

# A Bayesian Approach to Block-Term Tensor Decomposition Model Selection and Computation

Paris V. Giampouras

Mathematical Inst. for Data Science IAASARS, National Observatory of Athens Dept. of Statistics and Insurance Science  
Johns Hopkins University  
Baltimore, MD 21218, USA  
parisg@jhu.edu

Athanasios A. Rontogiannis

152 36 Penteli, Greece  
tronto@noa.gr

Eleftherios Kofidis

University of Piraeus  
185 34 Piraeus, Greece  
kofidis@unipi.gr

**Abstract**—The so-called block-term decomposition (BTD) tensor model, especially in its rank- $(L_r, L_r, 1)$  version, has been recently receiving increasing attention due to its enhanced ability of representing systems and signals that are composed of blocks of rank higher than one, a scenario encountered in numerous and diverse applications. Its uniqueness and approximation have thus been thoroughly studied. Nevertheless, the challenging problem of estimating the BTD model structure, namely the number of block terms and their individual ranks, has only recently started to attract significant attention. In this work, a Bayesian approach is taken to addressing the problem of rank- $(L_r, L_r, 1)$  BTD model selection and computation, based on the idea of imposing column sparsity jointly on the factors and in a hierarchical manner and estimating the ranks as the numbers of factor columns of non-negligible energy. Approximate posterior inference for the proposed sophisticated Bayesian model is based on variational inference giving rise to an iterative algorithm that comprises closed-form updates. Its Bayesian nature completely avoids the ubiquitous in regularization-based methods task of hyper-parameter tuning. Simulation results with synthetic data are reported, which demonstrate the effectiveness of the proposed scheme in terms of both rank estimation and model fitting.

**Index Terms**—Bayesian inference, block-term decomposition (BTD), hierarchical iterative reweighted least squares (HIRLS), rank, tensor, variational inference (VI)

## I. INTRODUCTION

Block-Term Decomposition (BTD) was introduced in [1] as a tensor model that combines the Canonical Polyadic Decomposition (CPD) and the Tucker decomposition (TD), in the sense that it decomposes a tensor in a sum of tensors (block terms) that have low multilinear rank (instead of rank one as in CPD). Hence a BTD can be seen as a constrained TD, with its core tensor being block diagonal (see [1, Fig. 2.3]). It can also be seen as a constrained CPD having factors with (some) colinear columns [1]. In a way, BTD lies between the two extremes (in terms of core tensor structure), CPD and TD, and it is interesting to recall the related remark made in [1], namely that “the rank of a higher-order tensor is actually a combination of the two aspects: one should specify the number of blocks and their size”. Accurately and efficiently estimating

these numbers for a given tensor is the main subject of this paper.

Although [1] introduced BTD as a sum of  $R$  rank- $(L_r, M_r, N_r)$  terms ( $r = 1, 2, \dots, R$ ) in general, the special case of rank- $(L_r, L_r, 1)$  BTD has attracted a lot more of attention, because of both its more frequent occurrence in applications and the existence of more concrete and easier to check uniqueness conditions. This work will also focus on this special yet very popular BTD model. Consider a 3rd-order tensor,  $\mathcal{X} \in \mathbb{C}^{I \times J \times K}$ . Then its rank- $(L_r, L_r, 1)$  decomposition is written as  $\mathcal{X} = \sum_{r=1}^R \mathbf{E}_r \circ \mathbf{c}_r$ , where  $\mathbf{E}_r$  is an  $I \times J$  matrix of rank  $L_r$ ,  $\mathbf{c}_r$  is a nonzero column  $K$ -vector and  $\circ$  denotes outer product. Clearly,  $\mathbf{E}_r$  can be written as a matrix product  $\mathbf{A}_r \mathbf{B}_r^T$  with the matrices  $\mathbf{A}_r \in \mathbb{C}^{I \times L_r}$  and  $\mathbf{B}_r \in \mathbb{C}^{J \times L_r}$  being of full column rank,  $L_r$ . Eq. (1) can thus be re-written as

$$\mathcal{X} = \sum_{r=1}^R \mathbf{A}_r \mathbf{B}_r^T \circ \mathbf{c}_r. \quad (1)$$

BTD has been successfully used in a wide range of application areas and its uniqueness and computation have been thoroughly studied (cf. [2] for an extensive review). In general,  $R$  and  $L_r$ ,  $r = 1, 2, \dots, R$  are assumed *a-priori* known (and it is commonly assumed that all  $L_r$  are all equal to  $L$ , for simplicity). However, unless external information is given (such as in a telecommunications [3] or a hyperspectral image unmixing application with given or estimated ground truth [4]), there is no way to know these values beforehand. Model selection for BTD, that is estimating the number of block terms,  $R$ , and their ranks,  $L_r$ ,  $r = 1, 2, \dots, R$ , is clearly more challenging than in CPD and TD models and has only recently started to be studied (cf. [2] and references therein). The most recent contribution of this kind can be found in our work [2], which relies on a regularization of the squared approximation error function with the sum of the Frobenius norms of the factors reweighted by a diagonal weighting which jointly depends on the factors in two levels: the reweighted norms of  $\mathbf{A} \triangleq [\mathbf{A}_1 \ \mathbf{A}_2 \ \dots \ \mathbf{A}_R]$  and  $\mathbf{B} \triangleq [\mathbf{B}_1 \ \mathbf{B}_2 \ \dots \ \mathbf{B}_R]$  are combined and then coupled with the reweighted norm of  $\mathbf{C} \triangleq [\mathbf{c}_1 \ \mathbf{c}_2 \ \dots \ \mathbf{c}_R]$ . This two-level coupling naturally matches the structure of the model

