

Accelerating Block Coordinate Descent for Nonnegative Tensor Factorization

Andersen Man Shun Ang[†] Jeremy E. Cohen[‡] Nicolas Gillis[†] Le Thi Khanh Hien^{†*}

[†] Department of Mathematics and Operational Research
Faculté Polytechnique, Université de Mons
Rue de Houdain 9, 7000 Mons, Belgium

[‡] Université de Rennes, INRIA, CNRS, IRISA, France

Abstract

This paper is concerned with improving the empirical convergence speed of block-coordinate descent algorithms for approximate nonnegative tensor factorization (NTF). We propose an extrapolation strategy in-between block updates, referred to as heuristic extrapolation with restarts (HER). HER significantly accelerates the empirical convergence speed of most existing block-coordinate algorithms for dense NTF, in particular for challenging computational scenarios, while requiring a negligible additional computational budget.

Keywords. nonnegative tensor factorization, nonconvex optimization, block-coordinate descent, Nesterov extrapolation

1 Introduction

In this paper, we consider the approximate nonnegative tensor canonical polyadic decomposition (CPD) problem, which we refer to as nonnegative tensor factorization (NTF). A N -way array or N -th order tensor \mathcal{T} is a multidimensional array in the product $\mathbb{R}^{I_1 \times \dots \times I_N}$ of the vector spaces \mathbb{R}^{I_i} for $i = 1, 2, \dots, N$. A vector $x \in \mathbb{R}^{I_1}$ is a first-order tensor, and a matrix $M \in \mathbb{R}^{I_1 \times I_2}$ is a second-order tensor. The goal of NTF is to approximate a tensor \mathcal{T} by a structured tensor \mathcal{X} . Using the squared Frobenius norm as a distance metric, defined as $\|\mathcal{X}\|_F^2 = \sum_{j_1, j_2, \dots, j_N} \mathcal{X}_{j_1 j_2 \dots j_N}^2$, NTF is the following optimization problem:

$$\min_{a_p^{(i)} \geq 0, 1 \leq i \leq N, 1 \leq p \leq r} \left\| \mathcal{T} - \sum_{p=1}^r \bigotimes_{i=1}^N a_p^{(i)} \right\|_F^2, \quad (1)$$

where \bigotimes is a tensor product over N real vector spaces $\mathbb{R}^{I_1}, \dots, \mathbb{R}^{I_N}$ defined as follows:

$$\left[\bigotimes_{i=1}^N a_p^{(i)} \right]_{j_1, j_2, \dots, j_N} := \prod_{i=1}^N a_p^{(i)}(j_i), \quad \text{where } a_p^{(i)} \in \mathbb{R}^{I_i} \text{ for } i = 1, 2, \dots, N \text{ and } p = 1, 2, \dots, r.$$

In other words, NTF is a low¹ nonnegative rank approximation problem, since by definition any nonnegative rank r tensor \mathcal{X} of order N can be parameterized as

$$\mathcal{X} = \sum_{p=1}^r \bigotimes_{i=1}^N a_p^{(i)} \text{ where } a_p^{(i)} \in \mathbb{R}_+^{I_i} \text{ for } i = 1, 2, \dots, N \text{ and } p = 1, 2, \dots, r.$$

*This work was supported by the Fonds de la Recherche Scientifique - FNRS and the Fonds Wetenschappelijk Onderzoek - Vlaanderen (FWO) under EOS Project no O005318F-RG47, and by the European Research Council (ERC starting grant no 679515). E-mails: {manshun.ang, nicolas.gillis, thikhanhien}@umons.ac.be, jeremy.cohen@irisa.fr

¹Low means much smaller than the generic rank of tensors in the considered tensor space [1, 2].

Intuitively, using NTF to approximate a tensor means using a part-based decomposition to summarize its content with a few “simple” rank-one tensors $a_p^{(1)} \otimes a_p^{(2)} \otimes \dots \otimes a_p^{(N)}$ where the components $a_p^{(i)}$ are entrywise nonnegative, with $1 \leq p \leq r$ and $1 \leq i \leq N$. This idea finds numerous applications in diverse areas, among which chemometrics or psychometrics are historical examples, see [3, 4, 5].

In this paper, we use the Frobenius norm to measure the error of approximation. It is arguably the most widely used measure, mostly because it has some nice properties (in particular, the subproblems in each block of variables is a convex quadratic problem; see below) and it corresponds to the maximum likelihood estimator in the presence of i.i.d. Gaussian noise. NTF is a non-convex optimization problem. Moreover, no closed-form solution is known to solve NTF; in fact, the problem is NP-hard already for the matrix case, that is, for $N = 2$; see [6]. Therefore, there has been a large amount of works dedicated to solving NTF using various optimization heuristics; see § 2 for a review of the state-of-the-art algorithms. However, note that unlike unconstrained approximate tensor factorization, NTF is well-posed in the sense that there always exists an optimal solution; see [7]. Moreover, a solution \mathcal{X}^* to NTF is almost always² unique for $N > 2$, and the solution to (1) also has exactly rank r ; see [8].

Outline and contribution This paper focuses on computing solutions to NTF as fast as possible. We derive new Block-Coordinate Descent (BCD) algorithms for NTF, that aim at being faster than existing BCD algorithms. To achieve this empirical speed-up in convergence speed, an extrapolation scheme “à la Nesterov” is used every time a block has been optimized, before switching to another block. The proposed Heuristic Extrapolation with Restarts (HER) algorithm consists of the following steps:

1. Initialize $A^{(i)} = [a_1^{(i)}, \dots, a_r^{(i)}]$ and pairing variables $\hat{A}^{(i)}$ for $1 \leq i \leq N$.
2. Loop over the blocks $A^{(i)}$ ($1 \leq i \leq N$):
 - (a) Update $A^{(i)}$ by minimizing (1) where the other blocks are fixed and take the value of the pairing variables $\hat{A}^{(j)}$ ($j \neq i$). For example, one can take a gradient step (see § 2.1 for more sophisticated strategies). Keep the previous value of $A^{(i)}$ in memory as $A_{old}^{(i)}$.
 - (b) Update the pairing variable using extrapolation:

$$\hat{A}^{(i)} = \max\left(0, A^{(i)} + \beta(A^{(i)} - A_{old}^{(i)})\right).$$

3. If the reconstruction error F has increased, reject the extrapolation and reset pairing variables $\hat{A}^{(i)} = A^{(i)}$ for $1 \leq i \leq N$; otherwise, update $A^{(i)} = \hat{A}^{(i)}$ for $1 \leq i \leq N$.
4. Update the parameter β ; see § 3.1.4 for the details.
5. If convergence criterion is not met, go back to 2.

This approach has been scarcely studied [9, 10, 11], while extrapolation is a rather well understood method to accelerate both convex and non-convex single-block descent algorithms; see for instance [12, 13]. The main novelty of this paper is to tackle a non-convex optimization problem using BCD with extrapolation between the block update, as opposed to inside each block update such as in [14] or after each outer loop as in [10]. This in-between extrapolation comes at almost no additional computational cost.

Extrapolated BCD algorithms are shown to be considerably faster than their standard counterparts in various difficult cases. These algorithms were observed to be slower than existing BCD algorithms only for extremely sparse tensors. Therefore, a contribution of this work is to experimentally show that using in-between block extrapolation allows to accelerate any BCD algorithm for dense NTF. This opens interesting questions for other optimization problems usually solved by BCD, for which such an extrapolation scheme may be applicable.

²The set of “bad” \mathcal{T} form an hypersurface.

Context Let us provide a brief historical note about tensor decomposition. The origin of tensor decomposition can be traced back to the work of Hitchcock [15, 16], whereas the idea of using multiway analysis is credited to the work of [17, 18]. Since then, especially after the work of Tucker in the field of psychometrics [19, 20], tensor decomposition has spread and become more and more popular in other fields such as chemometrics [21], signal processing [22, 23], data mining [24, 25], and many more. We refer the readers to [26, 27, 28, 29, 30] and references therein for comprehensive reviews of the applications of tensor decomposition. It is important to note that NTF is just one of many tensor decomposition models. Some other types of tensor decomposition or format include PARAFAC (that is, unconstrained approximate tensor factorization), Tucker/HOSVD [20, 31], and Tensor Train [32], to name a few. We focus on NTF in this paper.

Nonnegative matrix factorization (NMF), a key problem in machine learning and data analysis, is a special case of NTF when $N = 2$. First introduced in [33], it started becoming widely used after the seminal work of [34], and NMF has since then been deeply studied and well documented with variety of applications such as document classification, image processing, audio source separation and hyperspectral unmixing; see [35, 36, 37] and the references therein for more details. On the other hand, there are classes of data for which being represented by tensors is more natural. For example, a third-order tensor is preferably used to connect excitation-emission spectroscopy matrices in chemometrics [21], and RGB color images or 3D light field displays are generated as tensors [38]; see [35] for more examples. NTF was first introduced in [39] for fitting the latent class model in statistics. It has also been applied in model selection problem, sparse image coding in computer vision [40], sound source separation [41], image decomposition [42], text mining [43], among others; see [44, 35, 45, 46, 7] and reference therein for more applications of NTF.

Notation Below we recall some important notations in tensor algebra. First of all, the Kronecker product [47] of two matrices $A \in \mathbb{R}^{I_1 \times J_1}$ and $B \in \mathbb{R}^{I_2 \times J_2}$ is defined as follows:

$$A \boxtimes B = \left[\begin{array}{c|c|c} [A]_{11}B & \cdots & [A]_{1J_1}B \\ \hline \vdots & \ddots & \vdots \\ \hline [A]_{I_1 1}B & \cdots & [A]_{I_1 J_1}B \end{array} \right]. \quad (2)$$

Moreover, the Kronecker product of several matrices can be deduced from the above definition by associativity. The Khatri-Rao product $A \odot B$ is the columns-wise Kronecker product. Setting $A = [a_1, \dots, a_{J_1}]$ and $B = [b_1, \dots, b_{J_1}]$,

$$A \odot B = [a_1 \boxtimes b_1, \dots, a_{J_1} \boxtimes b_{J_1}]. \quad (3)$$

The Hadamard product (element-wise product) is denoted $A \otimes B$.

Compact decomposition notations: There exist several complementary notations to parameterize a low-rank tensor. In particular, grouping components $a_p^{(i)}$ as columns of factor matrices $A^{(i)} = [a_1^{(i)}, \dots, a_r^{(i)}]$, the following notations are equivalent:

$$\mathcal{X} = \sum_{p=1}^r \bigotimes_{i=1}^N a_p^{(i)} \quad (4)$$

$$= [A^{(1)}, \dots, A^{(N)}] \quad (5)$$

$$= \mathcal{I}_r \times_1 A^{(1)} \times_2 \dots \times_N A^{(N)} \quad (6)$$

$$\stackrel{\text{def}}{:=} \left(\bigotimes_{i=1}^N A^{(i)} \right)_a \mathcal{I}_r. \quad (7)$$

where \bigotimes_a is a tensor product of linear maps induced by the tensor product \otimes of vectors.

Equation (5) is the so-called Kruskal notation, equation (6) makes use of the n-mode product \times_p (see [26]), and equation (7) uses the fact that linear applications on tensor spaces of finite dimensions also form a tensor space with tensor product $(A \otimes_a B)(x \otimes y) := Ax \otimes By$. Because (7) exhibits this tensor product structure, we will make use of this compact formulation rather than the others.

Tensor unfoldings and useful formula: To derive partial derivatives of the NTF cost with respect to factors matrices, it is convenient to switch from a tensor formulation to a matrix description of the problem. More

precisely, the following relationships hold:

$$\mathcal{X} = \left(\bigotimes_{i=1}^N A^{(i)} \right) \mathcal{I}_r \quad \equiv \quad \forall i \in [1, N], X_{[i]} = A^{(i)} \begin{pmatrix} 1 \\ \circ \\ A^{(l)} \\ \text{\scriptsize } l \neq i \\ \text{\scriptsize } l = N \end{pmatrix}^T, \quad (8)$$

where unfoldings $X_{[i]}$ of a rank-one tensor \mathcal{X} are defined as follows:

$$X_{[i]} := a^{(i)} \otimes \begin{pmatrix} 1 \\ \boxtimes \\ \text{\scriptsize } l \neq i \\ \text{\scriptsize } l = N \end{pmatrix} a^{(l)} \in \mathbb{R}^{\mathcal{D}}, \quad \mathcal{D} = I_i \times \prod_{l \neq i} I_l. \quad (9)$$

Unfoldings of a general tensor are obtained by linearity of the unfolding maps. Note that several non-equivalent definitions are used in the tensor signal processing community; see [26] and [48].

2 The state-of-the-art algorithms for solving NTF

Below, we provide an overview of various techniques to solve NTF, which can be reformulated as follows

$$\min_{A^{(i)} \geq 0, 1 \leq i \leq N} F(A^{(1)}, \dots, A^{(N)}),$$

where

$$F(A^{(1)}, \dots, A^{(N)}) = \frac{1}{2} \left\| \mathcal{T} - \left(\bigotimes_{i=1}^N A^{(i)} \right) \mathcal{I}_r \right\|_F^2. \quad (10)$$

As a foreword, let us mention that there exist a wide range of algorithmic solutions for NTF (as for most of the tensor decomposition problems), that show different performances depending on the task at hand.

Algorithms for exact PARAFAC First of all, although the focus of this paper is approximate decompositions, it is interesting to point out that several algebraic techniques based on eigendecompositions have been proposed to deal with exact unconstrained tensor factorizations [49, 50, 51] whenever such a factorization exists. However, by design, these algorithms are typically not robust to noise, or may even be numerically unstable; see [52]. These techniques are sometimes used for initialization. Moreover, to the best of our knowledge, there does not exist an algorithm that computes NTF exactly for any rank using algebraic techniques. In fact, even in the matrix case, that is, for nonnegative matrix factorization (NMF), such techniques do not exist. Therefore, in this paper, we do not discuss exact NTF algorithms; this is a direction for further research.

Algorithms for approximate unconstrained tensor decomposition Because the Tensor Factorization (TF) model (that is, without nonnegativity constraints) has some interesting identifiability properties, it may occur that for well-conditioned tensors [53], approximate NTF can be computed with high precision by using a TF solver. Solving the TF problem is however harder in theory (the tensor low-rank unconstrained approximation problem is ill-posed, see [54]) and not really easier in practice than solving NTF. Actually, many algorithms that solve NTF are inspired from TF solvers and have similar complexity. Therefore we do not discuss TF solvers in what follows, and assume the reader is interested in solving NTF with specific algorithms that make use of the properties of the NTF problem.

All-at-once optimization A first class of widely used methods to solve NTF are all-at-once gradient-based methods. Indeed, it is quite straightforward to compute the gradients of F with respect to each matrix $A^{(i)}$. Let us denote

$$B^{(i)} = A^{(N)} \circ \dots \circ A^{(i+1)} \circ A^{(i-1)} \circ \dots \circ A^{(1)}. \quad (11)$$

Then the gradient of F with respect to $A^{(i)}$ is

$$\nabla_{A^{(i)}} F = \left(A^{(i)} (B^{(i)})^T - \mathcal{T}_{[i]} \right) B^{(i)}. \quad (12)$$

Algorithms	AO-AS[27]	AO-ADMM[58]	AO-Nesterov[62]	A-HALS[36]	APG[60]	iBPG[61]
Section	§ 2.1.1	§ 2.1.2	§ 2.1.3	§ 2.1.4	§ 2.2.1	§ 2.2.2

Table 1: Several block coordinate methods for solving NTF

Therefore, there is no obstacle to using any constrained gradient-based algorithm to (try to) find a stationary point of the non-convex NTF problem. To the best of our knowledge, the oldest all-at-once algorithm for NTF is a Gauss-Newton approach [55], but many approaches have been tested, including:

- Second-order optimization: using the fact that surrogates of the Hessian of F are heavily structured, second-order information can be used to solve NTF at a reasonable cost [56]. Limited-memory BFGS has also been employed when scalability is required [57]. To enforce the nonnegativity constraints, one can for instance square the variables, or use a variational approach (such a log-barrier).
- Primal-Dual optimization: the alternating direction method of multipliers has been tested for NTF, with however less promising results than its block-coordinate counterpart discussed below, see [58].
- Conjugate gradient: it has been reported that conjugate gradient can be used to solve NTF by squaring the variables; see [59].

Block Coordinate Descent (BCD) Methods Other than the above mentioned algorithms, BCD has become a standard and efficient scheme for solving NTF, mainly because (1) it essentially has cheap computation cost in each block update (BCD fixes all blocks except for one), (2) BCD can make use of recent developments in convex constrained optimization to efficiently solve NTF with respect to each block, and (3) under some suitable assumptions, many first-order BCDs and their accelerated versions have convergence guarantee in the context of general block-separable *non-convex composite optimization* problem that subsumes NTF as a special case, see for example [60] and [61] and the references therein. Below, we review several block coordinate methods for solving NTF; we list these algorithms in Table 1.

2.1 Alternating optimization (AO) framework

When solving NTF using BCDs, the blocks of variables that are alternatively updated must be chosen. It turns out that F is a quadratic function with respect to each matrix $A^{(i)}$ and therefore the optimization problem

$$\min_{A^{(i)} \geq 0} F(A^{(1)}, \dots, A^{(N)}) \quad (13)$$

is a linearly constrained quadratic programming problem referred to as Nonnegative Least Squares (NNLS). In particular, it is strictly convex if and only if $B^{(p)}$ is full column-rank. Therefore, it is quite natural to consider $A^{(i)}$ as the blocks in a BCD. The AO framework, which is a standard procedure to solve NTF, alternatively (exactly/inexactly) solves (13) for each block. We describe the AO framework in Algorithm 1. Note that the objective function of AO methods decreases after each block update. Depending on how the matrix-form NNLS problem (14) is solved, various implementations of AO algorithms can be obtained. Some of them are very efficient for solving NTF, they are surveyed below.

Algorithm 1: Alternating optimization framework

- 1: Input: a nonnegative N -way tensor
- 2: Output: nonnegative factors $A^{(1)}, A^{(2)}, \dots, A^{(N)}$.
- 3: Initialization: $(A_0^{(1)}, \dots, A_0^{(N)})$. Set $k = 1$.
- 4: **repeat**
- 5: **for** $i = 1, \dots, N$ **do**
- 6: Update $A_k^{(i)}$ as an exact/inexact solution of

$$\min_{A^{(i)} \geq 0} F\left(A_k^{(1)}, \dots, A_k^{(i-1)}, A^{(i)}, A_{k-1}^{(i+1)}, \dots, A_{k-1}^{(N)}\right). \quad (14)$$

- 7: $(A_{k-1}^{(i)})$ can be used as the initial point for the algorithm used to solve (14).
 - 7: **end for**
 - 8: Set $k = k + 1$.
 - 9: **until** some criteria is satisfied
-

2.1.1 AO-AS – solving NNLS with Active Set

When $A_k^{(i)}$ is updated by an exact solution of the NNLS problem (14), we obtain an alternating nonnegative least squares algorithm, usually referred to as ANLS in the literature. To solve exactly the NNLS subproblem (14), active set (AS) methods are usually rather effective and popular; see [27]. We will refer to AO-AS as the ANLS algorithm where the NNLS subproblems are solved with AS.

2.1.2 AO-ADMM – solving NNLS with ADMM

Designed to tackle a wide range of constrained tensor decomposition problems and various loss functions, AO-ADMM [58] applied to NTF boils down to using several steps of a primal-dual algorithm, the Alternating Direction Method of Multipliers (ADMM), to solve the cascaded nonnegative least squares problems. Therefore, AO-ADMM for NTF problem (1) is a variant of the AO framework that solves (14) inexactly.

