

Scalable Tensor Factorizations for Incomplete Data[☆]

Evrin Acar^{a,1}, Daniel M. Dunlavy^{b,1}, Tamara G. Kolda^{a,1,*}, Morten Mørup^c

^a*Sandia National Laboratories, Livermore, CA 94551-9159.*

^b*Sandia National Laboratories, Albuquerque, NM 87123-1318.*

^c*Technical University of Denmark, 2800 Kgs. Lyngby, Denmark.*

Abstract

The problem of incomplete data—i.e., data with missing or unknown values—in multi-way arrays is ubiquitous in biomedical signal processing, network traffic analysis, bibliometrics, social network analysis, chemometrics, computer vision, communication networks, etc. We consider the problem of how to factorize data sets with missing values with the goal of capturing the underlying latent structure of the data and possibly reconstructing missing values (i.e., tensor completion). We focus on one of the most well-known tensor factorizations that captures multi-linear structure, CANDECOMP/PARAFAC (CP). In the presence of missing data, CP can be formulated as a weighted least squares problem that models *only* the known entries. We develop an algorithm called CP-WOPT (CP Weighted OPTimization) that uses a first-order optimization approach to solve the weighted least squares problem. Based on extensive numerical experiments, our algorithm is shown to successfully factorize tensors with noise and up to 99% missing data. A unique aspect of our approach is that it scales to sparse large-scale data, e.g., $1000 \times 1000 \times 1000$ with five million known entries (0.5% dense). We further demonstrate the usefulness of CP-WOPT on two real-world applications: a novel EEG (electroencephalogram) application where missing data is frequently encountered due to disconnections of electrodes and the problem of modeling computer network traffic where data may be absent due to the expense of the data collection process.

Keywords: missing data, incomplete data, tensor factorization, CANDECOMP, PARAFAC, optimization

[☆]A preliminary version of this paper has appeared as [1].

*Corresponding author

Email addresses: eacarat@sandia.gov (Evrin Acar), dmdunla@sandia.gov (Daniel M. Dunlavy), tgtkolda@sandia.gov (Tamara G. Kolda), mm@imm.dtu.dk (Morten Mørup)

¹This work was funded by the Laboratory Directed Research & Development (LDRD) program at Sandia National Laboratories, a multiprogram laboratory operated by Sandia Corporation, a Lockheed Martin Company, for the United States Department of Energy's National Nuclear Security Administration under Contract DE-AC04-94AL85000.

1. Introduction

Missing data can arise in a variety of settings due to loss of information, errors in the data collection process, or costly experiments. For instance, in biomedical signal processing, missing data can be encountered during EEG analysis, where multiple electrodes are used to collect the electrical activity along the scalp. If one of the electrodes becomes loose or disconnected, the signal is either lost or discarded due to contamination with high amounts of mechanical noise. We also encounter the missing data problem in other areas of data mining, such as packet losses in network traffic analysis [2] and occlusions in images in computer vision [3]. Many real-world data with missing entries are ignored because they are deemed unsuitable for analysis, but this work contributes to the growing evidence that such data can be analyzed.

Unlike most previous studies on missing data which have only considered matrices, we focus here on the problem of missing data in *tensors* because it has been shown increasingly that data often have more than two modes of variation and are therefore best represented as multi-way arrays (i.e., tensors) [4, 5]. For instance, in EEG data each signal from an electrode can be represented as a time-frequency matrix; thus, data from multiple channels is three-dimensional (temporal, spectral, and spatial) and forms a three-way array [6]. Social network data, network traffic data, and bibliometric data are of interest to many applications such as community detection, link mining, and more; these data can have multiple dimensions/modalities, are often extremely large, and generally have at least some missing data. These are just a few of the many data analysis applications where one needs to deal with large multi-way arrays with missing entries. Other examples of multi-way arrays with missing entries from different disciplines have also been studied in the literature [7, 8, 9]. For instance, [7] shows that, in spectroscopy, intermittent machine failures or different sampling frequencies may result in tensors with missing fibers (i.e., the higher-order analogues of matrix rows or columns, see Figure 1). Similarly, missing fibers are encountered in multidimensional NMR (Nuclear Magnetic Resonance) analysis, where sparse sampling is used in order to reduce the experimental time [8].

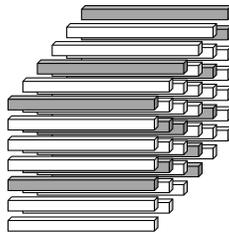


Figure 1: A 3-way tensor with missing row fibers (in gray).

Our goal is to capture the latent structure of the data via a higher-order factorization, even in the presence of missing data. Handling missing data in the context of matrix factorizations, e.g., the widely-used principal component

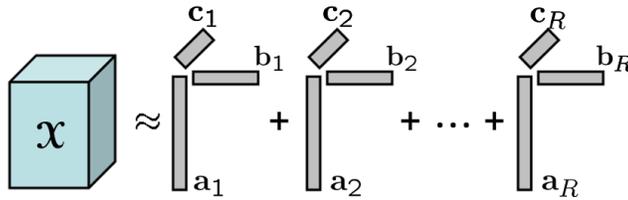


Figure 2: Illustration of an R -component CP model for a third-order tensor \mathcal{X} .

analysis, has long been studied [10, 11] (see [3] for a review). It is also closely related to the matrix completion problem, where the goal is to recover the missing entries [12, 13] (see §3 for more discussion). Higher-order factorizations, i.e., tensor factorizations, have emerged as an important method for information analysis [4, 5]. Instead of flattening (unfolding) multi-way arrays into matrices and using matrix factorization techniques, tensor models preserve the multi-way nature of the data and extract the underlying factors in each mode (dimension) of a higher-order array.

We focus here on the CANDECOMP/PARAFAC (CP) tensor decomposition [14, 15], which is a tensor model commonly used in various applications [6, 16, 17, 18, 19]. To illustrate differences between matrix and tensor factorizations, we introduce the CP decomposition for three-way tensors; discussion of the CP decomposition for general N -way tensors can be found in §4. Let \mathcal{X} be a three-way tensor of size $I \times J \times K$, and assume its rank is R (see [5] for a detailed discussion on tensor rank). With perfect data, the CP decomposition is defined by *factor matrices* \mathbf{A} , \mathbf{B} , and \mathbf{C} of sizes $I \times R$, $J \times R$, and $K \times R$, respectively, such that

$$x_{ijk} = \sum_{r=1}^R a_{ir} b_{jr} c_{kr}, \quad \text{for all } i = 1, \dots, I, j = 1, \dots, J, k = 1, \dots, K.$$

In the presence of noise, the true \mathcal{X} is not observable and we cannot expect equality. Instead, the CP decomposition should minimize the error function

$$f(\mathbf{A}, \mathbf{B}, \mathbf{C}) = \frac{1}{2} \sum_{i=1}^I \sum_{j=1}^J \sum_{k=1}^K \left(x_{ijk} - \sum_{r=1}^R a_{ir} b_{jr} c_{kr} \right)^2. \quad (1)$$

An illustration of CP for third-order tensors is given in Figure 2. The CP decomposition is extensible to N -way tensors for $N \geq 3$, and there are numerous methods for computing it [20].

In the case of incomplete data, a standard practice is to impute the missing values in some fashion (e.g., replacing the missing entries using average values along a particular mode). Imputation can be useful as long as the amount of missing data is small; however, performance degrades for large amounts of missing data [10] (also see §5.5). As a better alternative, factorizations of the data with imputed values for missing entries can be used to re-impute the missing

values and the procedure can be repeated to iteratively determine suitable values for the missing entries. Such a procedure is an example of the expectation maximization (EM) algorithm [21]. Computing CP decompositions by combining the alternating least squares method, which computes the factor matrices one at a time, and iterative imputation (denoted EM-ALS in this paper) has been shown to be quite effective and has the advantage of often being simple and fast. Nevertheless, as the amount of missing data increases, the performance of the algorithm may suffer since the initialization and the intermediate models used to impute the missing values will increase the risk of converging to a less than optimal factorization [7]. Also, the poor convergence of alternating methods due to their vulnerability to flatlining, i.e., stagnation, is noted in [3].

In this paper, though, we focus on using a weighted version of the error function to ignore missing data and model only the known entries. In that case, nonlinear optimization can be used to directly solve the weighted least squares problem for the CP model. The weighted version of (1) is

$$f_{\mathbf{W}}(\mathbf{A}, \mathbf{B}, \mathbf{C}) = \frac{1}{2} \sum_{i=1}^I \sum_{j=1}^J \sum_{k=1}^K \left\{ w_{ijk} \left(x_{ijk} - \sum_{r=1}^R a_{ir} b_{jr} c_{kr} \right) \right\}^2, \quad (2)$$

where \mathbf{W} , which is the same size as \mathbf{X} , is a nonnegative weight tensor defined as

$$w_{ijk} = \begin{cases} 1 & \text{if } x_{ijk} \text{ is known,} \\ 0 & \text{if } x_{ijk} \text{ is missing,} \end{cases} \quad \text{for all } i = 1, \dots, I, j = 1, \dots, J, k = 1, \dots, K.$$

Our contributions in this paper are summarized as follows. (a) We develop a scalable algorithm called CP-WOPT (CP Weighted OPTimization) for tensor factorizations in the presence of missing data. CP-WOPT uses first-order optimization to solve the weighted least squares objective function over all the factor matrices simultaneously. (b) We show that CP-WOPT can scale to sparse, large-scale data using specialized sparse data structures, significantly reducing the storage and computation costs. (c) Using extensive numerical experiments on simulated data sets, we show that CP-WOPT can successfully factor tensors with noise and up to 99% missing data. In many cases, CP-WOPT is significantly faster than the best published method in the literature [7]. (d) We demonstrate the applicability of the proposed algorithm on a real data set in a novel EEG application where data is incomplete due to failures of particular electrodes. This is a common occurrence in practice, and our experiments show that even if signals from almost half of the channels are missing, underlying brain activities can still be captured using the CP-WOPT algorithm, illustrating the usefulness of our proposed method. (e) In addition to tensor factorizations, we also show that CP-WOPT can be used to address the tensor completion problem in the context of network traffic analysis. We use the factors captured by the CP-WOPT algorithm to reconstruct the tensor and illustrate that even if there is a large amount of missing data, the algorithm is able to keep the relative error in the missing entries close to the modeling error.

