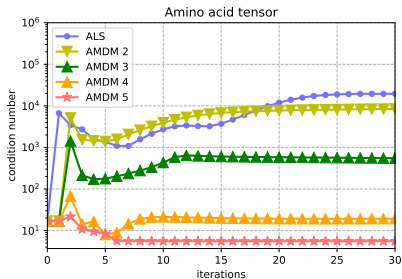
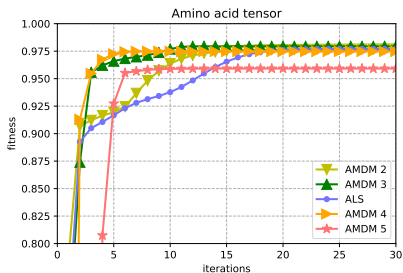
SIAM Conference on Parallel Processing for Scientific Computing
(PP22)

Outline

- 1 Overview
- 2 Motivation
- 3 New Alternating Update Scheme for CPD
- 4 Exact CP Decomposition
- 5 Approximate CP Decomposition
- 6 Conclusion and Future Work

Highlights

- Introduce computation of singular vals/vecs via considering multilinear function associated with the tensor with log barrier penalty
- Critical points of the above spectrally diagonalize an order N tensor
- Analyze local convergence of the algorithm for exact CPD of rank lesser than mode lengths
- Formulation that generalizes the algorithm to perform well conditioned¹ approximate CPD



¹P. Breiding and N. Vannieuwenhoven, SIMAX 2018

Overview

- Tensor: A multidimensional array \mathcal{X}
- Indices: $x_{i_1, i_2 \dots i_N}$ imply order = N
- CP tensor decomposition breaks down a tensor into sum of rank 1 components.
- CPD of an order 3 tensor \mathcal{X} with rank R and factor matrices $\mathbf{A}, \mathbf{B}, \mathbf{C}$:
 $\mathcal{X} = \llbracket \mathbf{A}, \mathbf{B}, \mathbf{C} \rrbracket$

$$x_{ijk} = \sum_{l=1}^R a_{il} b_{jl} c_{kl}$$

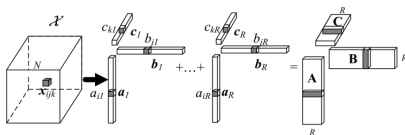


Figure: S.He et. al. Tensor Decomposition Based Electrical Data Recovery

Motivation: Singular Vectors via Variational Approach

Analogous to obtaining eigenvalues via critical points of $\mathbf{x}^T \mathbf{A} \mathbf{x}$ with unit l^2 -norm constraints,

L. Lim derives singular vectors and values ² of \mathbf{A} via

- critical points of $\frac{\mathbf{x}^T \mathbf{A} \mathbf{y}}{\|\mathbf{x}\|_2 \|\mathbf{y}\|_2}$ with unit norm constraints.
- Lagrangian is

$$L(\mathbf{x}, \mathbf{y}, \sigma) = \mathbf{x}^T \mathbf{A} \mathbf{y} - \sigma(\|\mathbf{x}\|_2 \|\mathbf{y}\|_2 - 1)$$

- First order conditions yields,

$$\mathbf{A} \frac{\mathbf{y}}{\|\mathbf{y}\|_2} = \sigma \frac{\mathbf{x}}{\|\mathbf{x}\|_2}, \quad \mathbf{A}^T \frac{\mathbf{x}}{\|\mathbf{x}\|_2} = \sigma \frac{\mathbf{y}}{\|\mathbf{y}\|_2}, \quad \|\mathbf{x}\|_2 \|\mathbf{y}\|_2 = 1$$

$$\mathbf{A} \mathbf{v} = \sigma \mathbf{u}, \quad \mathbf{A}^T \mathbf{u} = \sigma \mathbf{v}$$

Order 3 tensor eigen/singular values and vectors can be derived similarly

²Lek-Heng Lim Singular values and Eigenvalues of a tensor: A variational approach

