

OPTIMIZATION ON THE HIERARCHICAL TUCKER MANIFOLD - APPLICATIONS TO TENSOR COMPLETION

CURT DA SILVA¹ AND FELIX J. HERRMANN²

¹ DEPARTMENT OF MATHEMATICS, UNIVERSITY OF BRITISH COLUMBIA

² DEPARTMENT OF EARTH AND OCEAN SCIENCES, UNIVERSITY OF BRITISH COLUMBIA

ABSTRACT. In this work, we develop an optimization framework for problems whose solutions are well-approximated by *Hierarchical Tucker* (HT) tensors, an efficient structured tensor format based on recursive subspace factorizations. By exploiting the smooth manifold structure of these tensors, we construct standard optimization algorithms such as Steepest Descent and Conjugate Gradient for completing tensors from missing entries. Our algorithmic framework is fast and scalable to large problem sizes as we do not require SVDs on the ambient tensor space, as required by other methods. Moreover, we exploit the structure of the Gramian matrices associated with the HT format to regularize our problem, reducing overfitting for high subsampling ratios. We also find that the organization of the tensor can have a major impact on completion from realistic seismic acquisition geometries. These samplings are far from idealized randomized samplings that are usually considered in the literature but are realizable in practical scenarios. Using these algorithms, we successfully interpolate large-scale seismic data sets and demonstrate the competitive computational scaling of our algorithms as the problem sizes grow.

Keywords: Hierarchical Tucker tensors, tensor completion, Riemannian manifold optimization, Gauss-Newton, differential geometry, low-rank tensor.

2010 MSC: 15A69, 57N35, 65K10, 65Y20, 53B21, 90C90

1. INTRODUCTION

The matrix completion problem is concerned with interpolating a $m \times n$ matrix from a subset of its entries. The amount of recent successes in developing solution techniques to this problem is a result of assuming a *low-rank* model on the 2-D signal of interest and a uniform random sampling scheme [9], [8], [10]. The original signal is recovered by promoting low-rank structures subject to data constraints.

Using a similar approach, we consider the problem of interpolating a d -dimensional tensor from samples of its entries. That is, we aim to solve,

$$(1) \quad \min_{\mathbf{X} \in \mathcal{H}} \frac{1}{2} \|P_{\Omega} \mathbf{X} - b\|_2^2,$$

where P_{Ω} is a linear operator $P_{\Omega} : \mathbb{R}^{n_1 \times n_2 \times \dots \times n_d} \rightarrow \mathbb{R}^m$, $b \in \mathbb{R}^m$ is our subsampled data satisfying $b = P_{\Omega} \mathbf{X}^*$ for some “solution” tensor \mathbf{X}^* and \mathcal{H} is a *specific* class of low-rank tensors to be specified later. Under the assumption that \mathbf{X}^* is well approximated by an element in \mathcal{H} , our goal is to recover \mathbf{X}^* by solving (1). For concreteness, we concern ourselves with the case when P_{Ω} is a restriction operator, i.e.,

$$P_{\Omega} \mathbf{X} = \mathbf{X}_{i_1, i_2, \dots, i_d} \quad \text{if } (i_1, i_2, \dots, i_d) \in \Omega,$$

and $\Omega \subset [n_1] \times [n_2] \times \dots \times [n_d]$ is the so-called *sampling set*, where $[n] = \{1, \dots, n\}$. In the above equation, we suppose that $|\Omega| = m \ll n_1 n_2 \dots n_d$, so that P_{Ω} is a subsampling operator.

Unlike the matrix case, there is no unique notion of rank for tensors, as we shall see in Section 1.1, and there are multiple tensor formats that generalize a particular notion of *separability* from the matrix case—i.e, there is no unique extension of the SVD to tensors. Although each tensor format can lead to compressible representations of their respective class of low-rank signals, the truncation of a general signal to one of these formats requires access to the *fully* sampled tensor \mathbf{X} (or at the very least *query*-based access to the tensor [5]) in order to achieve reasonable accuracy, owing to the use of truncated SVDs acting on various *matricizations* of the tensor. As in matrix completion, randomized missing entries change the behavior of the singular values and vectors of these matricizations and hence of the final approximation. Moreover, when the tensor of interest is actually a discretized continuous signal, there can be a number of constraints,

physical or otherwise, that limit our ability to ideally sample it. For instance, in the seismic case, the tensor of interest is a multi-dimensional wavefield in the earth’s subsurface sampled at an array of receivers located at the surface. In real-world seismic experiments, budgetary constraints or environmental obstructions can limit both the total amount of time available for data acquisition as well as the number and placement of active sources and receivers. Since seismic processing, among other domains, relies on having fully sampled data for drawing accurate inferences, tensor completion is an important technique for a variety of scientific fields that acquire multidimensional data.

In this work, we consider the class of Hierarchical Tucker (abbreviated HT) tensors as our low-rank tensors of interest. The set of all such tensors is a smooth, embedded *submanifold* of $\mathbb{R}^{n_1 \times n_2 \times \dots \times n_d}$, first studied in [53], which we equip with a Riemannian metric. Using this Riemannian structure, we can construct optimization algorithms in order to solve (1) for d -dimensional tensors. We will also study some of the effects of higher dimensional sampling and extend ideas from compressive sensing and matrix completion to the HT tensor case for our specific seismic examples.

1.1. Previous Work. To provide the reader with some context on tensor representations, let us briefly detail some of the available structured tensor formats. We refer to [30, 31, 22] for a series of comprehensive overviews of structured tensor formats, and in particular to [24] for an algebraic and functional analytic point of view on the subject. In what follows, we let $N = \max_{i=1 \dots d} n_i$ be the maximum individual dimension size, $N^d := \prod_{i=1}^d n_i$ denote the dimension of the ambient space $\mathbb{R}^{n_1 \times n_2 \times \dots \times n_d}$, and, for each tensor format discussed, K denotes the maximum of all of the rank parameters associated to that format.

The so-called Candecomp/Parafac (CP) decomposition is a very straightforward application of the separation of variables technique. Very much like the SVD of a matrix, one stipulates that, for a tensor $\mathbf{f} \in \mathbb{R}^{n_1 \times n_2 \times \dots \times n_d}$, one can write it as

$$\mathbf{f} \approx \sum_{i=1}^K f_i^{(1)} \otimes f_i^{(2)} \otimes \dots \otimes f_i^{(d)}$$

where \otimes is the Kronecker product and $f_i^{(j)} \in \mathbb{R}^{n_j}$. In addition to its straightforward construction, the CP decomposition of rank K only requires dNK parameters versus the N^d of the full tensor and tensor-tensor operations can be performed efficiently on the underlying factors rather than the full tensors themselves (see [3] for a comprehensive set of MATLAB tools).

Unfortunately, despite the parsimoniousness of the CP construction, the approximation of an arbitrary (full) tensor by CP tensors has both theoretical and numerical difficulties. In particular, the set of all CP tensors of rank at most K is not closed, and thus a best K -rank approximation is difficult to compute in many cases [14]. Despite this shortcoming, various authors have proposed iterative and non-iterative algorithms in the CP format for approximating full tensors [31] as well as interpolating tensors with missing data, such as the Alternating Least Squares approach (a block Gauss-Seidel type method) proposed alongside the CP format in [11] and [26], with convergence analysis in [52], and a nonlinear least-squares optimization scheme in [2]. The authors in [55] extended the Alternating Least Squares analysis to ensure that it converges globally to a stationary point of a block-convex model, which encompasses a variety of matrix and tensor completion models including the CP format.

The CP format is a specific case of the more general *Tucker format*, which aims to write a tensor \mathbf{f} as a multilinear product

$$\mathbf{f} \approx U_1 \times_1 U_2 \times_2 \dots U_d \times_d \mathbf{C}$$

where $\mathbf{C} \in \mathbb{R}^{k_1 \times k_2 \times \dots \times k_d}$ is the so-called *core tensor* and the matrices $U_j \in \mathbb{R}^{n_j \times k_j}$, $j = 1, \dots, d$ are the *factors* of the decomposition. Here we use the notation of the multilinear product, that is, $U_i \times_i \mathbf{C}$ indicates that \mathbf{C} is multiplied by U_i in dimension i , e.g., see [14, 13]. We will elaborate on this construction in Section 2.2. The CP format follows from this formulation when the core tensor is *diagonal*, i.e., $\mathbf{C}_{i_1, i_2, \dots, i_d} = \mathbf{C}_{i_1, i_1, \dots, i_1} \delta_{i_1, i_2, \dots, i_d}$, where $\delta_{i_1, i_2, \dots, i_d} = 1$ when $i_1 = i_2 = \dots = i_d$ and 0 otherwise.

The Tucker format enjoys many benefits in terms of approximation properties over its CP counterpart. Namely, the set of all Tucker tensors of at most multilinear rank $\mathbf{k} = (k_1, k_2, \dots, k_d)$ is closed and as a result every tensor \mathbf{f} has a best at most multilinear rank- \mathbf{k} Tucker approximation. A near-optimal approximation can be computed efficiently by means of the Higher Order SVD [13]. For the tensor completion problem, the

authors in [19] consider the problem of recovering a Tucker tensor with missing entries using the Douglas-Rachford splitting technique, which decouples interpolation and regularization by nuclear norm penalization of different matricizations of the tensor into subproblems that are then solved via a particular proximal mapping. An application of this approach to seismic data is detailed in [32] for the interpolation problem and [33] for denoising. Depending on the size and ranks of the tensor to be recovered, there are theoretical and numerical indications that this approach is no better than penalizing the nuclear norm in a single matricization (see [44] for a theoretical justification in the Gaussian measurement case, as well as [50] for an experimental demonstration of this effect). Some preliminary results on theoretical guarantees for recovering low-rank Tucker tensors from subsampled measurements are given in [29] for pointwise measurements and a suitable, tensor-based incoherence condition and [40], which considers a nuclear norm penalty of the matricization of the first $d/2$ modes of \mathbf{X} as opposed to a sum of nuclear norms of each of its d modes, as is typically considered.

Aside from convex relaxations of the tensor rank minimization problem, the authors in [34] develop an alternative manifold-based approach to Tucker Tensor optimization similar to our considerations for the Hierarchical Tucker case and subsequently complete such tensors with missing entries. In this case, each evaluation of the objective and Riemannian gradient requires $O(d(N + |\Omega|)K^d + dK^{d+1})$ operations, whereas our method only requires $O(dNK^2 + d|\Omega|K^3 + dK^4)$ operations. As a result of using the Hierarchical Tucker format instead of the Tucker format, our method scales much better as d , N , and K grow. The differential geometric considerations for the Tucker format were first analyzed in [36]. It is from here where the phrase Dirac Frenkel variational principle arises in the context of dynamical systems, which corresponds to the Riemannian gradient vanishing at the optimum value of the objective in an optimization context.

Hierarchical Tucker (HT) tensors were originally introduced in [23, 20], with the subclass of Tensor Train (TT) tensors developed independently in [43, 41]. These so-called tensor network decompositions have also been previously explored in the quantum physics community, see for instance [25]. TT tensors are often considered over HT tensors owing to their explicit, non-recursive formulation and relative ease of implementation for numerical methods, see for instance [42], although many of the ideas developed for TT tensors extend to the HT tensor case.

Previous work in completing tensors in the Tensor Train format includes [21, 27], wherein the authors use an alternating least-squares approach for the tensor completion problem. The derivations of the smooth manifold structure of the set of TT tensors can be found in [28]. This work builds upon the manifold structure of Hierarchical Tucker tensors studied in [53]. The authors in [18] have considered the manifold geometry of tensor networks in Banach spaces, but we will not employ such general machinery here.

Owing to its extremely efficient storage requirements (which are *linear* in the dimension d as opposed to exponential in d), the Hierarchical Tucker format has enjoyed a recent surge in popularity for parametrizing high-dimensional problems. The hTucker toolbox [35] contains a suite of MATLAB tools for working with tensors in the HT format, including efficient vector space operations, matrix-tensor and tensor-tensor products, and truncations of full arrays to HT format. This truncation, the so-called Hierarchical SVD developed in [20], allows one to approximate a full tensor in HT format with a near-optimal approximation error. Even though the authors in [5] develop a HT truncation method that does not need access to every entry of the tensor in order to form the HT approximation, their approach requires algorithm-driven access to the entries, which does not apply for the seismic examples we consider below. A HT approach for solving dynamical systems is outlined in [37], which considers similar manifold structure as in this article applied in a different context. The authors in [45] also consider the smooth manifold properties of HT tensors to construct tensor completion algorithms using a Hierarchical SVD-based approach. As we shall see, since their methods rely on computing SVDs of large matrices, they will have difficulty scaling to tensors with large mode sizes N , unlike the methods discussed below.