The paper is organized as follows. We introduce the notation used throughout the paper in §2. In §3, we discuss related work in matrix and tensor factorizations. The computation of the function and gradient values for the general N -way weighted version of the error function and the presentation of the CP-WOPT method are given in §4. Numerical results on both simulated and real data are given in §5. Conclusions and future work are discussed in §6.

2. Notation

Tensors of order $N \geq 3$ are denoted by Euler script letters ($\mathbf{X}, \mathbf{Y}, \mathbf{Z}$), matrices are denoted by boldface capital letters ($\mathbf{A}, \mathbf{B}, \mathbf{C}$), vectors are denoted by boldface lowercase letters ($\mathbf{a}, \mathbf{b}, \mathbf{c}$), and scalars are denoted by lowercase letters (a, b, c). Columns of a matrix are denoted by boldface lower letters with a subscript ($\mathbf{a}_1, \mathbf{a}_2, \mathbf{a}_3$ are first three columns of \mathbf{A}). Entries of a matrix or a tensor are denoted by lowercase letters with subscripts, i.e., the (i_1, i_2, \dots, i_N) entry of an N -way tensor \mathbf{X} is denoted by $x_{i_1 i_2 \dots i_N}$.

An N -way tensor can be rearranged as a matrix; this is called *matricization*, also known as *unfolding* or *flattening*. The mode- n matricization of a tensor $\mathbf{X} \in \mathbb{R}^{I_1 \times I_2 \times \dots \times I_N}$ is denoted by $\mathbf{X}_{(n)}$ and arranges the mode- n one-dimensional “fibers” to be the columns of the resulting matrix. Specifically, tensor element (i_1, i_2, \dots, i_N) maps to matrix element (i_n, j) where

$$j = 1 + \sum_{\substack{k=1 \\ k \neq n}}^N (i_k - 1) J_k, \quad \text{with} \quad J_k = \begin{cases} 1, & \text{if } k = 1 \text{ or if } k = 2 \text{ and } n = 1, \\ \prod_{\substack{m=1 \\ m \neq n}}^{k-1} I_m, & \text{otherwise.} \end{cases}$$

Given two tensors \mathbf{X} and \mathbf{Y} of equal size $I_1 \times I_2 \times \dots \times I_N$, their Hadamard (elementwise) product is denoted by $\mathbf{X} * \mathbf{Y}$ and defined as

$$(\mathbf{X} * \mathbf{Y})_{i_1 i_2 \dots i_N} = x_{i_1 i_2 \dots i_N} y_{i_1 i_2 \dots i_N} \quad \text{for all } i_n \in \{1, \dots, I_n\} \text{ and } n \in \{1, \dots, N\}$$

The *inner product* of two same-sized tensors $\mathbf{X}, \mathbf{Y} \in \mathbb{R}^{I_1 \times I_2 \times \dots \times I_N}$ is the sum of the products of their entries, i.e.,

$$\langle \mathbf{X}, \mathbf{Y} \rangle = \sum_{i_1=1}^{I_1} \sum_{i_2=1}^{I_2} \dots \sum_{i_N=1}^{I_N} x_{i_1 i_2 \dots i_N} y_{i_1 i_2 \dots i_N}.$$

For a tensor \mathbf{X} of size $I_1 \times I_2 \times \dots \times I_N$, its *norm* is

$$\|\mathbf{X}\| = \sqrt{\langle \mathbf{X}, \mathbf{X} \rangle} = \sqrt{\sum_{i_1=1}^{I_1} \sum_{i_2=1}^{I_2} \dots \sum_{i_N=1}^{I_N} x_{i_1 i_2 \dots i_N}^2}.$$

For matrices (i.e., second-order tensors), $\|\cdot\|$ refers to the analogous Frobenius norm, and for vectors (i.e., first-order tensors), $\|\cdot\|$ refers to the analogous two-norm. We can also define a weighted norm as follows. Let \mathbf{X} and \mathbf{W} be two

tensors of size $I_1 \times I_2 \times \dots \times I_N$. Then the \mathbf{W} -weighted norm of \mathbf{X} is

$$\|\mathbf{X}\|_{\mathbf{W}} = \|\mathbf{W} * \mathbf{X}\|.$$

The weighted matrix norm is analogous.

Given a sequence of matrices $\mathbf{A}^{(n)}$ of size $I_n \times R$ for $n = 1, \dots, N$, the notation $\llbracket \mathbf{A}^{(1)}, \mathbf{A}^{(2)}, \dots, \mathbf{A}^{(N)} \rrbracket$ defines an $I_1 \times I_2 \times \dots \times I_N$ tensor whose elements are given by

$$\left(\llbracket \mathbf{A}^{(1)}, \mathbf{A}^{(2)}, \dots, \mathbf{A}^{(N)} \rrbracket \right)_{i_1 i_2 \dots i_N} = \sum_{r=1}^R \prod_{n=1}^N a_{i_n r}^{(n)},$$

for all $i_n \in \{1, \dots, I_n\}$ and $n \in \{1, \dots, N\}$.

For just two matrices, this reduces to familiar expressions: $\llbracket \mathbf{A}, \mathbf{B} \rrbracket = \mathbf{A}\mathbf{B}^\top$. Using the notation defined here, (2) can be rewritten as

$$f_{\mathbf{W}}(\mathbf{A}, \mathbf{B}, \mathbf{C}) = \frac{1}{2} \|\mathbf{X} - \llbracket \mathbf{A}, \mathbf{B}, \mathbf{C} \rrbracket\|_{\mathbf{W}}^2.$$

3. Related Work in Factorizations with Missing Data

In this section, we first review the approaches for handling missing data in matrix factorizations and then discuss how these techniques have been extended to tensor factorizations.

3.1. Matrix Factorizations

Matrix factorization in the presence of missing entries is a problem that has been studied for several decades; see, e.g., [10, 11]. The problem is typically formulated analogously to (2) as

$$f_{\mathbf{W}}(\mathbf{A}, \mathbf{B}) = \frac{1}{2} \left\| \mathbf{X} - \mathbf{A}\mathbf{B}^\top \right\|_{\mathbf{W}}^2. \quad (3)$$

A common procedure for solving this problem is EM which combines imputation and alternation [7, 22]. In this approach, the missing values of \mathbf{X} are imputed using the current model, $\hat{\mathbf{X}} = \mathbf{A}\mathbf{B}^\top$, as follows:

$$\bar{\mathbf{X}} = \mathbf{W} * \mathbf{X} + (\mathbf{1} - \mathbf{W}) * \hat{\mathbf{X}},$$

where $\mathbf{1}$ is the matrix of all ones. Once $\bar{\mathbf{X}}$ is generated, the matrices \mathbf{A} and \mathbf{B} can then be alternately updated according to the error function $\frac{1}{2} \|\bar{\mathbf{X}} - \mathbf{A}\mathbf{B}^\top\|^2$ (e.g., using the linear least squares method). See [22, 23] for further discussion of the EM method in the missing data and general weighted case.

Recently, a direct nonlinear optimization approach was proposed for matrix factorization with missing data [3]. In this case, (3) is solved directly using a 2nd-order damped Newton method. This new method is compared to other standard techniques based on some form of alternation and/or imputation as

well as hybrid techniques that combine both approaches. Overall, the conclusion is that nonlinear optimization strategies are key to successful matrix factorization. Moreover, the authors observe that the alternating methods tend to take much longer to converge to the solution even though they make faster progress initially. This work is theoretically the most related to what we propose—the main differences are 1) we focus on tensors rather than matrices, and 2) we use first-order rather than second-order optimization methods (we note that first order methods are mentioned as future work in [3]).

A major difference between matrix and tensor factorizations is worth noting here. In [22, 3], the lack of uniqueness in matrix decompositions is discussed. Given any invertible matrix \mathbf{G} , $[[\mathbf{A}, \mathbf{B}]] = [[\mathbf{A}\mathbf{G}, \mathbf{B}\mathbf{G}^{-\top}]]$. This means that there is an infinite family of equivalent solutions to (3). In [3], regularization is recommended as a partial solution; however regularization can only control scaling indeterminacies and not rotational freedom. In the case of the CP model, often a unique solution (including trivial indeterminacies of scaling and column permutation) can be recovered exactly; see, e.g., [5] for further discussion on uniqueness of the CP decomposition.

Factorization of matrices with missing entries is also closely related to the matrix completion problem. In matrix completion, one tries to recover the missing matrix entries using the low-rank structure of the matrix. Recent work in this area [12, 13] shows that even if a small amount of matrix entries are available and those are corrupted with noise, it is still possible to recover the missing entries up to the level of noise. In [13], it is also discussed how this problem relates to the field of compressive sensing, which exploits structures in data to generate more compact representations of the data. Practically speaking, the difference between completion and factorization is how success is measured. Factorization methods aim to increase accuracy in the factors; in other words, capture the underlying phenomena as well as possible. Completion methods, on the other hand, seek accuracy in filling in the missing data. Obviously, once a factorization has been computed, it can be used to reconstruct the missing entries. In fact, many completion methods use this procedure.

3.2. Tensor Factorizations

The EM procedure discussed for matrices has also been widely employed for tensor factorizations with missing data. If the current model is $[[\mathbf{A}, \mathbf{B}, \mathbf{C}]]$, then we fill in the missing entries of \mathcal{X} to produce a complete tensor according to

$$\bar{\mathcal{X}} = \mathcal{W} * \mathcal{X} + (\mathbf{1} - \mathcal{W}) * [[\mathbf{A}, \mathbf{B}, \mathbf{C}]],$$

where $\mathbf{1}$ is the tensor of all ones the same size as \mathcal{W} . The factor matrices are then updated using alternating least squares (ALS) as those that best fit $\bar{\mathcal{X}}$. See, e.g., [24, 25] for further details. As noted previously, we denote this method as EM-ALS.

Paatero [26] and Tomasi and Bro [7] have investigated direct nonlinear approaches based on Gauss-Newton (GN). The code from [26] is not widely available; therefore, we focus on [7] and its INDAFAC (INcomplete DAta paraFAC)

procedure which specifically uses the Levenberg-Marquardt version of GN for fitting the CP model to data with missing entries. The primary application in [7] is missing data in chemometrics experiments. This approach is compared to EM-ALS with the result being that INDAFAC and EM-ALS perform almost equally well in general with the exception that INDAFAC is more accurate for difficult problems, i.e., higher collinearity and systematically missing patterns of data. In terms of computational efficiency, EM-ALS is usually faster but becomes slower than INDAFAC as the percentage of missing entries increases and also depending on the missing entry patterns.

