# Stochastic Mirror Descent for Low-Rank Tensor Decomposition Under Non-Euclidean Losses

Wenqiang Pu<sup>®</sup>, Member, IEEE, Shahana Ibrahim<sup>®</sup>, Xiao Fu<sup>®</sup>, Senior Member, IEEE, and Mingyi Hong <sup>®</sup>, Senior Member, IEEE

Abstract—This work considers low-rank canonical polyadic decomposition (CPD) under a class of non-Euclidean loss functions that frequently arise in statistical machine learning and signal processing. These loss functions are often used for certain types of tensor data, e.g., count and binary tensors, where the least squares loss is considered unnatural. Compared to the least squares loss, the non-Euclidean losses are generally more challenging to handle. Non-Euclidean CPD has attracted considerable interests and a number of prior works exist. However, pressing computational and theoretical challenges, such as scalability and convergence issues, still remain. This work offers a unified stochastic algorithmic framework for large-scale CPD decomposition under a variety of non-Euclidean loss functions. Our key contribution lies in a tensor fiber sampling strategy-based flexible stochastic mirror descent framework. Leveraging the sampling scheme and the multilinear algebraic structure of low-rank tensors, the proposed lightweight algorithm ensures global convergence to a stationary point under reasonable conditions. Numerical results show that our framework attains promising non-Euclidean CPD performance. The proposed framework also exhibits substantial computational savings compared to state-of-the-art methods.

Index Terms—Tensor decomposition, stochastic optimization, mirror descent method,  $\Box$ -divergence, KL-divergence.

## I. INTRODUCTION

ANONICAL polyadic decomposition (CPD) has been used in many core tasks in signal processing and machine learning, such as neural signal analysis, video processing, array signal processing, text mining, social network analysis, link prediction, among others—see [1]—[4].

Manuscript received May 11, 2021; revised December 8, 2021 and March 10, 2022; accepted March 21, 2022. Date of publication March 31, 2022; date of current version April 26, 2022. The associate editor coordinating the review of this manuscript and approving it for publication was Dr. Evrim Acar. The work of Wenqiang Pu was supported by the NSFC, China under Grant 62101350. The work of Xiao Fu was supported in part by NSF ECCS under Grants 1808159 and IIS-1910118 and in part by ARO under Grant W911NF-19-1-0247. Howork of Mingyi Hong was supported in part by NSF under Grant CIF-1910385 and in part by ARO under Grant W911NF-19-1-0247. (Corresponding author: Wenqiang Pu.)

Wenqiang Pu is with Shenzhen Research Institute of Big Data, Shenzhen 518172, China (e-mail: wenqiangpu@outlook.com).

Shahana Ibrahim and Xiao Fu are with the Electrical Engineering and Computer Science Department, Oregon State University, Corvallis, OR 97331 USA (e-mail: ibrahish@oregonstate.edu; xiao.fu@oregonstate.edu).

Mingyi Hong is with the Electrical and Computer Engineering Department, University of Minnesota, Minneapolis, MN 55455 USA (e-mail: mhong@umn.edu).

This article has supplementary downloadable material available at https://doi.org/10.1109/TSP.2022.3163896, provided by the authors.

Digital Object Identifier 10.1109/TSP.2022.3163896

The majority of classic CPD models and algorithms were developed for least squares (LS) problems using the Euclidean distance-based fitting criterion; see [1]-[3], [5] and references therein. However, the Euclidean distance is unnatural for measuring the "distance" between many types of real-world data, e.g., stochastic, integer, and binary data. In principle, using certain "data geometry-aware" divergences to serve as the fitting criteria may greatly improve performance and robustness in practice [6]–[9]. For example, the "distance" between two probability distributions is typically measured by a proper divergence, such as the generalized Kullback-Leibler (KL) divergence [10]-[13] and Itakura-Saito (IS) divergence [14]. These divergences take into consideration that the data is constrained in the probabilistic simplex, and thus are often more effective relative to the LS criterion in analyzing data that are not generated over the entire Euclidean space. From a statistical estimation viewpoint, many non-Euclidean divergences are closely related to the maximum likelihood estimators (MLEs) under plausible data distributions. For example, the generalized KL divergence [15] and logistic loss [6]-[8] can be derived from the MLEs of count integer data and binary data that follow certain Poisson distributions and Bernoulli distributions, respectively.

However, computing CPD under non-Euclidean divergences is much more challenging compared with the case under Euclidean loss (or, the LS loss), especially when the data size becomes huge. Algorithms developed under the LS loss are often not easily extendable to handle these more complicated loss functions, due to the lack of "nice" properties that are possessed by the LS loss, e.g., the gradient Lipschitz continuity under relatively mild conditions. Below, we provide a brief review on existing developments for CPD models with specific loss function.

## A. Prior Works

Many existing non-Euclidean CPD approaches employ the block coordinate descent (BCD) paradigm [27] with divergence-specific strategies for block variable updating. For example, the work in [26] proposed a hierarchical alternating optimization algorithm for CPD with □ and □divergence. In [15], the generalized KL-divergence loss was considered, where a block majorization-minimization (MM) algorithm was developed. In [11], the exponential gradient algorithm was proposed for the KL-divergence. Similar strategies were developed for

1053-587X © 2022 IEEE. Personal use is permitted, but republication/redistribution requires IEEE permission. See https://www.ieee.org/publications/rights/index.html for more information.

Algorithm	Stochastic	Loss function	Data Type
First-order type algorithm [16]–[18]	No	LS	continuous
Primal-dual algorithm [19]	No	LS	continuous
(Quasi-)Second-order algorithm [20], [21]	No	LS	continuous
Stochastic optimization algorithm [22]-[25]	Yes	LS	continuous
Hierarchical alternating optimization [26]	No	$\alpha, \beta$ Div.	continuous
Majorization-minimization algorithm [15]	No	KL Div.	continuous and count
Multiplicative update algorithm [14]	No	LS, KL Div., IS Div.	binary
Exponential gradient algorithm [11]	No	KL Div.	continuous and count
Alternating optimization algorithm [8]	No	logistic loss and others	binary
Generalized Gaussian Newton algorithm [9]	No	$\beta$ -div.	continuous and count
Stochastic gradient descent algorithm [7]	Yes	general loss*	continuous, count, and binary
Stochastic mirror descent algorithm (this work)	Yes	general loss*	continuous, count, and binary

TABLE I
BRIEF REVIEW OF ALGORITHMS FOR CPD MODEL

the KL and IS divergences [14]. Recently in [8], several ML-based loss functions for binary data were considered and an alternating optimization algorithm with line search was proposed. Besides BCD, other optimization frameworks such as Gauss-Newton based methods [9], quasi-Newton methods [6], and stochastic gradient-based methods [7] were also developed for non-Euclidean CPD.

It is important to note that most of the algorithms mentioned above (such as [8], [9], [14], [15], [26], [28]), are batch algorithms, which utilize the entire data set to perform every update. They become increasingly slow when the size of the data increases. On the other hand, stochastic algorithms are effective in reducing per-iteration computational and memory burdens. Recently, a stochastic gradient descent (SGD) based algorithm [7] was proposed for CPD with a suite of statistical criterion based loss functions (including loss functions under missing values that are not considered in this work). The algorithm was developed based on randomly sampling the tensor entries. Hence, it is difficult to exploit some interesting multilinear algebraic properties of low-rank tensors to further improve computational efficiency. In addition, the SGD algorithm in [7] lacks convergence guarantees. Finally, in Table I, we summarize the properties of a number of recently developed algorithms for the CPD model.

## B. Contributions

In this paper, we develop a *unified stochastic mirror descent* (SMD) algorithmic framework for large-scale CPD under various non-Euclidean losses. Our major contributions are summarized as follows:

• Efficient fiber-sampled stochastic MD framework: We propose a block-randomized SMD algorithmic framework that is tailored for tensor decomposition. Both MD and SMD are known for its effectiveness in handling non-Euclidean losses [29], but directly applying generic SMD does not fully exploit the underlying CPD structure. We use a recently emerged tensor data sampling strategy (namely, fiber sampling [25], [30]) to assist designing SMD-type updates. The fiber sampling strategy judiciously uses the multilinear structure of low-rank tensors, which gives rise to structured (non-)convex subproblems. These

structures can often be exploited to come up with economical update rules for CPD.

• A suite of solutions for various losses and constraints: We carefully craft solutions for a series of non-Euclidean losses. The proposed algorithmic framework allows flexible choices of the local surrogate functions under the SMD framework to adapt to different loss functions. Such flexibility also helps offer lightweight updates when the latent factors are under a variety of constraints that are of interest in data analytics. In particular, we pay special attention to binary and integer data CPD problems, which find numerous applications across disciplines.

• Guaranteed convergence: We offer convergence characterizations for our block-randomized SMD-based non-Euclidean CPD framework. Establishing stationary-point convergence for generic SMD is already a challenging problem. The work in [31] on SMD requires its gradient estimation error converging to zero, which is unrealistic in many cases, especially under the context of CPD. In this work, we leverage the notion of relative smoothness and the tensor fiber sampling strategy to construct lightweight SMD updates for different losses. This design also helps circumvent stringent conditions (e.g., vanishing gradient estimation error) when establishing convergence. To our best knowledge, such convergence results for multi-block SMD under nonconvex settings have been elusive in the literature.

Part of the work appears in IEEE ICASSP 2021 [32]. The conference version considered algorithm design under the ☐ divergence loss. This journal version extends the ideas to handle more non-Euclidean losses, e.g., the logistic loss that is critical in binary data analysis. More importantly, this version provides unified convergence analysis for the proposed algorithmic structure. Some important practical considerations, e.g., stepsize scheduling, is also discussed and experimented. The journal version also contains substantially more simulations and real-data validation.

## C. Notation

We follow the conventional notation in signal processing, x, X, X, and X denote scalar, vector, matrix, and tensor, respectively. Given a matrix X,  $X^{c}$  and  $\exp(X)$  denote the entry-wise power and exponential operations respectively;  $X \pm c$  denote

<sup>\*</sup>The general loss in this table refers to many ML criterion motivated losses [6], [7] as well as statistical divergences such as KL div., IS div. and etc.

adding/subtracting a constant c for each entry of X;  $\operatorname{vec}(X)$  denote the vectorization operator that concatenates the columns of X. We use  $\mathbb{R}$ ,  $\mathcal{X}$  and  $\mathcal{X}$  to denote the Hadamard product, the Khatri-Rao product, and entry-wise division respectively. The denotes the transpose operation. Script letter is used to denote a discrete set and is the cardinality f. The denotes the Euclidean norm of vector, f. The denotes the Frobenius norm of matrix, and f, denotes the inner product. For a positive integer f, the set f, f, f, f is denoted as f for short. For a function f its domain and interior domain is denoted as domain and int dom f respectively. Other notation will be explained when it first appears.

## II. CPD UNDER NON-EUCLIDEAN LOSSES

Consider a data tensor  $\underline{X}$   $\mathbb{R}^{I_1 \times I_2 \times ... \times I_N}$ , where N is the order of the tensor  $\underline{X}$  and  $I_n > 0$  (n [N]) is the size of the nth mode of  $\underline{X}$ . Such multi-way data tensors arise in many applications. The entries of the data tensor  $\underline{X}$  could be continuous real numbers, non-negative integers or binaries. A general problem of interest is to approximate  $\underline{X}$  using a low rank tensor  $\underline{M}$ , defined as

$$\underline{\mathbf{M}} = \mathbf{L}_{R} \qquad \mathbf{A}_{\underline{f}}(:,r) \square \mathbf{A}_{\underline{f}}(:,r) \square \dots \square \mathbf{A} (:,r), \qquad (1)$$

where " $\square$ " denotes the outer product of vectors,  $\mathbf{A}_n = \mathbb{R}^{I_n \times R}$  is the mode-n latent factor matrix; R is the smallest positive integer such that (1) holds, and it is also known as the rank of  $\mathbf{M}$ .