P. V. Giampouras is supported by the European Union under the Horizon 2020 Marie Skłodowska-Curie Global Fellowship program: HyPPOCRATES-H2020-MSCA-IF-2018, Grant Agreement Number: 844290. The work of E. Kofidis has been partly supported by the University of Piraeus Research Center.

in (1), making explicit the different roles of  $\mathbf{A}$ ,  $\mathbf{B}$  and  $\mathbf{C}$ . This way, column sparsity is imposed *jointly* on the factors and in a *hierarchical* manner, which allows to estimate the ranks as the numbers of factor columns of non-negligible energy. Following a block coordinate descent solution approach, an alternating hierarchical iterative reweighted least squares (HIRLS) algorithm, called BT-D-HIRLS, was developed in [2] that manages to both reveal the ranks and compute the BT-D factors at a high convergence rate and low computational cost.

Nevertheless, BT-D-HIRLS, being a regularization-based method, faces the same challenge that all such methods have to address, namely to appropriately tune the regularization parameter so as to achieve the best possible performance. Although a rough guideline for the parameter selection has been given and utilized in [2] as a reference point for the trial-end-error search, this is still only a rule of thumb, not completely relieving the algorithm from the need to spend resources on searching for the most appropriate parameter value. To overcome this difficulty, we take in this paper an alternative, Bayesian approach, viewing the unknowns as random variables and tackling the problem as one of Bayesian modeling and inference [5]. The idea is again to impose column sparsity jointly on the factors in a hierarchical, two-level manner. This is achieved through a Bayesian hierarchy of priors with sparsity inducing effect, that realize the coupling of the columns of  $\mathbf{C}$  and the  $\mathbf{A}_r$ ,  $\mathbf{B}_r$  blocks at the outer level and that between the columns of corresponding blocks at the inner level. Thus,  $R$  is estimated as the number of columns of  $\mathbf{C}$  of non-negligible energy while the  $L_r$ 's are found similarly from the columns of the  $\mathbf{A}_r$ ,  $\mathbf{B}_r$  blocks. Approximate inference in the proposed probabilistic model is efficiently performed via variational inference [6] leading to an iterative algorithm that comprises closed-form updates and is fast converging. Its Bayesian nature completely avoids the need for parameter tuning. Simulation results with synthetic data are reported, which demonstrate the effectiveness of the proposed scheme in terms of both rank estimation and model fitting and in comparison with BT-D-HIRLS. Bayesian methods for tensor decomposition model selection and computation have been already reported in the literature, with [7], [8], and [9] being recent examples concerning TD, CPD, and Tensor Train (TT) decomposition, respectively. To the best of our knowledge, however, the present work is the first of its kind for BT-D in multilinear rank- $(L_r, L_r, 1)$  terms.

## II. PROBLEM STATEMENT

Given the  $I \times J \times K$  tensor

$$\mathbf{y} = \mathcal{X} + \sigma \mathcal{N}, \quad (2)$$

where  $\mathcal{X}$  is given by (1) and  $\mathcal{N}$  is the  $I \times J \times K$  noise tensor of zero-mean unit variance i.i.d. Gaussian entries, we aim at estimating  $R$ ,  $L_r$ ,  $r = 1, 2, \dots, R$  and the factor matrices  $\mathbf{A}_r = [\mathbf{a}_{r1} \ \mathbf{a}_{r2} \ \dots \ \mathbf{a}_{rL_r}] \in \mathbb{C}^{I \times L_r}$ ,  $\mathbf{B}_r = [\mathbf{b}_{r1} \ \mathbf{b}_{r2} \ \dots \ \mathbf{b}_{rL_r}] \in \mathbb{C}^{J \times L_r}$ ,  $\mathbf{C} \in \mathbb{C}^{K \times R}$ , subject of course to the inherent ambiguity resulting from the fact that only the product  $\mathbf{A}_r \mathbf{B}_r^T$  can be uniquely identified modulo

a scaling (with a counter-scaling of  $\mathbf{c}_r$ ) [1]. In terms of its mode unfoldings  $\mathbf{X}_{(1)} \in \mathbb{C}^{I \times JK}$ ,  $\mathbf{X}_{(2)} \in \mathbb{C}^{J \times IK}$  and  $\mathbf{X}_{(3)} \in \mathbb{C}^{K \times IJ}$ , the tensor  $\mathcal{X}$  can be written as [1]

$$\mathbf{X}_{(1)}^T = (\mathbf{B} \odot \mathbf{C}) \mathbf{A}^T \triangleq \mathbf{P} \mathbf{A}^T, \quad (3)$$

$$\mathbf{X}_{(2)}^T = (\mathbf{C} \odot \mathbf{A}) \mathbf{B}^T \triangleq \mathbf{Q} \mathbf{B}^T, \quad (4)$$

$$\mathbf{X}_{(3)}^T = [(\mathbf{A}_1 \odot_c \mathbf{B}_1) \mathbf{1}_{L_1} \ \dots \ (\mathbf{A}_R \odot_c \mathbf{B}_R) \mathbf{1}_{L_R}] \mathbf{C}^T \\ \triangleq \mathbf{T} \mathbf{C}^T, \quad (5)$$

where  $\odot$  denotes the Khatri-Rao product in its general (partition-wise) version and  $\odot_c$  is its column-wise version. In this paper, we follow a Bayesian approach to address the above problem. A fully Bayesian analysis is detailed next.