1.2. Contributions and Organization. In this paper, we extend the primarily theoretical results of [53] to practical algorithms for solving optimization algorithms on the HT manifold. In Section 3.1, we introduce the Hierarchical Tucker format. We restate some of the results of [53] in Section 3.1 to provide context for the Riemannian metric we introduce on the quotient manifold in Section 4. Equipped with this metric, we can now develop optimization methods on the HT manifold in Section 5 that are fast and SVD-free. For large-scale, high-dimensional problems, the computational costs of SVDs are prohibitive and affect the

scalability of tensor completion methods such as [19]. Since we are using the HT manifold rather than the Tucker manifold, we avoid an exponential dependence on the internal rank parameters as in [34]. We initially proposed the idea for a Riemannian metric on the HT manifold in [12] and in this paper, we have subsequently improved upon these results in this paper to reduce the overall computational overhead and speed up the convergence of the algorithm by using our Gauss-Newton method. In Section 5.4, we exploit the structure of HT tensors to regularize different matricizations of the tensor *without* having to compute SVDs of these matricizations, lessening the effects of overfitting when there are very few samples available. We conclude by demonstrating the effectiveness of our techniques on interpolating various seismic data volumes with missing data points in all dimensions as well as missing receivers, which is more realistic. Our numerical results are similar to those presented previously in [12], but much more extensive and include our regularization and Gauss-Newton based methods. In this paper, we also compare our method to a reference implementation of [34] and achieve very reasonable results for our seismic data volumes.

We note that the algorithmic results here generalize readily to complex tensor completion $\mathbb{C}^{n_1 \times n_2 \times \dots \times n_d}$ and more general subsampling operators P_Ω .

2. NOTATION

In this paper, we denote vectors by lower case letters x, y, z, \dots , matrices by upper case, plain letters A, B, C, \dots, X, Y, Z , and tensors by upper case, bold letters $\mathbf{X}, \mathbf{Y}, \mathbf{Z}$.

2.1. Matricization. We consider d -dimensional tensors \mathbf{X} of size $n_1 \times n_2 \times \dots \times n_d$. $t = (t_1, t_2, \dots, t_k) \subset \{1, \dots, d\}$ selects a subset of d dimensions and we denote $t^c := \{1, \dots, d\} \setminus t$ its complement. We let the *matricization* of a tensor \mathbf{X} along the modes $t \subset \{1, \dots, d\}$ be the matrix $X^{(t)}$ such that the indices in t are vectorized along the rows and the indices in t^c are vectorized along the columns, i.e., if we set $s = t^c$, then

$$X^{(t)} \in \mathbb{R}^{(n_{t_1} n_{t_2} \dots n_{t_k}) \times (n_{s_1} n_{s_2} \dots n_{s_{d-k}})}$$

$$(X^{(t)})_{(i_{t_1}, \dots, i_{t_k}), (i_{s_1}, \dots, i_{s_{d-k}})} := \mathbf{X}_{i_1, \dots, i_d}.$$

We also use the notation $(\cdot)_{(t)}$ for the *dematricization* operation, i.e., $(X^{(t)})_{(t)} = \mathbf{X}$, which reshapes the matricized version of \mathbf{X} along modes t back to its full tensor form.

2.2. Multilinear product. A natural operation to consider on tensors is that of the *multilinear product* [13, 53, 20, 31].

Definition 1. Given a d -tensor \mathbf{X} of size $n_1 \times n_2 \times \dots \times n_d$ and matrices $A_i \in \mathbb{R}^{m_i \times n_i}$, the *multilinear product* of $\{A_i\}_{i=1}^d$ with \mathbf{X} , is the $m_1 \times m_2 \times \dots \times m_d$ tensor $\mathbf{Y} = A_1 \times_1 A_2 \times_2 \dots A_d \times_d \mathbf{X}$, is defined in terms of the matricizations of \mathbf{Y} as

$$Y^{(i)} = A_i X^{(i)} A_d^T \otimes A_{d-1}^T \otimes \dots \otimes A_{i+1}^T \otimes A_{i-1}^T \dots \otimes A_1^T, \quad i = 1, 2, \dots, d.$$

Conceptually, we are applying operator each operator A_i to dimension i of the tensor \mathbf{X} , keeping all other coordinates fixed. For example, when A, X, B are matrices of appropriate sizes, the quantity AXB^T can be written as $AXB^T = A \times_1 B \times_2 X$. We remark in this instance that the ordering of the unfoldings matters and that this particular choice is compatible with the standard kronecker product. We refer to [31] for more details.

The standard Euclidean inner product between two d -dimensional tensors X and Y can be defined in terms of the standard Euclidean product for vectors, by letting

$$\langle \mathbf{X}, \mathbf{Y} \rangle := \text{vec}(\mathbf{X})^T \text{vec}(\mathbf{Y})$$

where $\text{vec}(\mathbf{X}) := X^{(1,2,\dots,d)}$ is the usual vectorization operator. This inner product induces a norm $\|\mathbf{X}\|_2$ on the set of all d -dimensional tensors in the usual way, $\|\mathbf{X}\|_2 = \sqrt{\langle \mathbf{X}, \mathbf{X} \rangle}$.

Here we state several properties of the multilinear product, which are straightforward to prove.

Proposition 1. Let $\{A_i\}_{i=1}^d, \{B_i\}_{i=1}^d$ be collections of linear operators and \mathbf{X}, \mathbf{Y} be tensors, all of appropriate sizes, so that the multilinear products below are well-defined. Then we have the following:

- (1) $(A_1 \times_1 \dots A_d \times_d) \circ (B_1 \times_1 \dots B_d \times_d \mathbf{X}) = (A_1 B_1) \times_1 \dots (A_d B_d) \times_d \mathbf{X}$ [14]
- (2) $\langle A_1 \times_1 \dots A_d \times_d \mathbf{X}, B_1 \times_1 \dots B_d \times_d \mathbf{Y} \rangle = \langle (B_1^T A_1) \times_1 \dots (B_d^T A_d) \times_d \mathbf{X}, \mathbf{Y} \rangle$

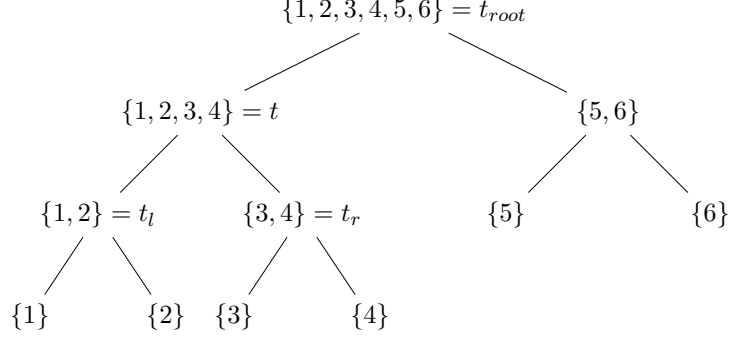


FIGURE 1. Complete dimension tree for $\{1, 2, 3, 4, 5, 6\}$.

2.3. Tensor-tensor contraction. Another natural operation to consider between two tensors is *tensor-tensor contraction*, a generalization of matrix-matrix multiplication. We define tensor-tensor contraction in terms of tensors of the same dimension for ease of presentation [17].

Definition 2. Given a d -tensor \mathbf{X} of size $n_1 \times \cdots \times n_d$ and a d -tensor \mathbf{Y} of size $m_1 \times \cdots \times m_d$, select $s, t \subset \{1, \dots, d\}$ such that $|s| = |t|$ and $n_{s_i} = m_{t_i}$ for $i = 1, \dots, |s|$. The tensor-tensor contraction of \mathbf{X} and \mathbf{Y} along modes s, t , denoted $\langle \mathbf{X}, \mathbf{Y} \rangle_{(s,t)}$, is defined as $(2d - (|s| + |t|))$ -tensor \mathbf{Z} of size (n_{s^c}, m_{t^c}) , satisfying

$$\mathbf{Z} = \langle \mathbf{X}, \mathbf{Y} \rangle_{(s,t)} = (X^{(s^c)} Y^{(t)})_{(s^c), (t^c)}.$$

Tensor tensor contraction over modes s and t merely sums over the dimensions specified by s, t in \mathbf{X} and \mathbf{Y} respectively, leaving the dimensions s^c and t^c free.

The inner product $\langle \mathbf{X}, \mathbf{Y} \rangle$ is a special case of tensor-tensor contraction when $s = t = \{1, \dots, d\}$.

We also make use of the fact that when the index sets s, t are $s, t = [d] \setminus i$ with \mathbf{X} , \mathbf{Y} , and A_i are appropriately sized for $i = 1, \dots, d$, then

$$(2) \quad \langle A_1 \times_1 A_2 \times_2 \dots A_d \times_d \mathbf{X}, \mathbf{Y} \rangle_{[d] \setminus i, [d] \setminus i} = A_i \langle A_1 \times_1 A_2 \times_2 \dots A_{i-1} \times_{i-1} A_{i+1} \times_{i+1} \dots A_d \times_d \mathbf{X}, \mathbf{Y} \rangle_{([d] \setminus i), ([d] \setminus i)}$$

i.e., applying A_i to dimension i commutes with contracting tensors over every dimension except the i th one.

3. SMOOTH MANIFOLD GEOMETRY OF THE HIERARCHICAL TUCKER FORMAT

In this section, we review the definition of the Hierarchical Tucker format (Section 3.1) as well as previous results [53] in the smooth manifold geometry of this format (Section 3.2). We extend these results in the next section by introducing a Riemannian metric on the space of HT parameters and subsequently derive the associated Riemannian gradient with respect to this metric. A reader familiar with the results in [53] can glance over this section quickly for a few instances of notation and move on to Section 4.

3.1. Hierarchical Tucker Format. The standard definition of the Hierarchical Tucker format relies on the notion of a *dimension tree*, chosen apriori, which specifies the format [20]. Intuitively, the dimension tree specifies which groups of dimensions are “separated” from other groups of dimensions, where “separation” is used in a similar sense to the SVD in two dimensions.

Definition 3. A dimension tree T is a non-trivial binary tree such that

- the root, t_{root} , has the label $t_{root} = \{1, 2, \dots, d\}$
- for every $t \notin L$, where L is the set of leaves of T , the labels of its left and right children, t_l, t_r , form a partition of the label for t , i.e., $t_l \cup t_r = t$ and $t_l \cap t_r = \emptyset$.

We set $N(T) := T \setminus L$. An example of a dimension tree when $d = 6$ is given in Figure 1.

Remark 1. For the following derivations, we take the point of view that each quantity with a subscript $(\cdot)_t$ is associated to the node $t \in T$. By Definition 3, for each $t \in T$, there is a corresponding subset of $\{1, \dots, d\}$ associated to t . If our HT tensor has dimensions $n_1 \times n_2 \times \dots \times n_d$, we let $n_t = \prod_{i \in t} n_i$ and, when $t \in N(T)$, n_t satisfies $n_t = n_{t_l} n_{t_r}$.

Definition 4. Given a dimension tree T and a vector of hierarchical ranks $(k_t)_{t \in T}$ with $k_t \in \mathbb{Z}^+$, $k_{t_{\text{root}}} = 1$, a tensor $\mathbf{X} \in \mathbb{R}^{n_1 \times n_2 \times \dots \times n_d}$ can be written in the Hierarchical Tucker format if there exist parameters $x = ((U_t)_{t \in L}, (\mathbf{B}_t)_{t \in N(T)})$ such that $\phi(x) = \mathbf{X}$, where

$$(3) \quad \begin{aligned} \text{vec}(\phi(x)) &= U_{t_l} \times_1 U_{t_r} \times_2 B_{t_{\text{root}}} & t &= t_{\text{root}} \\ U_t &= (U_{t_l} \times_1 U_{t_r} \times_2 \mathbf{B}_t)^{(1,2)} & t &\in N(T) \setminus t_{\text{root}} \end{aligned}$$

where $U_t \in \mathbb{R}_*^{n_t \times k_t}$, the set of full-rank $n_t \times k_t$ matrices, for $t \in L$ and $\mathbf{B}_t \in \mathbb{R}_*^{k_{t_l} \times k_{t_r} \times k_t}$, the set of 3-tensors of full multilinear rank, i.e.,

$$\text{rank}(B_t^{(1)}) = k_{t_l}, \quad \text{rank}(B_t^{(2)}) = k_{t_r}, \quad \text{rank}(B_t^{(3)}) = k_t.$$

We say the parameters $x = (U_t, \mathbf{B}_t)$ are in *Orthogonal Hierarchical Tucker* (OHT) format if, in addition to the above construction, we also have

$$(4) \quad \begin{aligned} U_t^T U_t &= I_{k_t} \quad \text{for } t \in L \\ (B_t^{(1,2)})^T B_t^{(1,2)} &= I_{k_t} \quad \text{for } t \in N(T) \setminus t_{\text{root}} \end{aligned}$$

We have made a slight modification of the definition of the HT format compared to [53] for ease of presentation. When $d = 2$, our construction is the same as the subspace decomposition introduced in [38] for low-rank matrices, but our approach is not limited to this case.