2.1.3 AO-Nesterov

When Nesterov’s accelerated gradient method is applied to solve the NNLS problem in (14), we obtain AO-Nesterov; see [63, 62].

2.1.4 A-HALS

The hierarchical alternating least square (HALS) algorithm was introduced for solving the nonnegative matrix factorization (NMF) problem $\min_{W \geq 0, H \geq 0} \|M - WH^T\|_F^2$ (that is, NTF when $N = 2$), and has been widely used for solving NMF as it performs extremely well in practice; see for example [35, 36].

HALS cyclically updates each column of the factor matrix $A^{(i)}$ by solving an NNLS problem with respect to that column while fixing the others. The optimal solution of this NNLS subproblem can be written in closed form. A-HALS, which is short for accelerated HALS, was proposed in [64] to accelerate HALS. A-HALS repeats updating each factor matrix several times before updating the other ones. Hence A-HALS can be considered as a variant of the AO framework where each NNLS is inexactly solved itself by a BCD with closed-form updates. Let us briefly describe A-HALS for solving NTF. The NNLS problem (14) of A-HALS is inexactly solved by repeating cyclically updating the columns of $A_{k-1}^{(i)}$. In particular, let $M = X_{(i)}$, $W = A_{k-1}^{(i)}$, and $H = A_{k-1}^{(N)} \odot \dots \odot A_{k-1}^{(i+1)} \odot A_k^{(i-1)} \dots \odot A_k^{(1)}$. The j -th column of $A_{k-1}^{(i)}$ is updated by

$$W_{:,j} = \frac{\max\left(0, MH_{j,:}^T - \sum_{l \neq j} W_{:,l} H_{l,:} H_{j,:}^T\right)}{\|H_{j,:}\|^2}.$$

It is worth noting that A-HALS for NTF has subsequential convergence guarantee (that is, every limit point is a stationary point of the objective function), see [65, Section 7].

2.2 Block proximal gradient type methods

The NNLS problem (14) does not have a closed-form solution. From Equation (12), we see that F , when restricted to $A^{(i)}$, is a $L^{(i)}$ -smooth function, that is, the gradient $\nabla_{A^{(i)}} F$ is Lipschitz continuous with the constant $L^{(i)} = \left\| (B^{(i)})^T B^{(i)} \right\|$, where $B^{(i)}$ is defined in (11). This property can be employed to replace the objective function in the NNLS problem (14) by its quadratic majorization function, that leads to a new minimization problem which has a closed-form solution. This minimization-majorization approach, in the literature of block-separable composite optimization problem with the block-wise L -smooth property, is known as proximal gradient block coordinate descent method (see e.g., [61]). Considering the NTF problem, the closed-form solution of minimizing the majorization function is a projected gradient step. Applying Nesterov-type acceleration for the proximal gradient step improves the convergence of the BCD algorithm. Below we review the two recent accelerated proximal gradient BCD methods that were proposed for solving the general composite optimization problem.

2.2.1 APG – An Alternating Proximal Gradient method for solving NTF

APG was proposed by Xu and Yin [60]; see the Appendix of [66] and [60, Section 3.2] for the algorithm pseudocode. APG *cyclically* update each block (a.k.a each factor matrix) $A^{(i)}$ by calculating an extrapolation point $\hat{A}_{k-1}^{(i)} = A_{k-1}^{(i)} + w_{k-1}^{(i)}(A_{k-1}^{(i)} - A_{k-2}^{(i)})$ (here $w_{k-1}^{(i)}$ is some extrapolation parameter) and embedding this point in a projected gradient step

$$A_k^{(i)} = \max \left(0, \hat{A}_{k-1}^{(i)} - \frac{1}{L_{k-1}^{(i)}} \left(\hat{A}_{k-1}^{(i)} (B_{k-1}^{(i)})^T - \mathcal{T}_{[i]} \right) B_{k-1}^{(i)} \right).$$

After all blocks are updated, APG needs a restarting step, that is, if the objective function has increased then the projected gradient step would be re-done by using the previous values of all blocks instead of using the extrapolation points.

2.2.2 iBPG – An inertial Block Proximal Gradient Method

Recently proposed in [61], iBPG computes two different extrapolation points $\hat{A}_{k-1}^{(i,1)}$ and $\hat{A}_{k-1}^{(i,2)}$: one is for evaluating the gradient and the other one for adding inertial force. iBPG updates one matrix factor using a projected gradient step

$$A_k^{(i)} = \max \left(0, \hat{A}_{k-1}^{(i,2)} - \frac{1}{L_{k-1}^{(i)}} \left(\hat{A}_{k-1}^{(i,1)} (B_{k-1}^{(i)})^T - \mathcal{T}_{[i]} \right) B_{k-1}^{(i)} \right),$$

see the Appendix of [66] for the algorithm pseudocode. Furthermore, similarly to A-HALS, iBPG allows updating each matrix factor some times before updating another one – this feature would help save some computational costs since some common expressions can be re-used when repeating updating the same block. iBPG does not require a restarting step which make it suitable for solving large-scale NTF problems where evaluating the objective functions is costly.

3 Making BCD significantly faster with HER

With modern machine learning applications of NTF in mind, for which input tensor sizes can be extremely large and NTF should be provided as a low-level routine, there would be a definite economical and scientific gain to speeding up NTF algorithms. Radically different approaches exist in the literature to speed up existing algorithms for solving NTF, such as parallel computing [67, 68], compression and sketching [69, 70]. The combinations and relationships between these methods is poorly understood. In this paper, we focus on the acceleration of BCD using extrapolation.

Extrapolation inside BCD algorithms such as the workhorse ALS algorithm has been studied in the tensor decomposition literature as an empirical trick to speed up computations and avoid swamps, see Section 3.3. However, in a recent work on rank-one approximations of rank-two tensors, Gong, Mohlenkamp and Young [71] used gradient flow to study transverse stability, and provided a deeper analysis of the optimization landscape of tensor low-rank approximation which further supports the use of extrapolation.

As reviewed in §2.2, we have seen that APG and iBPG accelerate block proximal gradient methods by using extrapolation points in the projected gradient step to update each factor matrix. In another line of works, AO (Algorithm 1) was accelerated by using extrapolation between each block update (rather than inside the block update as in APG and iBPG); in other words, each factor matrix is updated by the extrapolation between previous updated factors. In the literature of tensor decomposition, the second type of extrapolation has been used to accelerate alternating least squares algorithms for solving CPD. Those works will be reviewed in § 3.3. In the following, we introduce HER - a novel extrapolation scheme that can be categorized into the class of accelerated AO algorithms using extrapolation between block update.

3.1 Heuristic Extrapolation with Restarts (HER)

HER was first proposed for solving NMF in [11], and found to be extremely effective on NTF in a preliminary work [72]. The sketch of HER was given in the introduction and its pseudo-code is given in Algorithm 2. In the following, we elaborate on HER with more details.

Algorithm 2: HER

- 1: Input: a nonnegative N -way tensor
 - 2: Output: nonnegative factors $A^{(1)}, A^{(2)}, \dots, A^{(N)}$.
 - 3: Initialization: Choose $\beta_0 \in (0, 1)$, $\eta \geq \bar{\gamma} \geq \gamma \geq 1$ and 2 sets of initial factor matrices $(A_0^{(1)}, \dots, A_0^{(N)})$ and $(\hat{A}_0^{(1)}, \dots, \hat{A}_0^{(N)})$. Set $\bar{\beta}_0 = 1$ and $k = 1$.
 - 4: **repeat**
 - 5: **for** $i = 1, \dots, N$ **do**
 - 6: **Update step** Let $A_k^{(i)}$ be an exact/inexact solution of
$$\min_{A^{(i)} \geq 0} F\left(\hat{A}_k^{(1)}, \dots, \hat{A}_k^{(i-1)}, A^{(i)}, \hat{A}_{k-1}^{(i+1)}, \dots, \hat{A}_{k-1}^{(N)}\right). \quad (15)$$
 - 7: **Extrapolation step**
$$\hat{A}_k^{(i)} = \max\left(0, A_k^{(i)} + \beta_{k-1}(A_k^{(i)} - A_{k-1}^{(i)})\right). \quad (16)$$
 - 8: **end for**
 - 9: Compute $\hat{F}_k := F\left(\hat{A}_k^{(1)}, \hat{A}_k^{(2)}, \dots, \hat{A}_k^{(N-1)}, A_k^{(N)}\right)$.
 - 10: **if** $\hat{F}_k > \hat{F}_{k-1}$ **then**
 - 11: Set $\hat{A}_k^{(i)} = A_k^{(i)}$, $i = 1, \dots, N$ % abandon the sequence $\hat{A}_k^{(i)}$
 - 12: Set $\bar{\beta}_k = \beta_{k-1}$, $\beta_k = \beta_{k-1}/\eta$. % Update $\bar{\beta}$, decrease β
 - 13: **else**
 - 14: Set $A_k^{(i)} = \hat{A}_k^{(i)}$, $i = 1, \dots, N$. % keep the sequence $\hat{A}_k^{(i)}$
 - 15: Set $\bar{\beta}_k = \min\{1, \bar{\beta}_{k-1}\bar{\gamma}\}$, $\beta_k = \min\{\bar{\beta}_{k-1}, \beta_{k-1}\gamma\}$. % Increase $\bar{\beta}$ and β
 - 16: **end if**
 - 17: Set $k = k + 1$.
 - 18: **until** some criteria is satisfied
-

3.1.1 Update step – line 6

It is clear that Algorithm 2 has the form of an alternating optimization framework in which the key optimization sub-problem (15) is a NNLS problem. As reviewed in § 2, some efficient algorithms for the NNLS problem (15) include AS, ADMM, Nesterov’s accelerated gradient, or A-HALS. The main difference between AO and HER is that HER does not use the latest values of the other blocks $A^{(j)}$ ($j \neq i$) but employs the latest values of their extrapolation $\hat{A}^{(j)}$ ($j \neq i$). For convenience, we refer to $\{\hat{A}_k^{(i)}, i = 1, \dots, N\}_{k \geq 0}$ as the extrapolation sequence.

3.1.2 Extrapolation step – line 7

After the update of $A_k^{(i)}$, the same block of the extrapolation sequence $\hat{A}_k^{(i)}$ is updated by extrapolating $A_k^{(i)}$ along the direction $A_k^{(i)} - A_{k-1}^{(i)}$, see (16). Note that $\hat{A}_k^{(i)}$ produced by (16) is always feasible. It is possible to remove the projection in (16), but we do not consider such approach in this work. Note that, regarding feasibility, $A_k^{(i)}$ produced by line 6 of Algorithm 2 is always feasible regardless of the feasibility of $\hat{A}_k^{(i)}$.

3.1.3 The restart mechanism – lines 9-16

After the update-extrapolate process on all the blocks, a restart procedure is carried out to decide whether or not we replace $A^{(i)}$ ($1 \leq i \leq N$) with the extrapolation sequence. The command in line 14 of Algorithm 2 has the same spirit with the update in (24) where the factor matrices are updated by *the extrapolation between block update*.

It may raise a question why $F(A_{k-1}^{(1)}, \dots, A_{k-1}^{(N)})$ does not appear in the restarting condition – line 10. The answer is due to the practicality of the algorithm. As stated in [11], using F as the restart criterion is computationally much more expensive than using the approximate \hat{F} . When computing \hat{F} , no explicit computation is required; instead, one may reuse already computed components from the updates of $A^{(N)}$ and $\hat{A}^{(N)}$. This creates an important reduction of computational complexity. For example, let us consider an order- N NTF problem with factor matrices $\{A^{(i)}\}_{i=1,2,\dots,N}$ with size $I_1 \times r$, $I_2 \times r$, \dots up to $I_N \times r$. Reusing already computed components (such as gradient) in the update of the last block $(A^{(N)}, \hat{A}^{(N)})$, we can compute

$\hat{F}(\hat{A}^1, \dots, \hat{A}^{N-1}, A^N)$ under $I_N r^{N-1}$ flops. However, if we compute $F(A^1, \dots, A^N)$, it takes $\prod_{i=1}^N I_i$ flops. If $r^N \ll \prod_{i=1}^N I_i$, then such reduction in complexity from $\prod_{i=1}^N I_i$ to $I_N r^{N-1}$ is significant even when N is low. Furthermore, we can even rotate the tensor such that I_N is the mode with the smallest size among all the modes. In fact, computing the cost function naively can be as costly as one block update, and thus using \hat{F} instead of F as the restart criterion is important, since restart using F requires computing the cost function at each iteration, while restart using \hat{F} is much cheaper.

Moreover, note that if the iterates sequence is converging, then the extrapolated sequence also converges to the same limit point. Therefore, since F is a continuous map, if convergence of the iterates is observed then the surrogate cost \hat{F} will asymptotically converge to the same final value as F . Although we did not characterize how fast this convergence happens, this justifies to use \hat{F} as a surrogate at least near a stationary point.

3.1.4 The extrapolation parameters in lines 9-16

The extrapolation weight β_k is computed within the restart mechanism of lines 9-16 of Algorithm 2, and it is updated using four parameters; see Table 2.

Symbol	Name	Setting	Range	Requires tuning?
β_k	Extrapolation weight	update as (17)	$[0, 1]$	Yes for β_0
γ	Growth rate of β	constant	$[\bar{\gamma}, \eta]$	Yes
η	Decay rate of β	constant	$[\gamma, \infty)$	Yes
$\bar{\gamma}$	Growth rate of $\bar{\beta}$	constant	$[1, \gamma]$	Yes
$\bar{\beta}_k$	Upper bound for β	update as (18)	$[\beta_k, 1]$	No, $\bar{\beta}_0 = 1$

Table 2: Parameters in the HER scheme

In the initialization stage, we set the upper bound for β as $\bar{\beta}_0 = 1$, pick $\beta_0 \in (0, 1)$, and select η, γ and $\bar{\gamma}$ such that $1 < \bar{\gamma} \leq \gamma \leq \eta$. The parameter $\bar{\beta}$, which is initialized as 1, is called the upper bound parameter for β . This parameter is used to limit the growth of β ; see below for more details. The parameter γ is called the (multiplicative) growth rate of β : when the error decreases, β is updated with $\gamma\beta$. Similarly $\bar{\gamma}$ is the (multiplicative) growth rate of $\bar{\beta}$. Finally, η is called the decay rate of β . This value is used to update β with β/η when the error increases. The parameters $(\gamma, \bar{\gamma}, \eta)$ are fixed constants, while β and $\bar{\beta}$ are updated depending on the restart condition.

The update of β HER updates β_k as

$$\beta_{k+1} = \begin{cases} \beta_k/\eta & \text{if } \hat{F}_{k+1} > \hat{F}_k \\ \min\{\gamma\beta_k, \bar{\beta}_k\} & \text{if } \hat{F}_{k+1} \leq \hat{F}_k \end{cases}, \quad (17)$$

which is explained as follows :

- If restart occurs, that is, if $\hat{F}_{k+1} > \hat{F}_k$, we assume it is caused by an over-sized β_k (recall that, for $\beta_k = 0$, decrease is guaranteed by the update in line 6) and we shrink the value of β for the next iteration using the decay parameter η as in (17).
- Otherwise, $\hat{F}_{k+1} \leq \hat{F}_k$, and we assume β_k can safely be increased. We grow β for the next iteration as $\gamma\beta$. To prevent β grow indefinitely, we use an upper bound $\bar{\beta}$ as in (17).

The update of $\bar{\beta}$ HER updates $\bar{\beta}_k$ as follows

$$\bar{\beta}_{k+1} = \begin{cases} \beta_k & \text{if } \hat{F}_{k+1} > \hat{F}_k \\ \min\{\bar{\gamma}\bar{\beta}_k, 1\} & \text{if } \hat{F}_{k+1} \leq \hat{F}_k \end{cases}. \quad (18)$$

The explanations are as follows :

- If there is no restart, that is, if $\hat{F}_{k+1} \leq \hat{F}_k$, $\bar{\beta}$ is increased if $\bar{\beta}$ is smaller than 1.
- Otherwise $\hat{F}_{k+1} > \hat{F}_k$ and $\bar{\beta}_{k+1}$ is set to β_k to prevent β_{k+1} growing larger than β_k too fast in the future. In fact, β_k indicates a too large value for β since the error has increased.

Let us make a few remarks:

- The relationships between the parameters in HER is as follows:

$$0 < \beta_k \leq \bar{\beta}_k \leq 1 < \bar{\gamma} \leq \gamma \leq \eta < \infty. \quad (19)$$

By construction, $\beta_k \leq \bar{\beta}_k \leq 1$, while $\bar{\gamma} \leq \gamma$ ensures that $\bar{\beta}$ increases slower than β , while $\gamma \leq \eta$ ensures that β is decreased faster.

- We have observed that HER is more effective if the NNLS subproblems (15) are solved with relatively high precision. Empirically Fig. 6 suggests to use HER with repeated projected gradient steps rather than just a single step. The suffix 50 after the algorithms' name in Fig. 6 means that we run 50 iterations for the algorithms to solve (15).
- A drawback of the HER approach is the parameter tuning. There are 4 parameters to tune: $\beta_0, \gamma, \bar{\gamma}, \eta$. However HER is not too sensitive for reasonable values of the parameters; see Fig. 1 for an illustration. Therefore, all the experiments in this paper are executed with no parameter tuning, even in difficult cases when data are ill-conditioned or rank is very high; namely we will use $\beta_0 = 0.5, \gamma = 1.05, \bar{\gamma} = 1.01$ and $\eta = 1.5$.
- In the implementation, we initialize $\hat{A}_0^{(i)} = A_0^{(i)}, i = 1, \dots, N$.

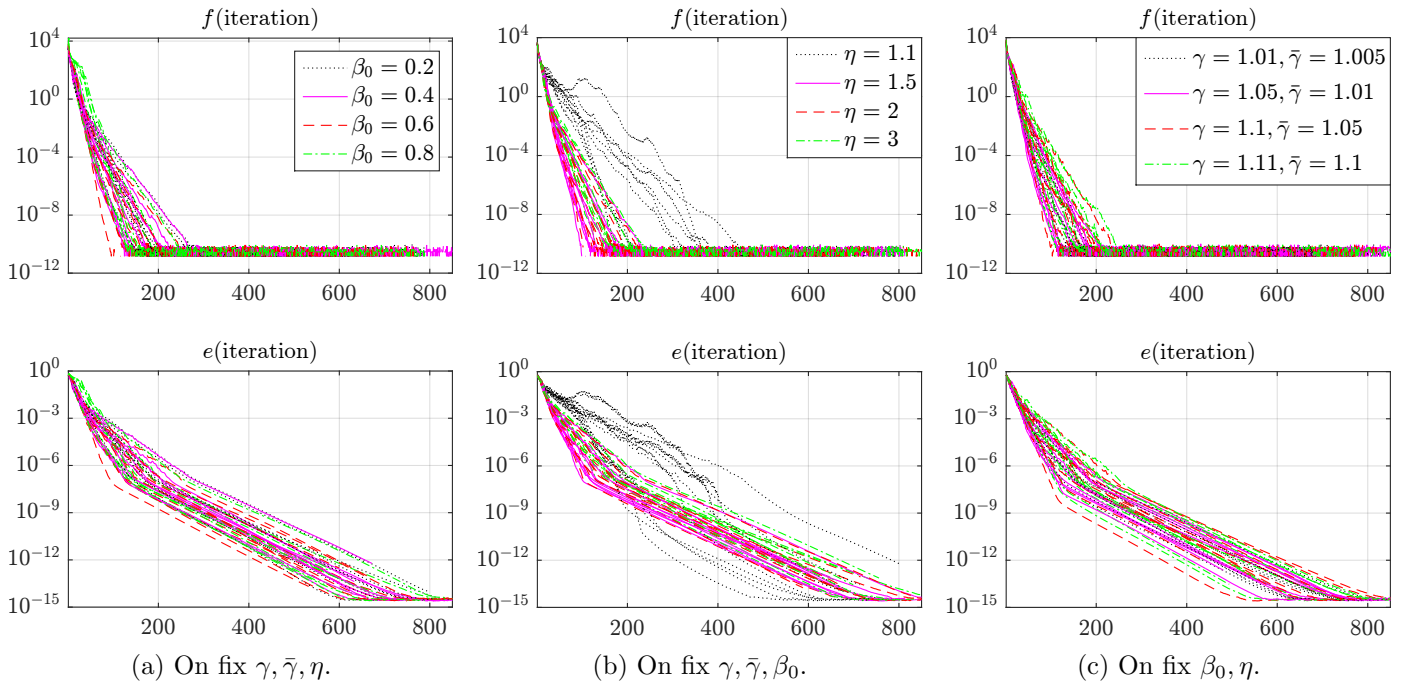


Figure 1: Comparison of HER with different parameters on the same NTF problems: a rank-10 factorization on noiseless tensors generated by random with size $50 \times 50 \times 50$. For each set of parameters, the decomposition is repeated 10 times over 10 different data tensors and initializations; see § 4.1 for more details. The top plots representing f display the error of the approximation, and the bottom plots representing e display the distance to the ground truth factors; see (26) and (27). The default set of hyper-parameters are $[\beta_0 = 0.5, \gamma = 1.05, \bar{\gamma} = 1.01, \eta = 1.5]$. The results here showed that HER is not very sensitive to its parameters as all the curves are not deviating away from each other, except for the case $\eta = 1.1$, suggesting that η should not be too small.