## III. PROPOSED BAYESIAN MODEL

Let  $R$  and the  $L_r$ 's be overestimated to  $R_{\text{ini}}$  and  $L_{\text{ini}}$ , respectively. Heavy-tailed multi-parameter Laplace distributions, known for their sparsity inducing effect, are placed over the columns of the  $\mathbf{A}_r$ 's,  $\mathbf{B}_r$ 's, and  $\mathbf{C}$  in a way that implicitly implements a regularization analogous to that of the BT-D-HIRLS method [2]. Namely, the number of block terms and the ranks of  $\mathbf{A}_r$ 's and  $\mathbf{B}_r$ 's are jointly penalized, while respecting the different role that these matrices play in the BT-D model. This results in the nulling of all but  $R$  columns of  $\mathbf{C}$ , and the nulling of all but  $L_r$  columns of the corresponding "surviving"  $\mathbf{A}_r$ ,  $\mathbf{B}_r$  blocks. Following the premise of the well-known automatic relevance determination (ARD) framework [5], [10], the priors are assigned via a 3-level hierarchy of prior distributions outlined in the following.

The likelihood function, which encodes the underlying causal relation between the data and the latent variables, can be written with respect to (w.r.t.) the mode-1 unfolding of  $\mathcal{Y}$  (cf. (3)) as follows:

$$p(\mathbf{Y}_{(1)} | \mathbf{A}, \mathbf{B}, \mathbf{C}, \beta) = \prod_{i=1}^I p(\mathbf{y}_{(1)i} | \mathbf{A}, \mathbf{B}, \mathbf{C}, \beta) \\ = \prod_{i=1}^I \mathcal{N}(\mathbf{y}_{(1)i} | \mathbf{P} \mathbf{a}_i, \beta^{-1} \mathbf{I}_{JK}) \quad (6)$$

where  $\beta$  is the noise precision and  $\mathbf{a}_i$ ,  $\mathbf{y}_{(1)i}$  are the  $i$ th rows<sup>1</sup> of  $\mathbf{A}$ ,  $\mathbf{Y}_{(1)}$ , respectively, in column form. In an analogous way the likelihood function can be written in terms of the rows  $\mathbf{b}_j$ ,  $\mathbf{c}_k$  and  $\mathbf{y}_{(2)j}$ ,  $\mathbf{y}_{(3)k}$  of  $\mathbf{B}$ ,  $\mathbf{C}$  and  $\mathbf{Y}_{(2)}$ ,  $\mathbf{Y}_{(3)}$ , respectively, if  $\mathcal{Y}$  is expressed w.r.t. its mode-2 and mode-3 unfoldings.  $\mathbf{A}$ ,  $\mathbf{B}$ ,  $\mathbf{C}$  are considered as unobserved random variables and are assigned 3-level hierarchical prior distributions. At the first level of the hierarchy, Gaussian distributions are placed over  $\mathbf{A}$ ,  $\mathbf{B}$ , and  $\mathbf{C}$ , namely,

$$p(\mathbf{A} | \mathbf{s}, \boldsymbol{\zeta}, \beta) = \prod_{i=1}^I \mathcal{N}(\mathbf{a}_i | \mathbf{0}, \beta^{-1} \mathbf{S}^{-1} (\mathbf{Z}^{-1} \otimes \mathbf{I}_{L_{\text{ini}}})) \quad (7)$$

$$p(\mathbf{B} | \mathbf{s}, \boldsymbol{\zeta}, \beta) = \prod_{j=1}^J \mathcal{N}(\mathbf{b}_j | \mathbf{0}, \beta^{-1} \mathbf{S}^{-1} (\mathbf{Z}^{-1} \otimes \mathbf{I}_{L_{\text{ini}}})) \quad (8)$$

<sup>1</sup>We denote matrix rows and columns with bold italic and roman letters, respectively.

and

$$p(\mathbf{C} \mid \zeta, \beta) = \prod_{k=1}^K \mathcal{N}(\mathbf{c}_k \mid \mathbf{0}, \beta^{-1} \mathbf{Z}^{-1}), \quad (9)$$

where  $\mathbf{S} = \text{diag}(\mathbf{s})$  with  $\mathbf{s} = [s_{rl}] \in \mathbb{R}^{L_{\text{ini}} R_{\text{ini}} \times 1}$  and  $\mathbf{Z} = \text{diag}(\zeta)$ ,  $\zeta \in \mathbb{R}^{R_{\text{ini}} \times 1}$ . Note that the priors of  $\mathbf{A}$  and  $\mathbf{B}$  are zero-mean with the same covariance matrix, which is formed from the diagonal matrices  $\mathbf{Z}$  and  $\mathbf{S}$ . This particular selection is of critical importance from an implicit regularization perspective, since it induces identical sparsity patterns over columns/sub-blocks of  $\mathbf{A}$  and  $\mathbf{B}$ . A sufficiently large value of  $\zeta_r$  will lead the  $r$ th column of  $\mathbf{C}$  (cf. (9)) and the entire set of the redundant  $L_{\text{ini}}$  columns of sub-matrices  $\mathbf{A}_r, \mathbf{B}_r$  (cf. (7), (8)) to zero. At the same time, the superfluous  $l$ th columns of the ‘‘surviving’’  $\mathbf{A}_r, \mathbf{B}_r$  are *jointly* forced to be zero when the value of  $s_{rl}$  becomes sufficiently large (cf. (7), (8)). Thus, the  $R_{\text{ini}}$  diagonal values of  $\mathbf{Z}$  act as weights that determine the number of nonzero block terms,  $R$ , while the  $L_{\text{ini}} R_{\text{ini}}$  diagonal values of  $\mathbf{S}$  determine the number of nonzero columns of  $\mathbf{A}$  and  $\mathbf{B}$ . The sophisticated selection of (7), (8) leads to *joint* sparsity imposition on the columns of the factors, in two levels. This way, ideas of ARD and sparse Bayesian learning (SBL) [10] are put into effect to address the challenging problem of BTD model selection.