Owing to the recursive construction (4), the intermediate matrices U_t for $t \in N(T)$ do not need to be stored. Instead, specifying U_t for $t \in L$ and \mathbf{B}_t for $t \in N(T)$ determines $\mathbf{X} = \phi(x)$ completely. Therefore, the overall number of parameters $x = ((U_t)_{t \in L}, (\mathbf{B}_t)_{t \in N(T)})$ is bounded above by $dNK + (d-2)K^3 + K^2$, where $N = \max_{i=1, \dots, d} n_i$ and $K = \max_{t \in T} k_t$. When $d \geq 4$ and $K \ll N$, this quantity is much less than the N^d parameters typically needed to represent \mathbf{X} .

Definition 5. The hierarchical rank of a tensor $\mathbf{X} \in \mathbb{R}^{n_1 \times n_2 \times \dots \times n_d}$ corresponding to a dimension tree T is the vector $\mathbf{k} = (k_t)_{t \in T}$ where $k_{t_{\text{root}}} = 1$ and for $t \in T \setminus t_{\text{root}}$,

$$k_t = \text{rank}(X^{(t)}).$$

We consider the set of *Hierarchical Tucker tensors of fixed rank* $\mathbf{k} = (k_t)_{t \in T}$, that is,

$$\mathcal{H} = \{\mathbf{X} \in \mathbb{R}^{n_1 \times n_2 \times \dots \times n_d} \mid \text{rank}(X^{(t)}) = k_t \quad \text{for } t \in T \setminus t_{\text{root}}\}.$$

We consider general HT tensors in the sequel, but we implement our algorithms with OHT parameters. In addition to significantly simplifying the resulting notation, this restriction allows us to avoid cumbersome and unnecessary matrix inversions, in particular for the resulting subspace projections in future sections. Moreover, the orthogonal format yields more accurate computations of inner products and subspace projections in finite arithmetic. This restriction does not reduce the expressibility of the HT format, however, since for any non orthogonalized parameters x such that $\mathbf{X} = \phi(x)$, there exists orthogonalized parameters x' with $\mathbf{X} = \phi(x')$ [20, Alg. 3].

We use the grouping $x = (U_t, \mathbf{B}_t)$ to denote $((U_t)_{t \in L}, (\mathbf{B}_t)_{t \in N(T)})$, as these are our “independent variables” of interest in this case. In order to avoid cumbersome notation, we also suppress the dependence on (T, \mathbf{k}) in the following, and presume a fixed dimension tree T and hierarchical ranks \mathbf{k} .

3.2. Quotient Manifold Geometry. In the interest of keeping this manuscript self-contained, we briefly summarize the key results of [53] that we will use for the following sections.

Given the full-rank constraints on each parameter, the space of possible $x = (U_t, \mathbf{B}_t)$ is written as

$$\mathcal{M} = \bigtimes_{t \in L} \mathbb{R}_*^{n_t \times k_t} \times \bigtimes_{t \in N(T)} \mathbb{R}_*^{k_{t_l} \times k_{t_r} \times k_t}.$$

\mathcal{M} is an open submanifold of $\mathbb{R}^{\sum_{t \in L} n_t k_t + \sum_{t \in N(T)} k_{t_l} k_{t_r} k_t}$ with corresponding tangent space

$$\mathcal{T}_x \mathcal{M} = \bigtimes_{t \in L} \mathbb{R}^{n_t \times k_t} \times \bigtimes_{t \in N(T)} \mathbb{R}^{k_{t_l} \times k_{t_r} \times k_t}.$$

Let $\phi : \mathcal{M} \rightarrow \mathcal{H}$ be the parameter to tensor map in (3). Then for each $\mathbf{X} \in \mathcal{H}$, then there is an inherent ambiguity in its representation by parameters x , i.e., $\mathbf{X} = \phi(x) = \phi(y)$ for distinct parameters x and y with the following relationship between them. Let \mathcal{G} be the Lie group

$$(5) \quad \mathcal{G} = \{(A_t)_{t \in T} : A_t \in GL(k_t), t \neq t_{\text{root}}, A_{t_{\text{root}}} = 1\}.$$

where $GL(p)$ is the matrix group of invertible $p \times p$ matrices and the group action of component-wise multiplication.

Let θ be the group action

$$(6) \quad \begin{aligned} \theta : \mathcal{M} \times \mathcal{G} &\rightarrow \mathcal{M} \\ (x, \mathcal{A}) &:= ((U_t, B_t), (A_t)) \mapsto \theta_x(\mathcal{A}) := (U_t A_t, A_{t_l}^{-1} \times_1 A_{t_r}^{-1} \times_2 A_t^T \times_3 \mathbf{B}_t). \end{aligned}$$

Then $\phi(x) = \phi(y)$ if and only if there exists a unique $\mathcal{A} = (A_t)_{t \in T} \in \mathcal{G}$ such that $x = \theta_{\mathcal{A}}(y)$ [53, Prop. 3]. Therefore these are the only types of ambiguities we must consider in this format.

It follows that the orbit of x ,

$$\mathcal{G}x = \{\theta_{\mathcal{A}}(x) : \mathcal{A} \in \mathcal{G}\},$$

is the set of all parameters that map to the same tensor $\mathbf{X} = \phi(x)$ under ϕ . This induces an equivalence relation on the set of parameters \mathcal{M} ,

$$x \sim y \text{ if and only if } y \in \mathcal{G}x.$$

If we let \mathcal{M}/\mathcal{G} be the corresponding quotient space of equivalence classes and $\pi : \mathcal{M} \rightarrow \mathcal{M}/\mathcal{G}$ denote the quotient map, then pushing ϕ down through π results in an injective function

$$\hat{\phi} : \mathcal{M}/\mathcal{G} \rightarrow \mathcal{H}$$

whose image is all of \mathcal{H} , and hence is an isomorphism (in fact, a diffeomorphism).

The vertical space, $\mathcal{V}_x \mathcal{M}$, is the subspace of $\mathcal{T}_x \mathcal{M}$ that is tangent to $\pi^{-1}(x)$. That is, $dx^v = (\delta U_t^v, \delta \mathbf{B}_t^v) \in \mathcal{V}_x \mathcal{M}$ when it is of the form [53, Eq. 26]

$$(7) \quad \begin{aligned} \delta U_t^v &= U_t D_t && \text{for } t \in L \\ \delta \mathbf{B}_t^v &= D_t \times_3 \mathbf{B}_t - D_{t_l} \times_1 \mathbf{B}_t - D_{t_r} \times_2 \mathbf{B}_t && \text{for } t \in N(T) \cup t_{\text{root}} \\ \delta B_{t_{\text{root}}}^v &= -D_{t_l} B_{t_{\text{root}}} - B_{t_{\text{root}}} D_{t_r}^T && \text{for } t = t_{\text{root}} \end{aligned}$$

where $D_t \in \mathbb{R}^{k_t \times k_t}$. A straightforward computation shows that $D\phi(x)|_{\mathcal{V}_x \mathcal{M}} \equiv 0$, and therefore for every $dx^v \in \mathcal{V}_x \mathcal{M}$, $\phi(x) = \phi(x + dx^v)$ to first order in dx^v . From an optimization point of view, moving from the point x to $x + dx^v$, for small dx^v , will not change the current tensor $\phi(x)$ and therefore for any search direction p , we must filter out the corresponding component in $\mathcal{V}_x \mathcal{M}$ in order to compute the gradient correctly. We accomplish this by projecting on to a *horizontal space*, which is any complementary subspace to $\mathcal{V}_x \mathcal{M}$. One such choice is [53, Eq. 26],

$$(8) \quad \mathcal{H}_x \mathcal{M} = \left\{ (\delta U_t^h, \delta \mathbf{B}_t^h) : \begin{aligned} &(\delta U_t^h)^T U_t = 0_{k_t} \text{ for } t \in L \\ &(\delta B_t^{(1,2)})^T (U_{t_r}^T U_{t_r} \otimes U_{t_l}^T U_{t_l}) B_t^{(1,2)} = 0_{k_t} \text{ for } t \in N(T) \setminus t_{\text{root}} \end{aligned} \right\}.$$

Note that there is no restriction on $B_{t_{\text{root}}}^h$, which is a matrix.

This choice has the convenient property that $\mathcal{H}_x \mathcal{M}$ is *invariant* under the action of θ , i.e., [53, Prop. 5]

$$(9) \quad D\theta(x, \mathcal{A})[\mathcal{H}_x \mathcal{M}, 0] = \mathcal{H}_{\theta_x(\mathcal{A})} \mathcal{M},$$

which we shall exploit for our upcoming discussion of a Riemannian metric. The horizontal space $\mathcal{H}_x \mathcal{M}$ allows us to uniquely represent abstract tangent vectors in $T_{\pi(x)} \mathcal{M}/\mathcal{G}$ with concrete vectors in $\mathcal{H}_x \mathcal{M}$.

4. RIEMANNIAN GEOMETRY OF THE HT FORMAT

In this section, we introduce a Riemannian metric on the parameter space \mathcal{M} that will allow us to use parameters x as representations for their equivalence class $\pi(x)$ in a well-defined manner when performing numerical optimization.

4.1. Riemannian metric. Since each distinct equivalence class $\pi(x)$ is uniquely identified with each distinct value of $\phi(x)$, the quotient manifold \mathcal{M}/\mathcal{G} is really our manifold of interest for the purpose of computations—i.e., we would like to formulate our optimization problem over the equivalence classes $\pi(x)$. By introducing a Riemannian metric on \mathcal{M} that respects its quotient structure, we can formulate concrete optimization algorithms in terms of the HT parameters without being affected by the non-uniqueness of the format—i.e., by optimizing over parameters x while implicitly performing optimization over equivalence classes $\pi(x)$. Below, we explain how to explicitly construct this Riemannian metric for the HT format.

Let $\eta_x = (\delta U_t, \delta \mathbf{B}_t)$, $\zeta_x = (\delta V_t, \delta \mathbf{C}_t) \in \mathcal{T}_x \mathcal{M}$ be tangent vectors at the point $x = (U_t, \mathbf{B}_t) \in \mathcal{M}$. We let $M_t = U_t^T U_t$. Then we define the inner product $g_x(\cdot, \cdot)$ at x as

$$(10) \quad \begin{aligned} g_x(\eta_x, \zeta_x) := & \sum_{t \in L} \text{tr}((M_t)^{-1} \delta U_t^T \delta V_t) \\ & + \sum_{t \in M(T) \setminus \mathbf{t}_{\text{root}}} \langle \delta \mathbf{B}_t, (M_{t_l}) \times_1 (M_{t_r}) \times_2 (M_t)^{-1} \times_3 \delta \mathbf{C}_t \rangle \\ & + \text{tr}(M_{(\mathbf{t}_{\text{root}})_r} \delta B_{\mathbf{t}_{\text{root}}}^T M_{(\mathbf{t}_{\text{root}})_l} \delta C_{\mathbf{t}_{\text{root}}}). \end{aligned}$$

By the full-rank conditions on U_t and \mathbf{B}_t at each node, by definition of the HT format, each M_t , for $t \in T$, is symmetric positive definite and varies smoothly with $x = (U_t, \mathbf{B}_t)$. As a result, g_x is a smooth, symmetric positive definite, bilinear form on $\mathcal{T}_x \mathcal{M}$, i.e., a Riemannian metric. Note that when x is in OHT, as in future sections, g_x reduces to the standard Euclidean product on the parameter space $\mathcal{T}_x \mathcal{M}$, making it straightforward to compute in this case.

Proposition 2. *On the Riemannian manifold (\mathcal{M}, g) , θ defined in (6) acts isometrically on \mathcal{M} , i.e., for every $\mathcal{A} \in \mathcal{G}$, $\xi_x, \zeta_x \in \mathcal{H}_x \mathcal{M}$*

$$g_x(\xi_x, \zeta_x) = g_{\theta_{\mathcal{A}}(x)}(\theta^* \xi_x, \theta^* \zeta_x)$$

where θ^* is the push-forward map, $\theta^* v = D\theta(x, \mathcal{A})[v]$.

Proof. Let $x = (U_t, \mathbf{B}_t) \in \mathcal{M}$, $y = (V_t, \mathbf{C}_t) = \theta_{\mathcal{A}}(x) = (U_t A_t, A_{t_l}^{-1} \times_1 A_{t_r}^{-1} \times_2 A_t^T \times_3 \mathbf{B}_t)$ for $\mathcal{A} \in \mathcal{G}$.

If we write $\eta_x = (\delta U_t, \delta \mathbf{B}_t)$, $\zeta_x = (\delta V_t, \delta \mathbf{C}_t)$ for $\eta_x, \zeta_x \in \mathcal{H}_x \mathcal{M}$, then, by (9), it follows that

$$\eta_y = \theta^* \eta_x = (\delta U_t A_t, A_{t_l}^{-1} \times_1 A_{t_r}^{-1} \times_2 A_t^T \times_3 \delta \mathbf{B}_t)$$

and similarly for ζ_y .

We will compare each component of the sum of (10) term by term. For ease of presentation, we only consider interior nodes $t \in N(T) \setminus \mathbf{t}_{\text{root}}$, as leaf nodes and the root node are handled in an analogous manner.