Both INDAFAC and CP-WOPT address the problem of fitting the CP model to incomplete data sets by solving (2). The difference is that INDAFAC is based on second-order optimization while CP-WOPT is first-order with a goal of scaling to larger problem sizes.

4. CP-WOPT Algorithm

We consider the general N -way CP factorization problem for tensors with missing entries. Let \mathcal{X} be a real-valued tensor of size $I_1 \times I_2 \times \dots \times I_N$ and assume its rank is known to be R .² Define a nonnegative weight tensor \mathcal{W} of the same size as \mathcal{X} such that

$$w_{i_1 i_2 \dots i_N} = \begin{cases} 1 & \text{if } x_{i_1 i_2 \dots i_N} \text{ is known,} \\ 0 & \text{if } x_{i_1 i_2 \dots i_N} \text{ is missing,} \end{cases}$$

for all $i_n \in \{1, \dots, I_n\}$ and $n \in \{1, \dots, N\}$. (4)

The N -way objective function is defined by

$$f_{\mathcal{W}}(\mathbf{A}^{(1)}, \mathbf{A}^{(2)}, \dots, \mathbf{A}^{(N)}) = \frac{1}{2} \left\| \left(\mathcal{X} - \llbracket \mathbf{A}^{(1)}, \dots, \mathbf{A}^{(N)} \rrbracket \right) \right\|_{\mathcal{W}}^2. \quad (5)$$

Due to the well-known indeterminacies of the CP model, it may also be desirable to add regularization to the objective function as in [20], but this has not been necessary thus far in our experiments.

Equation (5) is equivalent to

$$f_{\mathcal{W}}(\mathbf{A}^{(1)}, \mathbf{A}^{(2)}, \dots, \mathbf{A}^{(N)}) = \frac{1}{2} \|\mathcal{Y} - \mathcal{Z}\|^2,$$

where $\mathcal{Y} = \mathcal{W} * \mathcal{X}$ and $\mathcal{Z} = \mathcal{W} * \llbracket \mathbf{A}^{(1)}, \dots, \mathbf{A}^{(N)} \rrbracket$. (6)

The tensor \mathcal{Y} can be precomputed as neither \mathcal{W} nor \mathcal{X} change during the iterations. We will see later that \mathcal{Z} can be computed efficiently for sparse \mathcal{W} .

²In practice, the rank is generally not known and is not easily determined. Understanding the performance of the methods under consideration in that scenario is a topic of future work. Results in [20] indicate that direct optimization methods have an advantage over alternating least squares approaches when the rank is overestimated.

Our goal is to find matrices $\mathbf{A}^{(n)} \in \mathbb{R}^{I_n \times R}$ for $n = 1, \dots, N$ that minimize the weighted objective function in (5). This objective function can be considered as a mapping from the cross product of N two-dimensional vector spaces to \mathbb{R} , i.e.,

$$f_{\mathcal{W}} : \mathbb{R}^{I_1 \times R} \otimes \mathbb{R}^{I_2 \times R} \otimes \dots \otimes \mathbb{R}^{I_N \times R} \mapsto \mathbb{R}.$$

Although $f_{\mathcal{W}}$ is written as a function of matrices, it can be thought of as a vector function where the parameter vector contains the vectorized and stacked matrices $\mathbf{A}^{(1)}$ through $\mathbf{A}^{(N)}$, i.e.,

$$\left[\mathbf{a}_1^{(1)\top} \quad \dots \quad \mathbf{a}_R^{(1)\top} \quad \dots \quad \mathbf{a}_1^{(N)\top} \quad \dots \quad \mathbf{a}_R^{(N)\top} \right]^\top. \quad (7)$$

In this view, $f_{\mathcal{W}} : \mathbb{R}^P \mapsto \mathbb{R}$, where $P = R \sum_{n=1}^N I_n$. We derive the gradient of $f_{\mathcal{W}}$ in §4.1, and we show how to compute the function and gradient values efficiently when \mathcal{W} is dense in §4.2 and \mathcal{W} is sparse in §4.3. Once the function and gradient are known, any gradient-based optimization method [27] can be used to solve the optimization problem.

4.1. Gradient

We derive the gradient of (5) by computing the partial derivatives of $f_{\mathcal{W}}$ with respect to each element of the factor matrices, i.e., $a_{i_n r}^{(n)}$ for all $i_n \in \{1, \dots, I_n\}$, $n \in \{1, \dots, N\}$, and $r \in \{1, \dots, R\}$. We first rewrite the objective function in (5) as

$$f_{\mathcal{W}}(\mathbf{A}^{(1)}, \mathbf{A}^{(2)}, \dots, \mathbf{A}^{(N)}) = \frac{1}{2} \sum_{i_1=1}^{I_1} \sum_{i_2=1}^{I_2} \dots \sum_{i_N=1}^{I_N} w_{i_1 i_2 \dots i_N}^2 \left\{ x_{i_1 i_2 \dots i_N}^2 - 2x_{i_1 i_2 \dots i_N} \sum_{r=1}^R \prod_{n=1}^N a_{i_n r}^{(n)} + \left(\sum_{r=1}^R \prod_{n=1}^N a_{i_n r}^{(n)} \right)^2 \right\}.$$

Using this expression, it is clear to see that the partial derivatives are given by

$$\frac{\partial f_{\mathcal{W}}}{\partial a_{i_n r}^{(n)}} = \sum_{i_1=1}^{I_1} \dots \sum_{i_{n-1}=1}^{I_{n-1}} \sum_{i_{n+1}=1}^{I_{n+1}} \dots \sum_{i_N=1}^{I_N} w_{i_1 i_2 \dots i_N}^2 \left(-x_{i_1 i_2 \dots i_N} + \sum_{l=1}^R \prod_{m=1}^N a_{i_m l}^{(m)} \right) \prod_{\substack{m=1 \\ m \neq n}}^N a_{i_m r}^{(m)}$$

for all $i_n \in \{1, \dots, I_n\}, n \in \{1, \dots, N\}, r \in \{1, \dots, R\}$.

In matrix notation, using \mathbf{Y} and \mathbf{Z} from (6) and exploiting the fact that \mathcal{W} is binary, we can rewrite the gradient equation as

$$\frac{\partial f_{\mathcal{W}}}{\partial \mathbf{A}^{(n)}} = (\mathbf{Z}^{(n)} - \mathbf{Y}^{(n)}) \mathbf{A}^{(-n)}, \quad (8)$$

where

$$\mathbf{A}^{(-n)} = \mathbf{A}^{(N)} \odot \dots \odot \mathbf{A}^{(n+1)} \odot \mathbf{A}^{(n-1)} \odot \dots \odot \mathbf{A}^{(1)}$$

for $n = 1, \dots, N$. The symbol \odot denotes the Khatri-Rao product and is defined as follows for two matrices \mathbf{A} and \mathbf{B} of sizes $I \times K$ and $J \times K$ (both have the same number of columns):

$$\mathbf{A} \odot \mathbf{B} = [\mathbf{a}_1 \otimes \mathbf{b}_1 \quad \mathbf{a}_2 \otimes \mathbf{b}_2 \quad \dots \quad \mathbf{a}_N \otimes \mathbf{b}_N]$$

where \otimes denotes the vector Kronecker product.

4.2. Computations with Dense \mathcal{W}

Figure 3a shows the algorithmic steps for computing the function and gradient values. Recall that we are optimizing a function of P variables as defined by the factor matrices $\mathbf{A}^{(1)}$ through $\mathbf{A}^{(N)}$. From an optimization point of view, these are stacked to form a single vector of length P as in (7). The gradient will also be a vector of length P , but we compute it as a series of matrices $\mathbf{G}^{(n)} \equiv \frac{\partial f}{\partial \mathbf{A}^{(n)}}$ for $n = 1, \dots, N$. Once again, however, we can vectorize and stack the N matrices to get a vector of length P :

$$\left[\mathbf{g}_1^{(1)\top} \quad \dots \quad \mathbf{g}_R^{(1)\top} \quad \dots \quad \mathbf{g}_1^{(N)\top} \quad \dots \quad \mathbf{g}_R^{(N)\top} \right]^\top.$$

4.3. Computations with Sparse \mathcal{W}

If a large number of entries is missing, then \mathcal{W} is sparse. In this case, there is no need to allocate storage for every entry of the tensor \mathcal{X} . Instead, we can store and work with just the known values, making the method efficient in both storage and time. Figure 3b shows analogous computations to Figure 3a in the case that \mathcal{W} is sparse.

Let the ordered set $\mathcal{I} = \{\mathbf{i}^{(1)}, \mathbf{i}^{(2)}, \dots, \mathbf{i}^{(Q)}\}$ be the indices of the known values, i.e., all the locations where $\mathcal{W} = 1$. Each $\mathbf{i}^{(q)}$, $q \in \{1, \dots, Q\}$, is an N -tuple of indices whose n th entry is denoted by $i_n^{(q)}$. The known entries of \mathcal{X} can be stored in an array \mathbf{y} of length Q so that

$$y_q = x_{i_1^{(q)}, i_2^{(q)}, \dots, i_N^{(q)}} \quad \text{for } q = 1, \dots, Q.$$

Observe that $\|\mathcal{W} * \mathcal{X}\|^2 = \|\mathbf{y}\|^2$. Thus, the value of γ is the same in both Figure 3a and Figure 3b.

In the dense version (Figure 3a), we need to compute $\mathcal{Z} = \mathcal{W} * [\mathbf{A}^{(1)}, \dots, \mathbf{A}^{(N)}]$. In the sparse version, we observe that \mathcal{Z} is necessarily zero anywhere there is a missing entry in \mathcal{X} . Consequently, it is only necessary to compute the entries of \mathcal{Z} corresponding to known values in \mathcal{I} ; the vector \mathbf{z} corresponds to these known entries. The computations for Step 2 in Figure 3b can be done efficiently in MATLAB using the “expanded” vectors technique of [28, §3.2.4]. Let the r th summand be denoted by \mathbf{u} , i.e.,

$$u_q = \prod_{n=1}^N a_{i_n^{(q)} r}^{(n)} \quad \text{for } q = 1, \dots, Q.$$

1. Assume $\mathbf{y} = \mathbf{W} * \mathbf{X}$ and $\gamma = \|\mathbf{y}\|^2$ are precomputed.
2. Compute $\mathbf{Z} = \mathbf{W} * \llbracket \mathbf{A}^{(1)}, \dots, \mathbf{A}^{(N)} \rrbracket$.
3. Compute function value: $f = \frac{1}{2}\gamma - \langle \mathbf{y}, \mathbf{Z} \rangle + \frac{1}{2} \|\mathbf{Z}\|^2$.
4. Compute $\mathcal{J} = \mathbf{y} - \mathbf{Z}$.
5. Compute gradient matrices: $\mathbf{G}^{(n)} = -\mathbf{T}_{(n)}\mathbf{A}^{(-n)}$ for $n = 1, \dots, N$.