Denote an N -dimensional integer vector  $\mathbf{i}$  as the entry coordinate, i.e.,

*i* 
$$\square \square \{(i_1, i_2, ..., i_N) | i_n = 1, 2, ..., I_n, \square n\}.$$

Then the CPD problem can be formulated as the following minimization problem with a loss function of interest  $f(\cdot, \cdot)$ :

$$\begin{array}{cccc}
R \times R & \square R, \\
& \min & \frac{1}{I} & \mathbf{L} \\
\mathbf{A}_{1}, \mathbf{A}_{2}, \dots, \mathbf{A}_{N} & |I| & \mathbf{I} \\
& \text{s.t.} & \mathbf{M} = \mathbf{A} & (i, r), \square i & \mathbb{I}
\end{array}$$

$$\mathbf{A}_{n} \quad \mathbf{A}_{n} \quad \mathbf{I}_{n} & \mathbf{I}_{n} \qquad (2)$$

where  $X_i$  and  $M_i$  denote the entries of X and M indexed by i, respectively,  $A_n$  is a convex and closed set which captures the prior information about the structure of latent factors  $A_n$ , e.g., non-negativity, sparsity, and smoothness. By choosing proper loss functions f, Problem (2) is used for handling different types of data, e.g., continuous, count, and binary data. Several representative motivating examples are as follows,

• KL-divergence for count data: In many real-world scenarios, data is naturally recorded as nonnegative integers, e.g., crime numbers across locations and time¹ and email interactions recorded over months [33]. As an information-theoretic measure, the KL divergence was originally proposed for quantifying

<sup>1</sup>See official website of the city of Chicago, www.cityofchicago.org.

similarity between two probability distributions. The *generalized* KL-divergence that handles nonngeative quantities beyond distributions is also widely used in data analytics [10], [12], [15], [34]. The generalized KL divergence has an *f* defined as follows:

(KL-Div.) 
$$f(x, m) = m \square x \log(m)$$
, (3)

where x \ and m 0. Problem (2) with KL-divergence can also be interpreted as the MLE for estimating the Poisson parameter tensor  $\underline{\mathbf{M}}$  from independent samples  $\underline{\mathbf{X}}$ , where  $\underline{\mathbf{M}}$  has a low-rank structure [6], [7], [15]. To be more precise, the corresponding statistical model is

$$X_i \sim \text{Poisson}(\underline{M}_i), \Box i$$

where Poisson(m) denotes the Poisson distribution with a mean of m.

• Log loss for binary data: Binary data is also frequently encountered in data analytics, e.g., in adjacency matrix-based social network community detection [4], [35] and knowledge base analysis [6]–[8]. Binary data fitting is often handled using the following loss:

(Log Loss) 
$$f(x, m) = \log(1 + e^m) \square xm$$
, (4)

where  $x \in [0, 1]$  and m = 0. The log loss can be interpreted as MLE for finding the Bernoulli distribution parameter [8]. The associated binary data generation model of the independent samples is

$$\underline{X}_i \sim \text{Bernoulli}(\boxtimes_i), \ \boxtimes_i = e^{\underline{M}_i}/(1 + e^{\underline{M}_i}), \ \Box i,$$

where Bernoulli ( $\square$ ) denotes the Bernoulli distribution, is the probability for x taking 1, and  $\square \square$  has the same size  $\underline{\mathscr{K}}$ .

• □divergence: Non-Euclidean losses also find applications in continuous data CPD, especially under non-Gaussian and/or non-additive noise, e.g., multiplicative Gamma noise [9]. For example, the □divergence was found useful for neural signal analysis [26] and recently is studied as CPD fitting criterion [6], [7], [9], [32]. The □divergence is parametrized by a constant

The  $\square$ divergence subsumes the IS divergence ( $=\square$ ), the generalized KL divergence ( $\square$ = 1), and the Euclidean distance ( $\square$ = 2) as special cases. When  $\square$ = 0, it can also be interpreted as MLE corresponds to data with multiplicative Gamma noise. In music data analysis,  $\square$  < 2 was found useful, since such loss functions capture low intensity spectra components—but the Euclidean loss tends to focus on significant variations in data [34].

Remark 1: As one has seen in the examples, one way to select f is to take a statistical analysis viewpoint. Each entry of the data tensor is treated as a conditional independent random variable, whose generation is statistically modeled as follows:

$$X_i \sim p(X_i \mid \Box M_i)), \Box i \quad \Box$$
 (5)

Data Type	Distribution	Link Function	Loss function	Parameter Type	Name	
Continuous	Gaussian Gamma	$\theta(m) = m$ $\theta(m) = m$	$\frac{\frac{1}{2}(x-m)^2}{\frac{x}{m+\epsilon} + \log(m+\epsilon)}$	$x, m \in \mathbb{R}$ $x > 0, m \ge 0$	Euclidean Dis. IS Div.	
	-	-	$(m+\epsilon)^{\beta}/\beta - x(m+\epsilon)^{\beta-1}/(\beta-1)$	$z \geq 0, m \geq 0$	$\beta$ -Div. $(\beta \in \mathbb{R} \setminus \{0, 1, 2\})$	
Count	Poisson	$\theta(m) = m$	$m = x \log(m + \epsilon)$	$x \in \mathbb{N}, m \ge 0$	generalized KL Div.,	
		$\theta(m) = e^m$	$e^{m} - xm$	$x \in \mathbb{N}, m \in \mathbb{R}$	-	
Binary	Bernoulli	$\theta(m) = \frac{m}{1+m}$	$\log(m+1) - x\log(m+\epsilon)$	$x\in\{0,1\}, m\geq 0$	-	
		$\theta(m) = \frac{m}{1+m}$ $\theta(m) = \frac{e^m}{1+e^m}$	$\log(1+e^m) - xm$	$x \in \{0,1\}, m \in \mathbb{R}$	-	

TABLE II
DISTRIBUTIONS, LINK FUNCTIONS, AND LOSS FUNCTIONS FOR DIFFERENT TYPES OF DATA

where  $p(x; \square)$  is a distribution with natural parameter (e.g., the Poisson and Bernoulli distribution) and  $\square$ ():  $R_{\square}R$  is an invertible function whose inverse is often referred as the *link function* in statistics (e.g.,  $\square (n) = m$  and  $\square (n) = \frac{m}{1+m}$ ). A straightforward intuition behind model (5) is that, the observed data tensor  $\underline{X}$  is 'embedded' on a *latent* low rank tensor  $\underline{M}$ , whose generation procedure is characterized by  $p(x; \square)$  and  $\square$ (). To find  $\underline{M}$ , a statistically efficient way is choosing f() as the negative log-likelihood function associated with model (5), given as

$$f(\underline{X}_i, \underline{M}_i) \square \square \log p(\underline{X}_i \mid \square \underline{M}_i) + \text{constant},$$

which naturally leads to a non-Euclidean CPD problem (if the distribution is not Gaussian).

In Table II, some frequently used  $f(\cdot)$ , link function  $\square(0)$ , and distribution of our interests are illustrated. In the table, c > 0 is a sufficiently small number to avoid the function value or gradient being  $\square$  i.e.,  $c = 10^{\square}$ . This modification is often used in the literature [6].

Tackling Problem (2) at large scale is highly nontrivial. For example, a 5000  $\times$  5000 5000 tensor can be as large as 900 GB if the double precision arithmetic is used, which means that batch algorithms may not be a viable option. Instead, stochastic algorithms that sample 'partial data' per iteration become an attractive choice. In Euclidean loss CPD, it has been observed that stochastic algorithms can significantly reduce computational and memory cost per iteration; see [23], [25], [30]. Nonetheless, unlike Euclidean CPD, where various data sampling schemes and update rules may all offer competitive algorithms [5], non-Euclidean losses' complex structures may make stochastic algorithm design a more delicate process. In other words, the sampling schemes may affect the subsequent update rules' complexity and convergence properties of the algorithm.

Next, we offer a unified stochastic algorithmic structure that can efficiently tackle CPD under a variety of non-Euclidean losses. Our development is an integrated design of data sampling and *Lipschitz-like convexity* [36] based local surrogate construction, leveraging the underlying multilinear structure of low-rank tensors.

## III. PROPOSED APPROACH

## A. Preliminaries

A number of algorithmic frameworks have been considered for handling Problem (2) with non-Euclidean losses, e.g., stochastic gradient descent (SGD) [7], block coordinate descent (or, alternating optimization (AO)) [8], [15], and the Gauss-Newton (GN) method [9].

Let us denote  $\mathbf{A} := (\mathbf{A}_1, \mathbf{A}_2, ..., \mathbf{A}_N)$ ,  $\mathbf{A} := \mathbf{A}^1 \times \mathbf{A}^2$ ... We also use  $F(\mathbf{A})$  to represent the objective function of Problem (2). The updates of AO and SGD type algorithms can be summarized as follows:

(AO) 
$$\mathbf{A}_n^{t+1} \square \text{arg min } F(\mathbf{A}_n; \mathbf{A}_n^t),$$
 (6a)

(SGD) 
$$\mathbf{A}^{t+1} = \operatorname{Proj}_{\mathbf{A}}(\mathbf{A}^{t} \square \square \widehat{\mathbf{G}}^{t})$$
  
=  $\operatorname{arg min} (\mathbf{A}, \widehat{\mathbf{G}} \pm + \frac{1}{2} \square \mathbf{A}^{t} \square \mathbf{A}^{t} \square^{2}.$  (6b)

In (6),  $\mathbf{A}_{\square^n}$  corresponds to  $\mathbf{A}$  with  $\mathbf{A}_n$  being removed;  $F(\mathbf{A}^n; \mathbf{A}^{t\square})$  is the objective function with fixed  $\mathbf{A}^{t^n}$ ;  $\mathbf{G}_{is}$ 

represents the the step size; and Proj. () denotes the projection onto constraint set A Many deterministic non-Euclidean tensor decomposition algorithms take the AO route; see, e.g., [8], [15], whereas the recent work in [9] used a GN method to improve the iteration complexity (i.e., the number of iterations needed for reaching a satisfactory solution). However, the AO and GN methods face heavy per-iteration computational and memory complexities when handling large-scale tensors; see the "MTTKRP" challenge discussed in [1], [2], [5], [25]. The SGD approach in [7] is more lightweight in terms of the per-iteration resource consumption. Constructing  $\hat{\textbf{G}}$  can be fairly economical (c.f. [6, Theorem 3]) since only partial data is used. This makes the per-iteration complexity of the algorithm affordable, even if the tensor of interest is large.

However, simply using SGD for the non-Euclidean CPD problem may not be the most effective approach. One can see that from (6b), every iteration of SGD is equivalent to solving a quadratic program, which is used as a local surrogate of the original cost function. However, it is known that such quadratic functions may not be a good approximation for many non-Euclidean losses. In particular, using quadratic local surrogates may result in slow progresses since it largely ignores the geometry of the cost function [29], [36].

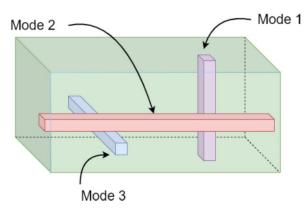


Fig. 1. Mode-*n* fibers of a third-order tensor, where n = 1, 2, 3.

In this section, we will propose a stochastic mirror descent (SMD) framework to handle the non-Euclidean CPD problem. Compared with SGD (6b), SMD replaces the quadratic term  $A_{1}$   $F^{t}$  by a general divergence metric named Bregman divergence  $D_{\square}(\mathbf{A}, \mathbf{A}^t)$  (see Definition 1),

(SMD) 
$$\mathbf{A}^{t+1} = \arg\min_{\mathbf{A}} (\mathbf{A}, \hat{\mathbf{G}} \pm + \frac{1}{\Box} D(\mathbf{A}, \mathbf{A}^t).$$
 (7)

If  $\square$  is properly designed, SMD can often take advantages of the problem geometry and attain substantial efficiency improvement compared to SGD, especially when non-Euclidean cost functions are used; see more on MD and SMD in the optimization literature [37]-[41]. The reason why using MD/SMD improves convergence over GD/SGD will become clearer later in the next subsections when specific examples of non-Euclidean losses and  $D_{\square}$ 's (e.g. the example in Fig. 2) are discussed.

Definition 1 (Bregman Divergence): Given a strongly convex function  $\square(\cdot)$ : dom  $\square(\square R^n)\square R$ , the Bregman divergence between **a** døm and **a**t Int dom is  $\square$ 

$$D_{\square}(\boldsymbol{a},\boldsymbol{a}) = \square(\boldsymbol{a}) \square \square(\boldsymbol{a}^t) \square \square \square(\boldsymbol{a}^t), \boldsymbol{a} \square \boldsymbol{a}^t \pm .$$

Remark 2: The Bregman divergence was originally defined in [42] using a Legendre function  $\square$  Here, we consider the case where  $\square$  is a strongly convex function; see more details in [36]. It

should be noted that  $D_{\square}(\boldsymbol{a}, \boldsymbol{a}^t) = 0$  of and  $D_{\square}(\boldsymbol{a}, \boldsymbol{a}^t) = 0$  if and only if  $\boldsymbol{a} = \boldsymbol{a}^t$  due to the strong convexity of  $\square$  in addition, if  $\square$  is a one-dimensional function,  $D_{\square}(\boldsymbol{a}, \boldsymbol{a}^t)$  can be defined in a coordinate-wise form as

$$D\square(\boldsymbol{a},\boldsymbol{a}) = j_{-1}\square(a_j)\square\square(a_j)\square$$

The general-purpose SMD scheme in (7) is not new in the optimization literature and its convergence behavior under the stochastic setting has been studied in recent works, e.g., [31], [38], [40]. However, the general-purpose SMD scheme does not exploit the structure of the CPD problem in (2). Next, together with a tensor fiber sampling strategy advocated in [5], [25], [30], we will develop a block randomized SMD algorithm to exploit

the multilinear structure of CPD, where only one block  $\mathbf{A}_n$  is randomly selected and updated per SMD iteration.