3.2 Discussion on convergence

Unfortunately, the HER framework to accelerate BCD methods for NTF cannot be guaranteed to converge using current existing proofs in the literature. The main distinction of HER that makes it very efficient but difficult to prove convergence is its *dynamic and flexibility* in the update of extrapolation parameters.

Specifically, HER performs a very aggressive extrapolation strategy: it extrapolates after each update of a block variable and uses that extrapolated point for the evaluation of the next block. Moreover, as long as the objective function decreases, it increases the extrapolation parameter β (as discussed in the previous section). The accelerated block proximal gradient methods using extrapolation inside the block update such as APG and iBPG do not have such strategy. Hence, although these methods have strong convergence guarantee, the way they choose extrapolation parameters is more conservative. Our extensive experiments strongly confirm the efficacy of dynamic and flexibility in choosing extrapolation parameters of HER when extrapolation between each block update is used and monotonicity of the objective function is taken into account. Studying convergence guarantees for the HER framework would be very challenging and is an important direction of further research.

3.3 Related works

Extrapolated AO algorithms can be traced back to a seminal work by Harshman [4], in which extrapolation was seen as a way to speed up the convergence of ALS. In this subsection we review two works on extrapolated AO algorithm: an older one [9] and a recent one [10]. To differentiate the extrapolation parameter used in HER (denoted as β_k), we denote the extrapolation parameter used in these works as ω_k .

General description Before we provide the details about the algorithms, let us point out a very important observation: unlike Algorithm 3.1, these algorithms use a global extrapolation parameter shared among all the blocks, and the extrapolation is conducted after all the blocks have been updated. That is, they first update all the blocks variables, then stack all the blocks together to form a vector x , and extrapolate it as

$$x_{k+1} = x_{k+\frac{1}{2}} + \omega_k(x_{k+\frac{1}{2}} - x_k), \quad (20)$$

where $x_{k+\frac{1}{2}}$ is the vector obtained by stacking all the blocks $A_k^{(1)}, \dots, A_k^{(N)}$ after all of them have been updated. The algorithms [9, 10] follow (20), and they differ in the way ω_k is computed, which we discuss below.

Extrapolated AO algorithms with Bro's sequence Bro revisited and optimized the seminal work by Harshman [4], and came up with an extrapolation scheme with convincing empirical speed-ups for computing CPD: the extrapolation parameter is tuned as

$$\omega_k = k^{\frac{1}{h(k)}} - 1,$$

where $h(k)$ is a recursive function so that $h(k+1) = h(k)$ if the error has not increased for more than four iterations, $h(k+1) = 1 + h(k)$ otherwise, and $h(1) = 3$. Moreover, no extrapolation is performed in the first few iterations because of stability issues. Furthermore, when the error increases, the extrapolation is not performed, that is, the extrapolated sequence is abandoned, as for HER.

Note that there is no particular modification of the Bro extrapolation scheme for aNCPD. In this paper we implement Bro-AHALS, Bro-ADMM and Bro-Nesterov – the three versions of Bro's accelerated methods in which we respectively use the same strategy using A-HALS (see § 2.1.4), ADMM and Nesterov's accelerated gradient method for solving the NNLS problem (14) inexactly.

Extrapolated AO algorithms with gradient ratio and line search Recently two heuristic approaches similar to Bro's have been proposed [10], and compute ω_k as follows

$$\text{Gradient ratio (GR):} \quad \omega_k = \nabla_x F_k(x) \Big|_{x=x_{k+\frac{1}{2}}} / \nabla_x F_{k-1}(x) \Big|_{x=x_{k-1}}, \quad (21)$$

$$\text{Linear search (LS):} \quad \omega_k = \underset{\omega}{\operatorname{argmin}} F\left(x_{k+\frac{1}{2}} + \omega(x_{k+\frac{1}{2}} - x_{k-1})\right), \quad (22)$$

where $\nabla_x F$ is the gradient of F with respect to x . As for Bro's accelerated algorithms, we implement in this paper GR-AHALS, GR-ADMM, GR-Nesterov, LS-AHALS, LS-ADMM and LS-Nesterov where we correspondingly use the same strategy as for A-HALS (see § 2.1.4), ADMM and Nesterov's accelerated gradient method for (14).

Note that GR and LS are designed for aCPD but not aNCPD, and similar to Bro's approach [9], there is no modification of the GR and LS schemes for the nonnegative decomposition case.

The modified implementation of Bro, GR and LS The algorithms Bro, GR and LS use a global extrapolation parameter shared among all the blocks, which is different from the extrapolation parameter β_k used in HER that is tuned independently for each block. Preliminary tests have showed that HER is always speeding up BCD algorithms much faster than Bro, GR and LS (see Fig. 7 in the next section). Such superiority can be explained in part by the fact that the block-wise tuning of β_k in HER gives HER much more degrees of freedom than Bro, GR and LS. Hence, to make a fair comparison between the different extrapolation strategies, we make the following modifications so that Bro, GR and LS have the same algorithmic structure as HER.

First, the update of x_k is performed block wise, that is, one $A^{(i)}$ at a time. Next, we extrapolate the blocks right after they have been updated, using the same extrapolation coefficients as described in Bro, GR and LS. It is important to note that, in the original algorithms, all the block variables are extrapolated with the same “global” extrapolation coefficient. That is, the extrapolation coefficients for every block $A^{(i)}$ in the original algorithm are the same. In the modification here, we “split” the global extrapolation coefficient into block-specific extrapolation coefficient. For example, in GR, the update-then-extrapolate step is performed for all i as

$$\text{Update: } A_{k+\frac{1}{2}}^{(i)} \text{ using (14),} \quad (23)$$

$$\text{Extrapolate: } A_{k+\frac{1}{2}}^{(i)} + \frac{\|\nabla_{A^{(i)}} F_k\|}{\|\nabla_{A^{(i)}} F_{k-1}\|} (A_{k+\frac{1}{2}}^{(i)} - A_k^{(i)}), \quad (24)$$

where $\nabla_{A^{(i)}} F_k$ is the gradient of F with respect to block $A^{(i)}$ at iteration k (see (12)), and $A_{k+\frac{1}{2}}^{(i)}$ is the block $A^{(i)}$ at iteration k just after the update, that is, we extrapolate the block right after it has been updated, as in HER. Moreover, (24) uses the ratio between the norm of the gradient of the current block A and the norm of the gradient of the same block in the last iteration. We use the same strategy on splitting the global extrapolation coefficient into block-specific extrapolation coefficient in Bro and LS. That is, the term $\frac{\|\nabla_{A^{(i)}} F_k\|}{\|\nabla_{A^{(i)}} F_{k-1}\|}$ in equation (24) is replaced by $(k^{\frac{1}{h(k)}} - 1)$ in Bro, and for LS, the extrapolation parameter ω_k is computed by solving a minimization subproblem: consider the update of the i th block at iteration k , we have

$$\omega_k = \underset{\omega}{\operatorname{argmin}} F\left(A_k^{(1)}, \dots, A_k^{(i-1)}, A_k^i + \omega(A_k^i - A_{k-1}^i), A_{k-1}^{(i+1)}, \dots, A_{k-1}^{(N)}\right). \quad (25)$$

By expanding F in terms of ω , (25) can be expressed as a second-order polynomial in ω , and hence a closed-form solution for ω exists.

Computational cost on Bro, GR and LS compared with HER The per-iteration cost in both GR and LS schemes is much larger than that of Bro. Both Bro, GR and LS have restart, but Bro’s extrapolation parameter is basically a scalar computation, while GR has multiple matrix-matrix multiplications and LS even has to solve a minimization sub-problem:

- Here we solve (25), which is a second-order polynomial in ω , exactly.
- For (22), for a 3rd-order NCPD problem, we need to minimize a sixth-order polynomial in ω .
- In general, for a N -order NCPD problem, we need to minimize a $2N$ -order polynomial in ω . Let $x = [A^{(1)}, \dots, A^{(N)}]$ denotes the stacking of the block into vector³, then

$$\begin{aligned} \beta_k &= \underset{\omega}{\operatorname{argmin}} F\left(x_{k+\frac{1}{2}} + \omega(x_{k+\frac{1}{2}} - x_{k-1})\right) \\ &= \underset{\omega}{\operatorname{argmin}} F\left([A_{k+\frac{1}{2}}^{(1)}, \dots, A_{k+\frac{1}{2}}^{(N)}] + \omega\left([A_{k+\frac{1}{2}}^{(1)}, \dots, A_{k+\frac{1}{2}}^{(N)}] - [A_{k-1}^{(1)}, \dots, A_{k-1}^{(N)}]\right)\right) \\ &\stackrel{(7)}{=} \underset{\omega}{\operatorname{argmin}} \frac{1}{2} \left\| \mathcal{T} - \left(\bigotimes_{a=1}^N \left(A_{k+\frac{1}{2}}^{(a)} + \omega(A_{k+\frac{1}{2}}^{(a)} - A_{k-1}^{(a)}) \right) \right) \mathcal{I}_r \right\|_F^2. \end{aligned}$$

We can see the cost of computing the coefficients for the polynomial can potentially be very high. Due to such reason, in the original paper, the ω_k in (22) is solved approximately using cubic line search in the Poblano toolbox.

³Here it is not the Kruskal notation (5).

Algorithms	Reference
HER-AS, HER-ADMM, HER-Nesterov, HER-AHALS	§ 3.1
AO-AS, AO-ADMM, AO-Nesterov, AHALS	§ 2.1
Bro-ADMM, Bro-Nesterov, Bro-AHALS	§ 3.3
GR-ADMM, GR-Nesterov, GR-AHALS	§ 3.3
LS-ADMM, LS-Nesterov, LS-AHALS	§ 3.3
APG, iBPG	§ 2.2

Table 3: Algorithms for solving NTF

- As pointed out in [11], exact line search has bad performance in NMF, so we believe this is the same for LS, for both the original form and the modified form.

In general, the per-iteration cost of the extrapolation step in Bro’s extrapolation is negligible, while for GR is higher than one ALS, and for LS it is much higher than one ALS. For this reason, we only run the modified LS in the numerical tests, i.e., (25), and we will see that, often LS performs the worst in the experiments, see for example Fig. 7.

Remarks on Bro, GR and LS compared with HER In the numerical experiments, we will compare the original form of Bro, GR, LS, as well as the modified version. There are several remarks on Bro, GR and LS.

- There are two sequences $A^{(i)}$ and $\hat{A}^{(i)}$ used in HER, while there is no auxiliary sequence in Bro, GR and LS: that is, the extrapolation in these approaches is conducted on the same block $A^{(i)}$.
- By splitting of the extrapolation coefficient into block extrapolation coefficients, the original GR and LS are improved as the ω_k in the new GR and LS are more adapted to each block variable.
- As Bro, GR and LS are designed for aCPD but not aNCPD, and similar to Bro’s approach [9], there is no modification of the GR and LS schemes for aNCPD. This means that there is no guarantee on feasibility of the iterates produced by these methods for aNCPD.

4 Experiments

In this section, we empirically prove the efficacy of HER by extensively test its performance on a rich set of synthetic data sets as well as real data sets. As presented in §3.1, HER is a scheme to accelerate AO algorithms by using extrapolation between block update; and as such, by using HER, we can derive several different algorithms corresponding to the solver we use for the NNLS problem (15). We stress out that HER can be used in combination with any BCD algorithm to make it faster. In this section, we combined it with the most well-known algorithms to tackle NTF, namely AS, ADMM, Nesterov accelerated gradient and AHALS for solving (15), which we denote by HER-AS, HER-ADMM, HER-Nesterov and HER-AHALS, respectively. We call HER-AO the set of these algorithms. Table 3 lists the algorithms that we implement and test in our experiments. Note that our goal is not to compare these various algorithms, but to show that HER can accelerate all of them significantly.

All experiments are run with MATLAB (v.2015a) on a laptop with 2.4GHz CPU and 16GB RAM. The code is available from <https://angms.science/research.html>.

Remark In this paper, we focus on dense NTF problems, for which the input tensor has mostly positive entries. However, the HER framework can also be applied to sparse tensors. This was in fact done for NMF in [11] with similar conclusions, that is, HER can accelerate algorithms significantly for sparse data sets as well. The problem with sparse data is not the algorithm itself, but rather its implementation. Handling sparse data also means dealing with extremely large dataset, which we cannot deal with our current implementation in Matlab. This could be fixed by integrating the proposed HER framework within a toolbox which features efficient sparse tensor manipulations and contractions, see for instance[73] and the references therein. We leave an efficient implementation of the HER framework for very sparse and large tensors for future works.

4.1 Set up

Performance measurement Two important factors in the evaluation of the performance of an algorithm are the data fitting error and the factor fitting error that are computed as follows. We use the value of the objective function

$$f_k := F \left(A_k^{(1)}, A_k^{(2)}, \dots, A_k^{(N-1)}, A_k^{(N)} \right) \quad (26)$$

to represent the data fitting error. Supposing the ground truth factor matrices $A_{\text{true}}^{(i)}$, $i = 1 \dots, N$ are available, then we compute the factor fitting error e_k as

$$e_k := \frac{1}{N} \sum_{i=1}^N \frac{\left\| \text{normalize}(A_{\text{true}}^{(i)}) - \text{normalize}(A_k^{(i)}) \Pi \right\|_F}{\left\| \text{normalize}(A_{\text{true}}^{(i)}) \right\|_F}. \quad (27)$$

Here $\text{normalize}(\cdot)$ is the column-wise normalization step (i.e., the i -th column of $\text{normalize}(A)$ is set to $\frac{A(:,j)}{\|A(:,j)\|_2}$), and Π is the permutation matrix computed through the Hungarian algorithm. The use of Π is to remove the permutation degree of freedom for matching the columns of $A^{(i)}$ to the column of $A_{\text{true}}^{(i)}$, and the use of normalization is to remove the scaling degree of freedom for matching the columns of $A^{(i)}$ to the column of $A_{\text{true}}^{(i)}$.

Generate a synthetic data To generate a synthetic tensor, we first generate ground truth factor matrices $A_{\text{true}}^{(i)} \in \mathbb{R}_+^{I_i \times r}$, $i = 1, \dots, N$ whose entries are sampled from i.i.d. uniform distributions in the interval $[0, 1]$. The tensor $\mathcal{T}^{\text{clean}} \in \mathbb{R}_+^{I_1 \times \dots \times I_N}$ is then constructed from $A_{\text{true}}^{(i)}$, $i = 1, \dots, N$. Finally, we form a synthetic data \mathcal{T} by adding some noise to $\mathcal{T}^{\text{clean}}$, $\mathcal{T} = \max(0, \mathcal{T}^{\text{clean}} + \sigma \mathcal{E})$, where $\sigma \geq 0$ is the noise level, and $\mathcal{E} \in \mathbb{R}^{I_1 \times \dots \times I_N}$ is a tensor whose entries are sampled from a unitary centered normal distribution.

Initialization, number of runs and plots For each run of an algorithm, we use a random initialization, i.e., the initial factor matrices $A_0^{(i)}$, $i = 1, \dots, N$, are generated by sampling uniform distributions in $[0, 1]$. Note that, testing a specific data tensor, we use the same initialization in one run of all algorithms. We run all the algorithms 20 times with 20 different initializations. We stop one run of an algorithm when the maximum time (which is chosen before running the algorithms) is reached.

In presenting the results, we plot $f - f_{\min}$; and if the ground truth is known, we also plot $e - e_{\min}$. Here f_{\min} and e_{\min} are respectively the minimal value of all the data fitting errors and the factor fitting errors obtained across all algorithms on all runs. In noiseless settings ($\sigma = 0$), exact factorization is possible, so we set $f_{\min} = 0$. In order to have a better observation of the performance of the algorithms, we plot the curves with respect to both time and iterations⁴. We remark that, ‘‘an iteration’’ for AO algorithms means the counter k of the outer loop after all blocks being updated. Regarding the time evaluation, we record the time stamp for each iteration, and then perform a linear interpolation to synchronize the time curves. Note that such linear interpolation does not reflect 100% truly the real convergence behaviour as it is just an linear estimate, but we consider such estimate to be accurate enough.

In our experiment, we emphasize on plotting the median curves of the 20 runs (which are the thick curves in the upcoming figures), because there may be large deviations between different runs.

Solving the NNLS problem (14) and (15) When using AS, ADMM, Nesterov’s accelerated gradient descent algorithm or AHALS to solve (14) or (15), in our implementation, we terminate the solver when the number of iterations reaches 50 or when $\|A_s^{(i)} - A_{s-1}^{(i)}\| \leq 10^{-2} \|A_0^{(i)} - A_1^{(i)}\|$, where s is the iteration counter of the solver.

Parameter set up for HER We use the following set of parameters for HER-AO (unless otherwise specified): $\beta_0 = 0.5$, $\gamma = 1.05$, $\bar{\gamma} = 1.01$, $\eta = 1.5$, and (16) is used for the extrapolation point.

⁴We do not report the number of MTTKRP (Matricized tensor times Khatri-Rao product) as all the algorithms in the experiments (except for AS) share the same number of MTTKRP (which is N for a tensor with order N), so the performance in terms of number of MTTKRP is contained implicitly in the plot with respect to the iterations.

List of experiments Table 4 lists the figures that report our diverse experiments on synthetic data and real data sets. All the experiments have $N = 3$ and the input tensor is dense.

Complete numerical experimental results We refer the interested reader to the appendix of the arXiv version of the current paper [66] for all the numerical experiments as reported in Table 4. The conclusions remain the same: HER significantly accelerate the convergence of BCD algorithms, while the HER-BCD outperforms both iBPG and APG.