(a) Dense \mathbf{W} .

1. Let $\mathcal{I} = \{\mathbf{i}^{(1)}, \mathbf{i}^{(2)}, \dots, \mathbf{i}^{(Q)}\}$ be an ordered set of all the locations where $\mathbf{W} = 1$, i.e., the indices of the known values. Let \mathbf{y} be the length- Q vector of the values of \mathbf{X} at the locations indicated by \mathcal{I} . Assume \mathbf{y} and $\gamma = \|\mathbf{y}\|^2$ are precomputed.

2. Compute the Q -vector \mathbf{z} as

$$z_q = \sum_{r=1}^R \prod_{n=1}^N a_{i_n^{(q)} r}^{(n)} \quad \text{for } q = 1, \dots, Q.$$

3. Compute function value: $f = \frac{1}{2}\gamma - \mathbf{y}^\top \mathbf{z} + \frac{1}{2} \|\mathbf{z}\|^2$.
4. Compute $\mathbf{t} = \mathbf{y} - \mathbf{z}$.
5. Compute gradient matrices $\mathbf{G}^{(n)}$ for $n = 1, \dots, N$ as follows:

$$g_{jr}^{(n)} = - \sum_{\substack{q=1 \\ q: i_n^{(q)}=j}}^Q \left(t_q \prod_{\substack{m=1 \\ m \neq n}}^N a_{i_m^{(q)} r}^{(m)} \right).$$

(b) Sparse \mathbf{W} .

Figure 3: CP-WOPT computation of function value and gradient.

Let $\mathbf{v}^{(n)}$ be “expanded” vectors of length Q for $n = 1, \dots, N$ defined by

$$v_q^{(n)} = a_{i_q^{(n)}r}^{(n)} \quad \text{for } q = 1, \dots, Q. \quad (9)$$

Then the vector \mathbf{u} can be calculated as

$$\mathbf{u} = \mathbf{v}^{(1)} * \mathbf{v}^{(2)} * \dots * \mathbf{v}^{(N)}.$$

This can naturally be done iteratively (i.e., only one $\mathbf{v}^{(n)}$ is calculated at a time) to reduce storage costs. Likewise, each summand \mathbf{u} can be iteratively added to compute the final answer \mathbf{z} .

In Step 3, the function values in Figure 3a and Figure 3b are clearly equivalent since \mathbf{y} and \mathbf{z} contain the nonzero values of \mathbf{Y} and \mathbf{Z} , respectively. Similarly, the vector \mathbf{t} in Step 4 of Figure 3b represents just the nonzero values of \mathcal{T} in Figure 3a.

The computation of the gradients in Step 6 of Figure 3b performs a matricized-tensor-times-Khatri-Rao-product (mttkrp) calculation, which has been described for the sparse data case in [28, §5.2.6]. Here we briefly summarize the methodology. The r th column of $\mathbf{G}^{(n)}$, $\mathbf{g}_r^{(n)}$ is calculated as follows. Let the vectors $\mathbf{v}^{(n)}$ be defined as above in (9), but define \mathbf{u} instead as

$$\mathbf{u} = \mathbf{t} * \mathbf{v}^{(1)} * \mathbf{v}^{(2)} * \dots * \mathbf{v}^{(n-1)} * \mathbf{v}^{(n+1)} * \dots * \mathbf{v}^{(N)}.$$

Then

$$\left(\mathbf{g}_r^{(n)}\right)_j = \sum_{i_q^{(q)}=j} u_q \quad \text{for } j = 1, \dots, I_n.$$

This can be computed efficiently using the `accumarray` function in MATLAB.

5. Experiments

On both real and simulated three-way data, we assess the performance of the CP-WOPT method. We compare CP-WOPT with other methods and also demonstrate its performance on two applications.

5.1. Computational environment

All experiments were performed using MATLAB 2009b on a Linux Workstation (RedHat 5.2) with 2 Quad-Core Intel Xeon 3.0GHz processors and 32GB RAM. Timings were performed using MATLAB’s `tic` and `toc` functions since `cputime` is known to produce inaccurate results for multi-CPU and/or multi-core systems.

CP-WOPT is implemented using the Tensor Toolbox [29]. We consider dense and sparse versions, based on the gradient and function computations shown in Figure 3a and Figure 3b, respectively. We use the nonlinear conjugate gradient (NCG) method with Hestenes-Stiefel updates [27] and the Moré-Thuente line search [30] provided in the Poblano Toolbox [31] as the optimization method.

We compare CP-WOPT to two other methods: EM-ALS (implemented in the N-way Toolbox for MATLAB, version 3.10 [32]) and INDAFAC [33], which is a damped Gauss-Newton method proposed by Tomasi and Bro [7]. INDAFAC is the more related method to CP-WOPT because it is also based on optimization over all factor matrices simultaneously. Previously, Tomasi and Bro showed that INDAFAC converged to solutions in many fewer iterations than EM-ALS. We also present several comparisons against EM-ALS as implemented in the `parafac` procedure in the N-way Toolbox.

The stopping conditions are set as follows. All algorithms use the relative change in the function value $f_{\mathcal{W}}$ in (2) as a stopping condition (set to 10^{-6}). In INDAFAC, the tolerance on the infinity norm of the gradient is set to 10^{-8} and the maximum number of iterations is set to 500. In CP-WOPT, the tolerance on the two-norm of the gradient divided by the number of entries in the gradient is set to 10^{-8} , the maximum number of iterations is set to 500, and the maximum number of function evaluations is set to 10000. In EM-ALS, the maximum number of iterations (equivalent to one function evaluation) is set to 10000.

All the methods under consideration are iterative methods. Starting points are generated using the left singular vectors of $\mathbf{X}_{(n)}$ (\mathbf{X} unfolded in mode n) with missing entries replaced by zero. In our preliminary experiments, this starting procedure produced significantly better results than random initializations, even with large amounts of missing data.

5.2. Validation metrics

If the true factors are known, then we can assess the recovery of the factors via the *factor match score* (FMS) defined as follows. Let the correct and computed factorizations be given by

$$\sum_{r=1}^R \lambda_r \mathbf{a}_r^{(1)} \circ \mathbf{a}_r^{(2)} \circ \dots \circ \mathbf{a}_r^{(N)} \quad \text{and} \quad \sum_{r=1}^{\bar{R}} \bar{\lambda}_r \bar{\mathbf{a}}_r^{(1)} \circ \bar{\mathbf{a}}_r^{(2)} \circ \dots \circ \bar{\mathbf{a}}_r^{(N)},$$

respectively. Without loss of generality, we assume that all the vectors have been scaled to unit length and that the scalars are positive. Further, we assume $\bar{R} \geq R$ so that the computed solution has at least as many components as the true solution. (One could add all-zero components to the computed solution if $\bar{R} < R$.) Recall that there is a permutation ambiguity, so all possible matchings of components between the two solutions must be considered. Under these conditions, the FMS is defined as

$$\text{FMS} = \max_{\sigma \in \Pi(R, \bar{R})} \frac{1}{R} \sum_{r=1}^R \left(1 - \frac{|\lambda_r - \bar{\lambda}_{\sigma(r)}|}{\max\{\lambda_r, \bar{\lambda}_{\sigma(r)}\}} \right) \prod_{n=1}^N |\mathbf{a}_r^{(n)\top} \bar{\mathbf{a}}_{\sigma(r)}^{(n)}|. \quad (10)$$

If $\bar{R} = R$, then the set $\Pi(R, \bar{R})$ is all permutations of 1 to R ; otherwise, it is all possible permutations of all $\binom{\bar{R}}{R}$ mappings of $\{1, \dots, \bar{R}\}$ to $\{1, \dots, R\}$. The FMS can be between 0 and 1, and the best possible FMS is 1. If $\bar{R} \geq R$, some components in the computed solution are completely ignored in the calculation of FMS, but this is not an issue for us because we use $\bar{R} = R$ throughout.

We also consider the problem of recovering missing data. Let \mathbf{X} be the original data and let $\hat{\mathbf{X}}$ be the tensor that is produced by the computed model. Then the *tensor completion score* (TCS) is

$$\text{TCS} = \frac{\|(\mathbf{1} - \mathbf{W}) * (\mathbf{X} - \hat{\mathbf{X}})\|}{\|(\mathbf{1} - \mathbf{W}) * \mathbf{X}\|}. \quad (11)$$

In other words, the TCS is the relative error in the missing entries. TCS is always nonnegative, and the best possible score is 0.

5.3. Simulated data

We consider the performance of the methods on moderately-sized problems of sizes $50 \times 40 \times 30$, $100 \times 80 \times 60$, and $150 \times 120 \times 90$. For all sizes, we consider the case of computing CP decompositions using $R = 5$. We consider 60%, 70%, 80%, 90% and 95% missing data. The experiments show that the underlying factors can be captured even if the CP model is fit to a tensor with significant amount of missing data. This is because the low-rank structure of the tensor is being exploited. A rank- R tensor of size $I \times J \times K$ has $R(I + J + K)$ degrees of freedom. The reason that the factors can be recovered even with 95% missing data is that there is still a lot more data than variables, i.e., the size of the data is equal to $0.05 IJK$ which is much greater than the $R(I + J + K)$ variables (see Table 1). Because it is a nonlinear problem, we do not know exactly how many data entries are needed in order to recover a CP model of a low-rank tensor. However, a lower bound for the number of entries needed to recover a low-rank matrix has been derived in [12].

5.3.1. Generating simulated data

We create the test problems as follows. Assume that the target size is $I \times J \times K$ and that the number of factors is R . We generate factor matrices \mathbf{A} , \mathbf{B} , and \mathbf{C} of sizes $I \times R$, $J \times R$, and $K \times R$, respectively, by randomly choosing each entry from $\mathcal{N}(0, 1)$ and then normalizing every column to unit length. Note that this method of generating the factor matrices ensures that the solution is unique with probability one for the sizes and number of components that we are considering because $R \ll \min\{I, J, K\}$. Consequently, each matrix has R linearly independent columns with probability one and the k-rank of each factor matrix is R . We therefore satisfy the necessary conditions defined by Kruskal [34] for uniqueness.