At the second level of the hierarchy of priors, inverse Gamma priors are assigned over  $\mathbf{s}$  and  $\zeta$ ,

$$p(\mathbf{s}) = \prod_{r=1}^{R_{\text{ini}}} \prod_{l=1}^{L_{\text{ini}}} \mathcal{IG} \left( s_{rl} \mid \frac{I+J+1}{2}, \frac{\delta_{rl}}{2} \right), \quad (10)$$

$$p(\zeta) = \prod_{r=1}^{R_{\text{ini}}} \mathcal{IG} \left( \zeta_r \mid \frac{(I+J)L_{\text{ini}} + K + 1}{2}, \frac{\rho_r}{2} \right), \quad (11)$$

where  $\delta_{rl}$  and  $\rho_r$  are the scale parameters of the distributions over  $s_{rl}$  and  $\zeta_r$ , respectively. The third level involves Gamma prior distributions over these variables, namely,

$$p(\delta_{rl}) = \mathcal{G}(\delta_{rl} \mid \psi, \tau), \quad (12)$$

$$p(\rho_r) = \mathcal{G}(\rho_r \mid \mu, \nu), \quad (13)$$

where  $\psi, \tau, \mu, \nu$  take very small positive values rendering the respective priors non-informative. Note that these priors are conjugate w.r.t. the likelihood functions and w.r.t. each other, which guarantees that the posterior distributions will belong to the same class of distributions with the priors [5]. Moreover, one can show that heavy-tailed multi-parameter Laplace marginal distributions show up over the columns of  $\mathbf{A}, \mathbf{B}$  and  $\mathbf{C}$  after integrating out the variables  $\mathbf{s}$  and  $\zeta$ . Finally, a non-informative Gamma prior is also placed on the noise precision variable  $\beta$ , namely,

$$p(\beta) = \mathcal{G}(\beta \mid \kappa, \theta). \quad (14)$$

The above Bayesian model is depicted in the form of a graphical model in Fig. 1.

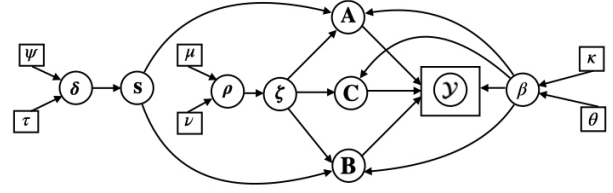


Fig. 1. The proposed Bayesian model.

#### IV. APPROXIMATE POSTERIOR INFERENCE

Let  $\Theta$  be the cell array which includes all unobserved variables, that is,  $\Theta \triangleq \{\mathbf{A}, \mathbf{B}, \mathbf{C}, \mathbf{s}, \zeta, \beta, \rho, \delta\}$ . The exact joint posterior of the variables of the adopted Bayesian model is given by

$$p(\Theta \mid \mathcal{Y}) = \frac{p(\mathcal{Y}, \Theta)}{\int p(\mathcal{Y}, \Theta) d\Theta}. \quad (15)$$

Due to the complexity of the model, the marginal distribution of  $\mathcal{Y}$  in the denominator is computationally intractable. Therefore, we follow a variational inference (VI) approach for approximating (15). The idea is to approximate the posterior by a distribution which is as close as possible to the exact posterior in terms of the Kullback-Leibler divergence [6]. VI allows for an efficient approximate inference process even in vastly complicated Bayesian models that involve high-dimensional variables. It is usually coupled with mean-field approximation, namely, the assumption that the posterior distribution can be factorized w.r.t. the involved variables, implying statistical independence among them. In our case, the approximate posterior  $q(\Theta)$  of  $p(\Theta \mid \mathcal{Y})$  is written in the form

$$q(\Theta) = q(\beta) \prod_{i=1}^I q(\mathbf{a}_i) \prod_{j=1}^J q(\mathbf{b}_j) \prod_{k=1}^K q(\mathbf{c}_k) \times \prod_{r=1}^{R_{\text{ini}}} \prod_{l=1}^{L_{\text{ini}}} q(s_{rl}) q(\delta_{rl}) \prod_{r=1}^{R_{\text{ini}}} q(\zeta_r) q(\rho_r) \quad (16)$$

Denoting the individual variables above by  $\theta_i$ , the corresponding VI-based posteriors are known to satisfy [6]

$$q(\theta_i) = \frac{\exp(\langle \ln(p(\mathcal{Y}, \Theta)) \rangle_{i \neq j})}{\int \exp(\langle \ln(p(\mathcal{Y}, \Theta)) \rangle_{i \neq j}) d\theta_i}, \quad (17)$$

where  $\langle \cdot \rangle_{i \neq j}$  denotes expectation w.r.t. all  $q(\theta_j)$ s but  $q(\theta_i)$ . To solve (17) a block coordinate ascent approach is taken, employing the cyclic update rule, namely solving for  $q(\theta_i)$  given  $q(\theta_j)$ ,  $j \neq i$ . It follows that the posterior distribution of  $\mathbf{a}_i$  is given by

$$q(\mathbf{a}_i) = \mathcal{N}(\langle \mathbf{a}_i \rangle, \Sigma_{\mathbf{a}}), \quad (18)$$

with<sup>2</sup>

$$\langle \mathbf{a}_i \rangle = \langle \beta \rangle \Sigma_{\mathbf{a}} \langle \mathbf{P} \rangle^T \mathbf{y}_{(1)i}, \quad (19)$$

$$\Sigma_{\mathbf{a}} = \langle \beta \rangle^{-1} (\langle \mathbf{P}^T \mathbf{P} \rangle + \langle \mathbf{S} \rangle (\langle \mathbf{Z} \rangle \otimes \mathbf{I}_{L_{\text{ini}}}))^{-1}, \quad (20)$$