For $t \notin L \cup \mathbf{t}_{\text{root}}$, let $\widetilde{\delta \mathbf{B}}_t$ be the component of η_y at the node t , i.e.,

$$\widetilde{\delta \mathbf{B}}_t = A_{t_l}^{-1} \times_1 A_{t_r}^{-1} \times_2 A_t^T \times_3 \delta \mathbf{B}_t$$

and similarly for $\widetilde{\delta \mathbf{C}}_t$.

The above inner product, evaluated at x , between $\delta \mathbf{B}_t$ and $\delta \mathbf{C}_t$, evaluated at x , can be written as

$$\text{vec}(\delta \mathbf{B}_t)^T (M_t^{-1} \otimes M_{t_r} \otimes M_{t_l}) \text{vec}(\delta \mathbf{C}_t),$$

and similarly for the inner product between $\widetilde{\delta \mathbf{B}}_t$ and $\widetilde{\delta \mathbf{C}}_t$ at y .

By setting $\tilde{M}_t := V_t^T V_t$, a quick computation shows that $\tilde{M}_t = A_t^T U_t^T U_t A_t = A_t^T M_t A_t$. Therefore, we have that the inner product between $\widetilde{\delta \mathbf{B}_t}$ and $\widetilde{\delta \mathbf{C}_t}$ at y is

$$\begin{aligned} & \text{vec}(\widetilde{\delta \mathbf{B}_t})^T (\tilde{M}_t^{-1} \otimes \tilde{M}_{t_r} \otimes \tilde{M}_{t_l}) \text{vec}(\widetilde{\delta \mathbf{C}_t}) \\ &= \text{vec}(\delta \mathbf{B}_t)^T (A_t \otimes A_{t_r}^{-T} \otimes A_{t_l}^{-T}) ((A_t^{-1} M_t^{-1} A_t^{-T}) \otimes (A_{t_r}^T M_{t_r} A_{t_r}) \otimes (A_{t_l}^T M_{t_l} A_{t_l})) \\ & \quad (A_t^T \otimes A_{t_r}^{-1} \otimes A_{t_l}^{-1}) \text{vec}(\delta \mathbf{C}_t) \\ &= \text{vec}(\delta \mathbf{B}_t)^T (M_t^{-1} \otimes M_{t_r} \otimes M_{t_l}) \text{vec}(\delta \mathbf{C}_t). \end{aligned}$$

Therefore, this term in the inner product is invariant under the group action θ and by adding the terms for each $t \in T$, we obtain that

$$g_x(\xi_x, \zeta_x) = g_{\theta_A(x)}(\theta^* \xi_x, \theta^* \zeta_x).$$

□

As we are interested in carrying out our optimization using the HT parameters x as proxies for their equivalence classes $\pi(x)$, this proposition states that if we measure inner products between two tangent vectors at the point x , we obtain the same result as if we had measured the inner product between two tangent vectors transformed by θ_A at the point $\theta_A(x)$. In this sense, once we have a unique association of tangent vectors in \mathcal{M}/\mathcal{G} with a subspace of $\mathcal{T}_x \mathcal{M}$, we can use the actual representatives, the parameters x , instead of the abstract equivalence class $\pi(x)$, in a well-defined way during our optimization. This shows that \mathcal{M}/\mathcal{G} , endowed with the Riemannian metric

$$g_{\pi(x)}(\xi, \zeta) := g_x(\xi_x^h, \zeta_x^h)$$

where ξ_x^h, ζ_x^h are the horizontal lifts at x of ξ, ζ , respectively, is a Riemannian quotient manifold of \mathcal{M} (i.e., $\pi : \mathcal{M} \rightarrow \mathcal{M}/\mathcal{G}$ is a Riemannian submersion) [1, Sec. 3.6.2].

In summary, by using this Riemannian metric and restricting our optimization to only consider horizontal tangent vectors, we can implicitly formulate our algorithms on the abstract quotient space by working with the concrete HT parameters. Below, we will derive the Riemannian gradient in this context.

Remark 2. *It should be noted that although the horizontal space (8) is complementary to the vertical space (7), it is demonstrably not perpendicular to $\mathcal{V}_x \mathcal{M}$ under the Riemannian metric (10). Choosing a horizontal space which is perpendicular to $\mathcal{V}_x \mathcal{M}$ under the standard Euclidean product (i.e., (10) when x is orthogonalized) is beyond the scope of this paper. Suffice to say, it can be done, as a generalization of the approach outlined in [38], resulting in a series of symmetry conditions on various multi-way combinations of parameters. The resulting projection operators involve solving a number of coupled Lyapunov equations, increasing with the depth of T . It remains to be seen whether such equations can be solved efficiently when d is large. We will not dwell on this point here, as we will not be needing orthogonal projections for our computations in the following.*

We restrict ourselves to orthogonal parameters from this point forward, for the reasons stated previously. In order to not overburden ourselves with notation, we use the notation \mathcal{M} to refer to orthogonal parameters and the corresponding group acting on \mathcal{M} as \mathcal{G} for the remainder of this paper. In particular, when restricting to orthogonal HT parameters, the general linear group $GL(k_t)$ in (5) is replaced by $O(k_t)$, the group of orthogonal $k_t \times k_t$ matrices. The expression for the horizontal space (8) and the Riemannian metric are also simplified since, for orthogonal parameters, $U_t^T U_t = I_{k_t}$ for all $t \in T \setminus t_{\text{root}}$.

4.2. Riemannian gradient. The problem we are interested in solving is

$$\min_{x \in \mathcal{M}} f(\phi(x))$$

for a smooth objective function $f : \mathbb{R}^{n_1 \times n_2 \times \dots \times n_d} \rightarrow \mathbb{R}$. We write $\hat{f} : \mathcal{M} \rightarrow \mathbb{R}$, where $\hat{f}(x) = f(\phi(x))$.

We need to derive expressions for the Riemannian gradient to update the HT parameters as part of local optimization procedures. Therefore, our primary quantity of interest is the *Riemannian gradient* of \hat{f} .

Definition 6. [1, Sec. 3.6] Given a smooth scalar function \hat{f} on a Riemannian manifold \mathcal{N} , the Riemannian gradient of \hat{f} at $x \in \mathcal{N}$, denoted $\nabla^R \hat{f}(x)$, is the unique element of $\mathcal{T}_x \mathcal{N}$ which satisfies

$$g_x(\nabla^R \hat{f}(x), \xi) = D\hat{f}(x)[\xi] \quad \forall \xi \in \mathcal{T}_x \mathcal{N}$$

with respect to the Riemannian metric $g_x(\cdot, \cdot)$.

Our manifold of interest in this case is $\mathcal{N} = \mathcal{M}/\mathcal{G}$, with the corresponding horizontal space $\mathcal{H}_x \mathcal{M}$ in lieu of the abstract tangent space $T_{\pi(x)} \mathcal{M}/\mathcal{G}$. Therefore, in the above equation, we can consider the horizontal lift ξ^h of the tangent vector ξ and instead write

$$g_x(\nabla^R \hat{f}(x), \xi^h) = D\hat{f}(x)[\xi^h].$$

Our derivation is similar to that of [53, Sec 6.2.2], except our derivations are more streamlined and cheaper computationally since we reduce the operations performed at the interior nodes $t \in N(T)$. By a slight abuse of notation in this section, we denote variational quantities associated to node t as $\delta Z_t \in \mathbb{R}^{n_{t_l} n_{t_r} \times k_t}$ and $\delta \mathbf{Z}_t \in \mathbb{R}^{n_{t_l} \times n_{t_r} \times k_t}$ where $(\delta Z_t)_{(1,2)} = \delta \mathbf{Z}_t$ is the reshaping of δZ_t into a 3-tensor. The Riemannian gradient will be denoted $(\delta U_t, \delta \mathbf{B}_t)$ and a general horizontal vector will be denoted by $(\delta V_t, \delta \mathbf{C}_t)$.

When $x = (U_t, \mathbf{B}_t)$ is orthogonalized, we use $\langle \cdot, \cdot \rangle$ to denote the Euclidean inner product. By the chain rule, we have that, for any $\xi = (\delta V_t, \delta \mathbf{C}_t) \in \mathcal{H}_x \mathcal{M}$,

$$\begin{aligned} D\hat{f}(x)[\xi] &= Df(\phi(x))[D\phi(x)[\xi]] \\ &= \langle \nabla_{\phi(x)} f(\phi(x)), D\phi(x)[\xi] \rangle. \end{aligned}$$

Then each tensor $\delta \mathbf{V}_t \in \mathbb{R}^{n_{t_l} \times n_{t_r} \times k_t}$, with $\delta V_{t_{\text{root}}} = D\phi(x)[\xi]$, satisfies the recursion

$$(11) \quad \delta \mathbf{V}_t = \delta V_{t_l} \times_1 U_{t_r} \times_2 \mathbf{B}_t + U_{t_l} \times_1 \delta V_{t_r} \times_2 \mathbf{B}_t + U_{t_l} \times_1 U_{t_r} \times_2 \delta \mathbf{C}_t,$$

for matrices $\delta V_{t_l} \in \mathbb{R}^{n_{t_l} \times k_{t_l}}$, $\delta V_{t_r} \in \mathbb{R}^{n_{t_r} \times k_{t_r}}$ and tensor $\delta \mathbf{C}_t \in \mathbb{R}^{k_{t_l} \times k_{t_r} \times k_t}$ satisfying [53, Lemma 2]

$$(12) \quad \delta V_{t_l}^T U_{t_l} = 0 \quad \delta V_{t_r}^T U_{t_r} = 0 \quad (\delta \mathbf{C}_t^{(1,2)})^T B_t^{(1,2)} = 0.$$

The third orthogonality condition is omitted when $t = t_{\text{root}}$.

Owing to this recursive structure, we compute $\langle \delta \mathbf{U}_t, \delta \mathbf{V}_t \rangle$, where $\delta \mathbf{U}_t$ is the component of the Riemannian gradient at the current node and recursively extract the components of the Riemannian gradient associated to the children, i.e., δU_{t_l} , δU_{t_r} , and $\delta \mathbf{B}_t$. Here we let $\delta U_{t_{\text{root}}} = \nabla_{\phi(x)} f(\phi(x))$ be the Euclidean gradient of $f(\phi(x))$ at $\phi(x)$, reshaped into a matrix of size $n_{(t_{\text{root}})_l} \times n_{(t_{\text{root}})_r}$.

We set

$$(13) \quad \begin{aligned} \delta U_{t_l} &= P_{U_{t_l}}^\perp \langle U_{t_r}^T \times_2 \delta \mathbf{U}_t, \mathbf{B}_t \rangle_{(2,3),(2,3)} \\ \delta U_{t_r} &= P_{U_{t_r}}^\perp \langle U_{t_l}^T \times_1 \delta \mathbf{U}_t, \mathbf{B}_t \rangle_{(1,3),(1,3)} \\ \delta \mathbf{B}_t &= U_{t_l}^T \times_1 U_{t_r}^T \times_2 \delta \mathbf{U}_t \end{aligned}$$

followed by a projection of $\delta \mathbf{B}_t^{(1,2)}$ on to $\text{span}(B_t^{(1,2)})^\perp$ if $t \in N(T) \setminus t_{\text{root}}$. Here $P_{U_t} = (I_{k_t} - U_t U_t^T)$ is the usual projection on to $\text{span}(U_t)^\perp$. After a straightforward calculation, it follows that $\langle \delta \mathbf{U}_t, \delta \mathbf{V}_t \rangle$ is equal to

$$\langle \delta U_{t_l}, \delta V_{t_l} \rangle + \langle \delta U_{t_r}, \delta V_{t_r} \rangle + \langle \delta \mathbf{B}_t, \delta \mathbf{C}_t \rangle$$

and δU_{t_l} , δU_{t_r} , and $\delta \mathbf{B}_t$ satisfy (12). Their recursively decomposed factors will therefore be in the horizontal space $\mathcal{H}_x \mathcal{M}$.

$\delta \mathbf{B}_t$ is the component of the Riemannian gradient at node t . If t_l is a leaf node, then we have extracted the component of the Riemannian gradient associated to t_l , namely δU_{t_l} . Otherwise, we set $\delta \mathbf{U}_{t_l} = (\delta U_{t_l})_{(1,2)}$ and apply the above recursion. We make the same considerations for the right children. In the above computations, the multilinear product operators are never formed explicitly and instead each operator is applied to various reshapings of the matrix or tensor of interest, see [16] for a reference Matlab implementation.