Fig.	Test description	$[I_1, I_2, I_3, r, \sigma]$
Synthetic data		
2	Cube size, low rank, noiseless	$[50, 50, 50, 10, 0]$
2	Unbalanced size, low rank, noiseless	$[150, 10^3, 50, 12, 0]$
3	Unbalanced size, larger rank, noiseless	$[150, 10^3, 50, 25, 0]$
4	Large cube size, low rank, noisy	$[500, 500, 500, 10, 0.01]$
5	Unbalanced size, low rank, noisy, ill-condition	$[150, 10^3, 50, 12, 0.001]$
6	HER-AO-gradients compared with APG and iBPG	$[150, 10^3, 50, 10, 0.01]$
7	Comparing {HER, Bro, GR, LS}-AHALS	$[50, 50, 50, 10, 0]$
		$[150, 10^3, 50, 12, 0.01]$
Real data		
8	Two HSI images : PaviaU and Indian Pine	$[610, 340, 103, 10]$
		$[145, 145, 200, 15]$
9	Big data : black-and-white video sequence	$[153, 238, 1.4 \times 10^4, \{10, 20, 30\}]$

Table 4: List of experiments on NTF.

4.2 Experiments on synthetic data sets

As listed in Table 4, the experiments on synthetic data sets are designed to simulate different kinds of situations that may occur in real applications, which includes : low rank, larger rank, noiseless, noisy, tensor with balanced size (cubic tensor), tensor with unbalanced size (rectangular tensor), and ill-conditioned tensor.

Figure 2, 3 4 and 5 strongly affirm that HER-ADMM and HER-AHALS significantly outperform their counterparts AO-ADMM and AHALS in term of both f_k and e_k . We stress that the improvement is often of several orders of magnitude (at least 10^4 in most cases). We observe the same result for HER-AS and HER-Nesterov vs AO-AS and AO-Nesterov. The full experiments can be viewed in Appendix C.

Compared with APG and iBPG, we observe from Fig. 6 that HER-Nesterov outperforms both APG and iBPG in term of f and significantly outperforms them in term of e . From extensive experiments (see Appendix C for more results), we observe that HER, the scheme that makes use of the extrapolation between block update scheme, shows much better performance than APG and iBPG, the accelerated block proximal gradient methods that use Nesterov-type extrapolation inside each block update.

Compared with Bro-AHALS, GR-AHALS and LS-AHALS, Fig. 7 shows that our HER-AHALS performs the best in the three experimental settings (only median is plotted here). Note that since the acceleration frameworks Bro, GR and LS are not designed for NTF, it is possible the iterates produced by these frameworks are infeasible. Here we only compare HER-AHALS with Bro- AHALS, GR- AHALS and LS- AHALS; the comparison of these methods where AHALS is replaced with AO-ADMM and AO-ADMM are available in Appendix C, and similar conclusions are drawn, namely that HER outperforms the other accelerations.

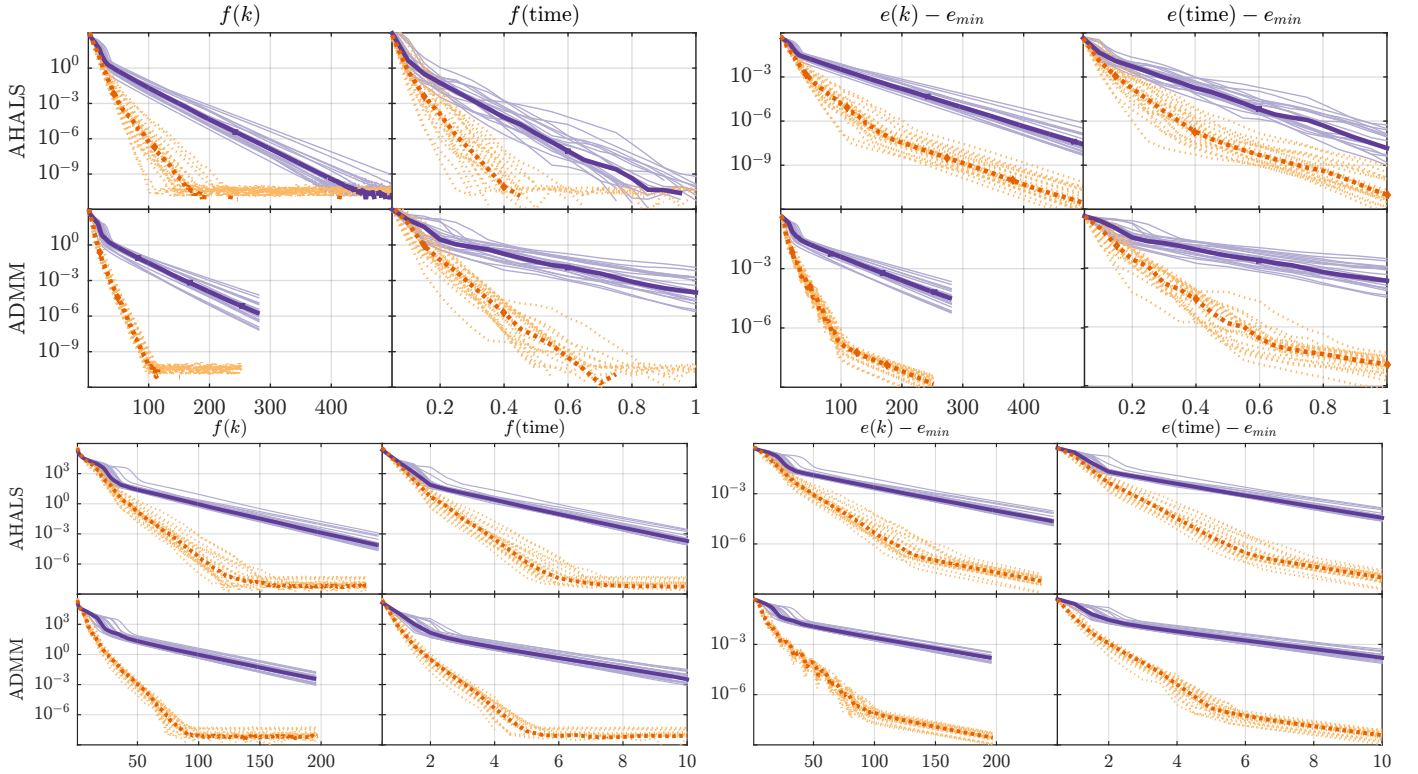


Figure 2: Convergence of algorithms : A-HALS and AO-ADMM without HER (solid purple) and with HER (dotted orange), on standard test case (top) : $[I_1, I_2, I_3, r] = [50, 50, 50, 10]$ and unbalanced sizes (bottom) $[I_1, I_2, I_3, r] = [150, 10^3, 50, 12]$. The results show that HER improves the convergence significantly, the convergence in both f and e for HER-accelerated methods are already multiple-order of magnitude better than the un-accelerated algorithms. Notice that due to a higher per-outer-iteration cost, ADMM-based algorithms (AO-ADMM and HER-AO-ADMM) run fewer number of outer-iteration than the AHALS-based algorithms. See Appendix C for the results on other algorithms where we observe a similar behaviour.

4.3 On real data

Two hyper-spectral images We test the performance of the algorithms on two hyperspectral images (HSI) PaviaU and Indian Pines⁵. They are nonnegative 3rd-order tensor; PaviaU has size $[610, 340, 103]$ with $r = 10$ and Indian Pines has size $[145, 145, 200]$ with $r = 15$. The r chosen are commonly used in practice.

We perform minimal pre-processing on the raw data : NaN or negative values (if any) are replaced by zero. Hence, it is possible the pre-processed data contains many zeros and being ill-conditioned. Figure 8 reports the performance of HER-AHALS, HER-ADMM and their counterparts AO-AHALS and AO-ADMM on the two data sets. As there are no ground truth factors, we only show f in the results.

We observe that there are multiple swamps, which are common for real datasets as the data are highly ill-conditioned (the condition numbers of the metricized pre-processed data tensor along all modes are $[593, 642, 1009]$ for Indian Pines and $[944, 462, 8083]$ for PaviaU). Nevertheless, considering the “best case” among the trials, HER-AHALS and HER-ADMM provide solutions with error $10^8 - 10^{10}$ times smaller than the best case of their un-accelerated counterparts. To compare with other algorithms, the readers can view the results in the Appendix of [66]. We observe that iBPG, APG and the AO (AO-AHALS and AO-ADMM) algorithms accelerated by GR, Bro and LS schemes are much slower than our AO (AO-AHALS and AO-ADMM) algorithms accelerated by HER. GR-AO and Bro-AO (for AO being AO-AHALS or AO-ADMM) even sometimes diverge.

On big data : video sequences We test HER-AHALS on the video data of the UCSD Anomaly Dataset [74]. Constructed by combining all the frame images of 70 surveillance video in the dataset, we have a tensor with sizes $153 \times 238 \times 14000$, where the first two modes are the screen resolution and the third mode is the number of frame. No pre-processing is performed on the raw data. Data of such size is too big to store in our computer memory, so we perform compression using Tucker decomposition, based on the built-in function from the Tensor toolbox [75]. We compare AHALS and HER-AHALS with $r \in \{10, 20, 30\}$. Results in Fig. 9

⁵Data available from http://www.ehu.es/ccwintco/index.php?title=Hyperspectral_Remote_Sensing_Scenes

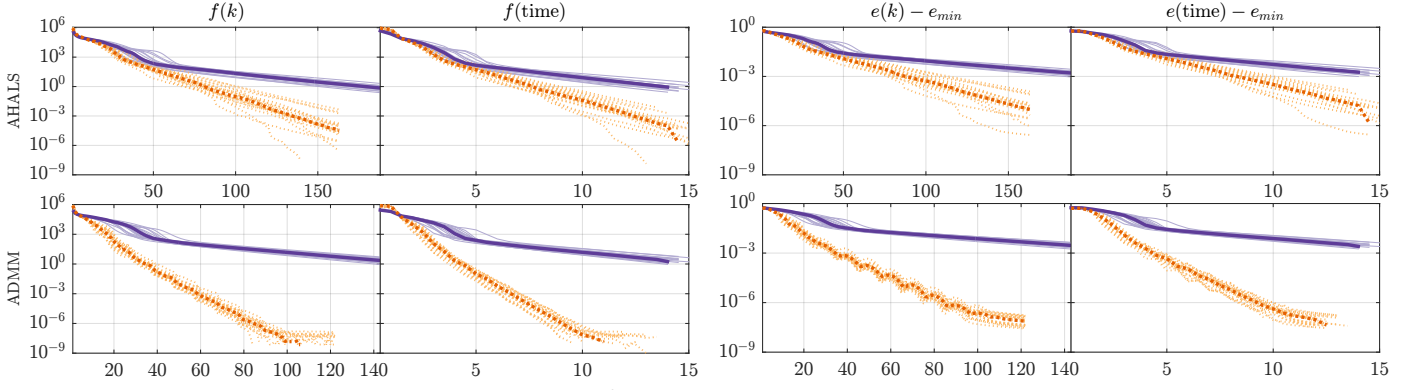


Figure 3: On large rank $[I_1, I_2, I_3, r] = [150, 10^3, 50, 25]$. For the plot set up, see Fig. 2. Results show HER improves the convergence speed significantly. See Fig. 2 for the plot set up, and Appendix C for the results on other algorithms.

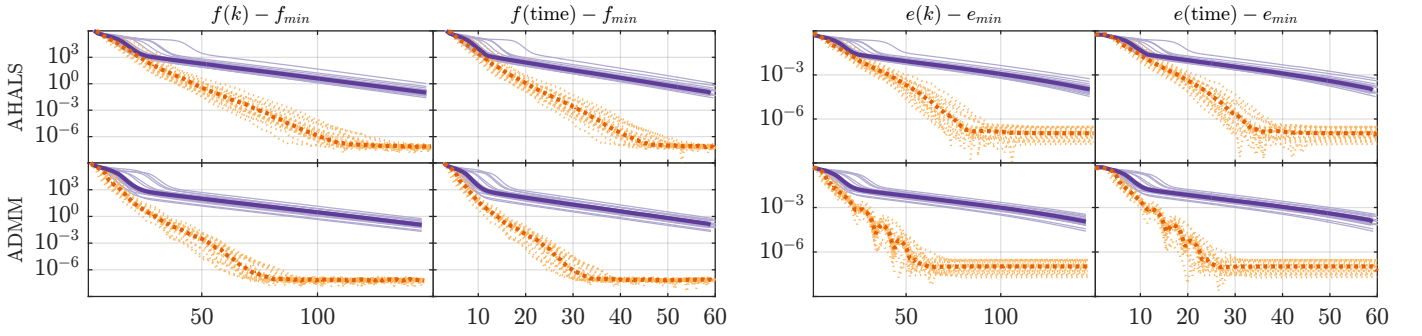


Figure 4: On big and noisy case $I_1 = I_2 = I_3 = 500, [r, \sigma] = [10, 0.01]$. Results show HER improves the convergence speed significantly. See Fig. 2 for the plot set up, and Appendix C for the results on other algorithms.

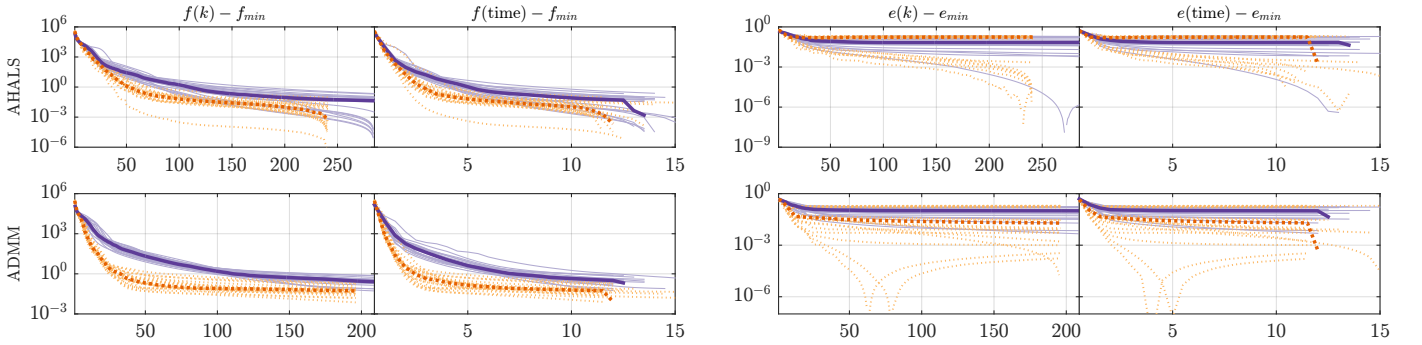


Figure 5: On ill-conditioned case $[I_1, I_2, I_3, r, \sigma] = [150, 10^3, 50, 12, 0.01]$, where $A_i(:, 1) = 0.99A_i(:, 2) + 0.01A_i(:, 1)$ for $i \in \{1, 2, 3\}$. For the plot set up, see Fig. 2. Results show HER improves the convergence speed. See Fig. 2 for the plot set up, and Appendix C for the results on other algorithms.

shows that HER-AHALS performs much better than AHALS. For the details on how HER works with Tucker compression, see the Appendix of [66].

As a conclusion for this section, we give some remarks on HER-AO. From our extensive experiments, we observe that HER-AS has inferior performance than others when the data is either big in size, high rank, or ill-conditioned. When the data has small size, all HER-AO algorithms have similar performance, and they all outperform their un-accelerated counterpart algorithms in term of both time and iteration. Among HER-AO algorithms, we highly recommend HER-AHALS for NTF as it shows good performance in all experiments.

5 Discussion and conclusion

In this paper, we have proposed an extrapolation strategy in-between block updates, referred to as heuristic extrapolation with restarts (HER), for improving the empirical convergence speed of block-coordinate descent algorithms for approximate nonnegative tensor factorization (NTF). HER significantly accelerates the empirical convergence speed of most existing block-coordinate algorithms for dense NTF, in particular for challenging computational scenarios, while requiring a negligible additional computational budget. The core of HER is to apply a special extrapolation-restart mechanism that aims to reduce the computational cost of restart while making sure the restart criterion follows the standard function restarts. The performance of HER was verified by the experiments reported in this paper. In all scenarios, HER-AHALS provides among the best results hence we recommend its use in practice.

Future works include deriving theoretical convergence for HER, and to apply it on other challenging applications.

References

- [1] Friedland S. On the generic and typical ranks of 3-tensors. *Linear Algebra and its Applications*. 2012;436(3):478–497.
- [2] Chiantini L, Ottaviani G. On generic identifiability of 3-tensors of small rank. *SIAM Journal on Matrix Analysis and Applications*. 2012;33(3):1018–1037.
- [3] Carroll JD, Chang JJ. Analysis of individual differences in multidimensional scaling via an n-way generalization of “Eckart-Young” decomposition. *Psychometrika*. 1970;35(3):283–319.
- [4] Harshman RA. Foundations of the PARAFAC procedure: Model and conditions for an “explanatory” multi-mode factor analysis. *UCLA Work Pap Phon*. 1970 11;16:1–84.
- [5] Kiers HAL. Towards a standardized notation and terminology in multiway analysis. *Journal of Chemometrics*. 2000;14(3):105–122.
- [6] Vavasis SA. On the complexity of nonnegative matrix factorization. *SIAM Journal on Optimization*. 2010;20(3):1364–1377.
- [7] Lim LH, Comon P. Nonnegative approximations of nonnegative tensors. *Journal of Chemometrics*. 2009;23(7-8):432–441. Available from: <https://onlinelibrary.wiley.com/doi/abs/10.1002/cem.1244>.
- [8] Qi Y, Comon P, Lim LH. Uniqueness of nonnegative tensor approximations. *IEEE Trans Inf Theory*. 2016 apr;62(4):2170–2183. ArXiv:1410.8129.
- [9] Bro R. *Multi-way Analysis in the Food Industry: Models, Algorithms, and Applications*. University of Amsterdam, The Netherlands; 1998.
- [10] Mitchell D, Ye N, De Sterck H. Nesterov acceleration of alternating least squares for canonical tensor decomposition: Momentum step size selection and restart mechanisms. *Numerical Linear Algebra with Applications*;27(4).
- [11] Ang AMS, Gillis N. Accelerating nonnegative matrix factorization algorithms using extrapolation. *Neural computation*. 2019;31(2):417–439.
- [12] Ghadimi S, Lan G. Accelerated gradient methods for nonconvex nonlinear and stochastic programming. *Mathematical Programming*. 2016;156(1-2):59–99.
- [13] Su W, Boyd S, Candès EJ. A differential equation for modeling Nesterov’s accelerated gradient method: theory and insights. *The Journal of Machine Learning Research*. 2016;17(1):5312–5354.
- [14] Liavas AP, Kostoulas G, Lourakis G, Huang K, Sidiropoulos ND. Nesterov-based alternating optimization for nonnegative tensor factorization: Algorithm and parallel implementation. *IEEE Transactions on Signal Processing*. 2017;66(4):944–953.
- [15] Hitchcock FL. The Expression of a Tensor or a Polyadic as a Sum of Products. *Journal of Mathematics and Physics*. 1927;6(1-4):164–189. Available from: <https://onlinelibrary.wiley.com/doi/abs/10.1002/sapm192761164>.