## B. Data Sampling

A key ingredient for stochastic algorithms lies in the data sampling strategy, which uses partial data to estimate the gradient in each iteration. Under the Euclidean loss, sub-tensor sampling [23], random entry sampling [22], and tensor fiber sampling [25], [30] were all considered—which all offered effective solutions. In principle, all the sampling strategies considered in the Euclidean case could still be used in the non-Euclidean cases. For example, the recent non-Euclidean CPD work in [7] used an entry sampling scheme. Nonetheless, since non-Euclidean losses are inherently more complex, different sampling strategies may lead to algorithms that admit drastically different updating rules and convergence properties.

In this work, we advocate the fiber sampling strategy that was used in [25], [30] for Euclidean loss CPD; see illustration in Fig. 1 for tensor fibers. We find this sampling strategy particularly handy in the non-Euclidean case, for a couple of reasons:

*Incorporating Prior Information on*  $\mathbf{A}_n$ : Randomly sampling some indexes i [22] or selecting a subtensor [23] faces an issue that the sampled data may relate to only some rows of  $A_n$ , while many important constraints under statistical non-Euclidean CPD are imposed on the columns of  $\mathbf{A}_n$ , e.g., the probability simplex constraints in [11], [43]-[45]. This may lead to difficulties in developing efficient updating rules while enforcing some types of prior information on  $\mathbf{A}_n$ . For example, the probability simplex constraints on  $\mathbf{A}_n(:, r)$ , norm constraints on  $\mathbf{A}_n(:, r)$ , and smoothness constraints on  $\mathbf{A}_n(:, r)$  all have important applications [5], but the information of  $\mathbf{A}_n(:, r)$  may not be included in the selected samples. With fiber sampling, every batch of sampled data contains information about one full  $\mathbf{A}_n$ , making dealing with such constraints easy.

• Convex Approximation to Optimize  $A_n$ : The fiber sampling strategy also provides a way to further exploit the blockwise structure of the loss functions. Even under complex non-Euclidean losses, with the notion of *Lipschitz-like convexity* [36], the loss function  $f(\cdot)$  with respect to each block  $\mathbf{A}_n$  can often be locally approximated (or, to be precise, majorized) by a strongly convex function—which will prove useful for deriving lightweight updates.

The fiber sampling scheme can be understood using the matrix unfolding representations of low-rank tensors. The mode-n matrix unfolding of  $\underline{\boldsymbol{X}}$  is a  $J_n \times I_n$  matrix with

 $J_n = \frac{\text{IT}_N}{I_m}$ , denoted as  $\mathbf{X}_n$ , and the entry-wise correspondence is  $\mathbf{X}_{1T} = \mathbf{X}_n(j, i_n)$ ,  $j = 1 + \sum_{k=1, k/=n}^{N} (i_k \square 1)S_k$ , where  $S_k = \sum_{m=1, m/=n}^{M} I_m$ . For the low-rank tensor in (1), its mode-n matrix unfolding can be expressed as  $\mathbf{M}_n = \mathbf{H}_n \mathbf{A}_n^T$  where  $\mathbf{H}_n = \mathbf{A}_N \times \mathbf{A}_{N_{\square}} \dots \times \mathbf{A}_{n+1} \times \mathbf{A}_{n+1}$  $\mathbf{A}_{n_{\square}} \times \cdots \times \mathbf{A}_{1}$  and denotes the Khatri-Rao product.

Based on the mode-n matrix unfolding for both X and M,

Problem (2) with 
$$\mathbf{A}_{\square^n}$$
 being fixed can be recast as:
$$\underbrace{\frac{1}{I}}_{J_n} \mathbf{I}_{J_n}^n f_j^n(\mathbf{A}_n; \mathbf{H}_n) \quad \text{s.t. } \mathbf{A}_n \quad \square_n,$$
(8)

where  $f^n(\mathbf{A}_n; \mathbf{H}_n) = \underline{1}^{\prod_{i=1}^n} f(\mathbf{X}_n(j, i), \mathbf{H}_n(j, :) \mathbf{A}_n(i, :)^T)$ . According to the reformulation in (8), fiber sampling can be

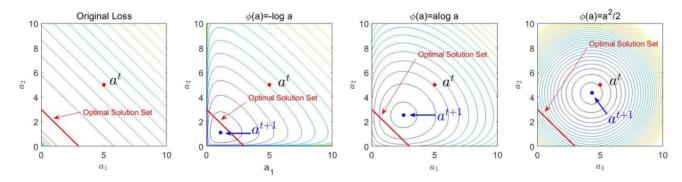


Fig. 2. Contour maps of upper bound functions over  $[0, 10]^2$  indicated by (11) with different  $\square a$ ). The loss function  $(x, m) = m \square x \log(m + E)$ , where x = 3,  $E = 10^{\square a}$ ,  $m = h^T a$  with h = (1, 1) and  $a = (a_1, a_2)$ , and L = 1 for all different  $\square a$ ). The solid red point is  $a^T = (5, 5)$  and the solid red line denotes the solution set,  $h^T a = x \square E$ .

understood as sampling the index j from  $1, \ldots, J_n$ . Note that  $f_j^n$  only relates to the jth row of  $X_n$ , where each row of  $X_n$  is a mode-n fiber. Thus, sampling  $f_j^n$  is equivalent to sampling rows of  $X_n$ . Denote  $f_n^n = 1, \dots, J_n$  as the index set of the sampled fibers. Then, the sampled version of Problem (8) becomes

$$\min_{\boldsymbol{A}_n \in \boldsymbol{A}_{n}} \frac{1}{|\mathcal{F}_n|} \prod_{j=1}^{|\mathcal{F}_n|} \prod_{i=1}^{|\mathcal{F}_n|} f(\hat{\boldsymbol{X}}_{n(j,i)}, \hat{\boldsymbol{H}}_{n(j,:)}, \hat{\boldsymbol{A}}_{n(i,:)^T}),$$

where  $\hat{\mathbf{X}}_n = \mathbf{X}_n(\boldsymbol{\mu},:)$  and  $\hat{\mathbf{H}}_n = \mathbf{H}_n(\boldsymbol{n},:)$ . Usually, the size of  $\boldsymbol{F}_n$  is chosen to satisfy  $\boldsymbol{F}_n = \boldsymbol{J}_n$ . This may help substantially reduce the amount of data samples worked with (and thus the computational workload) in each iteration. In addition, the loss term in the finite summation of (9) composes the nonlinear function f and the linear mapping  $\hat{\mathbf{H}}_n(j,:)\mathbf{A}_n(i,:)^T$ . Fiber sampling preserves such special composition form always taking summation over all i  $[1,2,\ldots,I_n]$ , i.e., one fiber relates to all entries of  $\mathbf{A}_n$ . This helps to derive a lightweight update of  $\mathbf{A}_n$  in Section III-C.

We also note that if f is the Euclidean loss, the subproblem (9) is a (constrained) least squares problem, which is convex (if f is a convex set) and can be relatively easily solved [25], [30]. However, when f is a non-Euclidean loss, the problem in (9) may still be nonconvex (e.g.,  $\Box$ divergence with  $\Box$ < 1) and challenging. One hopes to solve (9) using economical updates, which is often an art when non-Euclidean losses are considered.

## C. Block-Wise Approximation Via Bregman Divergence

Note that all the loss functions f(x, m)'s given in Table II can be decomposed into a convex part  $\hat{f}(x, m)$  plus a concave part  $\hat{f}(x, m)$  as  $f(x, m) = \hat{f}(x, m) + \hat{f}(x, m)$ . To see how to make use of such a convex-concave property together with the linear composition form to derive a lightweight update of  $\mathbf{A}_n$  from (9), we first give the definition of the Lipschitz-like convexity below.

Definition 2 (Lipschitz-like Convexity): The Lipschitz-like

convexity condition for a function pair ( ) if that

 $\Box L > 0$  such that  $L \Box f$  is convex on int dom where f and  $\Box$  are differentiable convex functions.

A more general definition of the Lipschitz-like convexity condition can be found in [36], which is defined for the so-called Legendre function  $\square$  and lower semicontinuous convex function f. This more general definition contains Definition 2 as a special case

Next, let us consider the (i, j)th component in (9). For notation simplicity, denote it as  $f(x, \mathbf{h}^T \mathbf{a})$ , where we have defined  $x = \hat{\mathbf{X}}_n(j, i)$ ,  $\mathbf{h}^T = \hat{\mathbf{H}}_n(j, :)$ , and  $\mathbf{a} = \mathbf{A}_n(i, :)^T$ . The following lemma constructs a strongly convex surrogate function of  $f(x, \mathbf{h}^T \mathbf{a})$  ( $a_r$  and  $h_r$  are the r-th components of  $\mathbf{a}$  and  $\mathbf{h}$  respectively):

Lemma 1: Let  $\square$ .): R R be a strongly convex function. Suppose that for the function pair (  $\not$ E) the following Lipschitz-like convexity condition holds:

$$\Box r, L (a_r) \Box \Box \not = x, \frac{h_r}{\Gamma_r} a_r \quad \text{is convex w.r.t. } a_r, \qquad (10)$$

where  $\Box := \frac{h_r \bar{a_r}}{h^T \bar{a_r}}$   $\Gamma$  and  $L, h_r, \bar{a_r} > 0$ , r, are constants. Then, the following holds:

$$f(x, \mathbf{h}^T \mathbf{a}) \square f(x, \mathbf{h}^T \bar{\mathbf{a}}) + (\square f(x, \mathbf{h}^T \bar{\mathbf{a}}), \mathbf{a} \square \bar{\mathbf{a}} \pm + LD_{\square}(\mathbf{a}, \bar{\mathbf{a}}),$$
(11)

where  $D_{\square}(\boldsymbol{a}, \bar{\boldsymbol{a}})$  is the Bregman divergence. The equality in (11) holds if  $\boldsymbol{a} = \bar{\boldsymbol{a}}$ .

*Proof:* See Appendix A in supplementary material.

The upper bound constructed in Lemma 1 is reminiscent of the majorization-minimization (MM) scheme developed for  $\square$  divergence [28]. The difference here is the choice for  $\square$  ). The MM scheme in [28] suggested to choose  $\square$  ) to be the convex part of  $f(\cdot)$  (up to a constant scaling) while Lemma 1 suggests to choose  $\square$  ) to 'fit' the geometry of the convex part of  $f(x, \mathbf{h}^T \mathbf{a})$ . Clearly, the condition in (10) is more general. The corresponding functions ( $\square$  ),  $\vec{F}(x, \cdot)$  ) are referred to as the Lipschitz-like convexity function pair [36].

Lemma 1 suggests that if one can find the appropriate  $\square$ ( ), then problem (9) can be approximately solved via the following SMD update:

$$\mathbf{A}_{n}^{t+1} = \arg\min_{\mathbf{A}} (\hat{\mathbf{G}}_{n}^{t} \mathbf{A} \square \mathbf{A}_{n}^{t} \pm \frac{1}{\square} D_{\square}(\mathbf{A}, \mathbf{A}_{n}^{t}), \quad (12)$$

where  $\mathbf{G}_n$  is the gradient estimation w.r.t.  $\mathbf{A}_n$  using sampled data at iteration t and  $D_{\square}(\mathbf{A}, \mathbf{A}^t_n)$  is the Bregman divergence

TABLE III Explicit Forms of  $\hat{\pmb{G}}_n^t$  for Different  $\pmb{\pounds}(\cdot)$ 

Loss Function $\ell(\cdot)$	$ \mathcal{F}_n  I_n \cdot \hat{G}_n^t$
$\frac{1}{2}(x-m)^2$	$\left(\hat{oldsymbol{X}}_{n}^{t}-\hat{oldsymbol{X}}_{n} ight)^{T}\hat{oldsymbol{H}}_{n}$
$\frac{\pi}{m+\epsilon} + \log(m+\epsilon)$	$\left[\left(\hat{X}_{n}^{t} - \hat{X}_{n} + \epsilon\right) \otimes \left(\hat{X}_{n}^{t} + \epsilon\right)^{-2}\right]^{T} \hat{H}_{n}$
$m - x \log(m + \epsilon)$	$\left[-\hat{oldsymbol{X}}_{n} \circledast \left(\hat{oldsymbol{X}}_{n}^{t} + \epsilon\right)^{-1} + 1\right]^{T} \hat{oldsymbol{H}}_{n}$
$e^m - xm$	$\left[\exp(\hat{X}_n^t) - \hat{X}_n\right]^T \hat{H}_n$
$\log(m+1) - \pmb{x} \log(m + \pmb{\epsilon})$	$\left[\left(\hat{\boldsymbol{X}}_{n}^{t}+1\right)^{1-1}-\hat{\boldsymbol{X}}_{n}\circledast\left(\hat{\boldsymbol{X}}_{n}^{t}+\epsilon\right)^{1-1}\right]^{T}\hat{\boldsymbol{H}}_{n}$
$\log(1+e^{m})-xm$	$\left[\exp(\hat{\boldsymbol{X}}_n^t) \circledast \left(\exp(\hat{\boldsymbol{X}}_n^t) + 1\right)^{t-1} - \hat{\boldsymbol{X}}_n\right]^T \hat{\boldsymbol{H}}_n$
$\beta$ -divergence	$\left[ (\hat{\boldsymbol{X}}_n^t + \epsilon)^{\cdot \beta - 2} \circledast (\hat{\boldsymbol{X}}_n^t - \hat{\boldsymbol{X}}_n + \epsilon) \right]^{T^*} \hat{\boldsymbol{H}}_n$

<sup>\*</sup>In this table,  $\hat{\mathbf{X}}_{n}^{t} = \hat{\mathbf{H}}_{n} \mathbf{A}_{n}^{t}$ 

between  $\mathbf{A}_n$  and  $\mathbf{A}_n^t$  For different f(x, m) in Table II, the expressions of  $\hat{\boldsymbol{G}}_n^t$  are summarized in Table III. The form of the gradient for general losses f(x, m) was recently presented in [6, Theorem 3].