We make the following observations in order to minimize the number of computations performed on intermediate tensors, which can be much larger than $\dim(\mathcal{M})$. In computing the terms

$$P_{U_{t_l}}^\perp \langle U_{t_r}^T \times_2 \delta \mathbf{U}_t, \mathbf{B}_t \rangle_{(2,3),(2,3)},$$

we have that $\delta \mathbf{U}_t = (P_{U_t}^\perp \delta \tilde{U}_t)_{(1,2)}$ for a matrix $\delta \tilde{U}_t \in \mathbb{R}^{n_{t_l} n_{t_r} \times k_t}$. Using (2), the above expression can be written as

$$(14) \quad \langle P_{U_{t_l}}^\perp \times_1 U_{t_r}^T \times_2 (P_{U_t}^\perp \delta \tilde{U}_t)_{(1,2)}, \mathbf{B}_t \rangle_{(2,3),(2,3)}.$$

We note that in the above, $P_{U_{t_l}}^\perp \times_1 U_{t_r}^T \times_2 (P_{U_t}^\perp \delta \tilde{U}_t)_{(1,2)} = (U_{t_r}^T \otimes P_{U_{t_l}}^\perp P_{U_t}^\perp \delta \tilde{U}_t)_{(1,2)}$, and the operator applied to $\delta \tilde{U}_t$ satisfies

$$\begin{aligned} U_{t_r}^T \otimes P_{U_{t_l}}^\perp P_{U_t}^\perp &= U_{t_r}^T \otimes P_{U_{t_l}}^\perp (I_{n_t} - U_t U_t^T) \\ &= U_{t_r}^T \otimes P_{U_{t_l}}^\perp (I_{n_t} - U_{t_r} \otimes U_{t_l} B_t^{(1,2)} (B_t^{(1,2)})^T U_{t_r}^T \otimes U_{t_l}^T) \\ &= U_{t_r}^T \otimes P_{U_{t_l}}^\perp. \end{aligned}$$

This means that, using (2), we can write (14) as

$$P_{U_{t_l}}^\perp \langle U_{t_r}^T \times_2 (\delta \tilde{U}_t)_{(1,2)}, \mathbf{B}_t \rangle_{(2,3),(2,3)}$$

i.e., we do not have to apply $P_{U_t}^\perp$ to the matrix $\delta \tilde{U}_t$ at the parent node of t . Applying this observation recursively and to the other terms in the Riemannian gradient, we merely need to orthogonally project the resulting extracted parameters $(\delta U_t, \delta \mathbf{B}_t)$ on to $\mathcal{H}_x \mathcal{M}$ after applying the formula (13) *without* applying the intermediate operators $P_{U_t}^\perp$, reducing the overall computational costs. We summarize our algorithm for computing the Riemannian gradient in Algorithm 1.

Algorithm 1 The Riemannian gradient $\nabla^R f$ at a point $x = (U_t, \mathbf{B}_t) \in \mathcal{M}$

Require: $x = (U_t, \mathbf{B}_t)$ parameter representation of the current point.

Compute $\mathbf{X} = \phi(x)$ and $\nabla_{\mathbf{X}} f(\mathbf{X})$, the Euclidean gradient of f , a $n_1 \times \dots \times n_d$ tensor.

$\delta U_{t_{\text{root}}} \leftarrow (\nabla_{\mathbf{X}} f(\mathbf{X}))_{(1,2)}$

for $t \in N(T)$, visiting parents before their children **do**

$\delta \mathbf{U}_t \leftarrow (\delta U_t)_{(1,2)}$

$\delta U_{t_l} \leftarrow \langle U_{t_r}^T \times_2 \delta \mathbf{U}_t, \mathbf{B}_t \rangle_{(2,3),(2,3)}$, $\delta U_{t_r} \leftarrow \langle U_{t_l}^T \times_1 \delta \mathbf{U}_t, \mathbf{B}_t \rangle_{(1,3),(1,3)}$

$\delta \mathbf{B}_t \leftarrow U_{t_l}^T \times_1 U_{t_r}^T \times_2 \delta \mathbf{U}_t$

if $t \neq t_{\text{root}}$ **then**

$\delta \mathbf{B}_t \leftarrow (P_{B_t^{(1,2)}}^\perp (\delta B_t)^{(1,2)})_{(1,2)}$

end if

end for

for $t \in L$ **do**

$\delta U_t \leftarrow P_{U_t}^\perp \delta U_t$

end for

return $\nabla^R f \leftarrow (\delta U_t, \delta \mathbf{B}_t)$

Algorithm 1 is computing the operator $D\phi(x)^* : \mathbb{R}^{n_1 \times n_2 \times \dots \times n_d} \rightarrow \mathcal{H}_x \mathcal{M}$ applied to the Euclidean gradient $\nabla_{\phi(x)} f(\phi(x))$. The forward operator $D\phi(x) : \mathcal{H}_x \mathcal{M} \rightarrow \mathcal{T}_{\phi(x)} \mathcal{H} \subset \mathbb{R}^{n_1 \times n_2 \times \dots \times n_d}$ can be computed using a component-wise orthogonal projection $P_x^H : \mathcal{T}_x \mathcal{M} \rightarrow \mathcal{H}_x \mathcal{M}$ followed by applying (11) recursively.

4.3. Tensor Completion Objective and Gradient. In this section, we specialize the computation of the objective and Riemannian gradient in the HT format to the case where the Euclidean gradient of the objective function is sparse, in particular for tensor completion. This will allow us to scale our method to high dimensions in a straightforward fashion as opposed to the inherently dense considerations in Algorithm 1. Here for simplicity, we suppose that our dimension tree T is *complete*, that is a full binary tree up to level depth(T) - 1 and all of the leaves at level depth(T) are on the leftmost side of T , as in Figure 1. This will ease the exposition as well as allow for a more efficient implementation compared to a noncomplete tree. In what follows below, we denote $\mathbf{i} = (i_1, i_2, \dots, i_d)$, with $1 \leq i_j \leq n_j$ for all $1 \leq j \leq d$, and let \mathbf{i}_t be the subindices of \mathbf{i} indexed by $t \in T$.

We consider a separable, smooth objective function on the HT manifold,

$$(15) \quad \hat{f}(x) = f(\phi(x)) = \sum_{\mathbf{i} \in \Omega} f_{\mathbf{i}}(\phi(x)_{\mathbf{i}}),$$

where $f_{\mathbf{i}} : \mathbb{R} \rightarrow \mathbb{R}$ is a smooth, single variable function. For the least-squares tensor completion problem, $f_{\mathbf{i}}(a) = \frac{1}{2}(a - b_{\mathbf{i}})^2$.

In this section, we also use the Matlab notation for indexing in to matrices, i.e., $A(m, n)$ is the (m, n) th entry of A , and similarly for tensors. Let $K = \max_{t \in T} k_t$.

4.3.1. *Objective function.* With this notation in mind, we write each entry of $P_{\Omega}\phi(x)$, indexed by $\mathbf{i} \in \Omega$, as

$$(P_{\Omega}\phi(x))(\mathbf{i}) = \sum_{r_l=1}^{k_{t_l}} \sum_{r_r=1}^{k_{t_r}} (U_{t_l})(\mathbf{i}_{t_l}, r_l) \cdot (U_{t_r})(\mathbf{i}_{t_r}, r_r) \cdot B_{t_{\text{root}}}(r_l, r_r), \quad \text{where } t = t_{\text{root}}.$$

Each entry of U_{t_l}, U_{t_r} can be computed by applying the recursive formula (3), i.e.,

$$U_t(\mathbf{i}_t, r) = \sum_{r_l=1}^{k_{t_l}} \sum_{r_r=1}^{k_{t_r}} (U_{t_l})(\mathbf{i}_{t_l}, r_l) \cdot (U_{t_r})(\mathbf{i}_{t_r}, r_r) \cdot \mathbf{B}_t(r_l, r_r, r)$$

with the substitutions of $t \rightarrow t_l, t_r$ as appropriate.

At each node $t \in T$, we perform at most K^3 operations and therefore the computation of $P_{\Omega}\phi(x)$ requires at most $2|\Omega|dK^3$ operations. The least squares objective, $\frac{1}{2}\|P_{\Omega}\phi(x) - b\|_2^2$, can be computed in $|\Omega|$ operations.

4.3.2. *Riemannian gradient.* The Riemannian gradient is more involved, notation-wise, to derive explicitly compared to the objective, so in the interest of brevity we only concentrate on the recursion for computing δU_1 below.

We let $\mathbf{Z} = \nabla_{\phi(x)} f(\phi(x))$ denote the Euclidean gradient of $f(\mathbf{X})$ evaluated at $\mathbf{X} = \phi(x)$, which has nonzero entries $\mathbf{Z}(\mathbf{i})$ indexed by $\mathbf{i} \in \Omega$. By expanding out (13), for each $\mathbf{i} \in \Omega$, δU_{t_l} evaluated at the root node with coordinates \mathbf{i}_{t_l}, r_l for $r_l = 1, \dots, k_{t_l}$ is

$$\delta U_{t_l}(\mathbf{i}_{t_l}, r_l) = \sum_{\mathbf{i}=(\mathbf{i}_{t_l}, \mathbf{i}_{t_r}) \in \Omega} \mathbf{Z}(\mathbf{i}) \sum_{r_r=1}^{k_{t_r}} U_{t_r}(\mathbf{i}_{t_r}, r_r) B_{t_{\text{root}}}(r_l, r_r), \quad \text{where } t = t_{\text{root}}.$$

For each $t \in N(T) \cup t_{\text{root}}$, we let $\widetilde{\delta U}_{t_l}$ denote the length k_{t_l} vector, which depends on \mathbf{i}_t , satisfying, for each $\mathbf{i} \in \Omega$, $r_l = 1, \dots, k_{t_l}$,

$$(\widetilde{\delta U}_{t_l})(\mathbf{i}_t, r_l) = \sum_{r_r=1}^{k_{t_r}} \sum_{r_t=1}^{k_t} U_{t_r}(\mathbf{i}_{t_r}, r_r) \mathbf{B}_t(r_l, r_r, r_t) \widetilde{\delta U}_t(\mathbf{i}_t, r_t).$$

This recursive construction above, as well as similar considerations for the right children for each node, yield Algorithm 2. For each node $t \in T \setminus t_{\text{root}}$, we perform $3|\Omega|K^3$ operations and at the root where we perform $3|\Omega|K^2$ operations. The overall computation of the Riemannian gradient requires at most $6d|\Omega|K^3$ operations, when T is complete, and a negligible $O(dK)$ additional storage to store the vectors $\widetilde{\delta U}_t$ for each fixed $\mathbf{i} \in \Omega$. The computations above are followed by componentwise orthogonal projection on to $\mathcal{H}_x \mathcal{M}$, which requires $O(d(NK^2 + K^4))$ operations and are dominated by the $O(d|\Omega|K^3)$ time complexity when $|\Omega|$ is large.

Therefore for large $|\Omega|$, each evaluation of the objective, with or without the Riemannian gradient, requires $O(d|\Omega|K^3)$ operations. Since $f(\mathbf{X})$ can be written as a sum of component functions $f_{\mathbf{i}}$, each depending only on the entry $\mathbf{X}(\mathbf{i})$, we can parallelize the computation of $f(\mathbf{X})$ and the Riemannian gradient in the following way. For a system setup with p processors, we first partition the sampling set Ω in to p disjoint subsets, with the subset, say Ω_j , being assigned to the j th processor. The objective (and possibly the Riemannian gradient) are computed at each processor independently using the set Ω_j and the results are added together afterwards. This "embarrassingly parallel" structure, as referred to in the parallel computing literature, allows us to scale these algorithms to large problems in a distributed environment.

In certain situations, when say $|\Omega| = pN^d$ for some $p \in [10^{-3}, 1]$ and d is sufficiently small, say $d = 4, 5$, it may be more efficient from a computer hardware point of view to use the dense linear algebra formulation in

Algorithm 1 together with an efficient dense linear algebra library such as BLAS, rather than Algorithm 2. The dense formulation requires $O(N^d K)$ operations when T is a balanced tree, which may be smaller than the $O(d|\Omega|K^3)$ operations needed in this case.

Remark 3. By comparison, the gradient in the Tucker tensor completion case [34] requires $O(d(|\Omega|+N)K^d+K^{d+1})$ operations, which scales much more poorly when $d \geq 4$ compared to using Algorithm 2. This discrepancy is a result of the structural differences between Tucker and Hierarchical Tucker tensors, the latter of which allows one to exploit additional low-rank behaviour of the core tensor in the Tucker format.