- [16] Hitchcock FL. Multiple Invariants and Generalized Rank of a P-Way Matrix or Tensor. *Journal of Mathematics and Physics*. 1928;7(1-4):39–79. Available from: <https://onlinelibrary.wiley.com/doi/abs/10.1002/sapm19287139>.
- [17] Cattell RB. “Parallel proportional profiles” and other principles for determining the choice of factors by rotation. *Psychometrika*. 1944;9(4):267–283.
- [18] Cattell RB. The three basic factor-analytic designs—Their interrelations and derivatives. *Psychological bulletin*. 1952 10;49:499–520.
- [19] Tucker LR. The extension of factor analysis to three-dimensional matrices. In: *Contributions to mathematical psychology*. Holt, Rinehart and Winston; 1964. p. 110–127.
- [20] Tucker LR. Some mathematical notes on three-mode factor analysis. *Psychometrika*. 1966;31(3):279–311.
- [21] Smilde AK, Bro R, Geladi P. *Multi Way Analysis — Applications in Chemical Sciences*. Wiley, West Sussex, England; 2004.
- [22] De Lathauwer L, De Moor B. From matrix to tensor : multilinear algebra and signal processing; 1996. p. 1–11.
- [23] Cichocki A, Mandic D, De Lathauwer L, Zhou G, Zhao Q, Caiafa C, et al. Tensor Decompositions for Signal Processing Applications: From two-way to multiway component analysis. *IEEE Signal Processing Magazine*. 2015;32(2):145–163.
- [24] Bader BW, Berry MW, Browne M. Discussion tracking in Enron email using PARAFAC. In: *Survey of Text Mining II*. Springer; 2008. p. 147–163.
- [25] Mørup M. Applications of tensor (multiway array) factorizations and decompositions in data mining. *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery*. 2011;1(1):24–40.
- [26] Kolda T, Bader B. *Tensor Decompositions and Applications*. *SIAM Review*. 2009;51(3):455–500.
- [27] Kim J, He Y, Park H. Algorithms for nonnegative matrix and tensor factorizations: a unified view based on block coordinate descent framework. *Journal of Global Optimization*. 2014;58(2):285–319.
- [28] Kroonenberg PM. *Applied Multiway Data Analysis*. vol. 702. Wiley, West Sussex, England; 2008.
- [29] Anandkumar A, Ge R, Hsu D, Kakade SM, Telgarsky M. Tensor decompositions for learning latent variable models. *The Journal of Machine Learning Research*. 2014;15(1):2773–2832.
- [30] Sidiropoulos ND, De Lathauwer L, Fu X, Huang K, Papalexakis EE, Faloutsos C. Tensor decomposition for signal processing and machine learning. *IEEE Transactions on Signal Processing*. 2017;65(13):3551–3582.
- [31] De Lathauwer L, De Moor B, Vandewalle J. A multilinear singular value decomposition. *SIAM journal on Matrix Analysis and Applications*. 2000;21(4):1253–1278.
- [32] Oseledets IV. Tensor-train decomposition. *SIAM Journal on Scientific Computing*. 2011;33(5):2295–2317.
- [33] Paatero P, Tapper U. Positive matrix factorization: A non-negative factor model with optimal utilization of error estimates of data values. *Environmetrics*. 1994;5(2):111–126. Available from: <https://onlinelibrary.wiley.com/doi/abs/10.1002/env.3170050203>.
- [34] Lee DD, Seung HS. Learning the parts of objects by non-negative matrix factorization. *Nature*. 1999;401(6755):788–791.
- [35] Cichocki A, Zdunek R, Phan AH, Amari S. *Nonnegative matrix and tensor factorizations: applications to exploratory multi-way data analysis and blind source separation*. John Wiley & Sons; 2009.
- [36] Gillis N. The why and how of nonnegative matrix factorization. *Regularization, Optimization, Kernels, and Support Vector Machines*. 2014;12(257):257–291.
- [37] Fu X, Huang K, Sidiropoulos ND, Ma W. Nonnegative matrix factorization for signal and data analytics: Identifiability, algorithms, and applications. *IEEE Signal Processing Magazine*. 2019;36(2):59–80.
- [38] Hirsch M, Lanman D, Wetzstein G, Raskar R. *Tensor Displays*. In: *ACM SIGGRAPH 2012 Emerging Technologies*; 2012. p. 24:1–24:1.
- [39] Carroll JD, De Soete G, Pruzansky S. *Multiway Data Analysis*. Amsterdam, The Netherlands, The Netherlands: North-Holland Publishing Co.; 1989. p. 463–472. Available from: <http://dl.acm.org/citation.cfm?id=120565.120614>.
- [40] Shashua A, Hazan T. Non-negative Tensor Factorization with Applications to Statistics and Computer Vision. In: *Proceedings of the 22nd International Conference on Machine Learning*. ICML; 2005. p. 792–799.
- [41] Fitzgerald D, Cranitch M, Coyle E. Non-negative tensor factorisation for sound source separation. *IET Conference Proceedings*. 2005 January:8–12(4). Available from: https://digital-library.theiet.org/content/conferences/10.1049/cp_20050279.
- [42] Welling M, Weber M. Positive Tensor Factorization. *Pattern Recogn Lett*. 2001;22(12):1255–1261.
- [43] Chi EC, Kolda TG. On tensors, sparsity, and nonnegative factorizations. *SIAM Journal on Matrix Analysis and Applications*. 2012;33(4):1272–1299.
- [44] Bro R, De Jong S. A fast non-negativity-constrained least squares algorithm. *Journal of Chemometrics*. 1997;11(5):393–401.
- [45] Friedlander MP, Hatz K. Computing non-negative tensor factorizations. *Optimization Methods and Software*. 2008;23(4):631–647.
- [46] Hazan T, Polak S, Shashua A. Sparse image coding using a 3D non-negative tensor factorization. In: *Tenth IEEE International Conference on Computer Vision (ICCV’05)*. vol. 1; 2005. p. 50–57.
- [47] Brewer J. Kronecker products and matrix calculus in system theory. *IEEE Transactions on circuits and systems*. 1978;25(9):772–781.

- [48] Cohen JE. About Notations in Multiway Array Processing. arXiv preprint. 2015. Available from: <https://arxiv.org/abs/1511.01306>.
- [49] De Lathauwer L, De Moor B, Vandewalle J. Computation of the canonical decomposition by means of a simultaneous generalized Schur decomposition. *SIAM journal on Matrix Analysis and Applications*. 2004;26(2):295–327.
- [50] Luciani X, Albera L. Semi-algebraic canonical decomposition of multi-way arrays and joint eigenvalue decomposition. In: 2011 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP). IEEE; 2011. p. 4104–4107.
- [51] Domanov I, De Lathauwer L. Canonical polyadic decomposition of third-order tensors: Reduction to generalized eigenvalue decomposition. *SIAM Journal on Matrix Analysis and Applications*. 2014;35(2):636–660.
- [52] Beltrán C, Breiding P, Vannieuwenhoven N. Pencil-based algorithms for tensor rank decomposition are not stable. *SIAM Journal on Matrix Analysis and Applications*. 2019;40(2):739–773.
- [53] Vannieuwenhoven N. Condition numbers for the tensor rank decomposition. *Linear Algebra and its Applications*. 2017;535:35–86.
- [54] Silva VD, Lim LH. Tensor rank and the ill-posedness of the best low-rank approximation problem. *SIAM J Matrix Analysis Appl*. 2008;30(3):1084–1127.
- [55] Paatero P. A weighted non-negative least squares algorithm for three-way “PARAFAC” factor analysis. *Chemometrics and Intelligent Laboratory Systems*. 1997;38(2):223–242.
- [56] Vervliet N. *Compressed Sensing Approaches to Large-scale Tensor Decompositions*. KU Leuven, Belgium; 2018.
- [57] Acar E, Dunlavy DM, Kolda TG. A scalable optimization approach for fitting canonical tensor decompositions. *Journal of Chemometrics*. 2011;25(2):67–86.
- [58] Huang K, Sidiropoulos ND, Liavas AP. A flexible and efficient algorithmic framework for constrained matrix and tensor factorization. *IEEE Transactions on Signal Processing*. 2016;64(19):5052–5065.
- [59] Royer JP, Thirion-Moreau N, Comon P. Nonnegative 3-way tensor factorization taking in to account possible missing data. In: 2012 Proceedings of the 20th European Signal Processing Conference (EUSIPCO). IEEE; 2012. p. 71–75.
- [60] Xu Y, Yin W. A block coordinate descent method for regularized multiconvex optimization with applications to nonnegative tensor factorization and completion. *SIAM Journal on Imaging Sciences*. 2013;6(3):1758–1789.
- [61] Hien LTK, Gillis N, Patrinos P. Inertial Block Proximal Method for Non-Convex Non-Smooth Optimization. In: Thirty-seventh International Conference on Machine Learning (ICML); 2020. p. 1–11.
- [62] Zhang Y, Zhou G, Zhao Q, Cichocki A, Wang X. Fast nonnegative tensor factorization based on accelerated proximal gradient and low-rank approximation. *Neurocomputing*. 2016;198:148–154.
- [63] Guan N, Tao D, Luo Z, Yuan B. NeNMF: An optimal gradient method for nonnegative matrix factorization. *IEEE Transactions on Signal Processing*. 2012;60(6):2882–2898.
- [64] Gillis N, Glineur F. Accelerated Multiplicative Updates and Hierarchical ALS Algorithms for Nonnegative Matrix Factorization. *Neural Computation*. 2012;24(4):1085–1105.
- [65] Razaviyayn M, Hong M, Luo Z. A Unified Convergence Analysis of Block Successive Minimization Methods for Nonsmooth Optimization. *SIAM Journal on Optimization*. 2013;23(2):1126–1153.
- [66] Ang AMS, Cohen JE, Gillis N, Khanh Hien LT. Accelerating Block Coordinate Descent for Nonnegative Tensor Factorization. arXiv e-prints. 2020 Jan:arXiv:2001.04321.
- [67] Ravindran N, Sidiropoulos ND, Smith S, Karypis G. Memory-efficient parallel computation of tensor and matrix products for big tensor decomposition. In: 2014 48th Asilomar Conference on Signals, Systems and Computers. IEEE; 2014. p. 581–585.
- [68] Ballard G, Hayashi K, Ramakrishnan K. Parallel nonnegative cp decomposition of dense tensors. In: 2018 IEEE 25th International Conference on High Performance Computing (HiPC); 2018. p. 22–31.
- [69] Vervliet N, Debals O, De Lathauwer L. Exploiting Efficient Representations in Large-Scale Tensor Decompositions. *SIAM Journal on Scientific Computing*. 2019;41(2):A789–A815.
- [70] Battaglino C, Ballard G, Kolda TG. A practical randomized CP tensor decomposition. *SIAM Journal on Matrix Analysis and Applications*. 2018;39(2):876–901.
- [71] Gong X, Mohlenkamp MJ, Young TR. The optimization landscape for fitting a rank-2 tensor with a rank-1 tensor. *SIAM Journal on Applied Dynamical Systems*. 2018;17(2):1432–1477.
- [72] Ang A, Cohen J, Gillis N. Accelerating Approximate Nonnegative Canonical Polyadic Decomposition using Extrapolation. In: GRETSI 2019 - XXVIIème Colloque francophone de traitement du signal et des images. Grets; 2019. p. 1–4.
- [73] Smith S, Karypis G. Tensor-matrix products with a compressed sparse tensor. In: Proceedings of the 5th Workshop on Irregular Applications: Architectures and Algorithms; 2015. p. 1–7.
- [74] Mahadevan V, Li W, Bhalodia V, Vasconcelos N. Anomaly detection in crowded scenes. In: 2010 IEEE Computer Society Conference on Computer Vision and Pattern Recognition. IEEE; 2010. p. 1975–1981.
- [75] Bader BW, Kolda TG. Efficient MATLAB computations with sparse and factored tensors. *SIAM Journal on Scientific Computing*. 2007;30(1):205–231.
- [76] Hackbusch W. *Tensor spaces and numerical tensor calculus*. vol. 42. Springer Science & Business Media; 2012.

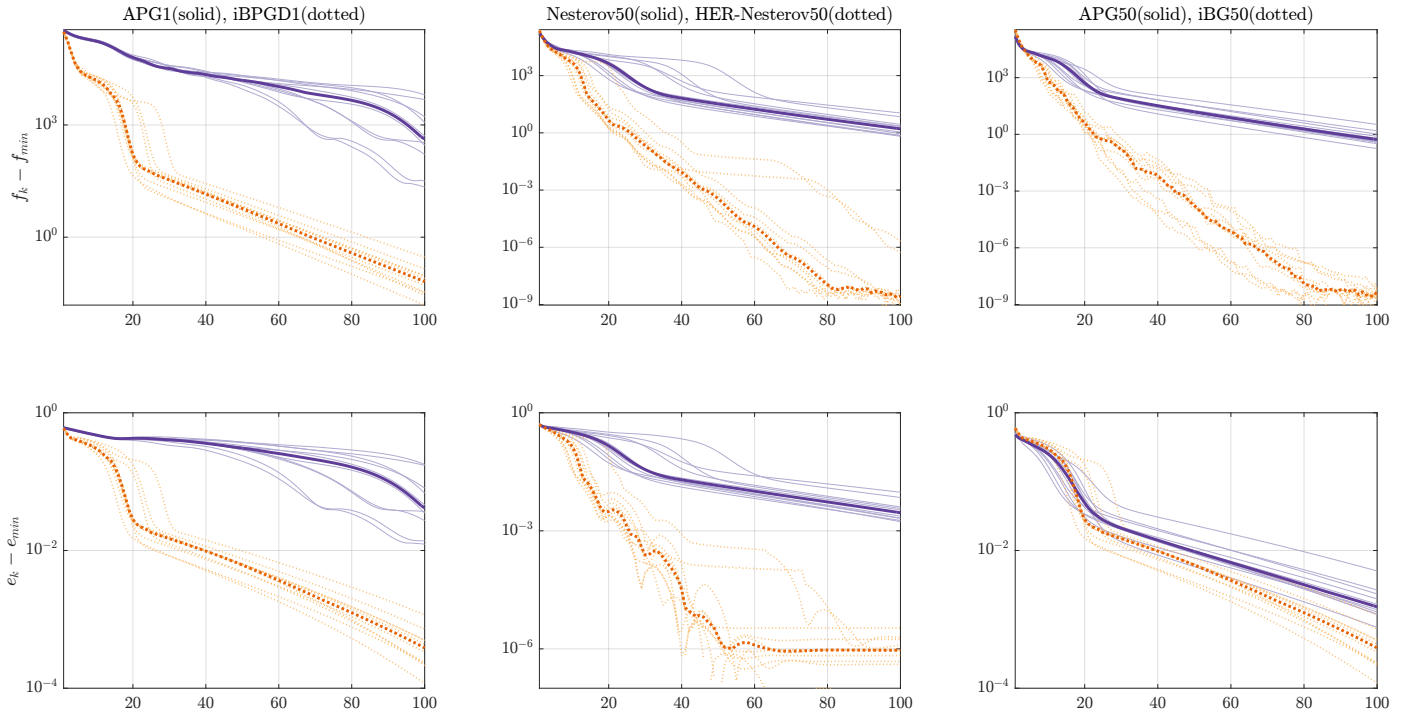


Figure 6: Comparing gradient algorithms on $[I_1, I_2, I_3, r, \sigma] = [150, 10^3, 50, 10, 0.01]$. Suffix number denotes the maximum number of inner iterations. Result shows HER works for both inexact and exact BCD using gradient. Here HER-Nesterov50 and HER-PGD50 are the best algorithms in both f and e . We do not plot the time plot here as they are similar to the iteration plot.

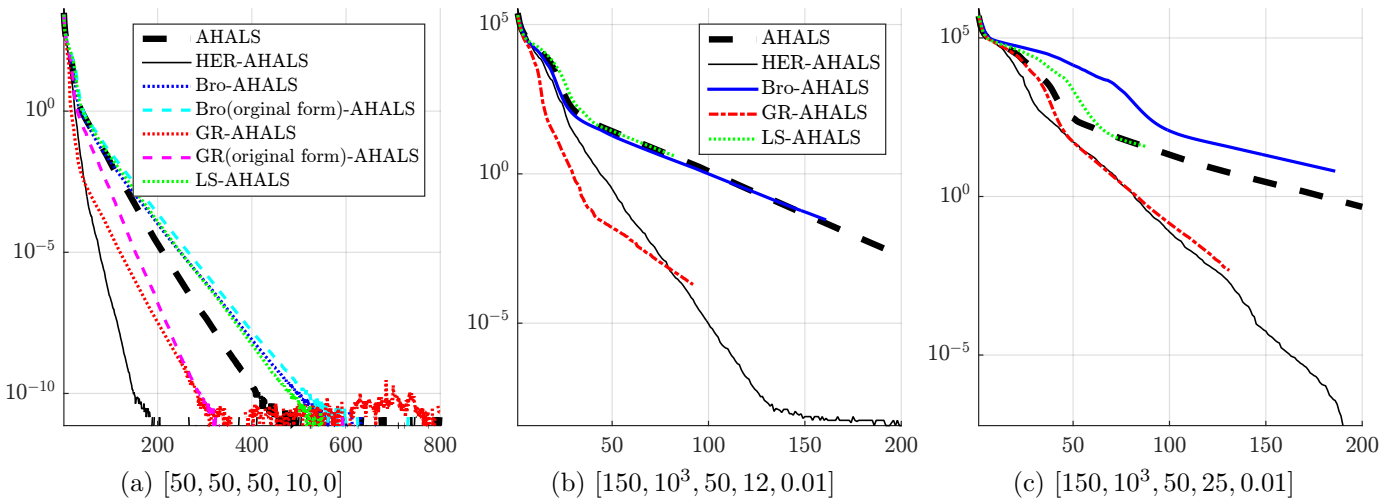


Figure 7: Comparing AHALS with different acceleration frameworks on synthetic datasets on 3 setting of $[I_1, I_2, I_3, r, \sigma]$. The curves are the median in $f(k) - f_{\min}$. The x-axis is the number of iteration, and all algorithm run with same run time limited. These results show that that (1) HER-AHALS performs better than all other algorithms, and (2) the modified Bro and GR algorithms perform better than their original counterpart; see subplot(a), while we do not show the result by Bro and GR in their original form in other subplots because they perform worse. Note that LS and GR run less number of iterations due to their larger per-iteration cost. Bro's approach has lower per-iteration cost, but it is even slower than vanilla AHALS. See the Appendix for more results.

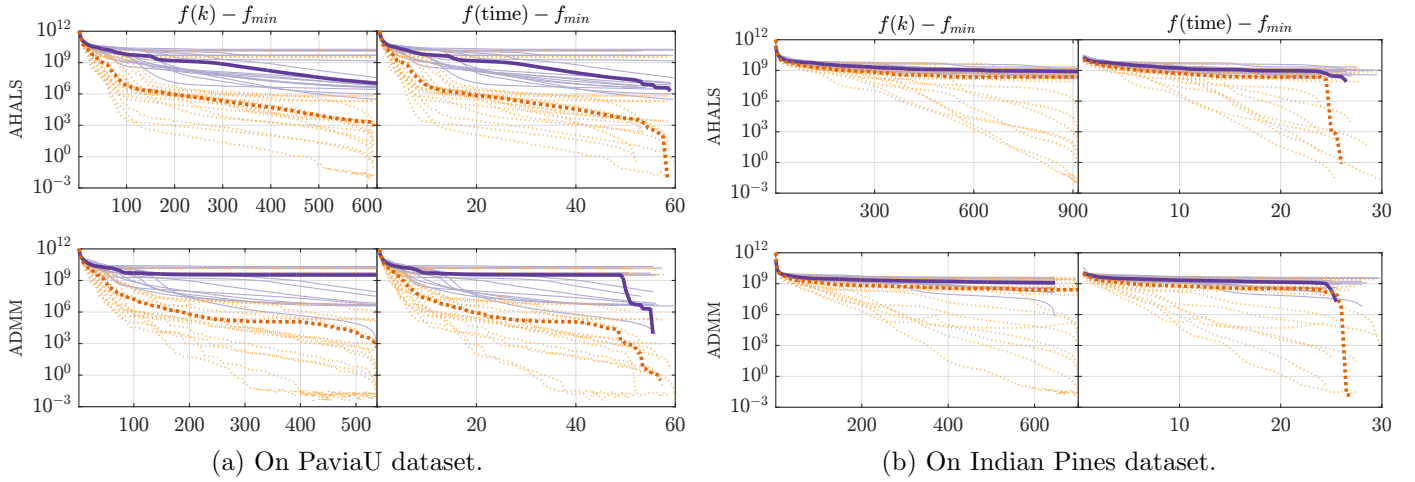


Figure 8: The results on HSI data. For the plot set up, see Fig. 2. Results show HER improve convergences. See the Appendix of [66] for more results.