We then create the data tensor as

$$\mathbf{X} = \llbracket \mathbf{A}, \mathbf{B}, \mathbf{C} \rrbracket + \eta \frac{\|\mathbf{X}\|}{\|\mathbf{N}\|} \mathbf{N}$$

where \mathbf{N} is a noise tensor of the same size as \mathbf{X} with entries randomly selected from $\mathcal{N}(0, 1)$ and η is the noise parameter. We use $\eta = 10\%$ in our experiments.

Finally, we set some entries of each generated tensor to be missing according to an indicator tensor \mathbf{W} . We consider two situations for determining the missing data: randomly missing entries and structured missing data in the form

of randomly missing fibers. In the first case, \mathbf{W} is a binary tensor such that exactly $\lfloor MIJK \rfloor$ randomly selected entries are set to zero, where $M \in (0, 1)$ defines the percentage of missing data. We require that every slice of \mathbf{W} (in every direction) have at least one nonzero because otherwise we have no chance of recovering the CP factors; this is related to the problem of *coherence* in the matrix completion problem where it is well-known that missing an entire row (or a column) of a matrix means that the true factorization can never be recovered. In the second case, without loss of generality, we only consider missing fibers in the third mode. In that case, \mathbf{W} is an $I \times J$ binary matrix with exactly $\lfloor MIJ \rfloor$ randomly selected entries are set to zero, and the binary tensor \mathbf{W} is created by stacking K copies of \mathbf{W} together. We again require that every slice of \mathbf{W} (in every direction) have at least one nonzero, which is equivalent to requiring that \mathbf{W} has no zero rows or columns.

For each problem size and missing data percentage, thirty independent test problems are generated.

5.3.2. Comparisons of CP-WOPT and INDAFAC for randomly missing entries

Perhaps contrary to intuition, we note that smaller problems are generally more difficult than larger problems for the same percentage of missing data. This is because the ratio of known entries to variables is higher for larger problems with the same percentage of missing data. Specifically, we define the ratio ρ as

$$\rho = \frac{\text{Number of known tensor entries}}{\text{Number of variables}} = \frac{(1 - M)IJK}{R(I + J + K)}. \quad (12)$$

In general, smaller values of ρ indicate more difficult problems. Table 1 lists the ratios for the problems under consideration in this section. Since this is a nonlinear problem, ρ does not tell the entire story, but it is at least a partial indicator of problem difficulty.

Table 1: Ratio (ρ) values for different problem sizes and percentages of missing values with $R = 5$.

M	$50 \times 40 \times 30$	$100 \times 80 \times 60$	$150 \times 120 \times 90$
60%	40	160	360
70%	30	120	270
80%	20	80	180
90%	10	40	90
95%	5	20	45

Table 2 reports the mean FMS results (across thirty samples) for CP-WOPT and INDAFAC on problems with randomly missing entries. Included in the table are the p -values computed for paired-sample t -tests of the FMS scores. We see that there is no appreciable difference in these results.

The advantage of CP-WOPT is that it can be significantly faster. We compare timings for the “dense” and “sparse” versions of CP-WOPT (both of which

get the same answers) with INDAFAC in Figure 4. When \mathcal{W} is dense (e.g., $M < 80\%$), the dense version of CP-WOPT is fastest; as the percentage of missing data increases, however, the sparse version of CP-WOPT has a clear speed advantage.

Table 2: FMS comparisons of CP-WOPT and INDAFAC for tensors with randomly missing entries. Mean FMS scores for each method along with p -values computed for paired-sample t -tests are presented.

Tensor Size	M	CP-WOPT	INDAFAC	p -value
$50 \times 40 \times 30$	60%	0.9878	0.9977	0.3256
	70%	0.9967	0.9762	0.0831
	80%	0.9815	0.9879	0.5944
	90%	0.9017	0.9142	0.6131
	95%	0.3907	0.3860	0.9183
$100 \times 80 \times 60$	60%	0.9908	0.9873	0.8197
	70%	0.9809	0.9990	0.1610
	80%	0.9984	0.9984	0.0708
	90%	0.9754	0.9957	0.0830
	95%	0.9263	0.9390	0.5611
$150 \times 120 \times 90$	60%	0.9724	0.9731	0.9753
	70%	0.9896	0.9995	0.3255
	80%	0.9888	0.9992	0.3255
	90%	0.9979	0.9913	0.3257
	95%	0.9742	0.9606	0.4184

5.3.3. Comparisons of CP-WOPT and INDAFAC for structured missing data

We also consider more structured missing data in the case of randomly missing fibers.

Table 3 reports the mean FMS results (across thirty samples) for CP-WOPT and INDAFAC on problems with randomly missing fibers. As for the results for randomly missing entries, the p -values computed for paired-sample t -tests of the FMS scores are included as well. Contrary to the results for randomly missing entries, though, there are several cases where there are significantly different FMS scores; specifically, the case of 90% missing data for tensors of size $50 \times 40 \times 30$ indicates a significantly higher FMS score for CP-WOPT versus INDAFAC. The 30 test-by-test results for this case are illustrated in Figure 5, where we see that none of the FMS scores reaches higher than 0.6 for either method and that more than half of the scores for INDAFAC are close to zero. However, as both methods fail to solve the problem in each test for this case, it is unclear what the significantly higher FMS scores indicate. Investigation of what can be learned about different factorization algorithms in cases with such low FMS scores is left for future work.

Once again, the advantage of CP-WOPT is that it can be significantly faster.

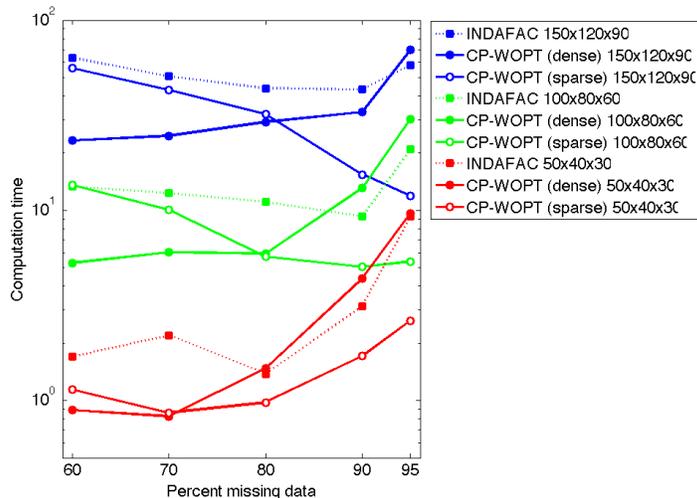


Figure 4: Time comparisons of both dense and sparse versions of CP-WOPT with INDAFAC for tensors with randomly missing entries.

We compare timings for the “dense” and “sparse” versions of CP-WOPT (both of which get the same answers) with INDAFAC applied to the data with randomly missing fibers in Figure 6. As with the factorizations involving randomly missing entries, when \mathcal{W} is dense (e.g., $M < 80\%$), the dense version of CP-WOPT is fastest; as the percentage of missing data increases, however, the sparse version of CP-WOPT has a clear speed advantage.

5.3.4. Comparisons of CP-WOPT and EM-ALS for randomly missing entries

It is well-known that the EM-ALS approach is fast but has some problems. We refer the reader to [7] for a detailed comparison of INDAFAC and EM-ALS. Here we just illustrate a few elements that are relevant to CP-WOPT.

Using the same experimental data as in §5.3.2, we used EM-ALS to also compute the factorizations. The FMS data for CP-WOPT and EM-ALS are statistically indistinguishable using paired-sample t -tests at the 5% confidence level except for two cases which are discussed in more detail below. We focus mainly on the ratio of EM-ALS time to CP-WOPT time (the best of the dense or sparse timings for each case), which is presented in Table 4. EM-ALS is faster for 60-80% missing data, with a speed-up of up to 5X. In the cases with 90-95% missing data, however, CP-WOPT (the sparse version in all cases) is faster, achieving more than a 10X speed-up in some cases. It may be possible to accelerate EM-ALS using some of the same sparse data techniques as were used for CP-WOPT, and that is a topic for future research.

There are two cases where the results were statistically distinguishable at the 5% confidence level: (a) $50 \times 40 \times 30$ with $M = 95\%$, and (b) $150 \times 120 \times 90$ with

Table 3: FMS comparisons of CP-WOPT and INDAFAC for tensors with randomly missing fibers. Mean FMS scores for each method along with p -values computed for paired-sample t -tests are presented.

Tensor Size	M	CP-WOPT	INDAFAC	p -value
$50 \times 40 \times 30$	60%	0.9845	0.9977	0.1609
	70%	0.9893	0.9760	0.4249
	80%	0.8261	0.7558	0.1144
	90%	0.2526	0.0944	0.0000
$100 \times 80 \times 60$	60%	0.9992	0.9992	0.0317
	70%	0.9988	0.9921	0.3256
	80%	0.9782	0.9577	0.0770
	90%	0.6301	0.6077	0.6209
$150 \times 120 \times 90$	60%	0.9997	0.9889	0.3256
	70%	0.9995	0.9869	0.3256
	80%	0.9990	0.9923	0.3259
	90%	0.8736	0.8903	0.4371

Table 4: Ratio of EM-ALS time to CP-WOPT time for varying sizes and missing data percentages.

M	$50 \times 40 \times 30$	$100 \times 80 \times 60$	$150 \times 120 \times 90$
60%	0.18	0.24	0.19
70%	0.32	0.26	0.21
80%	0.43	0.42	0.26
90%	2.81	1.45	1.05
95%	13.28	11.40	5.73

$M = 90\%$. The test-by-test results are illustrated in Figure 7. In the first case, both methods do poorly across the board, making comparisons difficult, though it might be argued that there are four places where EM-ALS crosses the 0.9 FMS threshold and so arguably yields interpretable results. In the second case, four of thirty examples are worse for EM-ALS than CP-WOPT. In contrast, Figure 8 shows two examples where the FMS data are statistically indistinguishable. It can be seen that sometimes both methods have trouble with the same problem while other times one or the other is able to solve it — this suggests that hybrid approaches may be a technique for future study.

5.4. Large-scale simulated data

A unique feature of CP-WOPT is that it can be applied to problems that are too large to fit in memory if dense storage were used. We consider two situations:

- (a) $500 \times 500 \times 500$ with 99% missing data (1.25 million known values), and
- (b) $1000 \times 1000 \times 1000$ with 99.5% missing data (5 million known values).