Algorithm 2 Objective & Riemannian gradient for separable objectives

Require: $x = (U_t, \mathbf{B}_t)$ parameter representation of the current point

```

 $f_x \leftarrow 0, \delta U_t, \delta \mathbf{B}_t \leftarrow 0, \widetilde{\delta U}_t \leftarrow 0$ 
for  $\mathbf{i} \in \Omega$  do
  for  $t \in N(T)$ , visiting children before their parents do
    for  $z = 1, 2, \dots, k_t$  do
       $U_t(\mathbf{i}_t, z) \leftarrow \sum_{w=1}^{k_l} \sum_{y=1}^{k_r} (U_{t_l})(\mathbf{i}_{t_l}, w) \cdot (U_{t_r})(\mathbf{i}_{t_r}, y) \cdot \mathbf{B}_t(w, y, z)$ 
    end for
  end for
   $f_x \leftarrow f_x + f_{\mathbf{i}}(U_{t_{\text{root}}}(\mathbf{i}))$ 
   $\delta U_{t_{\text{root}}} \leftarrow \nabla f_{\mathbf{i}}(U_{t_{\text{root}}}(\mathbf{i}))$ 
  for  $t \in N(T)$ , visiting parents before their children do
    for  $w = 1, \dots, k_{t_l}, y = 1, \dots, k_{t_r}, z = 1, \dots, k_t$  do
       $\delta \mathbf{B}_t(w, y, z) \leftarrow \delta \mathbf{B}_t(w, y, z) + \widetilde{\delta U}_t(z) \cdot (U_{t_l})(\mathbf{i}_{t_l}, w) \cdot (U_{t_r})(\mathbf{i}_{t_r}, y)$ 
    end for
    for  $w = 1, \dots, k_{t_l}$  do
       $\widetilde{\delta U}_{t_l}(w) \leftarrow \sum_{y=1}^{k_{t_r}} \sum_{z=1}^{k_t} (U_{t_r})(\mathbf{i}_{t_r}, y) \cdot \mathbf{B}_t(w, y, z) \cdot \widetilde{\delta U}_t(z)$ 
    end for
    for  $y = 1, \dots, k_{t_r}$  do
       $\widetilde{\delta U}_{t_r}(y) \leftarrow \sum_{w=1}^{k_{t_l}} \sum_{z=1}^{k_t} (U_{t_l})(\mathbf{i}_{t_l}, w) \cdot \mathbf{B}_t(w, y, z) \cdot \widetilde{\delta U}_t(z)$ 
    end for
  end for
for  $t \in L$  do
  for  $z = 1, 2, \dots, k_t$  do
     $\delta U_t(\mathbf{i}_t, z) \leftarrow \delta U_t(\mathbf{i}_t, z) + \widetilde{\delta U}_t(z)$ 
  end for
end for
  Project  $(\delta U_t, \delta \mathbf{B}_t)$  componentwise on to  $\mathcal{H}_x \mathcal{M}$ 
end for
return  $f(x) \leftarrow f_x, \nabla^R f(x) \leftarrow (\delta U_t, \delta \mathbf{B}_t)$ 

```

5. OPTIMIZATION

5.1. Reorthogonalization as a retraction. As is standard in manifold optimization, we employ a *retraction* on the tangent bundle in lieu of the exponential mapping, the former being much less computationally demanding compared to the latter.

Definition 7. A retraction on a manifold \mathcal{N} is a smooth mapping \mathcal{R} from the tangent bundle \mathcal{TN} onto \mathcal{N} with the following properties: Let \mathcal{R}_x denote the restriction of \mathcal{R} to $\mathcal{T}_x \mathcal{N}$.

- $\mathcal{R}_x(0_x) = x$, where 0_x denotes the zero element of $\mathcal{T}_x \mathcal{N}$
- With the canonical identification $\mathcal{T}_{0_x} \mathcal{T}_x \mathcal{N} \simeq \mathcal{T}_x \mathcal{N}$, \mathcal{R}_x satisfies

$$D\mathcal{R}_x(0_x) = \text{id}_{\mathcal{T}_x \mathcal{N}}$$

where $\text{id}_{\mathcal{T}_x\mathcal{N}}$ denotes the identity mapping on $\mathcal{T}_x\mathcal{N}$ (Local rigidity condition).

A retraction approximates the action of the exponential mapping to first order and hence much of the analysis for algorithms utilizing the exponential mapping can also be carried over by the usual modifications to those using retractions. For more details, we refer the reader to [1].

A computationally feasible retraction on the HT parameters is given by QR-based reorthogonalization. The QR-based orthogonalization of (potentially nonorthogonal) parameters $x = (U_t, \mathbf{B}_t)$, denoted $QR(x)$, is given in Algorithm 3. Alternative retractions, such as those based on the SVD, can also be considered in a similar manner but we do not explore those options here. Since we are implicitly performing optimization on the quotient space \mathcal{M}/\mathcal{G} , the choice of retraction should not significantly affect algorithmic performance.

Proposition 3. *Given $x \in \mathcal{M}$, $\eta \in \mathcal{T}_x\mathcal{M}$, let $QR(x)$ be the QR-based orthogonalization defined in Algorithm 3. Then $\mathcal{R}_x(\eta) := QR(x + \eta)$ is a retraction on \mathcal{M} .*

We refer to Appendix A for the proof of this proposition.

As before for the Riemannian metric, we can treat the retractions on the HT parameter space as implicitly being retractions on the quotient space as outlined below.

Since $\mathcal{R}_x(\eta)$ is a retraction on the parameter space \mathcal{M} , and our horizontal space is invariant under the Lie group action, by the discussion in [1, 4.1.2], we have the following

Proposition 4. *The mapping*

$$\tilde{\mathcal{R}}_{\pi(X)}(z) = \pi(\mathcal{R}_X(Z))$$

is a retraction on \mathcal{M}/\mathcal{G} , where $\mathcal{R}_X(Z)$ is the QR retraction previously defined on \mathcal{M} , $\pi : \mathcal{M} \rightarrow \mathcal{M}/\mathcal{G}$ is the projection operator, and Z is the horizontal lift at X of the tangent vector z at $\pi(X)$.

Algorithm 3 QR-based orthogonalization ([20, Alg. 3])

Require: HT parameters $x = (U_t, \mathbf{B}_t)$.

return $y = (V_t, \mathbf{C}_t)$ orthogonalized parameters such that $\phi(x) = \phi(y)$

for $t \in L$ **do**

$Q_t R_t = U_t$, where Q_t is orthogonal and R_t is upper triangular with positive diagonal elements

$V_t \leftarrow Q_t$

end for

for $t \in T \setminus (L \cup \text{t}_{\text{root}})$, visiting children before their parents **do**

$Z_t \leftarrow R_{t_l} \times_1 R_{t_r} \times_2 \mathbf{B}_t$

$Q_t R_t = Z_t^{(1,2)}$, where Q_t is orthogonal and R_t is upper triangular

$\mathbf{C}_t \leftarrow (Q_t)_{(1,2)}$

end for

$C_{\text{t}_{\text{root}}} \leftarrow R_{(\text{t}_{\text{root}})_l} B_{\text{t}_{\text{root}}} R_{(\text{t}_{\text{root}})_r}^T$

5.2. Vector transport. The notion of vector transport allows us to relax the isometry constraints of parallel transport between tangent spaces at differing points on a manifold and decrease computational complexity without sacrificing the convergence quality of our CG-based method. Since our parameter space \mathcal{M} is a subset of Euclidean space, given a point $x \in \mathcal{M}$ and a horizontal vector $\eta_x \in \mathcal{H}_x\mathcal{M}$, we take our vector transport $\mathcal{T}_{x,\eta_x} : \mathcal{H}_x\mathcal{M} \rightarrow \mathcal{H}_{R_x(\eta_x)}\mathcal{M}$ of the vector $\xi_x \in \mathcal{H}_x\mathcal{M}$ to be

$$\mathcal{T}_{x,\eta_x}\xi_x := P_{\mathcal{R}_x(\eta_x)}^h \xi_x$$

where P_x^h is the component-wise projection onto the horizontal space at x , given in (8). This mapping is well defined on \mathcal{M}/\mathcal{G} since $\mathcal{H}_x\mathcal{M}$ is invariant under θ , and induces a vector transport on the quotient space [1, Sec. 8.1.4].

5.3. Smooth optimization methods. Now that we have established the necessary components for manifold optimization, we present a number of concrete optimization algorithms for solving

$$\min_{x \in \mathcal{M}} f(\phi(x)).$$

5.3.1. *First-order methods.* Given the expressions for the Riemannian gradient and retraction, it is straightforward to implement the classical Steepest Descent algorithm with an Armijo line search on this Riemannian manifold, specialized from the general Riemannian manifold case [1] to the HT manifold in Algorithm 4. This algorithm consists of computing the Riemannian gradient, followed by a line search, HT parameter update, and a reorthogonalization. We outline a nonlinear conjugate gradient method as outlined in Algorithm 4.

Here g_i denotes the Riemannian gradient at iteration i of the algorithm, p_i is the search direction for the optimization method, and α_i is the step length.

We choose the Polak-Ribiere approach

$$\beta_i = \frac{\langle g_i, g_i - \mathcal{T}_{x_{i-1}, \alpha_{i-1} p_{i-1}}(g_{i-1}) \rangle}{\langle g_{i-1}, g_{i-1} \rangle}$$

to compute the CG-parameter β_i , so that the search direction p_i satisfies

$$p_i = -g_i + \beta_i \mathcal{T}_{x_{i-1}, \alpha_{i-1} p_{i-1}} p_{i-1}$$

and $p_1 = -g_1$.

Algorithm 4 General Nonlinear Conjugate Gradient method for minimizing a function f over \mathcal{H}

Require: Initial guess $x_0 = (U_t, \mathbf{B}_t)$, $0 < \sigma \leq 1$ sufficient decrease parameter for the Armijo line search, $0 < \theta < 1$ step size decrease parameter, $\gamma > 0$ CG restart parameter.

$p_{-1} \leftarrow 0$

$i \leftarrow 0$

for $i = 0, 1, 2, \dots$ **until** convergence **do**

$\mathbf{X}_i \leftarrow \phi(x_i)$

$f_i \leftarrow f(\mathbf{X}_i)$

$g_i \leftarrow \nabla^R \hat{f}(x_i)$

\triangleright Riemannian gradient of $\hat{f}(x)$ at x_i

$s_i \leftarrow \mathcal{T}_{x_{i-1}, \alpha_{i-1} p_{i-1}} \alpha_{i-1} p_{i-1}$

\triangleright Vector transport the previous search direction

$y_i \leftarrow g_i - \mathcal{T}_{x_{i-1}, \alpha_{i-1} p_{i-1}} g_{i-1}$

$L_i \leftarrow y_i^T s_i / \|s_i\|^2$

\triangleright Lipschitz constant estimate

$p_i \leftarrow -g_i + \beta_i \mathcal{T}_{x_{i-1}, \alpha_{i-1} p_{i-1}} p_{i-1}$

if $\langle p_i, g_i \rangle > -\gamma$ **then**

$p_i = -g_i$

\triangleright Restart CG direction

end if

if $y_i^T s_i > 0$ **then**

$\alpha \leftarrow -g_i^T p_i / (L_i \|p_i\|_2^2)$

else

$\alpha \leftarrow \alpha_{i-1}$

end if

 Find $m \in \mathbb{Z}$ such that $\alpha_i = \alpha \theta^m$ and

$f(x_i + \alpha_i p_i) - f_i \leq \sigma \alpha_i g_i^T p_i$

$f(x_i + \alpha_i p_i) \leq \min\{f(x_i + \alpha_i \theta p_i), f(x_k + \alpha_i \theta^{-1} p_i)\}$

\triangleright Find a quasi-optimal minimizer

$x_{i+1} \leftarrow \mathcal{R}_{x_i}(\alpha_i p_i)$

\triangleright Reorthogonalize

$i \leftarrow i + 1$

end for

5.3.2. *Line search.* As any gradient based optimization scheme, we need a good initial step size and a computationally efficient line search. Following [48], we use a variation of the limited-minimization line search approach to set the initial step length based on the previous search direction and gradient that are vector transported to the current point—i.e, we have

$$s_i = \mathcal{T}_{x_{i-1}, \alpha_{i-1} p_{i-1}} \alpha_{i-1} p_{i-1}$$

$$y_i = g_i - \mathcal{T}_{x_{i-1}, \alpha_{i-1} p_{i-1}} g_{i-1}.$$

In this context, s_i is the manifold analogue for the Euclidean difference between iterates, $x_i - x_{i-1}$ and y_i is the manifold analogue for the difference of gradients between iterates, $g_i - g_{i-1}$, which are standard optimization quantities in optimization algorithms set in \mathbb{R}^n .

Our initial step size for the direction p_i is given as

$$\alpha_0 = -g_i^T p_i / (L_i \|p_i\|_2^2)$$

where $L_i = y_i^T s_i / \|s_i\|_2^2$ is the estimate of the Lipschitz constant for the gradient [49, Eq. 16].

In this context, computing the gradient is much more expensive than evaluating the objective. For this reason, we use a simple Armijo-type back-/forward-tracking approach that only involves function evaluations and seeks to minimize the 1D function $f(x + \alpha p_i)$ quasi-optimally, i.e., to find $m \in \mathbb{Z}$ such that $\alpha = \theta^m \alpha_0$ for $\sigma > 0$

$$(16) \quad \begin{aligned} f(x_i + \alpha p_i) - f(x_i) &\leq \sigma \alpha g_i^T p_i \\ f(x_i + \alpha p_i) &\leq \min\{f(x_i + \theta \alpha p_i), f(x_i + \theta^{-1} \alpha p_i)\} \end{aligned}$$

so $\alpha \approx \alpha^* = \operatorname{argmin}_{\alpha} f(x_i + \alpha p_i)$ in the sense that increasing or decreasing α by a factor of θ will increase $f(x_i + \alpha p_i)$. After the first few iterations of our optimization procedure, we observe empirically that our line search only involves two or three additional function evaluations to verify the second inequality in (16), i.e., our initial step length α_0 is quasi-optimal.