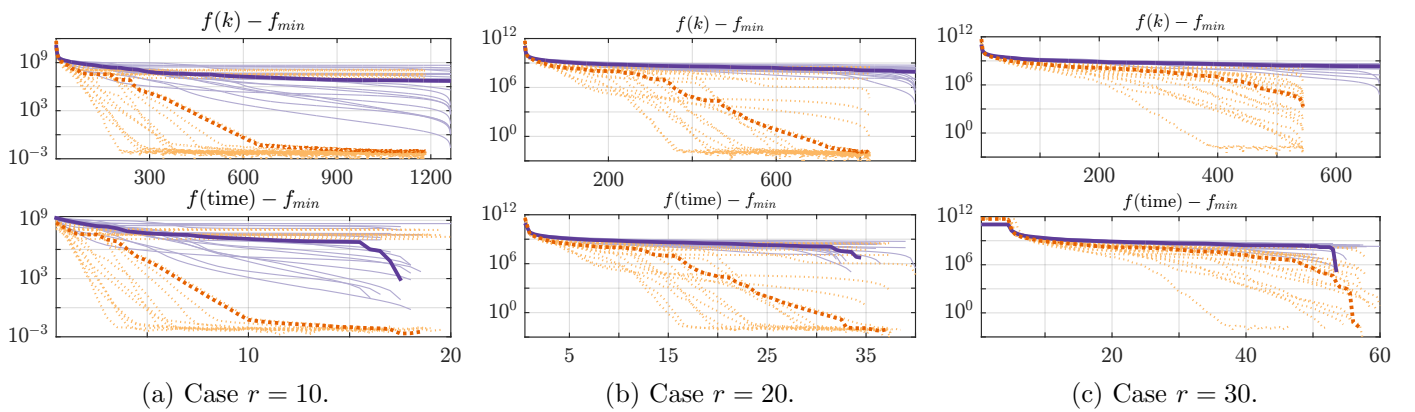


Figure 9: On video data [153, 238, 14000] for three values of r . Results show HER improve convergence and works well with Tucker-based compression.

Appendix

A Efficient tensor compression via Tucker format

Although there is a long history of using the Tucker model as a compression tool to pre-process big dataset, only recently has been formally discussed that compression does not actually imply transforming the large dataset into a smaller tensor [69]. Given the tensor \mathcal{T} , its Tucker format is expressed as :

$$\mathcal{T} = \left(\bigotimes_{p=1}^N U^{(p)} \right) \mathcal{G} \quad (28)$$

where $U^{(p)} \in \mathbb{R}^{n_p \times r_p}$, $\mathcal{G} \in \mathbb{R}^{r_1 \times \dots \times r_N}$ and $\{r_p\}_{p \leq N}$ are inputs integer parameters of the format, sometimes called Tucker ranks [76]. This representation is not unique but still offers a compressed expression of \mathcal{T} thus the name format rather than decomposition.

A typical situation is that of a tensor \mathcal{T} too big to fit in memory, since either too large and dense, or extremely large but sparse. Therefore, a third party may instead provide the data directly in a compact format such as the Tucker format. As Tucker format is in practice an approximation of the real data, the cost function of the aNCPD problem is modified as follows:

$$F_t(A^{(1)}, \dots, A^{(N)}) = \frac{1}{2} \left\| \left(\bigotimes_{p=1}^N U^{(p)} \right) \mathcal{G} - \left(\bigotimes_{p=1}^N A^{(p)} \right) \mathcal{I}_r \right\|_F^2. \quad (29)$$

On top of the storage gain, there is a huge computational burden ease in using structured representations of the data when computing the MTTKRP. Indeed, the gradient of F_t wrt say $A^{(1)}$ is obtained as follows:

$$\nabla_{A^{(1)}} F_t = -U^{(1)} G_{[1]} \left[\bigodot_{p=2}^N (U^{(p)})^T A^{(p)} \right] + A^{(1)} \left[\bigotimes_{p=2}^N (A^{(p)})^T A^{(p)} \right] \quad (30)$$

where \otimes is the Hadamard product. The equation involves only “cheap” products if Tucker ranks r_p are small compared to the data tensor dimensions n_p . In § 4.3, we check that indeed herBCD is compatible with accelerating the aNCPD using the Tucker format, and this actually opens the door to many problems that could not be tackled with simply herBCD, while enhancing at no cost the convergence speed of BCD algorithms for minimizing F_t . This contrasts with usual developments of fast techniques to solve aNCPD that typically do not consider other kind of acceleration in conjunction.

B Algorithm pseudocodes for APG and iBPG

In this section, we provide the details for the implementations of APG and iBPG.

Algorithm 3: APG

- 1: Input: nonnegative N -way tensor \mathcal{T}
- 2: Output: nonnegative factors $A^{(1)}, A^{(2)}, \dots, A^{(N)}$.
- 3: Initialization: Choose $\delta_w < 1$, $t_0 = 1$, and a set of initial factor matrices $(A_0^{(1)}, \dots, A_0^{(N)})$. Set $k = 1$.

4: **repeat**

5: **for** $i=1, \dots, N$ **do**

- 6: Compute $t_k = \frac{1}{2} (1 + \sqrt{1 + 4t_{k-1}^2})$, $\hat{w}_{k-1} = \frac{t_{k-1}-1}{t_k}$ and

$$w_{k-1}^{(i)} = \min \left(\hat{w}_{k-1}, \delta_w \sqrt{\frac{L_{k-2}^{(i)}}{L_{k-1}^{(i)}}} \right).$$

- 7: Compute an extrapolation point

$$\hat{A}_{k-1}^{(i)} = A_{k-1}^{(i)} + w_{k-1}^{(i)} (A_{k-1}^{(i)} - A_{k-2}^{(i)}).$$

- 8: Update $A_k^{(i)}$ by projected gradient step:

$$A_k^{(i)} = \max \left(0, \hat{A}_{k-1}^{(i)} - \frac{1}{L_{k-1}^{(i)}} \left(\hat{A}_{k-1}^{(i)} (B_{k-1}^{(i)})^T - \mathcal{T}_{[i]} \right) B_{k-1}^{(i)} \right) \quad (31)$$

9: **end for**

- 10: **if** $F(A_k) > F(A_{k-1})$ **then**

- 11: Update $A_k^{(i)}$ by the projected gradient step (31) with $\hat{A}_{k-1}^{(i)} = A_{k-1}^{(i)}$.

12: **end if**

- 13: Set $k = k + 1$.

14: **until** some criteria is satisfied

Algorithm 4: iBPG

- 1: Input: a nonnegative N -way tensor \mathcal{T}
 - 2: Output: nonnegative factors $A^{(1)}, A^{(2)}, \dots, A^{(N)}$.
 - 3: Initialization: Choose $\delta_w = 0.99$, $\beta = 1.01$, $t_0 = 1$, and 2 sets of initial factor matrices $(A_{-1}^{(1)}, \dots, A_{-1}^{(N)})$ and $(A_0^{(1)}, \dots, A_0^{(N)})$.
Set $k = 1$.
 - 4: Set $A_{\text{prev}}^{(i)} = A_{-1}^{(i)}$, $i = 1, \dots, N$. % $A_{\text{prev}}^{(i)}$ is to save the previous value of block i .
 - 5: Set $A_{\text{cur}}^{(i)} = A_0^{(i)}$, $i = 1, \dots, N$. % $A_{\text{cur}}^{(i)}$ is to save the current value of block i .
 - 6: **repeat**
 - 7: **for** $i=1, \dots, N$ **do**
 - 8: Compute $t_k = \frac{1}{2}(1 + \sqrt{1 + 4t_{k-1}^2})$, $\hat{w}_{k-1} = \frac{t_{k-1}-1}{t_k}$ and
$$w_{k-1}^{(i)} = \min \left(\hat{w}_{k-1}, \delta_w \sqrt{\frac{L_{k-2}^{(i)}}{L_{k-1}^{(i)}}} \right).$$
 - 9: **repeat**
 - 10: Compute two extrapolation points
$$\hat{A}^{(i,1)} = A_{\text{cur}}^{(i)} + w_{k-1}^{(i)} (A_{\text{cur}}^{(i)} - A_{\text{prev}}^{(i)}),$$
and
$$\hat{A}^{(i,2)} = A_{\text{cur}}^{(i)} + \beta w_{k-1}^{(i)} (A_{\text{cur}}^{(i)} - A_{\text{prev}}^{(i)})$$
 - 11: Set $A_{\text{prev}}^{(i)} = A_{\text{cur}}^{(i)}$.
 - 12: Update $A_{\text{cur}}^{(i)}$ by projected gradient step:
$$A_{\text{cur}}^{(i)} = \max \left(0, \hat{A}^{(i,2)} - \frac{1}{L_{k-1}^{(i)}} \left(\hat{A}^{(i,1)} (B_{k-1}^{(i)})^T - \mathcal{T}_{[i]} \right) B_{k-1}^{(i)} \right).$$
 - 13: **until** some criteria is satisfied
 - 14: Set $A_k^{(i)} = A_{\text{cur}}^{(i)}$.
 - 15: **end for**
 - 16: Set $k = k + 1$.
 - 17: **until** some criteria is satisfied
-

C Full experimental results

In this section, we provide more plots for the experimental set up presented in §4. This includes other algorithms, and other dimensions of the input tensors.

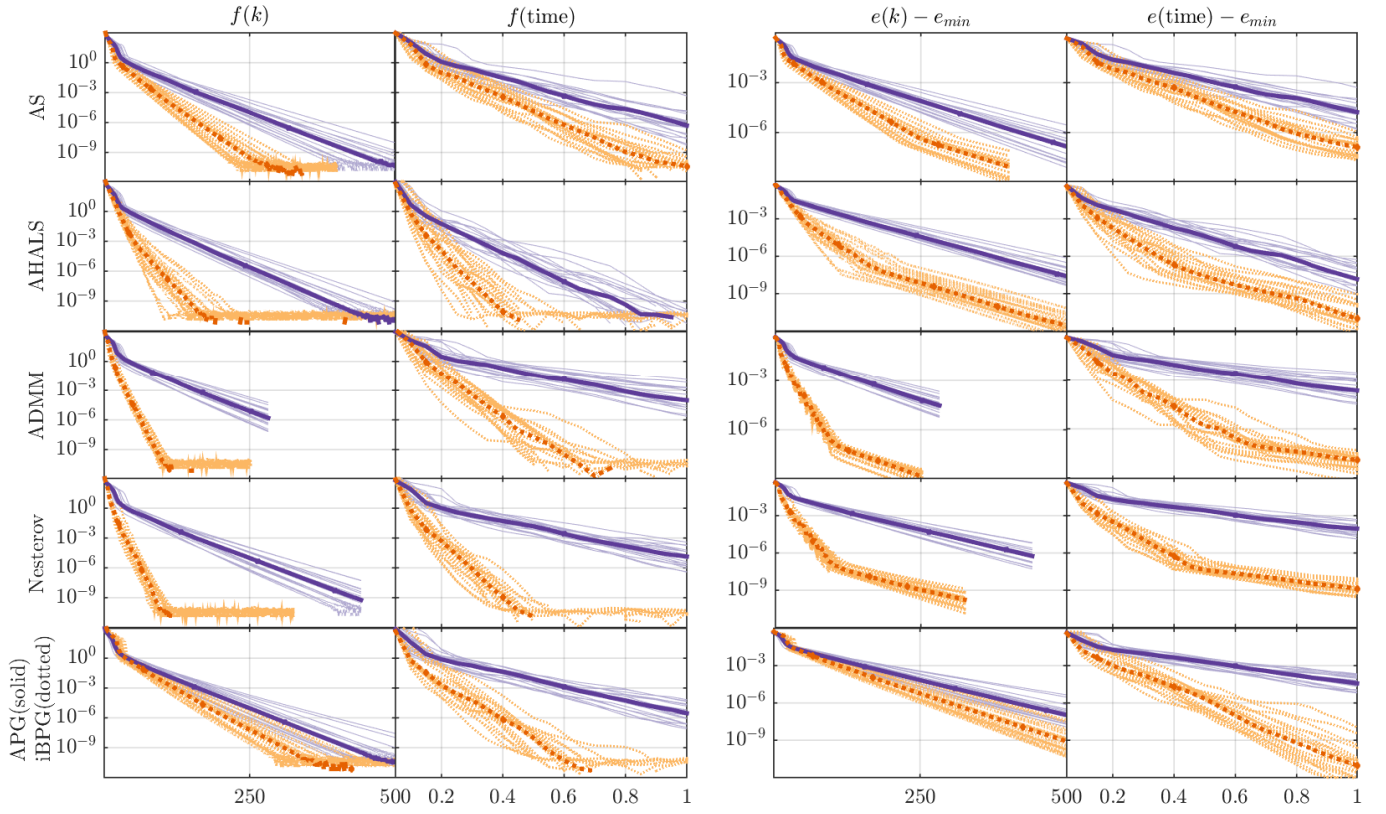


Figure 10: Experiments on synthetic data with $[I_1, I_2, I_3, r] = [50, 50, 50, 10]$ on algorithms without HER (solid purple curves) and with HER (dotted orange curves). Thick curves are median of the 20 thin curve with the same color. The results show that HER improves BCD algorithms, and that HER-BCD has a better performance than APG and iBPG.

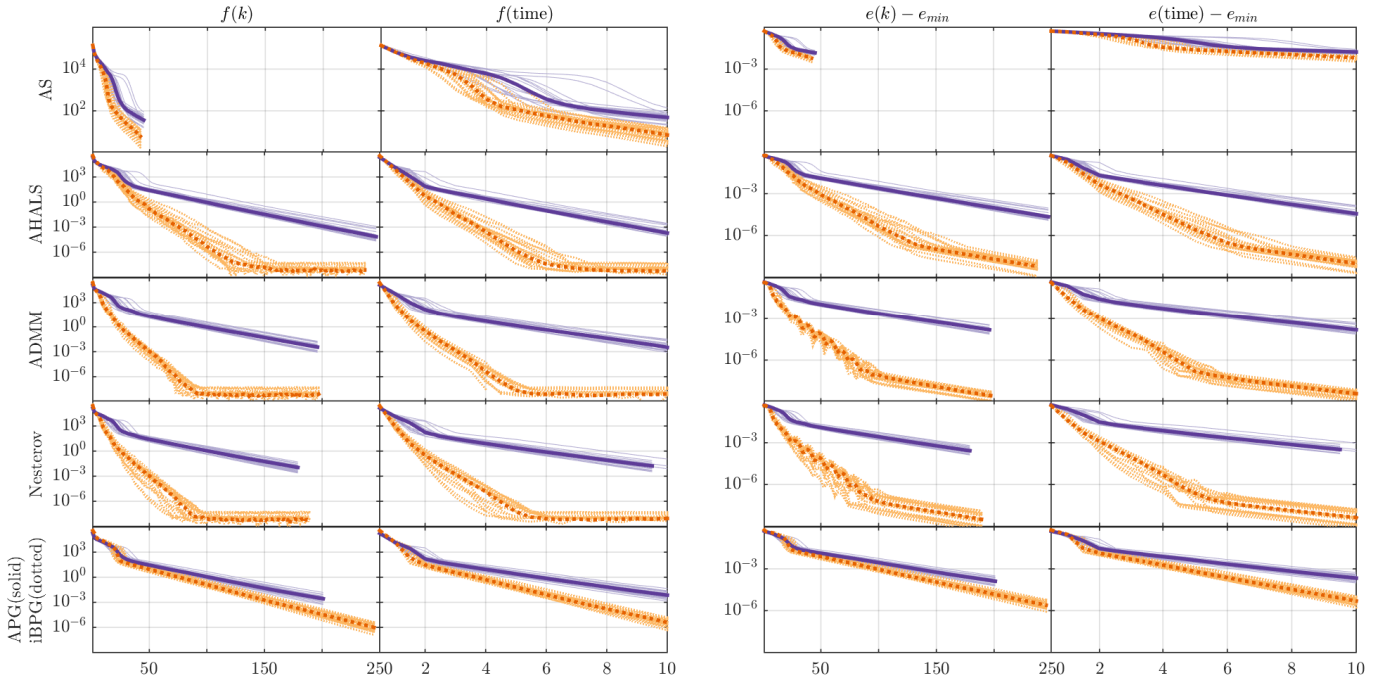


Figure 11: Experiments on synthetic data with unbalanced tensors with $[I_1, I_2, I_3, r] = [150, 10^3, 10, 10]$. The same setting as in Figure 10 is used. The results show that HER improves BCD algorithms, and that HER-BCD has a better performance than APG and iBPG. Note that here AS only ran approximately 40 iterations in 10 seconds due to high per-iteration cost.

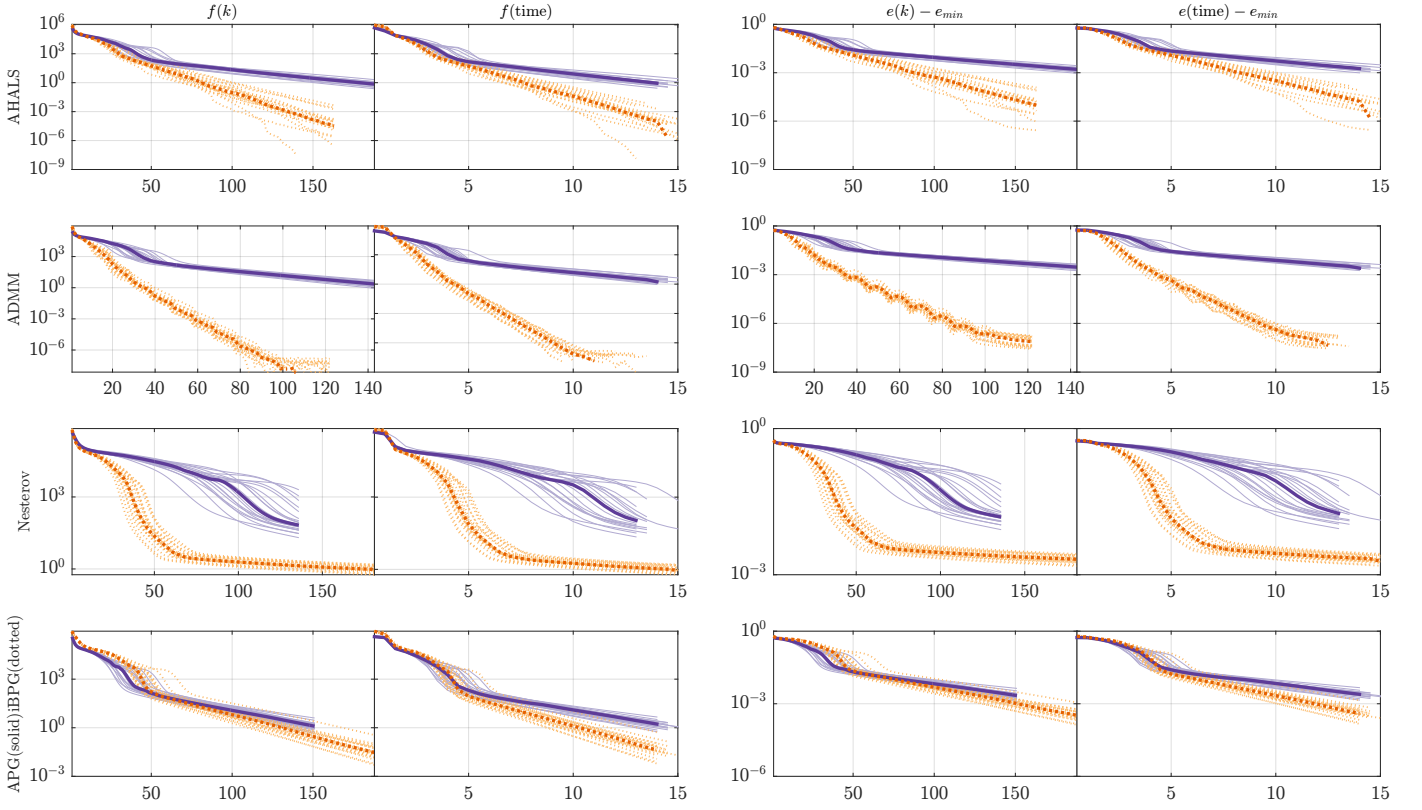


Figure 12: Experiments on synthetic data with $[I_1, I_2, I_3, r] = [150, 10^3, 50, 25]$. The same setting as in Figure 10 is used. The results show that HER improves BCD algorithms, and that HER-BCD has a better performance than APG and iBPG.