Compared with (7), the SMD step in (12) is only for the block variable  $A_n$ . The term  $A_n$  in  $A_n$  can use regarded order information and  $D_{\square}(A_n, A_n)$  in  $A_n$  can use regarded order

geometry-aware augmentation for properly approximating the loss in (9). If choosing L in (10) as  $\frac{1}{L}$  such that (10) holds, then (12) minimizes an upper bound function for the loss in (9). In addition, considering  $\Box(a) = \frac{1}{2} \int_a^a a^2$ , then (12) reduces to the well-known projected SGD.

Remark 3: For non-Euclidean losses, the choice of  $\square$  ) can heavily affect the behavior of the algorithm. On the other hand, the flexibility of using different  $\square$  )'s also entails opportunities of developing fast non-Euclidean CPD algorithms. An illustrative example using the generalized KL loss is shown in Fig. 2. One can see that using  $\Box(a) = a^2/2$ , the progress from  $a^t$  to  $a^{t+1}$ is very small. However, by using  $\Box(a)$ 's that are more adapted to the cost function's geometry (reflected by the contour of the cost function), i.e.,  $\square(a) = a \log a$  or  $\log a$ , the progress in one iteration can be much larger2.

# D. Stochastic Mirror Descent for CPD

The SMD step in (12) specifies the update for one latent matrix  $\mathbf{A}_n$ . Combining it with a random selection of block n, the

proposed algorithm is summarized in Algorithm 1. The major advantage of using uniformly at random block selection is that such a scheme leads to an unbiased gradient estimation [25] (up to a constant scaling), which simplifies the convergence analysis. In essence, the proposed algorithm is a block-randomized (inexact) coordinate descent method, which admits a similar structure as the algorithm in [25] for CPD under the Euclidean loss. The key difference is that Algorithm 1 employs SMD to solve each subproblem inexactly, while the algorithm in [25]

# Algorithm 1: Stochastic Mirror Descent (SMD) Algorithm.

Algorithm 1: Stochastic Mirror Descent (SMD) Algorithm.

Require: 
$$X, \mathbf{A}_1^0, \mathbf{A}_2^0, \dots, \mathbf{A}_N^0, \quad \{\Box_t^1 t = 0, 1, \dots \}$$

1: for  $t = 0, 1, \dots$ , until meet some convergence criteria do

2: Uniformly sample  $n \quad \{\Box_t, 2, \dots, N\}$ ;

3: Uniformly sample fibers  $F_n \quad \Box \{1, 2, \dots, J_n\}$ ;

4: Compute the sampled gradient  $\mathbf{G}_n^t$ ;

5:  $\mathbf{A}_n^{t+1} = \arg\min_{n \in \mathbb{N}} \mathbf{A}_n \mathbf{A}_n$ 

uses proximal gradient. For simplicity, we name it as Stochastic Mirror descent AlgoRiThm for CPD (SmartCPD).

Remark 4: If no data sampling or block sampling is considered, the full batch version of SmartCPD subsumes many existing non-Euclidean and Euclidean matrix/tensor decomposition algorithms as its special cases—see the algorithms in [11], [14], [15], [17], [25], [28], [46]. In particular, consider the KL divergence. If one chooses  $\square(a) = \square \log a$  and choose the step size  $\square$  properly, SmartCPD becomes the MM

algorithm [15]; when one uses  $\Box(a) = a \log a$ , then SmartCPD becomes the MD algorithm developed in [11]. This connection is not surprising, since MD includes many first-order approaches as its special cases. Nonetheless, using this connection, our convergence analysis (cf. Sec. IV) may also shed some light on the convergence behaviors of some existing algorithms whose convergence analyses were not considered at the time (e.g., [14]).

# E. Practical Implementation

To implement the SmartCPD algorithm in practice, several key aspects need to be considered carefully. In particular, as in all stochastic algorithms, the step size-related parameter selection (i.e.,  $\square$  and L in SmartCPD) needs to be carefully carried out. In addition, the  $\square$ () function should be chosen judiciously. In this subsection, we discuss these aspects in detail.

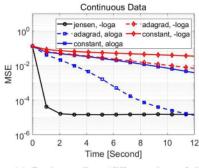
1) Choice of  $\square(\cdot)$ : As indicated by Lemma 1,  $\square(\cdot)$  should be chosen to adapt to the geometry of  $f(\cdot)$ —e.g., by setting  $\Box(\cdot)$  to be the convex part of  $f(\cdot)$ , if  $f(\cdot)$  has convex-concave structure. In addition, the update (12) needs to solve a subproblem which minimizes the constructed surrogate loss over constraint set  $A_n$ . Practically, it is desirable that this subproblem can be solved easily, preferably in closed form. Hence, the choice of  $\square$ should be an integrated consideration of the function geometry of f and the constraint  $\neg n$ . For example, consider f(x, m) = $x \log m$  with  $\mathbf{A}_n$  being the probability simplex constraint, i.e.,  $A \times 1 \boxtimes 1$ ,  $A_n(i,j)$  [f], i,j where  $1 \boxtimes$  enotes a vector (of proper length) with all elements being 1. Then, it is preferred to choose  $\square(a) = a \log a$  rather than  $\square(a) = -\log a$  since the former admits a closed-form solution of the MD update, which is also known as the exponential gradient descent or entropic

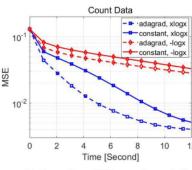
which admit closed-form solutions are summarized in Table IV.

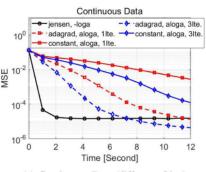
under various—n's

descent [29]. Several other examples of

<sup>&</sup>lt;sup>2</sup>In this illustrative example, we used L = 1. This is only for the ease of illustration. The L that satisfies the condition in (10) should be larger. Although the geometries are similar under L > 1 and L = 1, the former case is harder to visualize due to the condensed contours.







(a) Continuous Data (diff. step size and  $\phi$ )

(b) Count Data (diff. step size and  $\phi$ )

(c) Continuous Data (diff. no. of ite.)

Fig. 3. Numerical examples on  $100\ 100\ 100$  tensors with rank R=10, where the generalized KL divergence is selected as the loss function and non-negative constraints are considered. For SmartCPD, the number of fibers is 2R=20,  $b=10^{\square_3}$  for adagrad step size (see (13)) and  $\square=1$  for constant step size. The averaged mean squared error (MSE, see definition in (17)) of the latent matrix (averaged over 20 independent trials) is used as performance metric. The latent matrices are drawn from i.i.d. uniform distribution between 0 and 1. For continuous data case (Fig. 3(a) and (c)), Gaussian noise with SNR = 40 dB (see definition in Section V-A4) is added to the data; in count data case (Fig. 3(b)), each tensor entry is drawn from Poisson distribution with identity link function, i.e.,  $\square(m) = m$ .

TABLE IV  ${\it Examples for Pair ($A_{\it R}$) and Closed-Form Solution of (12) }$ 

$\phi(\cdot)$	$A_n$	Closed-form Solution
$-\log a$	non-negative	$egin{aligned} oldsymbol{A}_n^t \circledast \left[ \hat{oldsymbol{G}}_n^t \circledast \left( oldsymbol{A}_n^t \oslash L  ight) + 1  ight] \ oldsymbol{A}_n^t \circledast \exp(-\hat{oldsymbol{G}}_n^t \oslash L) \end{aligned}$
$a \log a$	non-negative	$m{A}_n^t \circledast \exp(-\hat{m{G}}_n^t \oslash L)$
$a^c$ $(c > 1 \text{ or } c < 0)$	non-negative	$\left[ (oldsymbol{A}_n^t)^{\cdot (c-1)} - rac{1}{c} \hat{oldsymbol{G}}_n^t \otimes L  ight]^{\cdot rac{1}{c-1}}$
$\frac{a \log a}{a^2}$	simplex many forms	$\operatorname{colnorm}(\boldsymbol{A}_n^t \circledast \exp(-\hat{\boldsymbol{G}}_n^t \otimes L))$ refer [25], [47]

\*In this table, colnorm( $\mathbf{A}$ ) =  $\mathbf{A}$  0 1( $\mathbf{A}^T$ 1) $^T$ , which denotes the column-wise normalization operation.

2) Choice of  $\square$ : Since the Bregman divergence  $D \sqsubseteq (A, A^t)$  is defined in an entry-wise summation form, scaling  $D \sqsubseteq (A, A^t)$  with a constant  $\frac{1}{1+t}$  for all entries may be less effective for approximating the loss in (9). In this work, we propose to use coordinate-dependent step size, i.e., different positive scaling factors  $\boxtimes_n^t (in, j)$ 's for different coordinates (in, j)'s, where  $in [In], j [R], and \boxtimes_n^t \mathbb{R}^{In \times R}$ . Consequently, the term  $\exists D \sqsubseteq (A, A_i^t)$  in (12) becomes  $\boxtimes_n^t (in, j)D \sqsubseteq (A_n(in, j), A^t(in, j))$ .

As illustrated in Fig. 3, such a coordinate-specific  $\boxtimes^t$  rempirically helps accelerate convergence. In what follows, we provide two

schemes on choosing  $\boxtimes_t^k(i_n,j)$ 's.

Jensen's inequality based choice: One important criteria is to choose  $\boxtimes_t^t(i_n,j)$  such that the inequality in (11) holds for each coordinate  $(i_n,j)$ . Such a choice makes the local approximation in (12) be an upper bound of (9). In the case where  $f(\cdot)$  enjoys the convex-concave property and the constraint sets  $A_n$  satisfy  $A_n = \mathbb{R}_+^{I_n \times R}$ ,  $I_{\square}$  Lemma 1 implies that one can choose  $\square$  to be the convex part of  $f(\cdot)$ . Consequentially,  $\boxtimes_t^t(i_n,j)$  can be derived based on the Jensen's inequality in each iteration. A number of examples of  $\boxtimes_t^t(i_n,j)$  with respect to different  $f(\cdot)$  are given in Table V, where the derivation are based on (18) and (19) in supplementary material. We note that the Jensen's inequality based  $\boxtimes_t^t(i_n,j)$  choice is often used by popular methods such as MM [48] and block successive upper bound minimization (BSUM) [49] to simplify the solution of the subproblems. The

TABLE V
(  $\Box_n^t$ ) Based on Jensen's Inequality

Loss Function	$\phi(a)$	$ \mathcal{F}_n I_n\cdot\Gamma_n^t$
Eucl. Dis.	$\frac{1}{2}a^{2}$	$rac{1}{2}(\hat{m{X}}_n^t)^T\hat{m{H}}_n\otimesm{A}_{n_c}^t$
IS Div.	$\frac{1}{a}$	$A_n^t \circledast \left[\hat{oldsymbol{X}}_n \oslash (\hat{oldsymbol{X}}_n^t)^{-2} ight]^T \hat{oldsymbol{H}}_n$
KL Div.	$-\log a$	$(oldsymbol{A}_n^t)^T \otimes \left[ \hat{oldsymbol{H}}_n^T (\hat{oldsymbol{X}}_n \otimes (\hat{oldsymbol{X}}_n^t))  ight]$
$\beta$ -Div. $(\beta > 1)$	$a^{\beta}$	$rac{1}{3}(\hat{m{X}}_n^t)^{ar{eta}-1})^T\hat{m{H}}_n\otimes (m{A}_n^t)^{ar{eta}-1}$
$\beta$ -Div. ( $\beta < 1$ )	$a^{\beta-1}$	$\frac{1}{(1-\beta)}[\hat{oldsymbol{X}}_n\otimes(\hat{oldsymbol{X}}_n^t)^{\cdot\beta-2}]\hat{oldsymbol{H}}_n\oslash(oldsymbol{A}_n^t)^{\cdot\beta-2}$

\*In this table,  $\hat{X}_{n}^{t} = \hat{H}_{n} A_{n}^{t}$ 

difference here is that our construction uses sampled fibers instead of the whole data. Such Jensen's inequality-based step-size choice also works well with certain non-Euclidean losses under our stochastic settings, as will be seen in the experiments.