<sup>2</sup>All  $\mathbf{a}_i$ s turn out to have the same covariance matrix,  $\Sigma_{\mathbf{a}}$ , and similarly for the  $\mathbf{b}_j$ s and the  $\mathbf{c}_k$ s.

where  $\langle \cdot \rangle$  denotes expectation w.r.t the posterior of the involved variable. The posterior of  $\mathbf{b}_j$  results in an analogous manner:

$$q(\mathbf{b}_j) = \mathcal{N}(\langle \mathbf{b}_j \rangle, \Sigma_{\mathbf{b}}), \quad (21)$$

with

$$\langle \mathbf{b}_j \rangle = \langle \beta \rangle \Sigma_{\mathbf{b}} \langle \mathbf{Q} \rangle^T \mathbf{y}_{(2)j} \quad (22)$$

$$\Sigma_{\mathbf{b}} = \langle \beta \rangle^{-1} (\langle \mathbf{Q}^T \mathbf{Q} \rangle + \langle \mathbf{S} \rangle (\langle \mathbf{Z} \rangle \otimes \mathbf{I}_{L_{\text{ini}}}))^{-1}. \quad (23)$$

The posterior of  $\mathbf{c}_k$  is

$$q(\mathbf{c}_k) = \mathcal{N}(\langle \mathbf{c}_k \rangle, \Sigma_{\mathbf{c}}) \quad (24)$$

with

$$\langle \mathbf{c}_k \rangle = \langle \beta \rangle \Sigma_{\mathbf{c}} \langle \mathbf{T} \rangle^T \mathbf{y}_{(3)k} \quad (25)$$

$$\Sigma_{\mathbf{c}} = \langle \beta \rangle^{-1} (\langle \mathbf{T}^T \mathbf{T} \rangle + \langle \mathbf{Z} \rangle)^{-1}. \quad (26)$$

Next, the approximate posteriors of the variables belonging to the second level of hierarchy are given. The posterior of  $s_{rl}$  turns out [11] to be a GIG pdf,

$$q(s_{rl}) = \mathcal{GIG} \left( s_{rl} \mid -\frac{1}{2}, \langle \beta \rangle \langle \zeta_r \rangle (\langle \mathbf{a}_{rl}^T \mathbf{a}_{rl} \rangle + \langle \mathbf{b}_{rl}^T \mathbf{b}_{rl} \rangle), \langle \delta_{rl} \rangle \right) \quad (27)$$

with mean

$$\langle s_{rl} \rangle = \sqrt{\frac{\langle \delta_{rl} \rangle}{\langle \beta \rangle \langle \zeta_r \rangle (\langle \mathbf{a}_{rl}^T \mathbf{a}_{rl} \rangle + \langle \mathbf{b}_{rl}^T \mathbf{b}_{rl} \rangle)}}, \quad (28)$$

where  $\langle \mathbf{a}_{rl}^T \mathbf{a}_{rl} \rangle$  and  $\langle \mathbf{b}_{rl}^T \mathbf{b}_{rl} \rangle$  are the  $((r-1)L_{\text{ini}} + l, (r-1)L_{\text{ini}} + l)$  entries of

$$\langle \mathbf{A}^T \mathbf{A} \rangle = \langle \mathbf{A} \rangle^T \langle \mathbf{A} \rangle + I \Sigma_{\mathbf{a}} \quad (29)$$

and

$$\langle \mathbf{B}^T \mathbf{B} \rangle = \langle \mathbf{B} \rangle^T \langle \mathbf{B} \rangle + J \Sigma_{\mathbf{b}}, \quad (30)$$

respectively. Similarly, the approximate posterior of  $\zeta_r$  is also GIG, with  $\langle \zeta_r \rangle$  given by

$$\langle \zeta_r \rangle = \sqrt{\frac{\langle \rho_r \rangle}{\langle \beta \rangle (\sum_{l=1}^{L_{\text{ini}}} \langle s_{rl} \rangle (\langle \mathbf{a}_{rl}^T \mathbf{a}_{rl} \rangle + \langle \mathbf{b}_{rl}^T \mathbf{b}_{rl} \rangle) + \langle \mathbf{c}_r^T \mathbf{c}_r \rangle)}} \quad (31)$$

and  $\langle \mathbf{c}_r^T \mathbf{c}_r \rangle$  denoting the  $(r, r)$  entry of

$$\langle \mathbf{C}^T \mathbf{C} \rangle = \langle \mathbf{C} \rangle^T \langle \mathbf{C} \rangle + K \Sigma_{\mathbf{c}}. \quad (32)$$

It can be shown (as in [11]) that, at the third level of hierarchy, the approximate posteriors of  $\delta_{rl}, \rho_r$  and  $\beta$  are Gamma distributions with  $\langle \delta_{rl} \rangle, \langle \rho_r \rangle$  and  $\langle \beta \rangle$  given in Table I, where the resulting *Bayesian-BTD (BBTD)* algorithm is summarized (cf. [12] for a more detailed presentation). The rest of the first- and second-order statistics that are required in the algorithm implementation are computed as in Table II, based on the assumption of statistically independent  $\mathbf{A}, \mathbf{B}, \mathbf{C}$  (cf. (16)) and making use of the identities for the Grammians of Khatri-Rao products proved in [2, Appendix C]. \* stands for the Hadamard product.