Motivation: Singular Vectors via Lagrangian Optimization

One can also obtain singular values and vectors by considering bilinear form $f(\mathbf{x}, \mathbf{y}) = \mathbf{x}^T \mathbf{A} \mathbf{y}$ with $\|\mathbf{x}\|_2 \neq 0$, $\|\mathbf{y}\|_2 \neq 0$,

$$\mathcal{L}_f(\mathbf{x}, \mathbf{y}) = \mathbf{x}^T \mathbf{A} \mathbf{y} - \log(\|\mathbf{x}\|_2 \|\mathbf{y}\|_2)$$

Critical points satisfy $\mathbf{A} \mathbf{v} = \sigma \mathbf{u}$ and $\mathbf{A}^T \mathbf{u} = \sigma \mathbf{v}$ for

$$\mathbf{u} = \mathbf{x} / \|\mathbf{x}\|_2, \quad \mathbf{v} = \mathbf{y} / \|\mathbf{y}\|_2, \quad \sigma = 1 / (\|\mathbf{x}\|_2 \|\mathbf{y}\|_2).$$

Similarly for an order 3 tensor \mathcal{T} , consider

$$\mathcal{L}_f(\mathbf{x}, \mathbf{y}, \mathbf{z}) = \sum_{i,j,k} t_{ijk} x_i y_j z_k - \log(\|\mathbf{x}\|_2 \|\mathbf{y}\|_2 \|\mathbf{z}\|_2).$$

Critical points satisfy equations

$$\sum_{j,k} t_{ijk} v_j w_k = \sigma \mathbf{u}, \quad \sum_{i,k} t_{ijk} u_i w_k = \sigma \mathbf{v}, \quad \sum_{i,j} t_{ijk} u_i v_j = \sigma \mathbf{w}$$

with $\mathbf{u} = \mathbf{x} / \|\mathbf{x}\|_2$, $\mathbf{v} = \mathbf{y} / \|\mathbf{y}\|_2$, $\mathbf{w} = \mathbf{z} / \|\mathbf{z}\|_2$, $\sigma = 1 / (\|\mathbf{x}\|_2 \|\mathbf{y}\|_2 \|\mathbf{z}\|_2)$

Motivation: Spectral Diagonalization

This notion can be generalized for $R > 1$ vectors, since

$$\mathbf{x}^T \mathbf{A} \mathbf{y} = \langle \mathbf{A}, \mathbf{x} \mathbf{y}^T \rangle$$

consider $\langle \mathbf{A}, \mathbf{B} \rangle = \langle \text{vec}(\mathbf{A}), \text{vec}(\mathbf{B}) \rangle$,

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

$$\begin{aligned} \mathcal{L}_f(\mathbf{X}, \mathbf{Y}) &= \langle \mathbf{A}, \mathbf{X} \mathbf{Y}^T \rangle - \frac{1}{2} (\log(\det(\mathbf{X}^T \mathbf{X})) - \log(\det(\mathbf{Y}^T \mathbf{Y}))) \\ &= \text{tr}(\mathbf{X}^T \mathbf{A} \mathbf{Y}) - \frac{1}{2} \text{tr}(\log(\mathbf{X}^T \mathbf{X} \mathbf{Y}^T \mathbf{Y})). \end{aligned}$$

The critical points of \mathcal{L}_f satisfy $\mathbf{A} \mathbf{Y} \mathbf{X}^T \cong \mathbf{I}$ and $\mathbf{A}^T \mathbf{X} \mathbf{Y}^T \cong \mathbf{I}$

- $\mathbf{X} \rightarrow$ invariant subspace of $\mathbf{A} \mathbf{A}^T$
- $\mathbf{Y} \rightarrow$ invariant subspace of $\mathbf{A}^T \mathbf{A}$

and diagonalize \mathbf{A} in the sense that

$$\mathbf{X}^T \mathbf{A} \mathbf{Y} = \mathbf{I}$$

Motivation: Spectral Diagonalization

Similarly for an order 3 tensor \mathcal{T} , $\langle \mathcal{T}, \mathcal{Y} \rangle = \langle \text{vec}(\mathcal{T}), \text{vec}(\mathcal{Y}) \rangle$

$$\mathcal{L}_f(\mathbf{X}, \mathbf{Y}, \mathbf{Z}) = \langle \mathcal{T}, \llbracket \mathbf{X}, \mathbf{Y}, \mathbf{Z} \rrbracket \rangle - \frac{1}{2} \text{tr}(\log(\mathbf{X}^T \mathbf{X} \mathbf{Y}^T \mathbf{Y} \mathbf{Z}^T \mathbf{Z}))$$