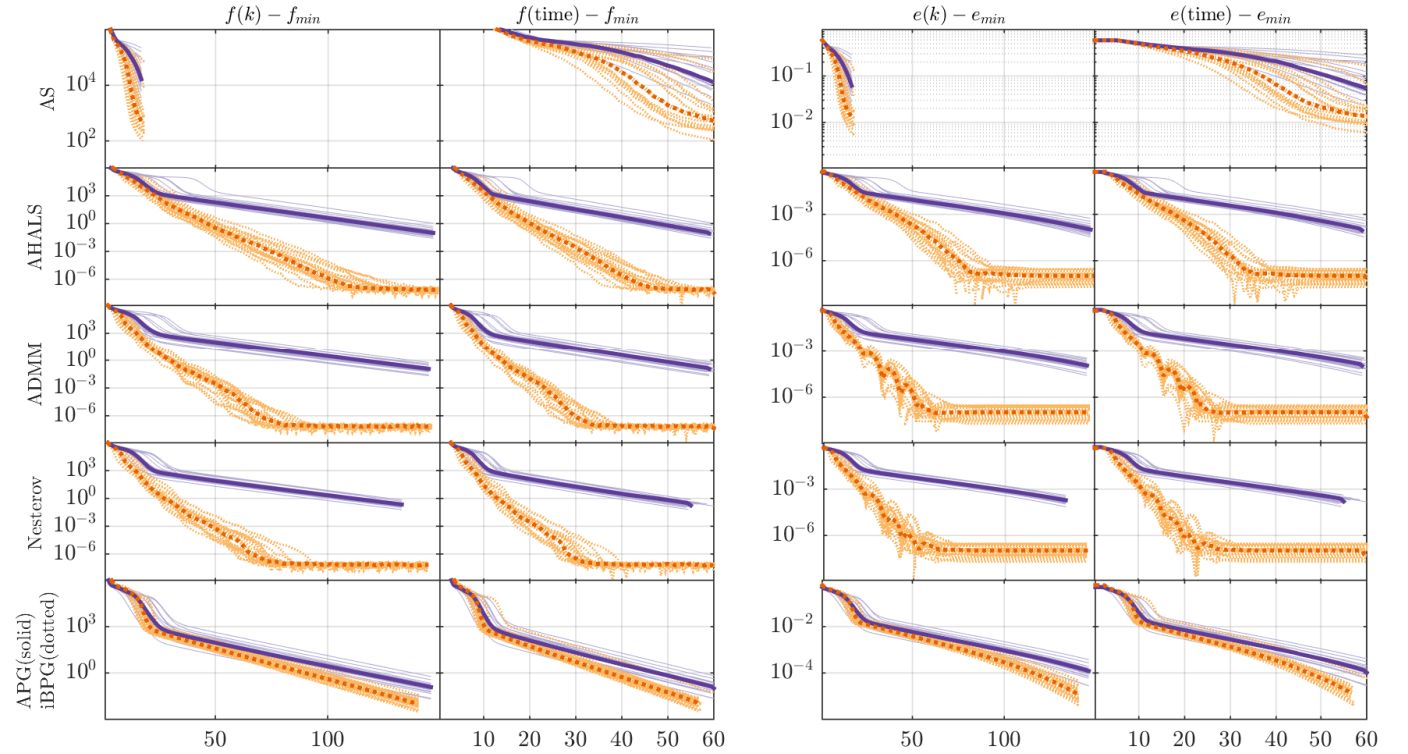


Figure 13: Experiments on synthetic data with $[I_1, I_2, I_3, r, \sigma] = [500, 500, 500, 10, 0.01]$ on algorithms without HER (solid purple curves) and with HER (dotted orange curves). The same setting as in Figure 10 is used. The results show that HER improves BCD algorithms, and that HER-BCD has a better performance than APG and iBPG. Here AS only ran approximately 40 iterations in 10 seconds due to high per-iteration cost.

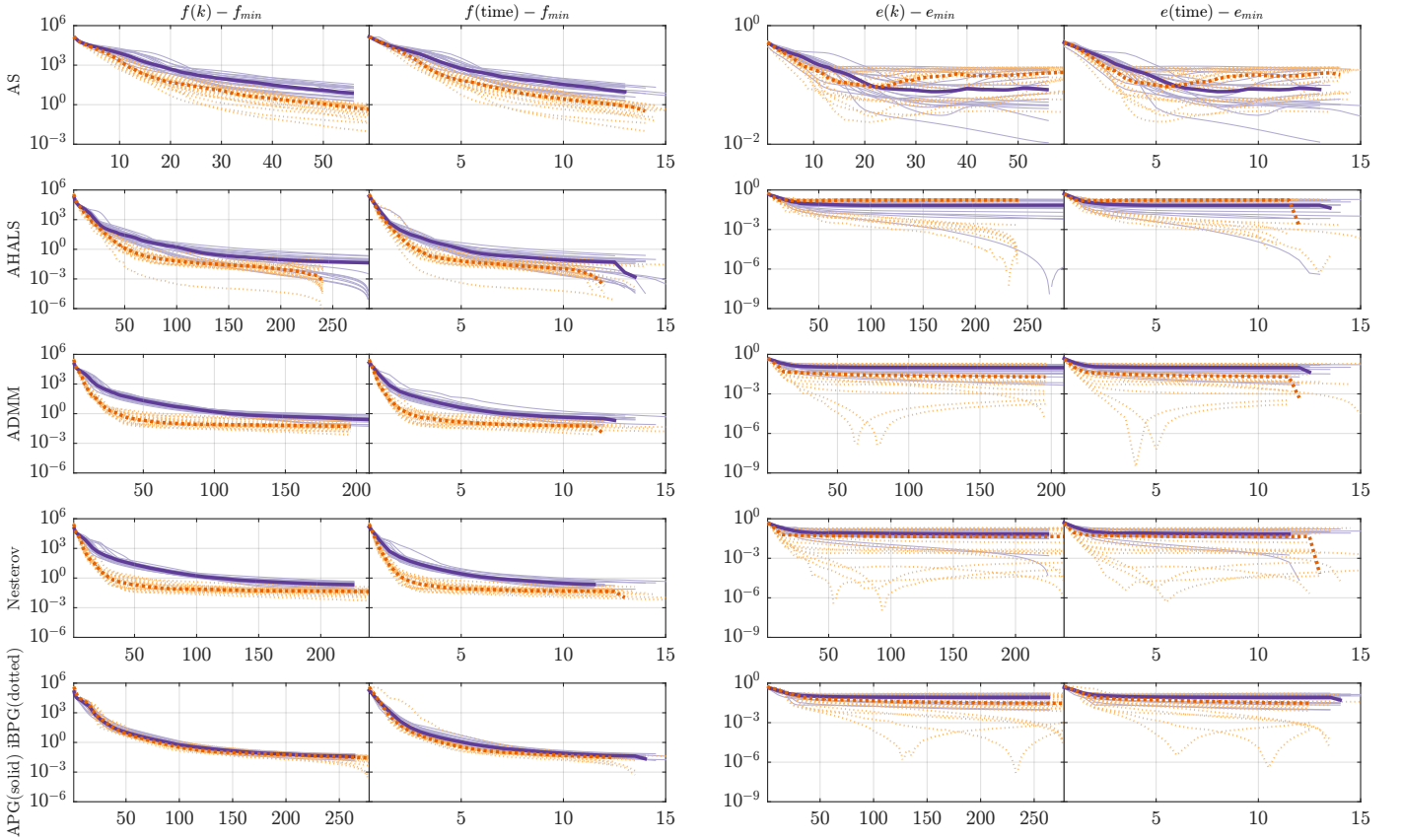


Figure 14: Experiments on synthetic data with $[I_1, I_2, I_3, r, \sigma] = [150, 10^3, 50, 12, 0.01]$ and with ill-conditioned tensors; ground truth mode factors $A_i(:, 1) = 0.99A_i(:, 2) + 0.01A_i(:, 1)$ for $i \in \{1, 2, 3\}$. In terms of best case performance, the results show that HER improves BCD algorithms, and that HER-BCD has a better performance than APG and iBPG on f , and a better performance than APG on e .

We now show the full results on comparing HER-BCD with Bro-BCD, GR-BCD (using Gradient ratio, in the modified form) and LS-BCD (line search, in the modified form), with BCD being AHALS, AO ADMM or AO Nesterov. While Fig.15 to Fig.23 will show the full result, the table below shows the summary of these results: here each curve represents the average over 20 trials. In general, HER-BCD has the best performance among all the extrapolated BCD methods.

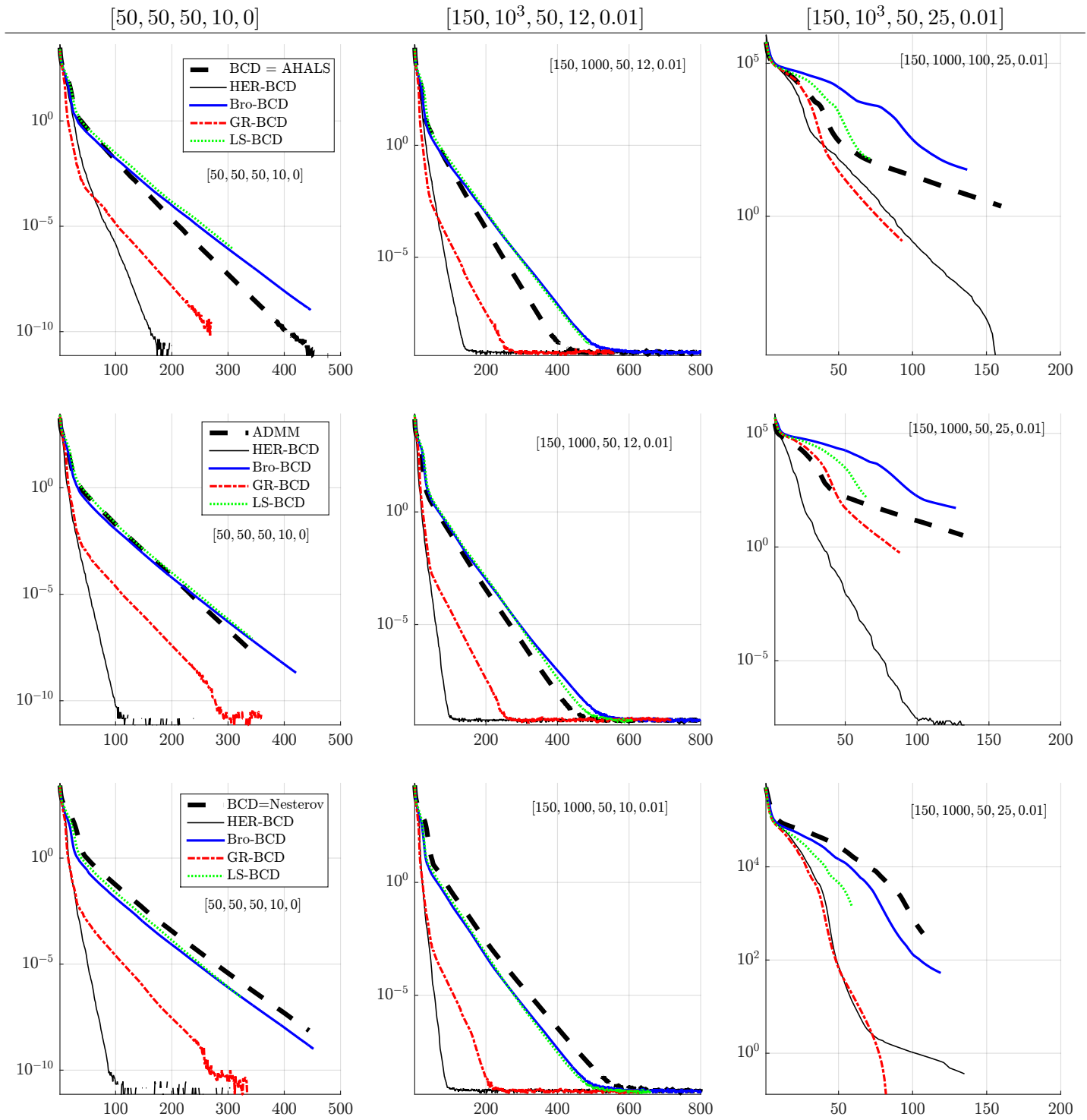


Table 5: Summary of the results of Fig.15 to Fig.23.

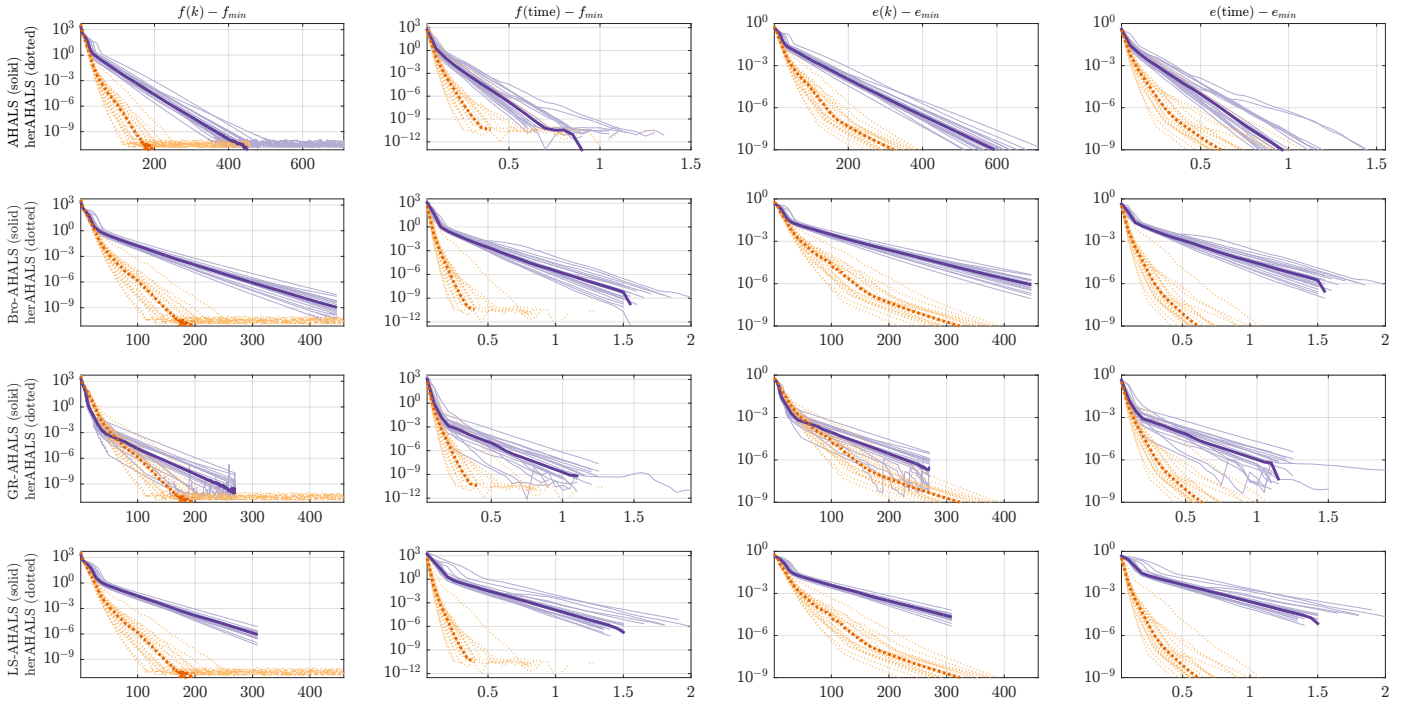


Figure 15: Experiments on synthetic data with $[I_1, I_2, I_3, r] = [50, 50, 50, 10]$ for different AHALS algorithms. Results show HER-AHALS has better performance than all other AHALS algorithms.

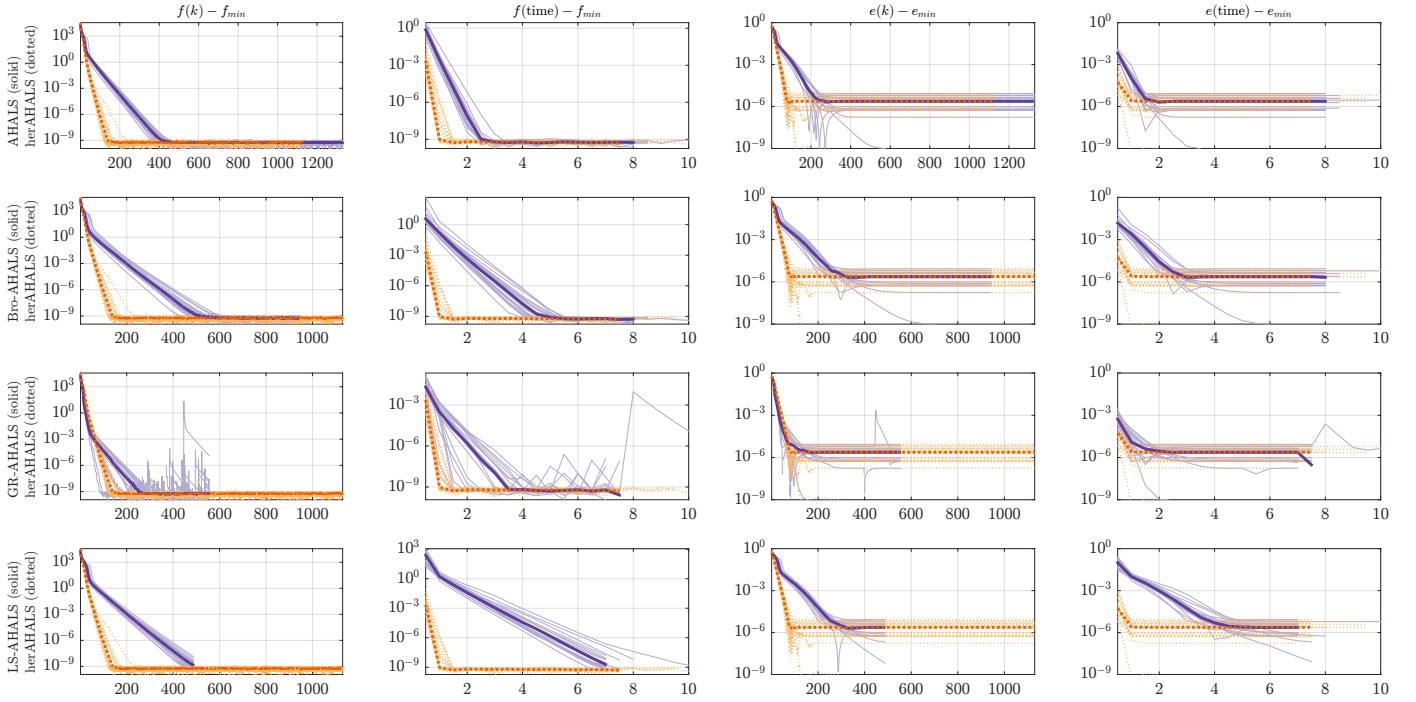


Figure 16: Experiments on synthetic data with $[I_1, I_2, I_3, r, \sigma] = [150, 1000, 50, 12, 0.01]$ for different AHALS algorithms. Results show HER-AHALS has better performance than all other AHALS algorithms.