Because $\phi(\mathcal{R}_x(\alpha\eta)) = \phi(x + \alpha\eta)$ for any $x \in \mathcal{M}$ and horizontal vector η , where \mathcal{R}_x is the QR retraction, Armijo linesearches do not require reorthogonalization at the intermediate steps, which further reduces computational costs. The vector $x + \alpha\eta$ is only orthogonalized once an appropriate step length α is found.

5.3.3. Gauss-Newton Method. Because of the least-squares structure of our tensor completion problem (1), we can approximate the Hessian by the Gauss-Newton Hessian

$$H_{GN} := D\phi^*(x)D\phi(x) : \mathcal{H}_x \mathcal{M} \rightarrow \mathcal{H}_x \mathcal{M},$$

which arises from linearizing the objective function

$$\min_{x \in \mathcal{M}} \frac{1}{2} \|\phi(x) - b\|_2^2$$

around the current point x . Note that we do not use the “true” Gauss-Newton Hessian for the tensor completion objective (1), which is $D\phi^*(x)P_{\Omega}^*P_{\Omega}D\phi(x)$, since for even moderate subsampling ratios, $P_{\Omega}^*P_{\Omega}$ is close to the zero operator and this Hessian is very poorly conditioned as a result. Moreover, as we shall see, this formulation will allow us to simplify the application of H_{GN} and its inverse.

Since $D\phi(x) : \mathcal{H}_x \mathcal{M} \rightarrow \mathcal{T}_{\phi(x)}\mathcal{H}$ is an isomorphism, it is easy to see that H_{GN} is symmetric and positive definite on $\mathcal{H}_x \mathcal{M}$. The solution to the Gauss-Newton equation is then

$$H_{GN}\xi = -\nabla^R f(x)$$

for $\xi \in \mathcal{H}_x \mathcal{M}$. We can simplify the computation of H_{GN} by exploiting the recursive structure of $D\phi^*(x)$ and $D\phi(x)$, thereby avoiding intermediate vectors of size $\mathbb{R}^{n_1 \times n_2 \times \dots \times n_d}$ in the process. At the root node, by (13), we have that

$$\delta U'_{t_l} = P_{U_{t_l}}^{\perp} \langle U_{t_r}^T \times_2 D\phi(x)[\xi], \mathbf{B}_t \rangle_{(2,3),(2,3)}, \quad \delta U'_{t_r} = P_{U_{t_r}}^{\perp} \langle U_{t_l}^T \times_1 D\phi(x)[\xi], \mathbf{B}_t \rangle_{(1,3),(1,3)}, \quad \delta B'_t = U_{t_l}^T D\phi(x)[\xi] U_{t_r}$$

where

$$D\phi(x)[\xi] = \delta U_{t_l} \times_1 U_{t_r} \times_2 B_t + U_{t_l} \times_1 \delta U_{t_r} \times_2 B_t + U_{t_l} \times_1 U_{t_r} \times_2 \delta B_t, \quad t = t_{\text{root}}.$$

In the above expression, $D\phi(x)$ is horizontal, so that for each $t \in T \setminus t_{\text{root}}$, δU_t is perpendicular to U_t (12). A straightforward computation simplifies the above expression to

$$\delta U'_{t_l} = \delta U_{t_l} G_{t_l}, \quad \delta U'_{t_r} = \delta U_{t_r} G_{t_r}, \quad \delta B'_t = \delta B_{t_{\text{root}}}.$$

where $G_{t_l} = B_{t_{\text{root}}} B_{t_{\text{root}}}^T$ and $G_{t_r} = B_{t_{\text{root}}}^T B_{t_{\text{root}}}$. This expression gives us the components of the horizontal vector $\delta U'_{t_l}, \delta U'_{t_r}$ sent to the left and right children, respectively, as well as the horizontal vector $\delta B'_t$.

We proceed recursively by considering a node $t \in N(T) \cup \mathbf{t}_{\text{root}}$ and let $\delta U_t G_t$ be the contribution from the parent node of t . By applying the adjoint partial derivatives, followed by an orthogonal projection on to $\mathcal{H}_x \mathcal{M}$, we arrive at a simplified form for the Gauss-Newton Hessian

$$\begin{aligned} P_{U_{t_l}}^\perp \langle U_{t_r}^T \times_2 \delta U_t G_t, \mathbf{B}_t \rangle_{(2,3),(2,3)} &= \langle \delta U_{t_l} \times_1 G_t \times_3 \mathbf{B}_t, \mathbf{B}_t \rangle_{(2,3),(2,3)} \\ &:= \delta U_{t_l} G_{t_l}, \\ P_{U_{t_r}}^\perp \langle U_{t_l}^T \times_1 \delta U_t G_t, \mathbf{B}_t \rangle_{(1,3),(1,3)} &= \delta U_{t_r} G_{t_r}, \\ P_{B_t^{(1,2)}}^\perp (U_{t_l}^T \times_1 U_{t_r}^T \times_2 G_t \times_3 \delta U_t) &= \delta G_t \times_3 \mathbf{B}_t. \end{aligned}$$

In these expressions, the matrices G_t are the Gramian matrices associated to the HT format, initially introduced in [20] and used for truncation of a general tensor to the HT format as in [51]. They satisfy, for $x = (U_t, \mathbf{B}_t) \in \mathcal{M}$, $G_{\mathbf{t}_{\text{root}}} = 1$ and

$$(17) \quad G_{t_l} = \langle G_t \times_3 \mathbf{B}_t, \mathbf{B}_t \rangle_{(2,3),(2,3)}, \quad G_{t_r} = \langle G_t \times_3 \mathbf{B}_t, \mathbf{B}_t \rangle_{(1,3),(1,3)},$$

i.e., the same recursion as G_t in the above derivations. Each G_t is a $k_t \times k_t$ symmetric positive definite matrix (owing to the full rank constraints of the HT format) and also satisfies

$$(18) \quad \lambda_j(G_t) = \sigma_j(X^{(t)})^2$$

where $\lambda_j(A)$ is the j th eigenvalue of the matrix A and $\sigma_j(A)$ is the j th singular value of A .

Assuming that each G_t is well conditioned, applying the inverse of H_{GN} follows directly, summarized in Algorithm 5. For the case where our solution HT tensor exhibits *quickly*-decaying singular values of

Algorithm 5 $H_{GN}^{-1}\zeta$

Require: Current point $x = (U_t, \mathbf{B}_t)$, horizontal vector $\zeta = (\delta U_t, \delta \mathbf{B}_t)$.

Compute $(G_t)_{t \in T}$ using (17)

for $t \in T \setminus \mathbf{t}_{\text{root}}$ **do**

if $t \in L$ **then**

$\widetilde{\delta U}_t \leftarrow \delta U_t G_t^{-1}$

else

$\widetilde{\delta \mathbf{B}}_t \leftarrow G_t^{-1} \times_3 \delta \mathbf{B}_t$

end if

end for

return $H_{GN}^{-1}\zeta \leftarrow (\widetilde{\delta U}_t, \widetilde{\delta \mathbf{B}}_t)$

the matricizations, as is typically the assumption on the underlying tensor, the Gauss-Newton Hessian becomes poorly conditioned as the iterates converge to the solution, owing to (18). This can be remedied by introducing a small $\epsilon > 0$ and applying $(G_t + \epsilon I)^{-1}$ instead of G_t^{-1} in Algorithm 5 or by applying H_{GN}^{-1} by applying H_{GN} in a truncated PCG method. For efficiency purposes, we find the former option preferable. Alternatively, we can also avoid ill-conditioning via regularization, as we will see in the next section.

Remark 4. *We note that applying the inverse Gauss-Newton Hessian to a tangent vector is akin to ensuring that the projection on to the horizontal space is orthogonal, as in [53, 6.2.2]. Using this method, however, is much faster than the previously proposed method, because applying Algorithm 5 only involves matrix-matrix operations on the small parameters, as opposed to operations on much larger intermediate matrices that live in the spaces between the full tensor space and the parameter space.*

5.4. Regularization. In the tensor completion case, when there is little data available, interpolating on the HT manifold is susceptible to overfitting if one chooses the ranks $(k_t)_{t \in T}$ for the interpolated tensor too high. In that case, one can converge to solutions in $\text{null}(P_\Omega)$ that try leave the current manifold, associated to the ranks $(k_t)_{t \in T}$, to another nearby manifold corresponding to higher ranks. This can lead to degraded results in practice, as the actual ranks for the solution tensor are almost always unknown. One can use cross-validation techniques to estimate the proper internal ranks of the tensor, but we still need to ensure that the solution tensor has the predicted ranks for this approach to be successful – i.e., the iterates x must stay away from the boundary of \mathcal{H} .

To avoid our HT iterates converging to the manifold boundary, we introduce a regularization term on the singular values of the HT tensor $\phi(x) = \mathbf{X}$. To accomplish this, we exploit the hierarchical structure of \mathbf{X} and specifically the property of the Gramian matrices G_t in (18) to ensure that *all* matricizations of \mathbf{X} remain well-conditioned *without* having to perform SVDs on each matricization $X^{(t)}$. The latter approach would be prohibitively expensive when d or N are even moderately large.

Instead, we penalize the growth of the Frobenius norm of $X^{(t)}$ and $(X^{(t)})^\dagger$, which indirectly controls the largest and smallest singular values of $X^{(t)}$. We implement this regularization via the Gramian matrices in the following way. From (18), it follows that $\text{tr}(G_t) = \|G_t\|_* = \|X^{(t)}\|_F^2$ and likewise $\text{tr}(G_t^{-1}) = \|G_t^{-1}\|_* = \|(X^{(t)})^\dagger\|_F^2$. Our regularizer is then

$$H((\mathbf{B}_{t'})_{t' \in T}) = \sum_{t \in T} \text{tr}(G_t) + \text{tr}(G_t^{-1}).$$

A straightforward calculation shows that for $\mathcal{A} \in \mathcal{G}$

$$(G_t)_{t \in T, x} = (A_t^T G_t A_t)_{t \in T, \theta_{\mathcal{A}}(x)}$$

for A_t orthogonal. Therefore, our regularizer R is well-defined on the quotient manifold in the sense that it is θ -invariant on the parameter space \mathcal{M} . This is the same regularization term considered in [34], used for (theoretically) preventing the iterate from approaching the boundary of \mathcal{H} . In our case, we can leverage the structure of the Gramian matrices to implement this regularizer in a computationally feasible way.

Since in the definition of the Gramian matrices (17), G_t is computed recursively via tensor-tensor contractions (which are smooth operations), it follows that the mapping $g : (\mathbf{B}_t)_{t \in N(T)} \rightarrow (G_t)_{t \in T}$ is smooth. In order to compute its derivatives, we consider a node $t \in T \setminus t_{\text{root}}$ and consider the variations of its left and right children, i.e.,

$$(19) \quad \delta G_{t_r} = \frac{\partial G_{t_r}}{\partial \mathbf{B}_t} \delta \mathbf{B}_t + \frac{\partial G_{t_r}}{\partial G_t} \delta G_t, \quad \delta G_{t_l} = \frac{\partial G_{t_l}}{\partial \mathbf{B}_t} \delta \mathbf{B}_t + \frac{\partial G_{t_l}}{\partial G_t} \delta G_t.$$

We can take the adjoint of this recursive formulation, and thus obtain the gradient of g , if we compute the adjoint partial derivatives in (19) as well as taking the adjoint of the recursion itself. To visualize this process, we consider the relationship between input variables and output variables in the recursion as a series of small directed graphs, shown in Figure 2.

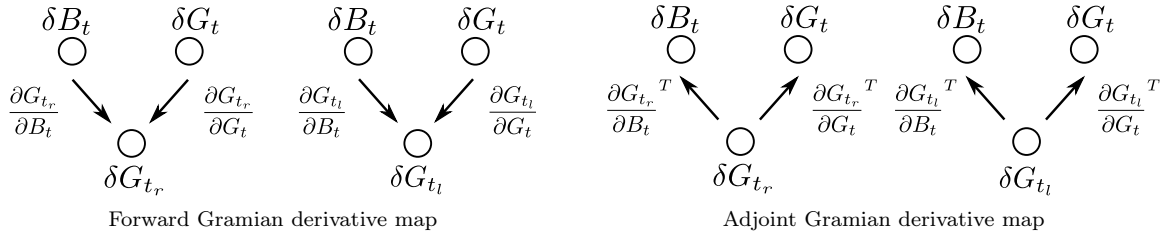


FIGURE 2. Adjoint Gramian derivative map

These graphs can be understood in the context of Algorithmic Differentiation, whereby the forward mode of this derivative map propagates variables up the tree and the adjoint mode propagates variables down the tree and adds (accumulates) the contributions of the relevant variables.

Since we only consider tangent vectors $\delta \mathbf{B}_t$ that are in the horizontal space at x , each extracted component is projected on to $(B_t^{(1,2)})^\perp$. We summarize our results in the following algorithms.