The critical points of \mathcal{L}_f diagonalize the \mathcal{T} such that

$$\mathcal{P} = \mathcal{T} \times_1 \mathbf{X} \times_2 \mathbf{Y} \times_3 \mathbf{Z},$$

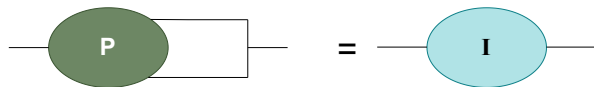


Figure: $p_{ijj} = p_{iij} = p_{iji} = \delta_{ij}$

implying \mathcal{P} has R elementary eigenvectors with unit eigenvalues, different from \mathcal{P} . Comon's idea of diagonalizing a tensor with orthogonal matrices ³

³P. Comon, M. Sorensen Tensor diagonalization with orthogonal transformation

New Alternating Update Scheme

Consider a rank R CP decomposition of a tensor $\mathcal{X} \in \mathbb{R}^{s \times s \times s}$,

$$\mathcal{X} = \llbracket \mathbf{A}, \mathbf{B}, \mathbf{C} \rrbracket, \text{ i.e. } x_{ijk} = \sum_{r=1}^R a_{ir} b_{jr} c_{kr},$$

which maybe obtained by ALS via minimizing

$$f(\mathbf{A}, \mathbf{B}, \mathbf{C}) = \frac{1}{2} \|\mathcal{X} - \llbracket \mathbf{A}, \mathbf{B}, \mathbf{C} \rrbracket\|_F^2$$

by alternating updates such as

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

We propose a different update, which for $R \leq s$ is,

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

Convergence to Exact Decomposition

When seeking an exact CP decomposition of rank $R \leq s$

- ALS achieves a linear convergence rate ⁴
- High order convergence possible via optimizing all factors, eg. using Gauss-Newton ^{5,6,7}, but is expensive
- The proposed algorithm achieves atleast quartic convergence per sweep of alternating updates
 - per subsweep, convergence order is α where α is the real positive root of $x^{N-1} - \sum_{i=0}^{N-2} x^i$ for order N tensor, i.e., $(1 + \sqrt{5})/2$ for order 3.
 - cost per iteration roughly the same as ALS (dominated by MTTKRP) and therefore easily parallelizable

⁴A. Uschmajew, SIMAX 2012

⁵P. Paatero, Chemometrics and Intelligent Laboratory Systems 1997

⁶A.H. Phan, P. Tichavsky, A. Cichocki, SIMAX 2013

⁷N. Singh, L. Ma, H. Yang, E.S., SISC 2021.

Exact Decomposition Error Analysis

The error in one factor scales with the product of errors in the other factors

Lemma

Suppose $\mathcal{X} = \llbracket \mathbf{A}^{(1)}, \dots, \mathbf{A}^{(N)} \rrbracket$, where each $\mathbf{A}^{(i)} \in \mathbb{R}^{s_i \times R}$ is full rank with $s_i \geq R$ and $\tilde{\mathbf{A}}^{(n)} = \mathbf{A}^{(n)} \mathbf{D}^{(n)} + \mathbf{\Delta}^{(n)}$ and satisfies $\|\mathbf{\Delta}^{(n)}\|_F = \epsilon_n$ for $n = 1, \dots, N - 1$, then $\exists \epsilon > 0$ such that if $\epsilon_n < \epsilon$ for $n = 1, \dots, N - 1$,

$$\tilde{\mathbf{A}}^{(N)} = \mathcal{X}_{(N)} (\tilde{\mathbf{A}}^{(1)\dagger T} \odot \dots \odot \bar{\mathbf{A}}^{(N-1)\dagger T})$$

satisfies

$$\|\tilde{\mathbf{A}}^{(N)} \mathbf{D}^{(N)} - \mathbf{A}^{(N)}\|_F = O\left(\prod_{n=1}^N \epsilon_n^{N-1}\right),$$

for some diagonal $\mathbf{D}^{(N)}$.