As BCD and LS-BCD consistently perform slower than HER-BCD, we will not show the results on comparing HER-BCD with them in the subsequent plots.

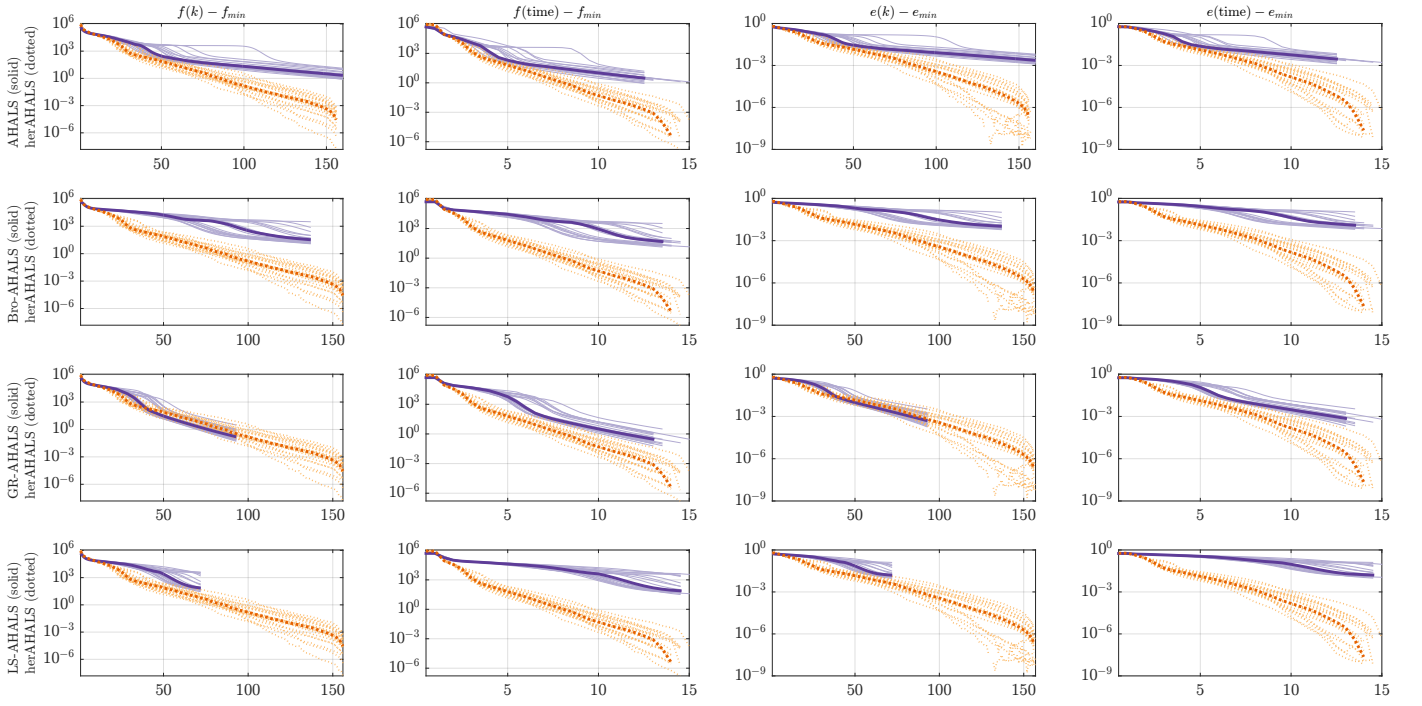


Figure 17: Experiments on synthetic data with $[I_1, I_2, I_3, r, \sigma] = [150, 1000, 50, 25, 0.01]$ for different AHALS algorithms. Results show HER-AHALS has better performance than all other AHALS algorithms.

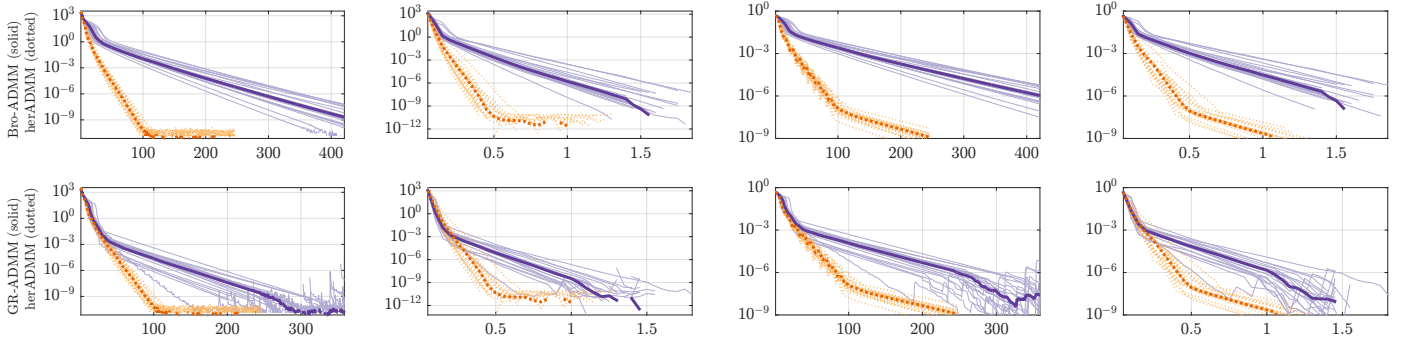


Figure 18: Experiments on synthetic data with $[I_1, I_2, I_3, r] = [50, 50, 50, 10]$ over different AO-ADMM algorithms. Results show HER-ADMM has better performance than all other AO-ADMM algorithms.

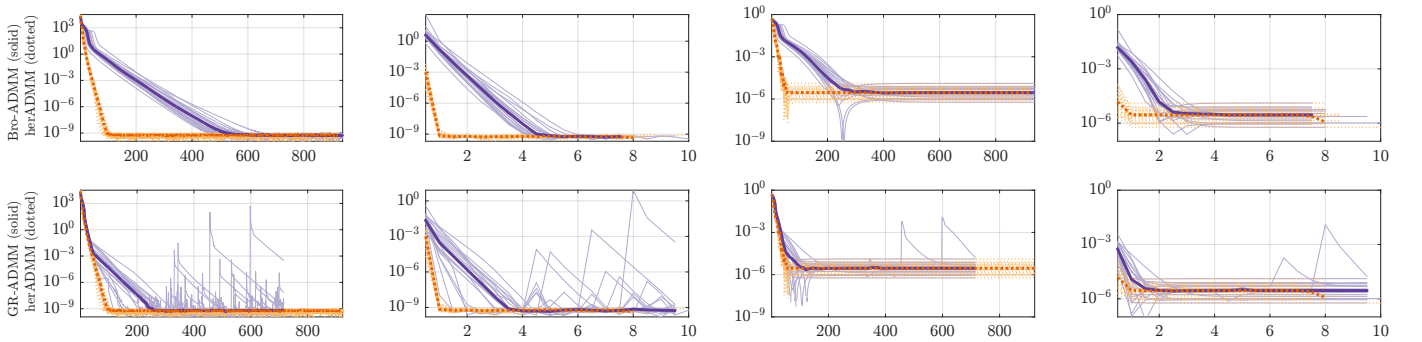


Figure 19: Experiments on synthetic data with $[I_1, I_2, I_3, r, \sigma] = [150, 1000, 50, 12, 0.01]$ for different AO-ADMM algorithms. Results show HER-ADMM has better performance than all other AO-ADMM algorithms.

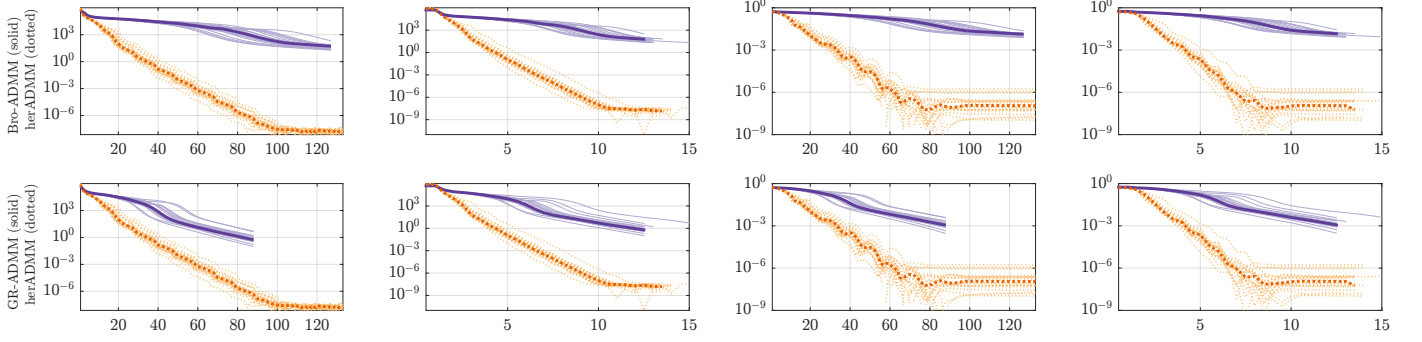


Figure 20: Experiments on synthetic data with $[I_1, I_2, I_3, r, \sigma] = [150, 1000, 50, 25, 0.01]$ for different AO-ADMM algorithms. Results show HER-ADMM has better performance than all other AO-ADMM algorithms.

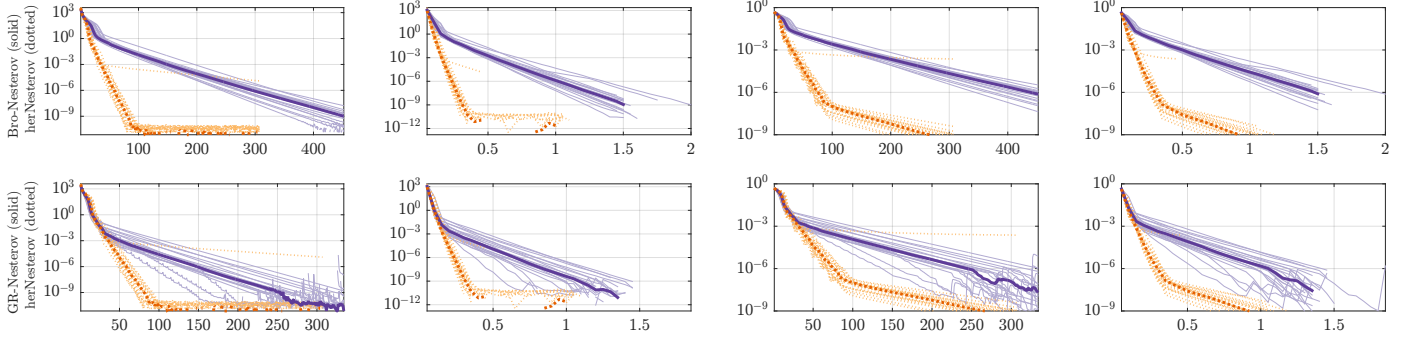


Figure 21: Experiments on synthetic data with $[I_1, I_2, I_3, r] = [50, 50, 50, 10]$ over different AO-Nesterov algorithms. Results show HER-Nesterov has better convergence performance than all other AO-Nesterov algorithms.

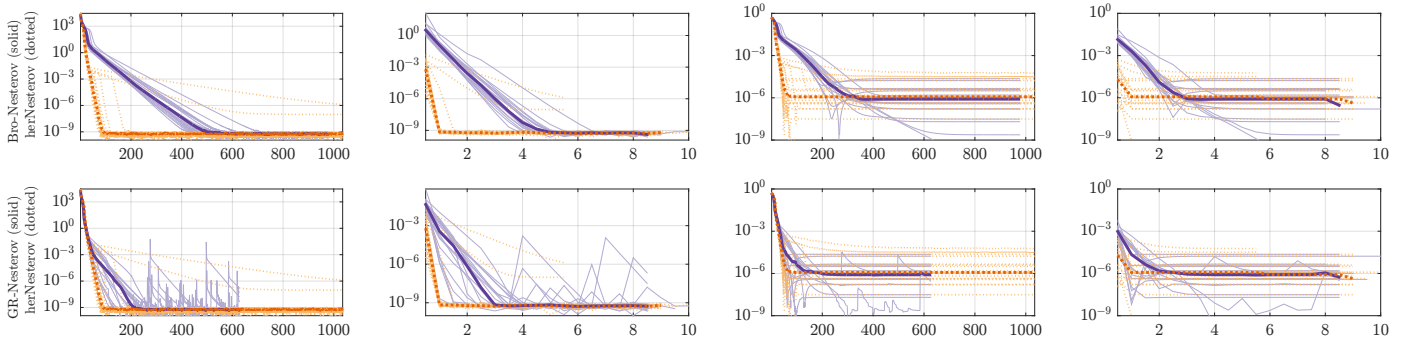


Figure 22: Experiments on synthetic data with $[I_1, I_2, I_3, r, \sigma] = [150, 1000, 50, 12, 0.01]$ for different AO-Nesterov algorithms. Results show HER-Nesterov has better performance than all other AO-Nesterov algorithms.

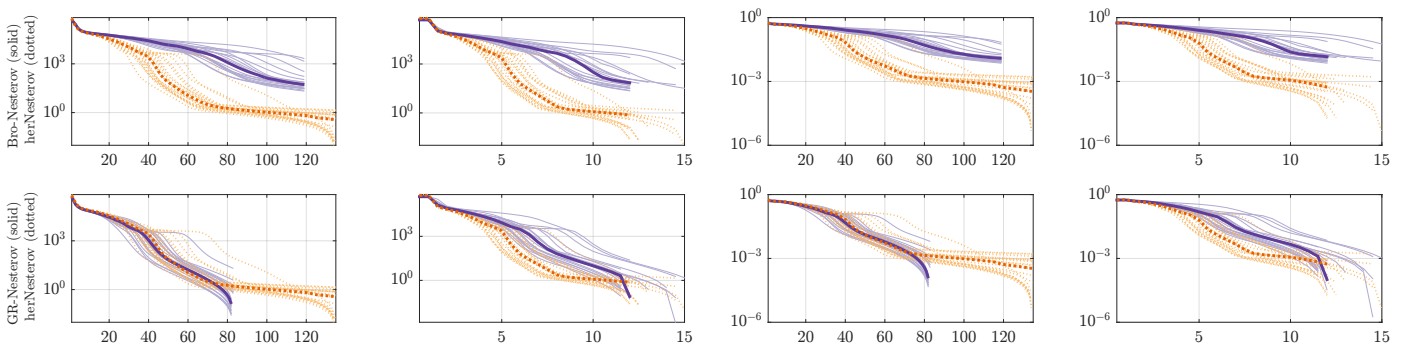


Figure 23: Experiments on synthetic data with $[I_1, I_2, I_3, r, \sigma] = [150, 1000, 50, 25, 0.01]$ for different AO-Nesterov algorithms. Results show HER-Nesterov has better performance than all other AO-Nesterov algorithms, except GR-Nesterov.

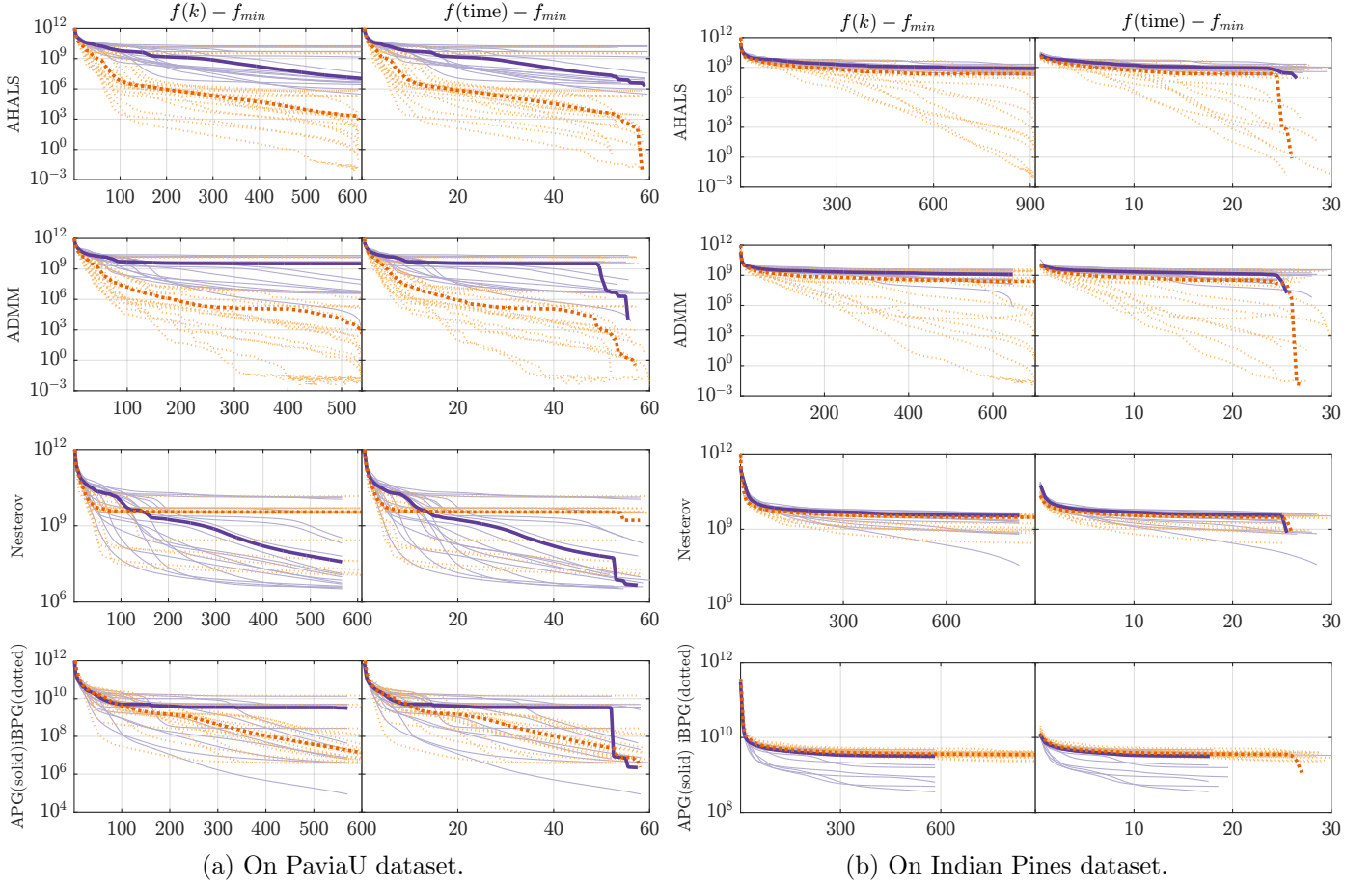


Figure 24: Experiments with HSI data. Gradient-based methods (Nesterov, HER-Nesterov, APG and iBPG) perform worse than AHALS and ADMM. In terms of f , the best run of HER-ADMM and HER-AHALS are about $10^7 - 10^9$ times better than AO-Nesterov, HER-Nesterov, APG and iBPG. In terms of e , the best run for HER-ADMM and HER-AHALS are about $10^{10} - 10^{12}$ times better than Nesterov, APG and iBPG.