• Adaptive step size based choice: The deep learning community has developed a number of effective adaptive step size scheduling methods, e.g., the Adagrad [50] and Adam type schemes [51]. These schemes typically exploit the past iterations' gradient information to scale the current sampled gradient in a coordinate-wise manner. The upshot of these methods is that they often require very small amount of step size tuning, yet offer

highly competitive empirical performance; also see theoretical understanding in [51]. Under the Bregman divergence, the adaptive step size schemes can be used to 'schedule'  $\sum_{i=1}^{t} (i_n, j)$  in each iteration. Specifically, we propose the following Adagrad step size rule for  $\sum_{i=1}^{t} h(i_n, j)$ ,

where  $\hat{\boldsymbol{G}}_{n}^{t}(i_{n}, j)$  denotes the  $(i_{n}, j)$ -th entry of  $\hat{\boldsymbol{G}}_{n}^{t}$  and b > 0 is a constant to ensure  $\bigotimes_{i=1}^{t} (i_{n}, j) > 0$ . Note that the adaptive step size was considered in fiber-sampling based stochastic Euclidean CPD in [25] and entry-sampling based stochastic non-Euclidean CPD [7], and encouraging results were observed in both cases.

Remark 5: It can be observed from numerical examples in Fig. 3 that, for continuous data (Fig. 3(a)), using different  $\sum n(in,j) > 0$  for different coordinate (in,j) exhibits much faster convergence behavior than using a constant step size

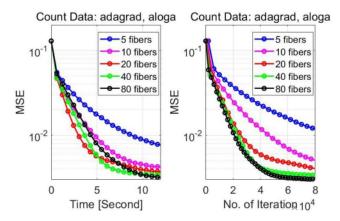


Fig. 4. Averaged MSE of count tensor with different number of fibers ( $I_n = 100$  and R = 10).

for all coordinates. For count data (Fig. 3(b)), we have similar observations. Another observation from Fig. 3 is that, for count or binary data, the Jensen's inequality based step size scheme may be less competitive, especially when the data contains many zeros. The zero entries may make  $\sum_{n}(i_n,j)$  very small, thereby causing numerical issues. The 'adagrad' scheme in (13) has empirically much more stable convergence behavior under such circumstances. In practice, a user could make a choice between a more conservative algorithm with better convergence understandings (constant step size  $\square$ ) or a more aggressive one with less theoretical guarantees (coordinate-wise step size  $\sum_{n=1}^{\infty} n$ ). This may be analogous to constant step size SGD or Adam/Adagrad type SGD, where the former has better convergence understanding while the latter works faster for some problems in practice.

- 3) Inner Iterations: The proposed Algorithm 1 only contains one iteration per block. Nonetheless, one can also extend it to multiple SMD updates (i.e., inner iterations) per block. We have observed that implementing with a few more inner iterations could improve the practical convergence behavior—as shown in Fig. 3(c). There are two ways of having multiple inner iterations. The first way is to repeat lines 3-5 in Algorithm 1 to update  $\mathbf{A}_n$  several times before moving to the next block, where the fibers are re-sampled for each inner iteration. The second way is that, for fixed sampled fibers, repeat lines 4-5 for multiple times. Both methods work reasonably well in practice. In Fig. 3(c), we use the latter because it uses less samples and has similar performance (in our experience) compared with re-sampling per inner iteration.
- 4) Number of Fibers: By regarding f(x, m) as a certain 'distance' measure between x and m, the SMD step (12) in SmartCPD can be understood as approximately solving the linear system  $\hat{X}_n = \hat{H}_n A_n^T$  under such 'distance' measure. From this viewpoint, the number of fibers should be chosen to be larger than rank R, e.g.,  $p^n = 2R$ . In Fig. 4, with the same problem setting as in Fig. 3(b), SmartCPD (adagrad step size,  $\Box(a) = a \log a$ , 1 inner iteration) with different numbers of fibers are compared. Clearly, increasing number of fibers improves per iteration efficiency but also needs more time to compute the sampled gradient. Using 2R fibers seems to strike a good balance between the gradient estimation's variance and the per-iteration complexity.

### IV. CONVERGENCE ANALYSIS

In general, convergence guarantees of stochastic tensor decomposition algorithms are difficult to establish, as nonconvex constrained optimization problems are intrinsically harder to analyze under stochastic settings. The non-Euclidean version is even more so, since the sampled subproblems may still be nonconvex. There has not been an analytical framework for SMD based non-Euclidean tensor/matrix factorization. The nonconvex block SMD in the optimization literature [31] is the most closely related to our algorithmic framework. However, the convergence analysis there does not cover the proposed algorithm. Specifically, the convergence of the algorithm in [31] hinges on some special incremental block averaging steps, which is not used in our algorithm. More importantly, the algorithm in [31] requires that the block-wise gradient estimation error vanishes to zero when the iterations progress. This may require implementing the algorithm with an active variance reduction technique, e.g., increasing the batch size in each iteration [25], which is not entirely realistic, and it is somewhat against the purpose of using stochastic algorithms.

Our goal is to offer tailored convergence analysis for SmartCPD that does not rely on conditions like incremental block averaging or vanishing gradient estimation error. We note that for constrained problems, convergence analysis for SMD with adaptive step size  $\boxtimes h$  is very challenging. Theoretical understanding of SGD with adaptive step size scheme was recently discussed in [51]—but the SMD case is still an open problem. In this work, we focus on the case where  $\boxtimes_{n}^{t}(i_{n}, j) = \square > 0$  are all identical.

By using Bregman divergence as a proximity measure, the recent notion of *relative smoothness* [41] (or called *L*-smooth adaptable [37]) can be defined as below:

Definition 3 (Relative Smoothness): A continuous differentiable function f is relatively smooth to a strongly convex funcwith  $L(0 < L <_{\square})$  if Lf is convex on int dom  $\square$ If f is convex, relative smoothness (Definition 3) becomes the Lipschitz-like convexity in Definition 2. Also, there are several different concepts which are closely related to relative smoothness, e.g., relative weak convexity [39], [40], and relative continuity [38], which were used in several recent works to analyze single-block SMD type algorithms' convergence; see [37]-[40]. Our analysis leverages the notion of relative smoothness, generalizing this prior work to cover multi-block SmartCPD. Particularly, the proof takes advantage of the multilinear low-rank tensor structure and the block-randomized fiber sampling strategy to brige the gap between the single-block and multi-block cases.

The objective function in (2) is denoted as  $F(\mathbf{A})$  with  $\mathbf{A} = (\mathbf{A}_1, \mathbf{A}_2, \dots, \mathbf{A}_N)$ . Then, Problem (2) can be re-expressed as:

$$\min_{\mathbf{A}} F(\mathbf{A}) + h(\mathbf{A}) \tag{14}$$

where  $h(\mathbf{A}) = \sum_{n=1}^{N} h_n(\mathbf{A}_n)$  and  $h_n(\mathbf{A}_n)$  is the indicator function of set n, i.e.,  $h_n(\mathbf{A}_n) = 0$  if  $\mathbf{A}_n$  and otherwise  $h_n(\mathbf{A}_n) = 0$  Our first observation is as follows:

Lemma 2: If  $\mathbf{A}^t$  nall reside in a compact set for all n, there exists  $0 < L < \frac{1}{2}$  such that for any  $\mathbf{A}^t$ ,  $\mathbf{A}$  in this compact set function F is relatively smooth to with L, i.e.,

$$|F(\mathbf{A}^t) \square F(\mathbf{A}) \square \square \square F(\mathbf{A}), \mathbf{A}^t \square \mathbf{A} \pm | \square LD \square (\mathbf{A}^t, \mathbf{A}),$$
(15)

where  $\square \cdot$ ) can be any strongly convex function. *Proof:* See Appendix B in the supplementary material.

It is easy to see Lemma 2 immediately implies inequality (11) in Lemma 1 summed over all the entries. Also, Lemma 2 has no restriction on the choice of  $\square(\cdot)$ . This brings much convenience for the practical usage of Algorithm 1. More importantly, Lemma 2 holds for the all latent factors simultaneously (as opposed to one latent factor as in Lemma 1). This is critical for establishing the convergence guarantee. Also, we note that inequality (15) is a generalization of standard Lipschitz-continuous gradient property of  $F(\mathbf{A})$  under Bregman divergence, which is also known as *relative smoothness* [52]. We remark that the assumption that  $\mathbf{A}^t_n$  live in a compact set is not easy to check in advance, yet it is not hard to satisfy. When the constraints  $\mathbf{A}$  are compact sets, then this assumption is naturally satisfied per the defined update in (12). Even if  $\mathbf{A}^n$  is not bounded, unbounded iterates are rarely (if not ever) observed in our extensive numerical experiments.

Based on Lemma 2, convergence of the proposed Algorithm 1 can be guaranteed for any strongly convex function  $\square$  ) (see Theorem 1). Our convergence analysis starts by using the following "reference" function in each iteration:

$$L^{(\boldsymbol{A};\boldsymbol{A}^t)} \coloneqq F(\boldsymbol{A}) + h(\boldsymbol{A}) + \frac{1}{2\square} D_{\square}(\boldsymbol{A},\boldsymbol{A}^t)$$

where  $0 < \square < \frac{1}{2L}$  is a constant such that  $L(\mathbf{A}; \mathbf{A}^t)$  is strongly convex in  $\mathbf{A}$ . The minimal value and minimizer of  $L(\mathbf{A}; \mathbf{A}^t)$  for a given  $\mathbf{A}^t$ , which are also known as Bregman Moreau envelope  $L(\mathbf{A}^t)$  and Bregman proximal mapping  $L(\mathbf{A}^t)$  [53], are defined as

$$M(\mathbf{A}^t) = \min_{\mathbf{A}} L(\mathbf{A}; \mathbf{A}^t), \ T(\mathbf{A}^t) = \arg\min_{\mathbf{A}} L(\mathbf{A}; \mathbf{A}^t).$$

Denote  $\hat{\mathbf{A}} = \{\mathbf{A}\}$  for a given  $\mathbf{A}$ , we use the following lemma to show that  $D_{\square}(\hat{\mathbf{A}}, \mathbf{A})$  is a measure for attaining stationary points of Problem (14).

Lemma 3:  $\mathbf{A}$  is a stationary point of Problem (14), i.e.,  $\mathbf{0} \boxtimes \mathbb{D} F(\mathbf{A}) + \partial h(\mathbf{A})$ , where  $\partial h(\mathbf{A})$  denotes the subgradient, if and only if  $D \square (\hat{\mathbf{A}}, \mathbf{A}) = 0$ .

*Proof:* see Appendix C in the supplementary material. ☐ The key step for analyzing Algorithm 1 is to quantify its one iteration behavior:

*Lemma 4:* Suppose that  $\{A^t\}$  all reside in a compact set for

all n, and that  $\square$  is a  $\square$  strongly convex function. Let  $\boldsymbol{A}$  t+1 be generated by Algorithm 1 at iteration t, then we have

$$\mathsf{E} M(\mathbf{A}^{t^{+1}})$$

$$\square M(\mathbf{A}^{t}) \square_{C_{1}} \square_{D_{\square}} (\hat{\mathbf{A}}^{t}, \mathbf{A}^{t}) + c_{2} \square \mathbb{E} I \hat{\mathbf{G}}^{t} I^{2}, \quad (16)$$

where  $c_1 = \frac{(1 \square 2)}{4 \square^2} \stackrel{I}{\text{and}} c_2 = \frac{1}{4 \square}$  The expectation in (16) is taken over the random variable responsible for fiber sampling in iteration t.

Proof: See Appendix D in the supplementary material. ☐
The lemma implies that, in expectation, if the step size ☐
is properly chosen, then (X\*) decreases after every iteration.
Using Lemma 4 as a stepping stone, we show our main result:
Theorem 1: Suppose that the assumptions in Lemma 4 hold.
1) for diminishing step size ☐ satisfying ☐ ☐ ☐ ☐ and
☐ ☐ ☐ ☐ ☐ ☐ Algorithm 1 converge to a stationary point in expectation,

$$\liminf_{t \cap \square} \mathsf{E} \ D_{\underline{\mathsf{I}}}(\hat{\boldsymbol{A}}^t, \boldsymbol{A}^t) = 0;$$

2) for a constant step size  $\Box = \Box_{\overline{T}}$ ,  $\frac{1}{1 \cdot \overline{T}}$  stationary solution of Problem (14) can be obtained by Algorithm 1 within T iterations,

$$\min_{1 \square \Gamma \cup T} \mathsf{E} \ D_{\square}(\hat{\boldsymbol{A}}^t, \boldsymbol{A}^t) \ \square C / \ T,$$

where C > 0 is a constant. The expectations are taken over the fiber and block sampling in all iterations jointly.