Exact Decomposition Error Analysis

A rough sketch of proof of the above Lemma follows from substituting true decomposition in the update rule

$$\begin{aligned}\tilde{\mathbf{A}}^{(N)} &= \mathbf{A}^{(N)} \left((\tilde{\mathbf{A}}^{(1)\dagger} \mathbf{A}^{(1)}) * \dots * (\tilde{\mathbf{A}}^{(N-1)\dagger} \mathbf{A}^{(N-1)}) \right)^T \\ &= \mathbf{A}^{(N)} \left((\mathbf{D}^{(1)} - \tilde{\mathbf{A}}^{(1)\dagger} \mathbf{\Delta}^{(1)}) * \dots * (\mathbf{D}^{(N-1)} - \tilde{\mathbf{A}}^{(N-1)\dagger} \mathbf{\Delta}^{(N-1)}) \right)^T \\ &= \mathbf{A}^{(N)} \left(\mathbf{D} + (-1)^{N-1} \tilde{\mathbf{A}}^{(1)\dagger} \mathbf{\Delta}^{(1)} * \dots * \tilde{\mathbf{A}}^{(N-1)\dagger} \mathbf{\Delta}^{(N-1)} \right)^T,\end{aligned}$$

where \mathbf{D} is diagonal matrix.

Exact Decomposition Experimental Performance

Rate of convergence of AMDM only depends on the (matrix) rank of underlying factors

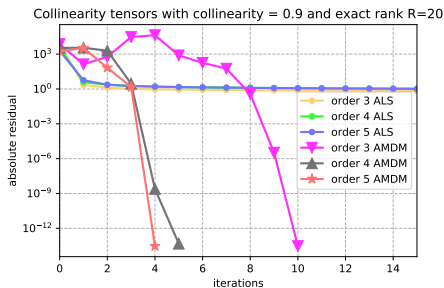


Figure: CP Decomposition of synthetic tensors with rank 20 and 100^3 entries

Approximate CP Decomposition

- The proposed update for \mathbf{A} minimizes

$$\frac{1}{2} \|(\mathbf{X} - \llbracket \mathbf{A}, \mathbf{B}, \mathbf{C} \rrbracket)_{(1)} (\mathbf{C}^{\dagger T} \otimes \mathbf{B}^{\dagger T})\|_F^2.$$

- The residual being

$$\mathbf{X}_{(1)} (\mathbf{C}^{\dagger T} \odot \mathbf{B}^{\dagger T}) - \mathbf{A} (\mathbf{I} \odot \mathbf{I})$$

- Residual transformation tends to equalize the weight of contribution of the error associated with different rank-1 parts of the CP decompositions.
- Similar property observed when Mahalanobis distance metric is considered

Mahalanobis Distance Objective

- Original motivation for the method came from optimizing CPD with general distance metrics with Ardavan Afshar, C. Qian, and J. Sun ⁸.
- Consider an order 3 tensor \mathcal{X} and Mahalanobis distance objective

$$f(\mathbf{A}, \mathbf{B}, \mathbf{C}) = \frac{1}{2} \|\mathcal{X} - \mathcal{Y}\|_M = \frac{1}{2} \text{vec}(\mathcal{X} - \mathcal{Y})^T \mathbf{M} \text{vec}(\mathcal{X} - \mathcal{Y}),$$

$$\text{where } \mathcal{Y} = \llbracket \mathbf{A}, \mathbf{B}, \mathbf{C} \rrbracket,$$

$$\text{with } \mathbf{M} = \bigotimes_{k=1}^3 \mathbf{M}^{(k)-1} \text{ being SPD.}$$

Minimization with respect to \mathbf{A} results in the following update

$$\mathbf{AZ} = \mathbf{X}_{(1)} \mathbf{L},$$

$$\text{where } \mathbf{L} = \left(\mathbf{M}^{(3)-1} \mathbf{C} \right) \odot \left(\mathbf{M}^{(2)-1} \mathbf{B} \right),$$

$$\text{and } \mathbf{Z} = \left(\mathbf{B}^T \mathbf{M}^{(2)-1} \mathbf{B} \right) * \left(\mathbf{C}^T \mathbf{M}^{(3)-1} \mathbf{C} \right).$$