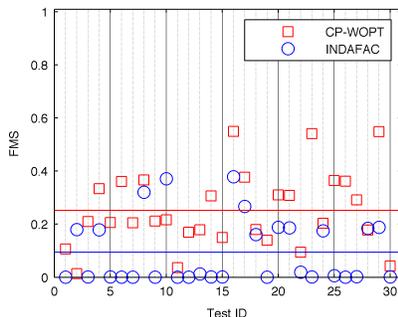


Figure 5: Test-by-test results comparing CP-WOPT and INDAFAC with significantly different FMS data for the case of tensors of size $50 \times 40 \times 30$ with 90% randomly missing fibers. Means are shown as solid lines.

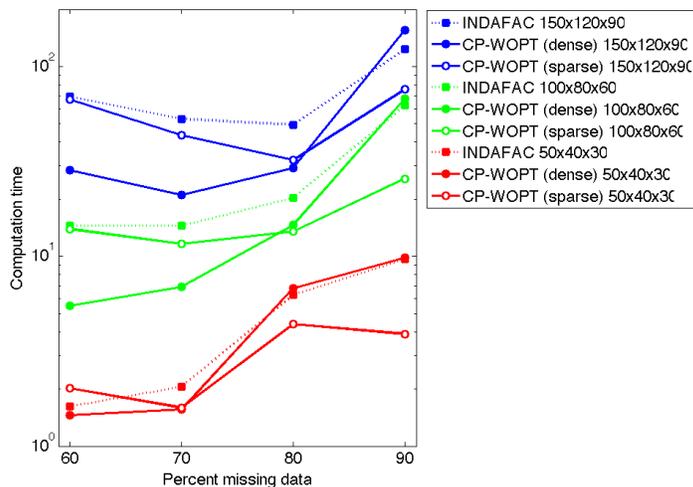
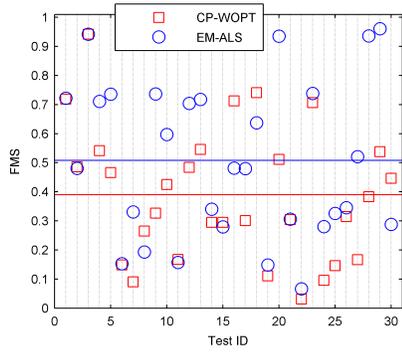


Figure 6: Time comparisons of both dense and sparse versions of CP-WOPT with INDAFAC for tensors with randomly missing fibers.

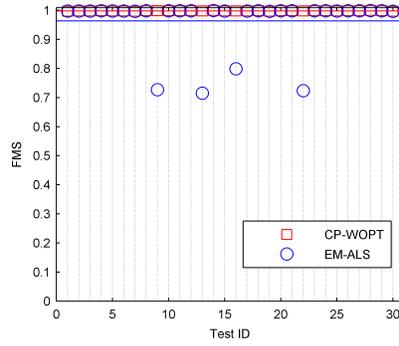
5.4.1. Generating large-scale simulated data

The method for generating test problems of this size are necessarily slightly different than in the dense case because we cannot, for example, generate a full noise tensor. In fact, the tensors $\mathbf{Y} = \llbracket \mathbf{A}, \mathbf{B}, \mathbf{C} \rrbracket$ (the noise-free data) and \mathcal{N} (the noise) are never explicitly fully formed.

Let the tensor size be $I \times J \times K$ and the missing value rate be M . We generate the factor matrices as described in §5.3, using $R = 5$ components. Next, we create the set \mathcal{I} with $(1 - M)IJK$ randomly generated indices; this set represents the indices of the known (non-missing) values in our tests. The

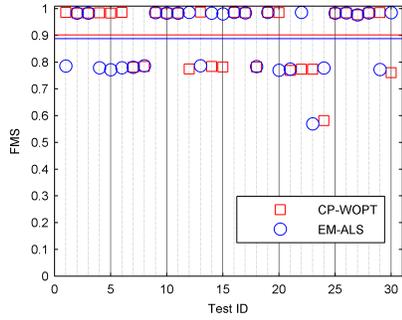


(a) $50 \times 40 \times 30$ with $M = 95\%$

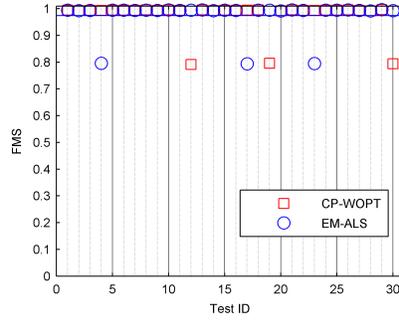


(b) $150 \times 120 \times 90$ with $M = 90\%$

Figure 7: Test-by-test results comparing CP-WOPT and EM-ALS with significantly different FMS data. Means are shown as solid lines.



(a) $50 \times 40 \times 30$ with $M = 90\%$



(b) $150 \times 120 \times 90$ with $M = 95\%$

Figure 8: Test-by-test results comparing CP-WOPT and EM-ALS with significantly indistinguishable FMS data. Means are shown as solid lines.

binary indicator tensor \mathcal{W} is stored as a sparse tensor [28] and defined by

$$w_{ijk} = \begin{cases} 1 & \text{if } (i, j, k) \in \mathcal{I}, \\ 0 & \text{otherwise.} \end{cases}$$

Rather than explicitly forming all of $\mathcal{Y} = \llbracket \mathbf{A}, \mathbf{B}, \mathbf{C} \rrbracket$, we only calculate its values for those indices in \mathcal{I} . This is analogous to the calculations described for Step 2 of Figure 3b. All missing entries of \mathcal{Y} are set to zero, and \mathcal{Y} is stored as a sparse tensor. Finally, we set

$$\mathcal{X} = \mathcal{Y} + \eta \frac{\|\mathcal{X}\|}{\|\mathcal{N}\|} \mathcal{N},$$

where \mathcal{N} is a *sparse* noise tensor such that

$$n_{ijk} = \begin{cases} \mathcal{N}(0, 1) & \text{if } (i, j, k) \in \mathcal{I}, \\ 0 & \text{otherwise.} \end{cases}$$

For each problem size, ten independent test problems are generated.

5.4.2. Computational results with CP-WOPT on large-scale data

The computational set-up was the same as that described in §5.1 except that the tolerance of the two-norm of the gradient divided by the number of entries in the gradient was set to 10^{-10} (previously 10^{-8}). The results across all ten runs for each problem size are shown in Figure 9.

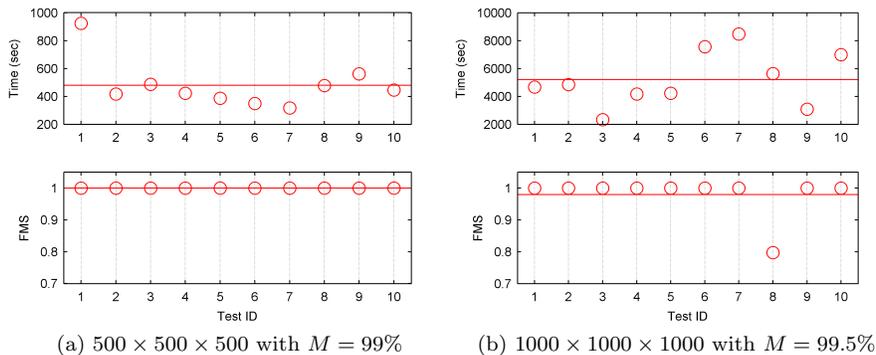


Figure 9: Results for large-scale problems with at least 99% missing data. The means are shown as solid lines.

In the $500 \times 500 \times 500$ case, storing a dense version of the tensor (assuming 8 bytes per entry) would require exactly 1GB. Storing just 1% of the data (1.25M entries) in sparse format (assuming 32 bytes per entry, one for the value and three for the indices) requires only 40MB. In our randomly generated experiments, all ten problems were solved with an FMS score greater than 0.99, with solve times ranging between 200 and 1000 seconds.

In the $1000 \times 1000 \times 1000$ case, storing a dense version of the tensor would require 8GB. Storing just 0.5% of the data as a sparse tensor (5M entries) requires 160MB. In our randomly generated experiments, only nine of the ten problems were solved with an FMS score greater than 0.99. The solve times ranged from 2000 to 10000 seconds, approximately 10 times slower than the $500 \times 500 \times 500$ case, which had half as many variables and 1/4 the nonzero entries. In the failed case, the optimization method exited because the relative change in the function tolerance was smaller than the tolerance. We restarted the optimization method with no changes except that we use the solution that had been computed previously as the initial guess. After 985 seconds of additional work, an answer was found with an FMS score of 0.9999.

5.5. EEG data

In this section, we demonstrate the use of CP-WOPT algorithm in multi-channel EEG analysis by capturing the underlying brain dynamics even in the presence of missing signals. We use an EEG data set collected to observe the gamma activation during proprioceptive stimuli of left and right hand [19]. The data set contains multi-channel signals (64 channels) recorded from 14 subjects during stimulation of left and right hand (i.e., 28 measurements in total). For each measurement, the signal from each channel is represented in both time and frequency domains using continuous wavelet transform and vectorized (forming a vector of length 4392); in other words, each measurement can be represented by a *channels* by *time-frequency* matrix. The data for all measurements can then be arranged as a *channels* by *time-frequency* by *measurements* tensor of size $64 \times 4392 \times 28$. For details about the data, see [19].

We model the data using a CP model with $R = 3$, denoting \mathbf{A} , \mathbf{B} , and \mathbf{C} as the extracted factor matrices corresponding to the channels, time-frequency, and measurements modes, respectively. We demonstrate the columns of the factor matrices in each mode in Figure 10a. The 3-D head plots correspond to the columns of \mathbf{A} , i.e., coefficients corresponding to the channels ranging from low values in blue to high values in red. The time-frequency domain representations correspond to the columns of \mathbf{B} rearranged as a matrix and again ranging from low values in blue to high values in red. The bar plots represent the columns of \mathbf{C} . Note that 3 rows of images in Figure 10a (3-D head plot, matrix plot, bar plot) correspond to columns $r = 1, 2, 3$ of the factor matrices (\mathbf{A} , \mathbf{B} , \mathbf{C}), respectively. Observe that the first row of images highlights the differences between left and right hand stimulation while the second and third rows of images pertain to frontal and parietal activations that are shared by the stimuli. Unlike [19], we do not use nonnegativity constraints; we convert tensor entries from complex to real by using the absolute values of the entries and center the data across the channels mode before the analysis.