*Proof:* See Appendix E in the supplementary material.  $\Box$ 

## V. SIMULATIONS

We use synthetic and real data experiments to showcase the effectiveness of SmartCPD. The MATLAB code is available at https://github.com/WQPu/SmartCPD.git.

## A. Synthetic Data

We evaluate the numerical performance on different types of synthetic data, i.e., continuous, count, and binary data, under various non-Euclidean losses in Table II (c in Table II is set to be  $10^{\square b}$ ).

1) Baselines and Performance Metric: We use two recent competitive algorithms as our main baselines. The first one is an entry-sampling based stochastic non-Euclidean CPD optimization algorithm, namely, GCP-OPT, proposed in [7]. The second baseline is the generalized Gauss-Newton (GGN) method for non-Euclidean CPD [9]. The GCP-OPT method is implemented in Tensor Toolbox [54] and 'Adam' is selected as the optimization solver. The sampling rule of GCP-OPT is the default 'uniform' setting for dense tensors unless we specify it for particular examples. The GGN method [9] is implemented by nlsb\_gndl and the 'preconditioner' is set as 'block-Jacobi'; see more details in the Tensorlab toolbox [55]. For the proposed SmartCPD, different choices of step size \(\sigma\_n\) and function \(\sigma\_n\) will be specified according to data type.

We generate third-order tensors with different sizes and ranks, each dimension of the tensor keeps the same as  $I = I_1 = I_2 = I$ . For the two stochastic algorithms, in each iteration, SmartCPD samples 2 R fibers<sup>3</sup> and GCP-OPT samples  $2I_nR$  entries. This makes the two algorithms use the same number of samples per iteration, though it differs from the recommendation in [7, Section 5.1]. The other hyperparameters of the baselines follow their default settings. The MSE of the latent matrices is used as a performance metric [25], which is

<sup>&</sup>lt;sup>3</sup>Note that that using 2R is not the only choice, but based on our experience.

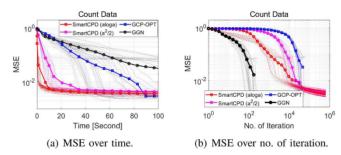


Fig. 5. MSE of count tensor (I = 100 and R = 20). defined as  $MSE = \frac{1}{N} MSE_n$  with N = 1

$$MSE_{n} = \min_{\square(r) \mid \mathbb{M} \mid_{r=1}}^{\mathbf{I}_{R}} \frac{1}{1} \underbrace{\mathbf{A}_{(:, \square(r))}}_{n} \qquad \frac{\hat{\mathbf{A}}_{(:, r)}}{I \hat{\mathbf{A}}_{n}(:, r)I^{1}} \stackrel{1}{\longrightarrow} \frac{1}{I \hat{\mathbf{A}}_{n}(:, r)I^{1}},$$
(17)

where  $\hat{A}_n$  denotes the estimated  $A_n$  and [1], ..., [R] represents a permutation of the set [R] = [1], ..., R, which is used to fix the intrinsic column permutation in CPD.

2) Count Data: We first evaluate the performance on count data tensors. For the proposed SmartCPD,  $\sum_n^t f_n$  is scheduled following (13). We use  $\sum_n^t f_n (a) = a \log a$  and  $\sum_n^t f_n (a) = a \log a$  in the SMD update. The constant b in (13) is  $10^{-15}$  for all simulation trials

We use the loss function  $f(x, m) = m_{\perp} x \log(m + c)$  as in [15] and draw the latent matrices  $\mathbf{A}_1$ ,  $\mathbf{A}_2$ , and  $\mathbf{A}_3$  from i.i.d. uniform distribution between 0 and  $A_{\max}$ , where  $A_{\max} = 0.5$  is a positive constant. For each column of the latent matrices, 5% elements are randomly selected and replaced by i.i.d. samples from uniform distribution between 0 and  $10A_{\max}$ . This way, the elements have more diverse scales. The observed count data tensor  $\mathbf{X}$  is generated following the Poisson distribution, i.e.,  $\mathbf{X}_i$  Poisson( $\mathbf{M}_i$ ). We set  $\mathbf{A}_n$  for all n as the nonnegative orthant.

Fig. 5(a) shows the performance of the algorithms under I = 100 and R = 20, where the solid lines correspond to the average MSEs and dashed lines are for individual trials. One can see that SmartCPD improves the MSE quickly in all trials. On average, it brings the MSE below 10<sup>□</sup> using less than 10 seconds, while the best baseline, i.e., GCP-OPT takes at least around 40 seconds to reach the same MSE level. Fig. 5(b) shows the MSEs of the algorithms against the number of iteration. Clearly, GGN has the best iteration complexity since it exploits second-order information based on all samples. Both SmartCPDs enjoy better iteration complexity than GCP-OPT. In addition, SmartCPD ( $a \log a$ ) performs better than SmartCPD  $(a^2/2)$  (which is equivalent to the BrasCPD algorithm in [25]). GCP-OPT achieves smaller MSE for some trials. Note that both SmartCPD  $(a^2/2)$  and GCP-OPT can be viewed as SGD type algorithms with fiber sampling and entry sampling, respectively.

In many trials, SmartCPD uses at least one order of magnitude fewer data entries to reach MSE=  $10^{\square}$ . Fig. 6 shows a closer look at the runtime breakouts of the two stochastic algorithms. It can be found that both use similar time for sampling but GCP-OPT requires slightly more time for computing the gradient as well as updating all latent matrices. These

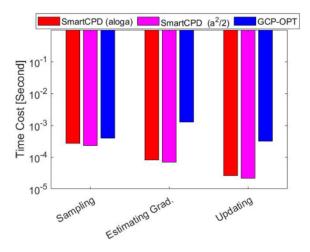


Fig. 6. Detailed time cost per iteration of SmartCPD and GCP-OPT algorithms (I = 100 and R = 20).

TABLE VI
MSE AFTER 100 SECONDS (I = 100) AND DIFF. R)

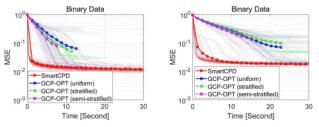
Rank $R$		10	20	50
SmartCPD $(a \log a)$	mean	2.7E-3	3.7E-3	6.6E-3
	median	2.7E-3	3.7E-3	6.5E-3
GCP-OPT	mean	9.6E-3	0.028	0.408
	median	2.3E-3	3.3E-3	0.406
GGN	mean	0.023	0.031	0.288
	median	1.9E-3	6.5E-3	0.281

results clearly show that using non-Euclidean proximal term (e.g.,  $a \log a$ ) and fiber sampling both play important roles in accelerating KL-divergence CPD.

In Table VI, the achieved MSEs after 100 seconds when I=100 and R=10, 20, 50 are compared. It can be observed that SmartCPD ( $a \log a$ ) outputs MSEs that are smaller than  $10^{\square}$  for all cases, while GGN and GCP-OPT could not reach MSE 0.1 in 100 seconds when R=50, in terms of both mean and median. Also, GCP-OPT and GGN can often attain a slightly smaller MSE value relative to SmartCPD when the rank is not large, e.g., R=10.

3) Binary Data: Next, we evaluate the performance on binary tensors. We use the loss function that is related to the MLE of Bernoulli tensors, i.e.,  $f(x, m) = \log(m + 1) - x \log(m + 1)$ c) (cf. Table II). We set  $\Box(a) = a \log a$  and non-negativity constraints are considered. Note that GGN is developed for □divergence, and thus cannot be used for this loss function. Hence, we only use GCP-OPT to benchmark our algorithm in the binary case. We set I = 100 as before and each entry of the binary tensor is generated from the Bernoulli distribution, i.e.,  $\underline{X}_i = 1$  with probability  $\underline{M}_i/(1 + \underline{M}_i)$ . The latent matrices are generated as in Sec. V-A2. First, we set R = 5, 10 and the constant  $A_{\text{max}} = 0.3$ . This results in the generated binary tensor having about 4% (R = 5) and 8% (R = 10) non-zero entries, respectively. For such sparse binary tensors, 'stratified' and 'semi-stratified' sampling rules in GCP-OPT are also included in the comparison.

Fig. 7 shows the convergence curves (averaged over 20 independent trials) and SmartCPD achieves the fastest convergence speed. Also, GCP-OPT with 'semi-stratified' shows better



(a)  $R=5, \approx 4\%$  non-zero entries. (b)  $R=10, \approx 8\%$  non-zero entries.

Fig. 7. MSE of binary tensor  $(I_n = 100)$  and R = 5, 10).

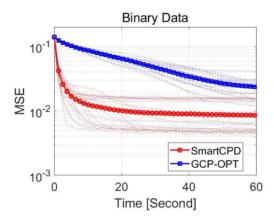


Fig. 8. MSE of binary tensor ( $I_n = 100$ , R = 20, and 15% non-zero entries).

convergence behavior than the other two sampling rules (i.e., 'uniform' and 'stratified').

In Fig. 8, we increase rank to R=20 and set the constant  $A_{\text{max}}$  0.2, 0.3, 0.4, 0.5 to generate four different binary data tensors, with about 5\%, 15\%, 28\%, 40\% nonzero entries, respectively. For R=20, we observe GCP-OPT with 'uniform' sampling rule performs more consistent than the other two sampling rules and hence, in this figure we only present the result of GCP-OPT with 'uniform' sampling rule. Fig. 8 shows the MSEs against time of 20 independent trials when the tensor has 15% nonzero entries. Similar as before, SmartCPD requires much shorter time to achieve MSF-  $10^{\text{CP}}$ . In addition, the histograms of MSEs after 60 seconds are presented in Fig. 9. One can see from there that that SmartCPD consistently outperforms the baseline, and the advantage is more articulated when the data becomes denser.

4) Continuous Data: We also evaluate the performance of SmartCPD on continuous tensors under Edivergence. We consider the multiplicative Gamma noise, i.e.,

Gamma: 
$$X_i = M_i \cdot N_i$$
,

in which  $\underline{N}_i$  is i.i.d. Gamma noise. The SNR for Gamma noise is defined the same way as in [9], which is determined by the shape parameter of Gamma distribution. The latent matrices  $A_1$ ,  $A_2$ , and  $A_3$  are drawn from the i.i.d. uniform distribution between 0 and 1 and non-negativity constraints are considered. Since the  $\square$ div. loss functions satisfy the convex-concave property, function  $\square$  and  $\square$  are used based on the Jensen's inequality for SmartCPD; see details in Table V. We find such setup

TABLE VII

MSE of Continuous Data Tensor With Gaussian Noise After 100 Seconds (I = 300, R = 20, and Different  $\square$ 

β		-1	-0.5	0	0.5	1	2	3
SmartCPD								
GCP-OPT	mean med.	9.2E 4 4.8E-5	1.5E 4 4.1E-5	1.4E 4 2.2E-5	5.8E 4 3.8E-5	4.1E 5 1.8E-5	1.7E 4 2.9E-5	1.7E 4 3.3E-5
GGN							<1E 6 <1E-6	

particularly efficient for dealing with dense and continuous tensors under \( \subseteq \divergence. \)

Fig. 10 shows simulations where the tensor has a size of 300  $\times$  300  $\times$ 300 and R=20. The average MSEs over 20 independent trials for = (Dwith SNR = 20 dB are compared. Clearly, the two stochastic algorithms, i.e., SmartCPD and GCP-OPT have faster convergence than GGN. Further, SmartCPD is even faster than GCP-OPT, especially in the beginning. We also observe that the original SmartCPD's MSE saturates at a certain level after few iterations.

Next, we test the proposed algorithm on CPD under various  $\square$ , i.e.,  $\square$  105, 0, 0.5, 1, 2, 3. The data is generated using an additive noise model, i.e.,  $X_i = M_i + N_i$ , where  $N_i$  denotes zero-mean i.i.d. Gaussian noise and M is generated as in the previous simulation.

In Table VII, the mean and median MSE over 20 independent trials after 100 seconds are compared with SNR = 40 dB. SmartCPD achieves smaller mean MSE than other two algorithms and both SmartCPD and GCP-OPT have similar median MSEs. This is because 100 seconds is not enough for GCP-OPT to ensure all trials converge; see the convergence curves for = 1 in Fig. 11. As a deterministic second-order algorithm, GGN attains better MSEs for = 2 This shows effectiveness of exploiting the second-order information for handling the Euclidean loss.

5) Column Constraints: In Fig. 12, we evaluate the performance for problems with simplex constraints, i.e.,  $\mathbf{A}^{T} \mathbf{1} \boxtimes$ **1**,  $\mathbf{A}_n(i,j)$   $\bigcirc$ , i,j and the loss function is f(x,m) = m $x \log(m+c)$  (KL div.). Tensors that represent joint probability distributions with I = 100 and R = 10 are considered; see [11], [43]–[45], [56] for details. The latent matrices  $\mathbf{A}_1$ ,  $\mathbf{A}_2$ , and  $\mathbf{A}_3$ are drawn from i.i.d. uniform distribution between 0 and 1, then each column is normalized so that it represents a probability mass function. For SmartCPD, the step size scheme in (13) with  $\Box(a) = a \log a$  is used. Note that both the baselines are not able to handle such simplex constraints, while this problem frequently arises in probabilistic tensor decomposition; see [11], [43], [44], [57]. To benchmark our algorithm, we use the deterministic block MD algorithm in [11]. In Fig. 12, one can see that SmartCPD has much faster convergence behavior than the MD algorithm due to stochastic sampling.