⁸A. Ardavan, K. Yin, S. Yan, C. Qian, J.C. Ho, H. Park, and J. Sun, AAAI 2021

Generalizing AMDM to Hybrid Algorithms

Decompose factors into sum of two matrices and using first θ singular values and vectors for each factor to construct $\mathbf{M}^{(k)}$,

$$\mathbf{M}^{(1)} = \mathbf{A}_1 \mathbf{A}_1^T + (\mathbf{I} - \mathbf{A}_1 \mathbf{A}_1^\dagger),$$

$$\mathbf{M}^{(2)} = \mathbf{B}_1 \mathbf{B}_1^T + (\mathbf{I} - \mathbf{B}_1 \mathbf{B}_1^\dagger),$$

$$\mathbf{M}^{(3)} = \mathbf{C}_1 \mathbf{C}_1^T + (\mathbf{I} - \mathbf{C}_1 \mathbf{C}_1^\dagger).$$

leads to an update that is a hybrid of AMDM and ALS, since

$$\mathbf{AZ} = \mathbf{X}_{(1)} \mathbf{L},$$

$$\text{where } \mathbf{L} = \left((\mathbf{C}_1^{\dagger T} + \mathbf{C}_2) \odot (\mathbf{B}_1^{\dagger T} + \mathbf{B}_2) \right),$$

$$\text{and } \mathbf{Z} = \left((\mathbf{C}_1^\dagger \mathbf{C}_1 + \mathbf{C}_2^T \mathbf{C}_2) * (\mathbf{B}_1^\dagger \mathbf{B}_1 + \mathbf{B}_2^T \mathbf{B}_2) \right).$$

Generalizing AMDM for All CP Ranks

It can be theoretically shown that AMDM converges linearly for CP rank $R > s$

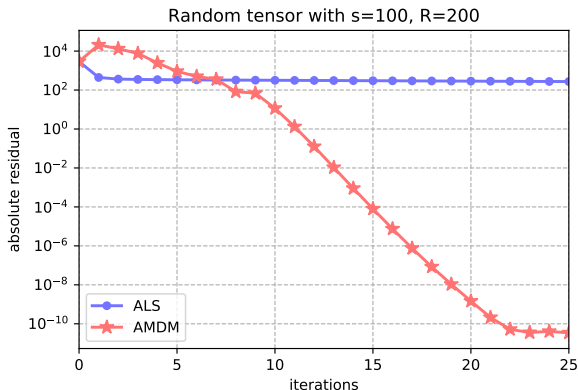
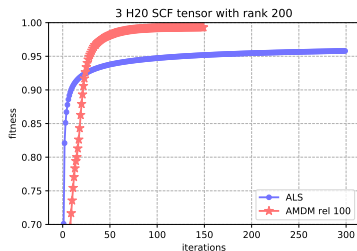
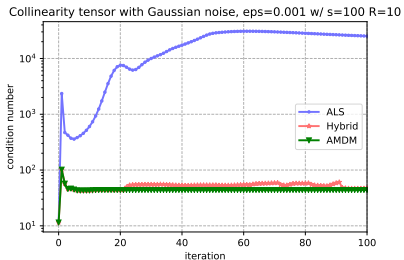
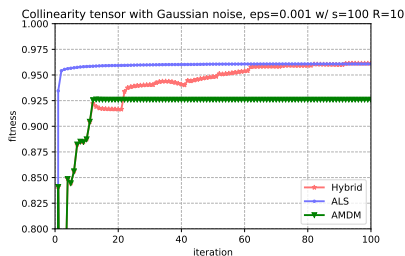


Figure: Linear convergence for exact CPD of a $100 \times 100 \times 100$ tensor with CP rank $R = 200$

Approximate Decomposition Results with AMDM

Using Hybrid algorithms leads to better conditioned and accurate decompositions.



Open Questions about AMDM

- Relation of AMDM with eigenvectors or singular vectors of a tensor
- Other views of the method (other than Mahalanobis Distance minimization)
- Existence of stationary points of AMDM for rank lesser than mode lengths case
- Quantifying conditioning of the alternating update in AMDM