TABLE I  
THE BBTD ALGORITHM

<b>Input:</b> $\mathcal{Y}, R_{\text{ini}}, L_{\text{ini}}$
<b>Output:</b> $\hat{R}, \hat{L}_r, r = 1, 2, \dots, \hat{R}, \hat{\mathbf{A}}, \hat{\mathbf{B}}, \hat{\mathbf{C}}$
<b>Initialize</b> $\langle \mathbf{B} \rangle, \langle \mathbf{C} \rangle, \langle \beta \rangle, \langle \mathbf{S} \rangle, \langle \mathbf{Z} \rangle, \langle \delta \rangle, \langle \rho \rangle, \Sigma_{\mathbf{b}}, \Sigma_{\mathbf{c}}$
<b>repeat</b>
$\Sigma_{\mathbf{a}} = \langle \beta \rangle^{-1} (\langle \mathbf{P}^T \mathbf{P} \rangle + \langle \mathbf{S} \rangle (\langle \mathbf{Z} \rangle \otimes \mathbf{I}_{L_{\text{ini}}}))^{-1}$
$\langle \mathbf{A} \rangle = \langle \beta \rangle \mathbf{Y}_{(1)} \langle \mathbf{P} \rangle \Sigma_{\mathbf{a}}$
$\Sigma_{\mathbf{b}} = \langle \beta \rangle^{-1} (\langle \mathbf{Q}^T \mathbf{Q} \rangle + \langle \mathbf{S} \rangle (\langle \mathbf{Z} \rangle \otimes \mathbf{I}_{L_{\text{ini}}}))^{-1}$
$\langle \mathbf{B} \rangle = \langle \beta \rangle \mathbf{Y}_{(2)} \langle \mathbf{Q} \rangle \Sigma_{\mathbf{b}}$
$\Sigma_{\mathbf{c}} = \langle \beta \rangle^{-1} (\langle \mathbf{T}^T \mathbf{T} \rangle + \langle \mathbf{Z} \rangle)^{-1}$
$\langle \mathbf{C} \rangle = \langle \beta \rangle \mathbf{Y}_{(3)} \langle \mathbf{T} \rangle \Sigma_{\mathbf{c}}$
$r = 1, 2, \dots, R_{\text{ini}}, l = 1, 2, \dots, L_{\text{ini}}$
$\langle s_{rl} \rangle = \sqrt{\frac{\langle \delta_{rl} \rangle}{\langle \beta \rangle \langle \zeta_r \rangle (\langle \mathbf{a}_{rl}^T \mathbf{a}_{rl} \rangle + \langle \mathbf{b}_{rl}^T \mathbf{b}_{rl} \rangle)}}$
$\langle \frac{1}{s_{rl}} \rangle = \frac{1}{\langle \delta_{rl} \rangle} + \frac{1}{\langle s_{rl} \rangle}$
$\langle \delta_{rl} \rangle = \frac{2\psi + I + J + 1}{2\tau + (\frac{1}{s_{rl}})}$
$r = 1, 2, \dots, R_{\text{ini}}$
$\langle \zeta_r \rangle = \sqrt{\frac{\langle \rho_r \rangle}{\langle \beta \rangle (\sum_{l=1}^{L_{\text{ini}}} \langle s_{rl} \rangle (\langle \mathbf{a}_{rl}^T \mathbf{a}_{rl} \rangle + \langle \mathbf{b}_{rl}^T \mathbf{b}_{rl} \rangle) + \langle \mathbf{c}_r^T \mathbf{c}_r \rangle)}}$
$\langle \frac{1}{\zeta_r} \rangle = \frac{1}{\langle \rho_r \rangle} + \frac{1}{\langle \zeta_r \rangle}$
$\langle \rho_r \rangle = \frac{2\mu + (I+J)L_{\text{ini}} + K + 1}{2\nu + (\frac{1}{\zeta_r})}$
$\langle \beta \rangle = (2\kappa + (I+J)L_{\text{ini}}R_{\text{ini}} + KR_{\text{ini}} + IJK + 1) / (2\theta +$
$\langle \ \mathbf{Y}_{(1)}^T - \mathbf{P}\mathbf{A}^T\ _{\mathbb{F}}^2 \rangle + \sum_{r=1}^{R_{\text{ini}}} \langle \zeta_r \rangle [\sum_{l=1}^{L_{\text{ini}}} \langle s_{rl} \rangle (\langle \mathbf{a}_{rl}^T \mathbf{a}_{rl} \rangle +$
$\langle \mathbf{b}_{rl}^T \mathbf{b}_{rl} \rangle) + \langle \mathbf{c}_r^T \mathbf{c}_r \rangle])$
<b>until convergence</b>