## B. Real Data

 Chicago Crime Data: We apply the algorithms to the Chicago Crime dataset. The dataset records crime reports in the

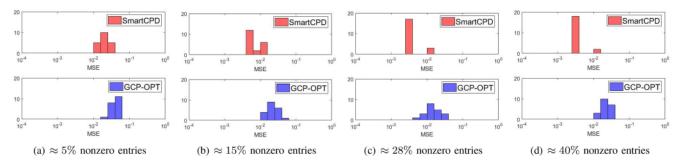


Fig. 9. Histogram of MSE after 60 seconds of  $100 \times 100 \times 100$  binary tensor (rank 20) with different level of sparsity.

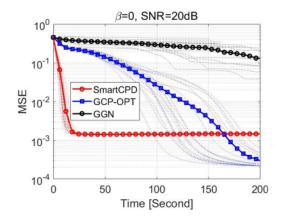


Fig. 10. MSE of continuous data tensor with multiplicative Gamma noise  $(I_n = 300, R = 20, \text{ and } = 0)$ 

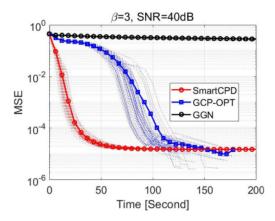


Fig. 11. MSE of continuous data tensor with additive Gaussian noise ( $I_n = 300$ , R = 20, and  $= 3\boxed{1}$ 

city of Chicago, Illinois, United States, between January 1, 2001 to December 11, 2017. The original dataset is published in the official website of the city of Chicago (www.cityofchicago.org). Here, we use the version in [58]. The data is in the form of a fourth order tensor (day<sub>x</sub> hour<sub>x</sub>crime community) with integer entries representing the number of crimes reported. The size of the tensor is 6186<sub>x</sub> 24 x77 32. It has 5,330,678 (1.5%) nonzero entries.

We choose the loss function corresponding to the Poisson distribution, i.e.,  $f(x, m) = m_{\perp} x \log(m + c)$  (see Table II), and non-negativity constraints are considered for the latent matrices. In every iteration, 40 fibers are used by SmartCPD. For

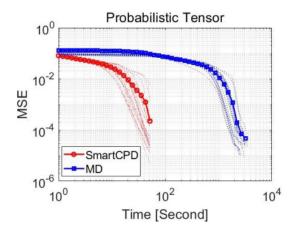


Fig. 12. MSE of continuous data tensor with simplex constraint ( $I_n = 100$ , R = 10, and  $\square = 1$ ).

GCP-OPT, 40 (6186 + 24 + 77 + 32)/4 tensor entries are sampled in every iteration, such that we sample an equal amount of data on average for both the algorithms, though it differs from the recommendation in [7, Section 5.1]. The uniform sampler of GCP-OPT is used, but note that this choice is not the GCP-OPT default for sparse tensors. All algorithms under test are stopped when the relative change in the  $\square$ divergence (in this case,  $\square$ = 1) is less than  $10^{\square}$ .

Fig. 14 shows the cost change against time of the algorithms under R=5 and R=10, respectively. Each algorithm is run for 20 trials and in each trial, the factor matrices are initialized by randomly sampling its entries from uniform distribution between 0 and 1. One can see that the proposed algorithm SmartCPD exhibits a fast runtime performance in this case. For each of the cases, the SmartCPD takes only about 5 seconds to reach a low cost value, whereas the baselines take much more time but still could not attain the same cost value. In some trials, especially when R=5, there exist some cases where GGN did not converge. More details on the algorithm-output latent factors can be found in Fig. 13.

2) Plant-Pollinator Network Data: We also consider the plant-pollinator network dataset published by [59]. The dataset consists of plant-pollinator interactions collected over 12 meadows in the Oregon Cascade Mountains, USA; also see [12] for detailed data descriptions. We extract the number of interactions between 562 pollinator species and 124 plant species over 123 days during 2011 and 2015. Hence, we form a count-type

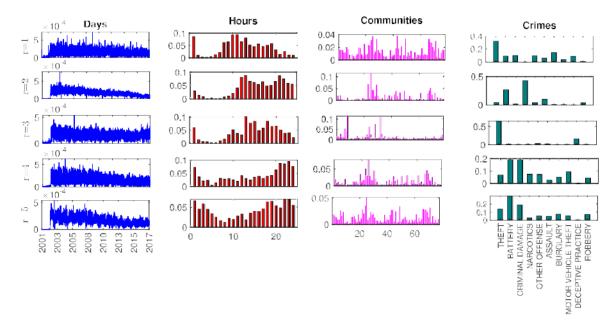


Fig. 13. Learned latent factors by SmartCPD when it reached the stopping criterion with cost value 0.043 (time= 35.81 sec.). Out of the 32 crimes reported. we chose the most frequent 10 crimes for visualization.

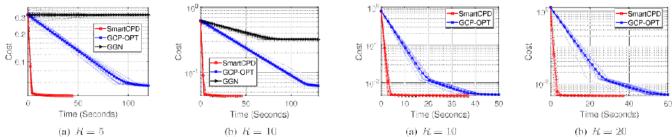


Fig. 14. Convergence of the algorithms (Count data, Chicago Crime, KL div.,  $size = 6186 \times 24 \times 77 \times 32$ 

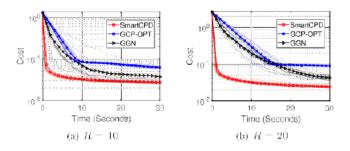


Fig. 15. Convergence of the algorithms (Count data, Plant-Pollinator network, gen. KL div., size =  $562 \times 124 \times 123$ ).

third-order tensor with a size of 562, 124, 123, where each entry represents the number of interactions between a particular plant and a pollinator on a specific day. The data is fairly sparse with only 8,370 (-0.1%) nonzero entries. We choose the loss function corresponding to the Poisson distribution, i.e.,  $f(x, m) = m_{\perp} x \log(m + c)$  along with non-negativity constraints. Other algorithm settings are also same as those used in the Chicago Crime data.

Fig. 15 plots the cost values of the algorithms against time, for different values of R. Both the baselines GCP-OPT and GGN

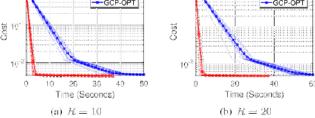


Fig. 16. Convergence of the algorithms (binary data, UCI chat network, log loss, size =  $400 \times 400 \times 196$ ).

work reasonably well in this dataset. However, the proposed SmartCPD has better performance—it converges to lower cost values compared to the other baselines. Also, one can observe that the SmartCPD is about 15 times faster in reaching low cost values compared to GCP-OPT and GGN.

3) UCI Chat Network Data: We test the algorithms on realworld binary data using a social network which contains the online interactions of the students from the University of California, Irvine, USA. The original dataset was published by [60], which includes 59,835 online messages sent between 1,899 students over 196 days from March 2004 to October 2004. We select 400 most prolific senders and form a third-order binary tensor of size 400, 400, 196 having 18862 (0.06%) nonzero entries. Each entry of the binary tensor indicates if sender i has sent a message to receiver j on the kth day. We choose the loss function corresponding to the Bernoulli distribution, i.e.,  $f(m, x) = \log(m + 1) - x \log(m + c)$  (see Table II). Other settings and parameters are as before. Since GGN method is not designed for the loss function considered in this case, it is not included.

Fig. 16 shows the cost value change against time in seconds for different values of R. Similar to the previous datasets, the proposed SmartCPD shows considerable runtime advantages over GCP-OPT. In all the trials, SmartCPD is at least 40 times faster for converging to a cost value that is later attained by GCP-OPT.

## VI. CONCLUSION

In this work, we proposed a unified SMD algorithmic framework for low-rank CPD under non-Euclidean losses. By integrating a fiber-sampling strategy within the SMD optimization technique, the proposed framework is flexible in dealing with a variety of loss functions and constraints that are of interest in real-world data analytics. By its stochastic nature, the proposed algorithm enjoys low computational and memory costs. In addition, under different data types and loss functions, we discussed a number of "best practices," e.g., step size scheduling and local surrogate function construction, which were shown critical for effective implementation. We also provided rigorous convergence analysis that is tailored for the non-Euclidean CPD, since generic SMD proofs do not cover the proposed algorithm. We tested the algorithm over various types of simulated and real data. Substantial computational savings relative to state-ofthe-art methods were observed. These results show encouraging and promising performance of using geometry-aware algorithm design for large-scale tensor decomposition.

### REFERENCES

- N. D. Sidiropoulos, L. De Lathauwer, X. Fu, K. Huang, E. E. Papalexakis, and C. Faloutsos, "Tensor decomposition for signal processing and machine learning," *IEEE Trans. Signal Process.*, vol. 65, no. 13, pp. 3551–3582, Jul. 2017.
- [2] T. G. Kolda and B. W. Bader, "Tensor decompositions and applications," SIAM Rev., vol. 51, no. 3, pp. 455–500, 2009.
- [3] A. Cichocki et al., "Tensor decompositions for signal processing applications: From two-way to multiway component analysis," IEEE Signal Process. Mag., vol. 32, no. 2, pp. 145–163, Mar. 2015.
- [4] A. Anandkumar, R. Ge, D. Hsu, S. M. Kakade, and M. Telgarsky, "Tensor decompositions for learning latent variable models," *J. Mach. Learn. Res.*, vol. 15, pp. 2773–2832, 2014.
- [5] X. Fu, N. Vervliet, L. De Lathauwer, K. Huang, and N. Gillis, "Computing large-scale matrix and tensor decomposition with structured factors: A unified nonconvex optimization perspective," *IEEE Signal Process. Mag.*, vol. 37, no. 5, pp. 78–94, Sep. 2020.
- [6] D. Hong, T. G. Kolda, and J. A. Duersch, "Generalized canonical polyadic tensor decomposition," SLAM Rev., vol. 62, no. 1, pp. 133–163, 2020.
- [7] T. G. Kolda and D. Hong, "Stochastic gradients for large-scale tensor decomposition," SIAM J. Math. Data Sci., vol. 2, no. 4, pp. 1066–1095, Oct. 2020.
- [8] M. Wang and L. Li, "Learning from binary multiway data: Probabilistic tensor decomposition and its statistical optimality," J. Mach. Learn. Res., vol. 21, no. 154, pp. 1–38, 2020.
- [9] M. Vandecappelle, N. Vervliet, and L. De Lathauwer, "A second-order method for fitting the canonical polyadic decomposition with non-leastsquares cost," *IEEE Trans. Signal Process.*, vol. 68, pp. 4454–4465, Jul. 2020, doi: 10.1109/TSP.2020.3010719.
- [10] K. Huang and N. D. Sidiropoulos, "Kullback-Leibler principal component for tensors is not NP-hard," in *Proc. IEEE 51st Asilomar Conf. Signals,* Syst, Comput., 2017, pp. 693–697.
- [11] N. Kargas and N. D. Sidiropoulos, "Learning mixtures of smooth product distributions: Identifiability and algorithm," in *Proc. 22nd Int. Conf. Artif. Intell. Statist.*, 2019, pp. 388–396.