It is not uncommon in EEG analysis that the signals from some channels are ignored due to malfunctioning of the electrodes. This will create missing fibers in the tensor (as in Figure 1) when we arrange the data as described above. To reflect such cases of missing data, we randomly set data for one or more of the 64 channels for each measurement to be missing, center the tensor

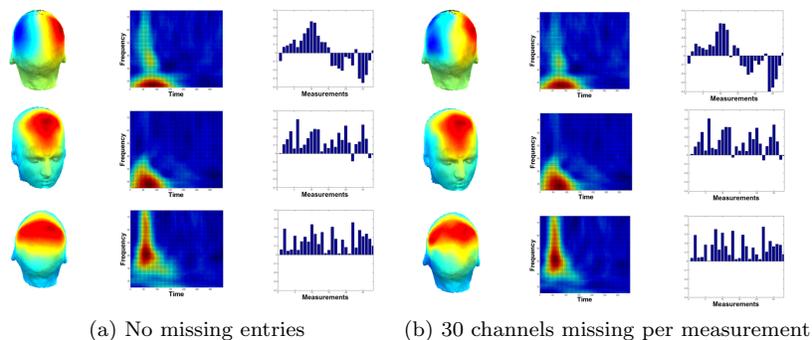


Figure 10: Columns of the CP factor matrices (\mathbf{A} , \mathbf{B} and \mathbf{C} with $R = 3$) extracted from the EEG data arranged as a *channels* by *time-frequency* by *measurements* tensor with. The 3-D head images were drawn using EEGLab [35].

across the channels mode ignoring the missing entries and then fit a CP model with $R = 3$ to the resulting data using the CP-WOPT algorithm. Let $\bar{\mathbf{A}}, \bar{\mathbf{B}}, \bar{\mathbf{C}}$ be the factor matrices extracted from a tensor with missing entries using the CP-WOPT algorithm. Table 5 illustrates how the number of missing channels per measurement affects the similarity between the columns of factor matrices extracted from missing data, i.e., $\bar{\mathbf{A}}, \bar{\mathbf{B}}, \bar{\mathbf{C}}$, and the columns of factor matrices extracted from the original data with no missing entries, i.e., $\mathbf{A}, \mathbf{B}, \mathbf{C}$. The similarity is defined in terms of FMS given in (10). For each number of missing channels, we generate 50 tensors with randomly missing channels and extract the corresponding 50 sets of $\bar{\mathbf{A}}, \bar{\mathbf{B}}$ and $\bar{\mathbf{C}}$. The values given in the second column of Table 5 are the average similarities between $\mathbf{A}, \mathbf{B}, \mathbf{C}$ and those 50 sets of $\bar{\mathbf{A}}, \bar{\mathbf{B}}$ and $\bar{\mathbf{C}}$. We observe that as the number of missing channels increases, the similarities decrease as expected. However, even up to 30 missing channels per measurement, or about 47% of the data, the extracted factor matrices match with the original factor matrices well, with similarity measures still around 0.90. Furthermore, Figure 10b images for $\bar{\mathbf{A}}, \bar{\mathbf{B}}$ and $\bar{\mathbf{C}}$ analogous to those for $\mathbf{A}, \mathbf{B}, \mathbf{C}$ in Figure 10a, illustrates that the underlying brain dynamics are still captured even when 30 channels per measurement are missing. Note that only slight local distortions can be observed with respect to the corresponding images for the original factor matrices in Figure 10a.

It can be argued that the activations of the electrodes are highly correlated and even if some of the electrodes are removed, the underlying brain dynamics can still be captured. However, in these experiments we do not set the same channels to missing for each measurement; the channels are randomly missing from each measurement. On the other hand, CP decomposition may not be able to recover factors when data has certain patterns of missing values, e.g., missing the signals from the same side of the brain for all measurements. However, it is not very likely to have such data. We still note that the success of the proposed

approach depends on the patterns of missing entries, which is the case for any factorization approach proposed for handling missing data.

Table 5: *EEG Analysis with Incomplete Data*: The similarity between the columns of $\mathbf{A}, \mathbf{B}, \mathbf{C}$ and the columns of factor matrices extracted by different approaches, i.e., CP-WOPT and Imputation. The similarity is measured in terms of FMS defined in (10).

Number of Missing Channels	CP-WOPT	Imputation
1	0.9959	0.9843
10	0.9780	0.8274
20	0.9478	0.6703
30	0.8949	0.4529
40	0.6459	0.2365

Finally, it may also look reasonable to impute missing entries simply with the mean rather than ignoring them. However, this is not a valid approach especially as the percentage of missing entries increases [10], which we also observe in Table 5. Let $\hat{\mathbf{A}}, \hat{\mathbf{B}}, \hat{\mathbf{C}}$ be the factor matrices extracted from data with missing entries when missing entries are replaced by the mean across the channel mode. Since the data is centered across the channels mode, missing entries are replaced with zeros. The third column of Table 5 shows how the similarities (again defined in terms of FMS given in (10)) between the columns of $\hat{\mathbf{A}}, \hat{\mathbf{B}}, \hat{\mathbf{C}}$ and the columns of the factor matrices extracted from the original data with no missing entries, i.e., $\mathbf{A}, \mathbf{B}, \mathbf{C}$, change as the amount of missing data increases. We can see that imputation may only work for small amounts of missing data and as the amount of missing data increases, the structure can be better captured by ignoring the missing entries rather than replacing them with the means.

5.6. Network traffic data

In the previous section, we were interested in tensor factorizations and capturing the underlying factors in the presence of missing data. In this section, we address the problem of recovering missing entries of a tensor; in other words, the tensor completion problem. One domain where this problem is frequently encountered is network traffic analysis. Network traffic data consists of traffic matrices (TM), which record the amount of network data exchanged between source and destinations pairs. Since TMs evolve over time, network data can be represented as a tensor. For instance, Géant data [36] records the traffic exchanged between 23 routers over a period of 4 months collected using 15-minute intervals. This data set forms a third-order tensor with modes: *source routers*, *destination routers*, and *time* and each entry indicates the amount of traffic sent

from a source to a destination during a certain time interval. Missing data arises due to the expense of the data collection process.

In our study, we have used the Géant data collected in April³; that corresponds to a tensor of size $23 \times 23 \times 2756$. Let \mathcal{T} represent this raw data tensor. After preprocessing \mathcal{T} as $\mathcal{X} = \log(\mathcal{T} + 1)$, we model \mathcal{X} using a 2-component CP model and extract factor matrices for each mode as illustrated in Figure 11. The first and second row correspond the first and second column of the factor matrices, respectively. This CP model fits the data well but not perfectly and there is some unexplained variation left in the residuals. Let $\hat{\mathcal{X}}$ be the tensor constructed using the computed factor matrices. If we define the modeling error as $\frac{\|\mathcal{X} - \hat{\mathcal{X}}\|}{\|\mathcal{X}\|}$, then the error is approximately 0.31 for a 2-component CP model. Even though extracting more components slightly lowers the modeling error, models with more components do not look appropriate for the data.

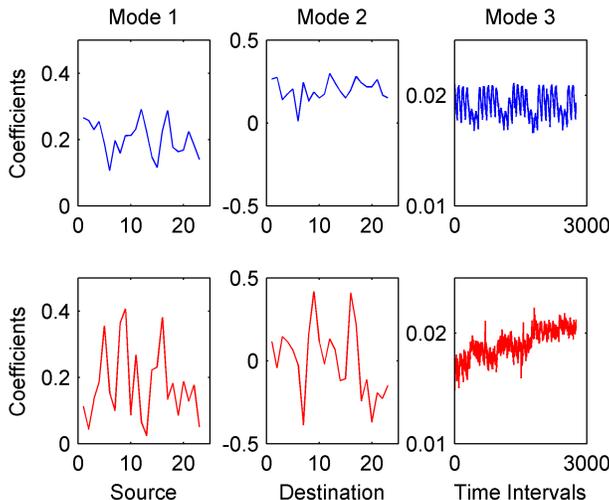


Figure 11: Factor matrices extracted from Géant data using a 2-component CP model. The first row illustrates the first column of the factor matrices in each mode and the second row shows the second column of the factor matrices.

In order to assess the performance of CP-WOPT in terms of recovering missing data, some entries of \mathcal{X} are randomly set to missing and a 2-component CP model is fit to the data with missing entries using the CP-WOPT algorithm. The extracted factor matrices are then used to reconstruct the data and fill in the missing entries. These recovered values are compared with the actual values based on the tensor completion score (TCS) defined in (11). Figure 12 shows how TCS differs for various amounts of missing data. These are the average

³The data was incomplete for the full four-month period, e.g., 6 days of data was missing at the end of February. For the analysis, we picked a period with no missing time slices.

TCS values across thirty runs, where random entries are set to missing. We observe that the relative error is around 0.31 when there is little amount of missing data and very slowly increases as we increase the amount of missing data. The error is only slightly higher, approximately 0.33, even when 95% of the entries are missing. The error increases sharply when we finally raise the amount of missing data to 99%. Note that even if there is no missing data, there will be error due to the modeling error, i.e., 0.31. Figure 12 demonstrates the robustness of the algorithm with respect to large amounts of missing data such that the relative error in the missing entries can be kept close to the level of the modeling error even when the amount of missing data is high.

We have addressed the tensor completion problem using a CP model, which gives easily-interpretable factors in each mode. However, for the tensor completion problem, the recovery of the missing entries is more important than the underlying factors and their interpretability. Therefore, modeling the data using a restricted model like CP may not be the best approach in this case. It may be possible to achieve lower modeling error using a more flexible model such as a Tucker3 model, or using matrix factorizations on the unfolded tensor [2]. If these models are fit to the data using algorithms as robust as CP-WOPT, missing entries may be recovered more accurately. This is a topic of future research.

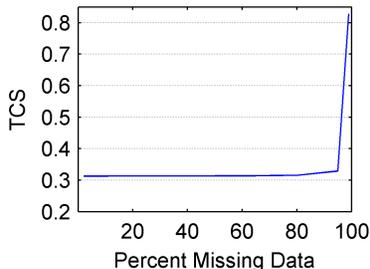


Figure 12: Tensor Completion Score for different amounts of missing data for Géant data when a 2-component CP model fitted using CP-WOPT is used to recover the missing entries.