TABLE II  
FIRST- AND SECOND-ORDER STATISTICS REQUIRED IN BBTD

$\langle \mathbf{P} \rangle = \langle \mathbf{B} \rangle \odot \langle \mathbf{C} \rangle$
$\langle \mathbf{Q} \rangle = \langle \mathbf{C} \rangle \odot \langle \mathbf{A} \rangle$
$\langle \mathbf{T} \rangle = [(\langle \mathbf{A}_1 \rangle \odot \langle \mathbf{B}_1 \rangle) \mathbf{1}_{L_{\text{ini}}} \cdots (\langle \mathbf{A}_{R_{\text{ini}}} \rangle \odot \langle \mathbf{B}_{R_{\text{ini}}} \rangle) \mathbf{1}_{L_{\text{ini}}}]$
$\langle \mathbf{P}^T \mathbf{P} \rangle = \langle \mathbf{B}^T \mathbf{B} \rangle * (\langle \mathbf{C}^T \mathbf{C} \rangle \otimes \mathbf{1}_{L_{\text{ini}} \times L_{\text{ini}}})$
$\langle \mathbf{Q}^T \mathbf{Q} \rangle = \langle \mathbf{A}^T \mathbf{A} \rangle * (\langle \mathbf{C}^T \mathbf{C} \rangle \otimes \mathbf{1}_{L_{\text{ini}} \times L_{\text{ini}}})$
$\langle \mathbf{T}^T \mathbf{T} \rangle = (\mathbf{I}_{R_{\text{ini}}} \otimes \mathbf{1}_{L_{\text{ini}}}) (\langle \mathbf{A}^T \mathbf{A} \rangle * \langle \mathbf{B}^T \mathbf{B} \rangle) (\mathbf{I}_{R_{\text{ini}}} \otimes \mathbf{1}_{L_{\text{ini}}})$
$\langle \ \mathbf{Y}_{(1)} - \mathbf{P}\mathbf{A}^T\ _{\mathbb{F}}^2 \rangle = \ \mathbf{Y}_{(1)}\ _{\mathbb{F}}^2 - 2\text{tr}\{\langle \mathbf{A} \rangle^T \mathbf{Y}_{(1)}^T \langle \mathbf{P} \rangle\} + \text{tr}\{\langle \mathbf{A}^T \mathbf{A} \rangle \langle \mathbf{P}^T \mathbf{P} \rangle\}$

$R$  is estimated as the number of columns of  $\langle \mathbf{C} \rangle$  of non-negligible energy and similarly for the  $L_r$ s and the corresponding blocks of  $\langle \mathbf{A} \rangle, \langle \mathbf{B} \rangle$ . The iterations stop when a convergence criterion is met (e.g., the relative difference of the tensor reconstruction errors in two consecutive iterations becomes less than a user-defined threshold) or the maximum number of iterations is reached.

## V. SIMULATION RESULTS

In this section, we evaluate the effectiveness of the BBTD algorithm in revealing the ranks and computing the model parameters, in comparison with the regularization-based BTD-HIRLS scheme with its regularization parameter being selected so as to achieve the minimum approximation error. A  $55 \times 55 \times 20$  tensor  $\mathcal{Y}$  is generated as in (2), with  $R = 5$  and the  $L_r$ s being set as  $L_1 = 8, L_2 = 6, L_3 = 4, L_4 = 5$  and  $L_5 = 3$ . The entries of  $\mathbf{A}_r, \mathbf{B}_r$  and  $\mathbf{C}$  are i.i.d., sampled from the standard Gaussian distribution. The noise power is set so as to result in a signal-to-noise ratio  $\text{SNR} = 10 \log_{10} \|\mathcal{X}\|_{\mathbb{F}}^2 / (\sigma^2 \|\mathcal{N}\|_{\mathbb{F}}^2)$  of only 5 dB. Both  $R$

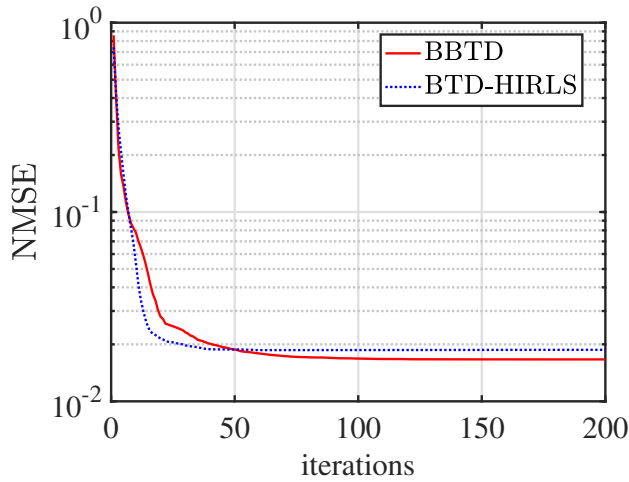


Fig. 2. NMSE vs. iterations.

and all  $L_r$ s are overestimated as 10. The evolution of the median of the  $\text{NMSE} = \sum_{r=1}^R \frac{\|\mathbf{A}_r \mathbf{B}_r^T \circ \mathbf{c}_r - \hat{\mathbf{A}}_r \hat{\mathbf{B}}_r^T \circ \hat{\mathbf{c}}_r\|_F^2}{\|\mathbf{A}_r \mathbf{B}_r^T \circ \mathbf{c}_r\|_F^2}$  values<sup>3</sup> obtained for 100 independent realizations of  $\mathcal{Y}$  is plotted in Fig. 2. It should be stressed that both BTD-HIRLS (as demonstrated in [2]) and BBTD appear to be insensitive to different initializations, hence a single random initialization is used here. Observe that both algorithms converge in less than 100 iterations.