- [12] X. Fu, E. Seo, J. Clarke, and R. A. Hutchinson, "Link prediction under imperfect detection: Collaborative filtering for ecological networks," *IEEE Trans. Knowl. Data Eng.*, vol. 33, no. 8, pp. 3117–3128, Aug. 2021.
- [13] L. Cheng, X. Tong, S. Wang, Y.-C. Wu, and H. V. Poor, "Learning nonnegative factors from tensor data: Probabilistic modeling and inference algorithm," *IEEE Trans. Signal Process.*, vol. 68, pp. 1792–1806, Feb. 2020, doi: 10.1109/TSP.2020.2975353.
- [14] B. Ermis, E. Acar, and A. T. Cemgil, "Link prediction in heterogeneous data via generalized coupled tensor factorization," *Data Mining Knowl. Discov.*, vol. 29, no. 1, pp. 203–236, 2015.
- [15] E. C. Chi and T. G. Kolda, "On tensors, sparsity, and nonnegative factorizations," SLAM J. Matrix Anal. Appl., vol. 33, no. 4, pp. 1272–1299, 2012
- [16] P. Comon, X. Luciani, and A. L. De Almeida, "Tensor decompositions, alternating least squares and other tales," J. Chemometrics: A. J. Chemometrics Soc., vol. 23, no. 7-8, pp. 393–405, 2009.
- [17] A.-H. Phan, P. Tichavský, and A. Cichocki, "Fast alternating LS algorithms for high order candecomp/parafac tensor factorizations," *IEEE Trans. Signal Process.*, vol. 61, no. 19, pp. 4834–4846, Oct. 2013.
- [18] Y. Xu and W. Yin, "A block coordinate descent method for regularized multiconvex optimization with applications to nonnegative tensor factorization and completion," SIAM J. Imag. Sci., vol. 6, no. 3, pp. 1758–1789, 2013
- [19] K. Huang, N. D. Sidiropoulos, and A. P. Liavas, "A flexible and efficient algorithmic framework for constrained matrix and tensor factorization," *IEEE Trans. Signal Process.*, vol. 64, no. 19, pp. 5052–5065, Oct. 2016.
- [20] L. Sorber, M. Van Barel, and L. De Lathauwer, "Optimization-based algorithms for tensor decompositions: Canonical polyadic decomposition, decomposition in rank-(l<sub>r</sub>, l<sub>r</sub>, 1) terms, and a new generalization," SIAM J. Optim., vol. 23, no. 2, pp. 695–720, 2013.
- [21] A.-H. Phan, P. Tichavsky, and A. Cichocki, "Low complexity damped gauss-newton algorithms for candecomp/parafac," SIAM J. Matrix Anal. Appl., vol. 34, no. 1, pp. 126–147, 2013.
- [22] A. Beutel, P. P. Talukdar, A. Kumar, C. Faloutsos, E. E. Papalexakis, and E. P. Xing, "Flexifact: Scalable flexible factorization of coupled tensors on hadoop," in *Proc. SIAM Int. Conf. Data Mining*, 2014, pp. 109–117.
- [23] N. Vervliet and L. De Lathauwer, "A randomized block sampling approach to canonical polyadic decomposition of large-scale tensors," *IEEE J. Sel. Topics Signal Process.*, vol. 10, no. 2, pp. 284–295, Mar. 2016.
- [24] M. Sorensen and L. De Lathauwer, "Fiber sampling approach to canonical polyadic decomposition and application to tensor completion," SIAM J. Matrix Anal. Appl., vol. 40, no. 3, pp. 888-917, 2019.
- [25] X. Fu, S. Ibrahim, H.-T. Wai, C. Gao, and K. Huang, "Block-randomized stochastic proximal gradient for low-rank tensor factorization," *IEEE Trans. Signal Process.*, vol. 68, pp. 2170–2185, Mar. 2020, doi: 10.1109/TSP.2020.2982321.
- [26] A. Cichocki and A.-H. Phan, "Fast local algorithms for large scale non-negative matrix and tensor factorizations," *IEICE Trans. Fundamentals Electron., Commun. Comput. Sci.*, vol. 92, no. 3, pp. 708–721, 2009.
- [27] M. Razaviyayn, M. Hong, and Z.-Q. Luo, "A unified convergence analysis of block successive minimization methods for nonsmooth optimization," SIAM J. Optim., vol. 23, no. 2, pp. 1126–1153, 2013.
- [28] C. Févotte and J. Idier, "Algorithms for nonnegative matrix factorization with the □divergence," *Neural Computation*, vol. 23, no. 9, pp. 2421–2456, 2011.
- [29] A. Beck and M. Teboulle, "Mirror descent and nonlinear projected subgradient methods for convex optimization," *Operations Res. Lett.*, vol. 31, no. 3, pp. 167–175, 2003.
- [30] C. Battaglino, G. Ballard, and T. G. Kolda, "A practical randomized cp tensor decomposition," SIAM J. Matrix Anal. Appl., vol. 39, no. 2, pp. 876–901, 2018.
- [31] C. D. Dang and G. Lan, "Stochastic block mirror descent methods for nonsmooth and stochastic optimization," SIAM J. Optim., vol. 25, no. 2, pp. 856–881, 2015.
- [32] W. Pu, S. Ibrahim, X. Fu, and M. Hong, "Fiber-sampled stochastic mirror descent for tensor decomposition with □divergence," in *Proc. IEEE Int. Conf. Acoust.*, Speech Signal Process., 2021, pp. 2925–2929, doi: 10.1109/ICASSP39728.2021.9413830.
- [33] J. Shetty and J. Adibi, "The enron email dataset database schema and brief statistical report," Inf. Sci. Inst. Tech. Rep., Univ. Southern California, vol. 4, no. 1, pp. 120–128, 2004.
- [34] C. Févotte, N. Bertin, and J.-L. Durrieu, "Nonnegative matrix factorization with the Itakura-Saito divergence: With application to music analysis," *Neural Computation*, vol. 21, no. 3, pp. 793–830, 2009.

- [35] K. Huang and X. Fu, "Detecting overlapping and correlated communities without pure nodes: Identifiability and algorithm," in *Proc. Int. Conf. Mach. Learn.*, 2019, pp. 2859–2868.
- [36] H. H. Bauschke, J. Bolte, and M. Teboulle, "A descent lemma beyond Lipschitz gradient continuity: First-order methods revisited and applications," *Math. Operations Res.*, vol. 42, no. 2, pp. 330–348, 2017.
- [37] J. Bolte, S. Sabach, M. Teboulle, and Y. Vaisbourd, "First order methods beyond convexity and Lipschitz gradient continuity with applications to quadratic inverse problems," SIAM J. Optim., vol. 28, no. 3, pp. 2131–2151, 2018.
- [38] H. Lu, "Relative continuity for non-Lipschitz nonsmooth convex optimization using stochastic (or deterministic) mirror descent," *INFORMS J. Optim.*, vol. 1, no. 4, pp. 288–303, 2019.
- [39] D. Davis, D. Drusvyatskiy, and K. J. MacPhee, "Stochastic model-based minimization under high-order growth," 2018, arXiv:1807.00255.
- [40] S. Zhang and N. He, "On the convergence rate of stochastic mirror descent for nonsmooth nonconvex optimization," 2018, arXiv:1806.04781.
- [41] Q. Li, Z. Zhu, G. Tang, and M. B. Wakin, "Provable Bregman-Divergence based methods for nonconvex and non-Lipschitz problems," 2019, arXiv:1904.09712.
- [42] L. M. Bregman, "The relaxation method of finding the common point of convex sets and its application to the solution of problems in convex programming," USSR Comput. Math. Math. Phys., vol. 7, no. 3, pp. 200–217, 1067
- [43] A. Yeredor and M. Haardt, "Estimation of a low-rank probability-tensor from sample sub-tensors via joint factorization minimizing the Kullback-Leibler divergence," in *Proc. IEEE 27th Eur. Signal Process. Conf.*, 2019, pp. 1–5.
- [44] A. Yeredor and M. Haardt, "Maximum likelihood estimation of a low-rank probability mass tensor from partial observations," *IEEE Signal Process. Lett.*, vol. 26, no. 10, pp. 1551–1555, Oct. 2019.
- [45] S. Ibrahim and X. Fu, "Recovering joint probability of discrete random variables from pairwise marginals," *IEEE Trans. Signal Process.*, vol. 69, pp. 4116–4131, Jun. 2021, doi: 10.1109/TSP.2021.3090960.
- [46] L. T. K. Hien and N. Gillis, "Algorithms for nonnegative matrix factorization with the kullback-leibler divergence," 2020, arXiv:2010.01935.
- [47] N. Parikh and S. Boyd, "Proximal algorithms," Foundations Trends Optim., vol. 1, no. 3, pp. 127–239, 2014.
- [48] Y. Sun, P. Babu, and D. P. Palomar, "Majorization-minimization algorithms in signal processing, communications, and machine learning," *IEEE Trans. Signal Process.*, vol. 65, no. 3, pp. 794–816, Feb. 2017.
- [49] M. Hong, M. Razaviyayn, Z.-Q. Luo, and J.-S. Pang, "A unified algorithmic framework for block-structured optimization involving Big Data: With applications in machine learning and signal processing," *IEEE Signal Process. Mag.*, vol. 33, no. 1, pp. 57–77, Jan. 2016.
- [50] J. Duchi, E. Hazan, and Y. Singer, "Adaptive subgradient methods for online learning and stochastic optimization," J. Mach. Learn. Res., vol. 12, no. 61, pp. 2121–2159, 2011. [Online]. Available: http://jmlr.org/papers/ v12/duchi11a.html
- [51] X. Chen, S. Liu, R. Sun, and M. Hong, "On the convergence of a class of adam-type algorithms for non-convex optimization," in *Proc. Int. Conf. Learn. Representations*, 2019.
- [52] H. Lu, R. M. Freund, and Y. Nesterov, "Relatively smooth convex optimization by first-order methods, and applications," *SIAM J. Optim.*, vol. 28, no. 1, pp. 333–354, 2018.
- [53] M. Teboulle, "A simplified view of first order methods for optimization," Math. Program., vol. 170, no. 1, pp. 67–96, 2018.
- [54] T. Kola et al., "Tensor toolbox for MATLAB v. 3.0," Sandia National Lab. (SNL-NM), Albuquerque, NM (United States), Oct. 2017. [Online]. Available: https://www.tensortoolbox.org
- [55] N. Vervliet, O. Debals, L. Sorber, M. Van Barel, and L. De Lathauwer, "Tensorlab 3.0," Mar. 2016. [Online]. Available: www.tensorlab.net
- [56] N. Kargas, N. D. Sidiropoulos, and X. Fu, "Tensors, learning, and 'Kolmogorov extension' for finite-alphabet random vectors," *IEEE Trans. Signal Process.*, vol. 66, no. 18, pp. 4854–4868, Sep. 2018.
- [57] S. Ibrahim, X. Fu, N. Kargas, and K. Huang, "Crowdsourcing via pairwise co-occurrences: Identifiability and algorithms," in *Proc. 33rd Int. Conf. Neural Inf. Process. Syst.*, 2019, pp. 7847–7857.
- [58] S. Smith et al., "FROSTT: The formidable repository of open sparse tensors and tools," 2017. [Online]. Available: http://frostt.io/
- [59] J. A. Jones, R. A. Hutchinson, and V. W. Pfeiffer, "Plant pollinator data at hj andrews experimental forest, 2011 to 2015," *Environmental Data Initiative, ver. 5*.
- [60] T. Opsahl and P. Panzarasa, "Clustering in weighted networks," Social Netw., vol. 31, no. 2, pp. 155–163, 2009.



Wenqiang Pu (Member, IEEE) received the B.S. and Ph.D. degrees in electrical engineering from Xidian University, Xi'an, China, in 2013 and 2018, respectively. From January 2019 to October 2020, he was a Postdoc Researcher with the School of Science and Engineering, The Chinese University of Hong Kong (Shenzhen), Shenzhen, China. He is currently a Research Scientist with the Shenzhen Research Institute of Big Data, Shenzhen, China. His research interests include signal processing and optimization algorithms.



Shahana Ibrahim received the B.Tech. degree in electronics and communication engineering from the National Institute of Technology, Calicut, Kozhikode, India, in 2012, and the M.S. degree in electrical engineering, in 2019, from Oregon State University, Corvallis, OR, USA, where she is currently working toward the Ph.D. degree. From 2012 to 2017, she was a System Validation Engineer with Texas Instruments, Bengaluru, India. Her research interests include statistical machine learning and signal processing.



Xiao Fu (Senior Member, IEEE) received the B.Eng. and M.Sc. degrees from the University of Electronic Science and Technology of China, Chengdu, China, in 2005 and 2010, respectively, and the Ph.D. degree in electronic engineering from The Chinese University of Hong Kong, Hong Kong, in 2014. From 2014 to 2017, he was a Postdoctoral Associate with the Department of Electrical and Computer Engineering, University of Minnesota, Minneapolis, MN, USA. Since 2017, he has been an Assistant Professor with the School of Electrical Engineering and Computer

Science, Oregon State University, Corvallis, OR, USA. His research interests include the broad area of signal processing and machine learning.

Dr. Fu was the recipient of the Best Student Paper Award at ICASSP 2014, and Outstanding Postdoctoral Scholar Award at the University of Minnesota in 2016. His coauthored papers was the recipient of the Best Student Paper Awards from IEEE CAMSAP 2015 and IEEE MLSP 2019, respectively. He is a Member of the Sensor Array and Multichannel Technical Committee of the IEEE Signal Processing Society. He is also a Member of the Signal Processing for Multisensor Systems Technical Area Committee of EURASIP. He is the Treasurer of the IEEE SPS Oregon Chapter. He is the Editor of Signal Processing and is also an Associate Editor for the IEEE TRANSACTIONS ON SIGNAL PROCESSING. He was a Tutorial Speaker at ICASSP 2017 and SIAM Conference on Applied Linear Algebra 2021.



Mingyi Hong (Senior Member, IEEE) received the Ph.D. degree from the University of Virginia, Charlottesville, VA, USA, in 2011. He is currently an Associate Professor with the Department of Electrical and Computer Engineering, University of Minnesota, Minneapolis, MN, USA. His research interests include optimization theory and applications in signal processing and machine learning. He is on the IEEE Signal Processing for Communications and Networking Technical Committee. He is also an Associate Editor for the IEEE Transactions on Signal Processing.