6. Conclusions

The closely related problems of matrix factorizations with missing data and matrix completion have recently been receiving a lot of attention. In this paper, we consider the more general problem of tensor factorization in the presence of missing data, formulating the canonical tensor decomposition for incomplete tensors as a weighted least squares problem. Unlike imputation-based techniques, this formulation ignores the missing entries and models only the known data entries. We develop a scalable algorithm called CP-WOPT using gradient-

based optimization to solve the weighted least squares formulation of the CP problem.

Our numerical studies suggest that the proposed CP-WOPT approach is accurate and scalable. CP-WOPT can recover the underlying factors successfully with large amounts of missing data, e.g., 90% missing entries for tensors of size $50 \times 40 \times 30$. We have also studied how CP-WOPT can scale to problems of larger sizes, e.g., $1000 \times 1000 \times 1000$, and recover CP factors from large, sparse tensors with 99.5% missing data. Moreover, through the use of both dense and sparse implementations of the algorithm, we have showed that CP-WOPT was always faster in our studies when compared to the best alternative approach based on second-order optimization.

We have demonstrated the practical use of CP-WOPT algorithm in two different applications. In multi-channel EEG analysis, the factors extracted by the CP-WOPT algorithm can capture brain dynamics even if signals from some channels are missing, suggesting that practitioners can now make better use of incomplete data in their analyses. In network traffic analysis, CP-WOPT algorithm can be used in the context of tensor completion and recover the missing network traffic data.

In future studies, we plan to extend our results in several directions. We will include constraints such as non-negativity and penalties to encourage sparsity, which enable us to find more meaningful latent factors from large-scale sparse data. Finally, we will consider the problem of collective factorizations with missing data, where we are jointly factoring multiple tensors with shared factors.

- [1] E. Acar, D. M. Dunlavy, M. Mørup, T. G. Kolda, Scalable tensor factorizations with missing data, Tech. Rep. SAND2009-6764, Sandia National Laboratories, Albuquerque, New Mexico and Livermore, California, submitted for publication (Oct. 2009).
- [2] Y. Zhang, M. Roughan, W. Willinger, L. Qiu, Spatio-temporal compressive sensing and internet traffic matrices, in: SIGCOMM '09: Proceedings of the ACM SIGCOMM 2009 conference on Data communication, ACM, New York, NY, USA, 2009, pp. 267–278. doi:10.1145/1592568.1592600.
- [3] A. M. Buchanan, A. W. Fitzgibbon, Damped Newton algorithms for matrix factorization with missing data, in: CVPR'05: 2005 IEEE Computer Society Conference on Computer Vision and Pattern Recognition, Vol. 2, IEEE Computer Society, 2005, pp. 316–322. doi:http://doi.ieeecomputersociety.org/10.1109/CVPR.2005.118.
- [4] E. Acar, B. Yener, Unsupervised multiway data analysis: A literature survey, IEEE Transactions on Knowledge and Data Engineering 21 (1) (2009) 6–20. doi:10.1109/TKDE.2008.112.
- [5] T. G. Kolda, B. W. Bader, Tensor decompositions and applications, SIAM Rev. 51 (3) (2009) 455–500. doi:http://dx.doi.org/10.1137/07070111X.

- [6] F. Miwakeichi, E. Martínez-Montes, P. A. Valds-Sosa, N. Nishiyama, H. Mizuhara, Y. Yamaguchi, Decomposing EEG data into space-time-frequency components using parallel factor analysis, *NeuroImage* 22 (3) (2004) 1035–1045. doi:10.1016/j.neuroimage.2004.03.039.
- [7] G. Tomasi, R. Bro, PARAFAC and missing values, *Chemometrics and Intelligent Laboratory Systems* 75 (2) (2005) 163–180. doi:10.1016/j.chemolab.2004.07.003.
- [8] V. Y. Orekhov, I. Ibraghimov, M. Billeter, Optimizing resolution in multidimensional nmr by three-way decomposition, *Journal of Biomolecular NMR* 27 (2003) 165–173.
- [9] X. Geng, K. Smith-Miles, Z.-H. Zhou, L. Wang, Face image modeling by multilinear subspace analysis with missing values, in: *MM '09: Proceedings of the seventeen ACM international conference on Multimedia*, ACM, 2009, pp. 629–632. doi:10.1145/1631272.1631373.
- [10] A. Ruhe, Numerical computation of principal components when several observations are missing, Tech. Rep. UMINF-48-74, Department of Information Processing, Institute of Mathematics and Statistics, University of Umea, Umea, Sweden (1974).
- [11] K. R. Gabriel, S. Zamir, Lower rank approximation of matrices by least squares approximation with any choice of weights, *Technometrics* 21 (4) (1979) 489–498.
URL [http://links.jstor.org/sici?sici=0040-1706\(197911\)21:4<489:LRAOMB>2.0.CO;2-Q](http://links.jstor.org/sici?sici=0040-1706(197911)21:4<489:LRAOMB>2.0.CO;2-Q)
- [12] E. J. Candes, T. Tao, The power of convex relaxation: Near-optimal matrix completion, arXiv:0903.1476v1 (2009).
URL <http://arxiv.org/abs/0903.1476>
- [13] E. J. Candès, Y. Plan, Matrix completion with noise, arXiv:0903.3131v1 [cs.IT] (Mar. 2009).
URL <http://arxiv.org/abs/0903.3131>
- [14] J. D. Carroll, J. J. Chang, Analysis of individual differences in multidimensional scaling via an N-way generalization of “Eckart-Young” decomposition, *Psychometrika* 35 (1970) 283–319. doi:10.1007/BF02310791.
- [15] R. A. Harshman, Foundations of the PARAFAC procedure: Models and conditions for an “explanatory” multi-modal factor analysis, *UCLA working papers in phonetics* 16 (1970) 1–84, available at <http://www.psychology.uwo.ca/faculty/harshman/wpppfac0.pdf>.
- [16] T. G. Kolda, B. W. Bader, J. P. Kenny, Higher-order web link analysis using multilinear algebra, in: *ICDM 2005: Proceedings of the 5th IEEE International Conference on Data Mining*, IEEE Computer Society, 2005, pp. 242–249. doi:10.1109/ICDM.2005.77.

- [17] R. Bro, Review on multiway analysis in chemistry—2000–2005, *Critical Reviews in Analytical Chemistry* 36 (3–4) (2006) 279–293. doi:10.1080/10408340600969965.
- [18] E. Acar, C. A. Bingol, H. Bingol, R. Bro, B. Yener, Multiway analysis of epilepsy tensors, *Bioinformatics* 23 (13) (2007) i10–i18. doi:10.1093/bioinformatics/btm210.
- [19] M. Mørup, L. K. Hansen, S. M. Arnfred, ERPWAVELAB a toolbox for multi-channel analysis of time-frequency transformed event related potentials, *Journal of Neuroscience Methods* 161 (2) (2007) 361–368. doi:10.1016/j.jneumeth.2006.11.008.
- [20] E. Acar, T. Kolda, D. Dunlavy, An optimization approach for fitting canonical tensor decompositions, Tech. Rep. SAND2009-0857, Sandia National Laboratories (2009).
- [21] A. P. Dempster, N. M. Laird, D. B. Rubin, Maximum likelihood from incomplete data via the em algorithm, *Journal of the Royal Statistical Society, Series B* 39 (1) (1977) 1–38.
- [22] N. Srebro, T. Jaakkola, Weighted low-rank approximations, in: *IMCL-2003: Proceedings of the Twentieth International Conference on Machine Learning, 2003*, pp. 720–727.
- [23] H. A. L. Kiers, Weighted least squares fitting using ordinary least squares algorithms, *Psychometrika* 62 (2) (1997) 215–266. doi:10.1007/BF02295279.
- [24] R. Bro, Multi-way analysis in the food industry: Models, algorithms, and applications, Ph.D. thesis, University of Amsterdam, available at <http://www.models.kvl.dk/research/theses/> (1998).
- [25] B. Walczak, D. L. Massart, Dealing with missing data: Part I, *Chemometrics and Intelligent Laboratory Systems* 58 (1) (2001) 15–27. doi:10.1016/S0169-7439(01)00131-9.
- [26] P. Paatero, A weighted non-negative least squares algorithm for three-way “PARAFAC” factor analysis, *Chemometrics and Intelligent Laboratory Systems* 38 (2) (1997) 223–242. doi:10.1016/S0169-7439(97)00031-2.
- [27] J. Nocedal, S. J. Wright, *Numerical Optimization*, Springer, 1999.
- [28] B. W. Bader, T. G. Kolda, Efficient MATLAB computations with sparse and factored tensors, *SIAM Journal on Scientific Computing* 30 (1) (2007) 205–231. doi:10.1137/060676489.
- [29] B. W. Bader, T. G. Kolda, Tensor toolbox for matlab, version 2.4beta, <http://csmr.ca.sandia.gov/~tgkolda/TensorToolbox/> (last accessed November, 2009).

- [30] J. J. Moré, D. J. Thuente, Line search algorithms with guaranteed sufficient decrease, *ACM Transactions on Mathematical Software* 20 (3) (1994) 286–307. doi:10.1145/192115.192132.
- [31] D. M. Dunlavy, E. Acar, T. G. Kolda, Poblano optimization toolbox for matlab, version 1.0beta, <http://csmr.ca.sandia.gov/~tgkolda/TensorToolbox/> (last accessed November, 2009).
- [32] C. A. Andersson, R. Bro, The N-way toolbox for MATLAB, *Chemometrics and Intelligent Laboratory Systems* 52 (1) (2000) 1–4, see also <http://www.models.kvl.dk/source/nwaytoolbox/>. doi:10.1016/S0169-7439(00)00071-X.
- [33] G. Tomasi, Incomplete data PARAFAC (INDAFAC), <http://www.models.kvl.dk/source/indafac/index.asp> (last accessed May, 2009).
- [34] J. B. Kruskal, Three-way arrays: rank and uniqueness of trilinear decompositions, with application to arithmetic complexity and statistics, *Linear Algebra and its Applications* 18 (2) (1977) 95–138. doi:10.1016/0024-3795(77)90069-6.
- [35] A. Delorme, S. Makeig, EEGLAB: An open source toolbox for analysis of single-trial EEG dynamics, *J. Neurosci. Meth.* 134 (2004) 9–21. doi:10.1016/j.jneumeth.2003.10.009.
- [36] S. Uhlig, B. Quoitin, J. Lepropre, S. Balon, Providing public intradomain traffic matrices to the research community, *ACM SIGCOMM Computer Communication Review* 36 (1) (2006) 83–86. doi:10.1145/1111322.1111341.