Fig. 3 depicts the success rates in the recovery of the true  $R$  and  $L_r$ s. The Bayesian algorithm is seen to recover the correct number of block terms in this scenario with a slightly lower probability than the state-of-the-art BTD-HIRLS algorithm. It should be noted however that in the rare cases that BBTD fails, this happens only with a small deviation from the true  $R = 5$ , namely 6. As it can be observed from Figs. 2(b)–(f), where the success rates in recovering the individual ranks given a correctly estimated  $R$  are shown, BBTD performs similarly with or better than BTD-HIRLS, achieving almost 100% accuracy in all terms. Instead BTD-HIRLS exhibits lower accuracy in estimating  $L_1 = 8$ . That BBTD performs *overall* (in estimating *both*  $R$  and the  $L_r$ s) better than BTD-HIRLS can be attributed to the probabilistic model that is adopted, which allows BBTD to better capture the structure of BTD when it comes to the low-rankness of the  $\mathbf{A}_r, \mathbf{B}_r$ . Namely, BBTD implicitly weighs differently each pair  $\mathbf{A}_r, \mathbf{B}_r$ , as opposed to BTD-HIRLS, which uniformly weighs all block terms through a single regularization parameter. This also explains the fact that the BTD-HIRLS estimation for  $R$  is relatively more accurate.

## REFERENCES

- [1] L. De Lathauwer, “Decompositions of a higher-order tensor in block terms — Part II: Definitions and uniqueness,” *SIAM J. Matrix Anal. Appl.*, vol. 30, no. 3, pp. 1033–1066, 2008.

<sup>3</sup>The Hungarian algorithm is employed to match the  $\hat{R}$  estimated non-zero block terms with the true ones.

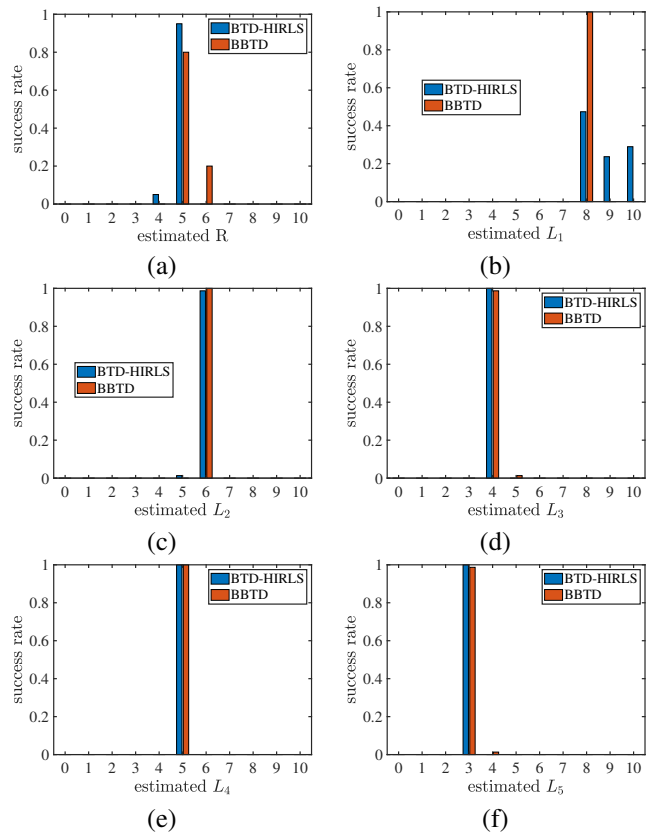


Fig. 3. Success rates of recovering (a)  $R$  and (b)–(f)  $L_r$ s via the BBTD and BTD-HIRLS algorithms.

- [2] A. A. Rontogiannis, E. Kofidis, and P. V. Giampouras, “Block-term tensor decomposition: Model selection and computation,” *IEEE J. Sel. Topics Signal Process.*, Jan. 2021.
- [3] L. De Lathauwer, “Block component analysis: a new concept for blind source separation,” in *Proc. LVA/ICA-2012*, Tel Aviv, Israel, Mar. 2012.
- [4] Y. Qian *et al.*, “Matrix-vector nonnegative tensor factorization for blind unmixing of hyperspectral imagery,” *IEEE Trans. Geosci. Remote Sens.*, vol. 55, no. 3, pp. 1776–1792, Mar. 2017.
- [5] S. Theodoridis, *Machine Learning — A Bayesian and Optimization Perspective*, 2nd ed. Academic Press, 2020.
- [6] D. G. Tzikas, A. C. Likas, and N. P. Galatsanos, “The variational approximation for Bayesian inference,” *IEEE Signal Process. Mag.*, pp. 131–146, Nov. 2008.
- [7] Q. Zhao, L. Zhang, and A. Cichocki, “Bayesian sparse Tucker models for dimension reduction and tensor completion,” arXiv:1505.02343v1 [cs.LG], May 2015.
- [8] L. Cheng *et al.*, “Towards probabilistic tensor canonical polyadic decomposition 2.0: Automatic tensor rank learning using generalized hyperbolic prior,” arXiv:2009.02472v1 [cs.LG], Sep. 2020.
- [9] L. Xu *et al.*, “Learning tensor train representation with automatic rank determination from incomplete noisy data,” arXiv:2010.06564v1 [eess.SP], Oct. 2020.
- [10] M. E. Tipping, “Sparse Bayesian learning and the relevance vector machine,” *J. Mach. Learn. Res.*, vol. 1, pp. 211–244, Jun. 2001.
- [11] P. V. Giampouras *et al.*, “Online sparse and low-rank subspace learning from incomplete data: A Bayesian view,” *Signal Process.*, vol. 137, pp. 199–212, 2017.
- [12] P. V. Giampouras, A. A. Rontogiannis, and E. Kofidis, “A Bayesian approach to block-term tensor decomposition model selection and computation,” arXiv:2101.02931v1 [stat.ME], Jan. 2021.