Applying Algorithm 7 to the gradient of $H(\mathbf{B}_t)$,

$$\nabla H(\mathbf{B}_t) = (V_t(I_{k_t} - S_t^{-2})V_t^T),$$

where $G_t = V_t S_t V_t^T$ is the eigenvalue decomposition of G_t , yields the Riemannian gradient of the regularizer. Note that here, we avoid having to compute SVDs of *any* matricizations of the full data $\phi(x)$, resulting in a method which is much faster than other tensor completion methods that require the SVDs on tensors

Algorithm 6 $Dg[\delta\mathbf{B}_t]$

Require: Current point $x = (U_t, \mathbf{B}_t)$, horizontal vector $dx = (\delta U_t, \delta \mathbf{B}_t)$.

Compute $(G_t)_{t \in T}$ using (17)

$\delta G_{t_{\text{root}}} \leftarrow 0$

for $t \in N(T)$, visiting parents before children **do**

$\delta G_{t_l} \leftarrow \langle \delta G_t \times_3 \mathbf{B}_t, \mathbf{B}_t \rangle_{(2,3),(2,3)} + 2\langle G_t \times_3 \delta \mathbf{B}_t, \mathbf{B}_t \rangle_{(2,3),(2,3)}$

$\delta G_{t_r} \leftarrow \langle \delta G_t \times_3 \mathbf{B}_t, \mathbf{B}_t \rangle_{(1,3),(1,3)} + 2\langle G_t \times_3 \delta \mathbf{B}_t, \mathbf{B}_t \rangle_{(1,3),(1,3)}$

end for

return $Dg[\delta\mathbf{B}_t] \leftarrow (\delta G_t)_{t \in T}$

Algorithm 7 $Dg^*[\delta G_t]$

Require: Current point $x = (U_t, \mathbf{B}_t)$, Gramian variations $(\delta G_t)_{t \in T}$, $\delta G_{t_{\text{root}}} = 0$.

Compute $(G_t)_{t \in T}$ using (17)

for $t \in T$ **do**

$\widetilde{\delta G}_t \leftarrow \delta G_t$

end for

for $t \in N(T)$, visiting children before parents **do**

$\delta \mathbf{B}_t \leftarrow (\widetilde{\delta G}_{t_l} + \widetilde{\delta G}_{t_l}^T) \times_1 G_t \times_3 \mathbf{B}_t + (\widetilde{\delta G}_{t_r} + \widetilde{\delta G}_{t_r}^T) \times_2 G_t \times_3 \mathbf{B}_t$

if $t \neq t_{\text{root}}$ **then**

$\delta \mathbf{B}_t \leftarrow (P_{B(1,2)}^\perp \delta \mathbf{B}_t^{(1,2)})_{(1,2)}$

$\widetilde{\delta G}_t \leftarrow \delta G_t + \langle \widetilde{G}_{t_l} \times_1 \mathbf{B}_t, \mathbf{B}_t \rangle_{(1,2),(1,2)} + \langle \widetilde{G}_{t_r} \times_2 \mathbf{B}_t, \mathbf{B}_t \rangle_{(1,2),(1,2)}$

end if

end for

return $Dg^*[\delta G_t] \leftarrow (\delta \mathbf{B}_t)_{t \in T}$

in $\mathbb{R}^{n_1 \times n_2 \times \dots \times n_d}$ [19]. Note that the cost of computing this regularizer $H(\mathbf{B}_t)$ and its gradient are almost negligible compared to the cost of computing the objective and its Riemannian gradient.

Finally, we should also note that the use of this regularizer is not designed to improve the recovery quality of problem instances with a relatively large amount of data and is useful primarily in the case where there is very little data so as to prevent overfitting, as we shall see in the numerical results section.

5.5. Convergence analysis. Our analysis here follows from similar considerations in [34, Sec. 3.6].

Theorem 5. *Let $\{x_i\}$ be an infinite sequence of iterates, with x_i generated at iteration i , generated from Algorithm 4 for the Gramian-regularized objective with $\lambda > 0$*

$$f(x) = \frac{1}{2} \|P_\Omega \phi(x) - b\|_2^2 + \lambda^2 \sum_{t \in T \setminus t_{\text{root}}} \text{tr}(G_t(x)) + \text{tr}(G_t^{-1}(x)).$$

Then $\lim_{i \rightarrow \infty} \|\nabla^R f(x_i)\| = 0$.

Proof. To show convergence, we merely need to show that the iterates remain in a sequentially compact set, since any accumulation point of $\{x_i\}$ is a critical point of f , by [1, Thm 4.3.1]. But this follows because by construction, since $f(x_i) \leq f(x_0) := C^2$ for all i . Letting $\mathbf{X}_i := \phi(x_i)$

$$\begin{aligned} & \frac{1}{2} \|P_\Omega \phi(x_i) - b\|_2^2 + \lambda^2 \sum_{t \in T \setminus t_{\text{root}}} \text{tr}(G_t(x_i)) + \text{tr}(G_t^{-1}(x_i)) = \\ & \frac{1}{2} \|P_\Omega \mathbf{X}_i - b\|_2^2 + \lambda^2 \sum_{t \in T \setminus t_{\text{root}}} \|X_i^{(t)}\|_F^2 + \|(X_i^{(t)})^\dagger\|_F^2 \leq C^2 \end{aligned}$$

This shows, in particular, that

$$\lambda^2 \sum_{t \in T \setminus t_{\text{root}}} \|X_i^{(t)}\|_F^2 \leq C^2 \quad \lambda^2 \sum_{t \in T \setminus t_{\text{root}}} \|(X_i^{(t)})^\dagger\|_F^2 \leq C^2$$

and therefore we have upper and lower bounds on the maximum and minimum singular values of $X_i^{(t)}$

$$\sigma_{\max}(X_i^{(t)}) \leq \|X_i^{(t)}\|_F \leq C/\lambda \quad \sigma_{\min}^{-1}(X_i^{(t)}) \leq \|(X_i^{(t)})^\dagger\|_F \leq C/\lambda$$

and therefore the iterates X_i stay within the compact set

$$\mathcal{C} = \{\mathbf{X} \in \mathcal{H} : \sigma_{\min}(X_k^{(t)}) \geq \lambda/C, \sigma_{\max}(X_k^{(t)}) \leq C/\lambda, t \in T \setminus \mathbf{t}_{\text{root}}\}.$$

One can show, as a modification of the proof in [53], that $\hat{\phi} : \mathcal{M}/\mathcal{G} \rightarrow \mathcal{H}$ is a homeomorphism on to its image, so that $\hat{\phi}^{-1}(C)$ is compact in \mathcal{M}/\mathcal{G} . We let $\|\cdot\|_T$ be the natural metric on \mathcal{M} , that is for $x = (U_t, \mathbf{B}_t)$, $y = (V_t, \mathbf{C}_t) \in \mathcal{M}$,

$$d(x, y) = \|x - y\|_T = \sum_{t \in L} \|U_t - V_t\|_F + \sum_{t \in N(T)} \|\mathbf{B}_t - \mathbf{C}_t\|_F.$$

A metric on \mathcal{M}/\mathcal{G} that generates the topology on the quotient space, is

$$(20) \quad \tilde{d}(\pi(x), \pi(y)) = \inf_{\mathcal{A}, \mathcal{B} \in \mathcal{G}} d(\theta_{\mathcal{A}}(x), \theta_{\mathcal{B}}(y)).$$

Note that this pseudo-metric is a metric which generates the topology on \mathcal{M}/\mathcal{G} by [7, Thm 2.1] since $\{\theta_{\mathcal{A}}\}_{\mathcal{A} \in \mathcal{G}}$ is a group of isometries acting on \mathcal{M} and the orbits of the action are closed by [53, Lemma 1]. Note that this metric is equivalent to

$$(21) \quad \tilde{d}(\pi(x), \pi(y)) = \inf_{\mathcal{A} \in \mathcal{G}} d(x, \theta_{\mathcal{A}}(y)),$$

which is well-defined and equal to (20) since $\|\theta_{\mathcal{A}}(x) - \theta_{\mathcal{B}}(y)\|_T = \|x - \theta_{\mathcal{A}^{-1}\mathcal{B}}(y)\|_T$ and \mathcal{A}, \mathcal{B} vary over \mathcal{G} .

Let $\{x_i\}$ be a sequence in $\pi^{-1}(\hat{\phi}^{-1}(C))$. By compactness of $\phi^{-1}(C)$ in \mathcal{M}/\mathcal{G} , we have that there is some set of HT parameters $y \in \phi^{-1}(C)$ such that, without loss of generality,

$$\tilde{d}(\pi(x_i), \pi(y)) \rightarrow 0 \text{ as } i \rightarrow \infty.$$

Then, by the characterization (21), since \mathcal{G} is compact, there exists a sequence $\mathcal{A}_i \subset \mathcal{G}$ such that

$$d(x_i, \theta_{\mathcal{A}_i}(y)) \rightarrow 0$$

Also by the compactness of \mathcal{G} , there exists a subsequence $\{\mathcal{A}_{i_j}\}$ that converges to $\mathcal{A} \in \mathcal{G}$, so $d(\theta_{\mathcal{A}_{i_j}}(y), \theta_{\mathcal{A}}(y))$ converges to 0 by continuity of the Lie group action θ . It then follows that

$$d(x_{i_j}, \theta_{\mathcal{A}}(y)) \leq d(x_{i_j}, \theta_{\mathcal{A}_{i_j}}(y)) + d(\theta_{\mathcal{A}_{i_j}}(y), \theta_{\mathcal{A}}(y)) \rightarrow 0 \quad \text{as } j \rightarrow \infty$$

And so any sequence $\{x_i\} \in \pi^{-1}(\hat{\phi}^{-1}(C))$ has a convergent subsequence and so $\pi^{-1}(\hat{\phi}^{-1}(C))$ is sequentially compact in \mathcal{M} . Therefore since the sequence x_k generated by Algorithm 4 remains in $\pi^{-1}(\hat{\phi}^{-1}(C))$ for all i , a subsequence of x_i converges to some $x \in \pi^{-1}(\hat{\phi}^{-1}(C))$, and so x is a critical point of f . \square

Although we have only shown convergence to a stationary point here, the authors in [47] consider the convergence of general line search methods that operate on the affine variety of rank at most k matrices. It would be an interesting extension of these results to the tensor case, but we leave this for future work. We also remark that this proof also holds for the Gauss-Newton method as well, since the results in [1, Thm 4.3.1] simply require that the search direction has a negative inner product with the gradient in the limit, which can be shown in this case.

6. NUMERICAL EXAMPLES

To address the challenges of large-scale tensor completion problems, as encountered in exploration seismology, we implemented the approach outlined in this paper in a highly optimized parallel Matlab toolbox entitled HTOpt (available at <http://www.math.ubc.ca/~curtd/software.html> for academic use). Contrary to the HT toolbox [35], whose primary function is performing operations on known HT tensors, our toolbox is designed to solve optimization problems in the HT format such as the seismic tensor completion problem. Our package includes the general optimization on HT manifolds detailed in Algorithm 1 as well as sparsity-exploiting objective & Riemannian gradient in Algorithm 2, implemented in Matlab. We also include a parallel implementation using the Parallel Matlab toolbox for both of these algorithms. All of the

following experiments were run on a single IBM x3550 workstation with 2 quad-core Intel 2.6Ghz processors with 16GB of RAM running Linux 2.6.18.

A simplified variation of the experiments below using the seismic data were presented previously in [12]. The experiments in our conference proceedings subsamples sources and use a Conjugate-Gradient method to solve the interpolation problem. In this paper, we use receiver subsampling and our subsequently developed Gauss-Newton and regularization methods to improve the recovery substantially while simultaneously greatly reducing the number of iterations required.

In this section, we compare our Gauss-Newton method with the interpolation scheme detailed in [34], denoted geomCG, for interpolating tensors with missing entries on the Tucker manifold. We have implemented a completely Matlab-based version of geomCG, which does not take advantage of the sparsity of the residual when computing the objective and Riemannian gradient, but uses Matlab’s internal calls to LAPACK libraries to compute matrix-matrix products and is much more efficient for this problem. To demonstrate this speedup, we compare our method to the reference mex implementation for geomCG from [34] on a randomly generated $200 \times 200 \times 200$ with multilinear rank 40 in each dimension, a training set comprised of 10% of the data points chosen randomly, and run for 20 iterations. The results of this comparison are shown in Table 1. We restrict our comparison to 3D tensors since the implementation of [34] available on the author’s website is strictly for 3D tensors.

Method	Time (s)	Training error	Test error	Difference from [34]
Original [34]	4701s	$1.091 \cdot 10^{-3}$	$2.792 \cdot 10^{-3}$	N/A
Our implementation	96.1s	$1.129 \cdot 10^{-3}$	$2.384 \cdot 10^{-3}$	$7.956 \cdot 10^{-4}$

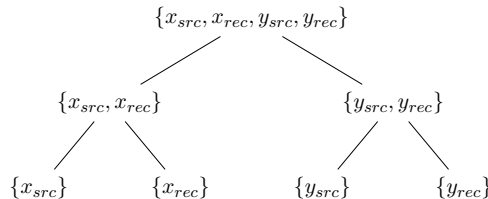
TABLE 1. Comparison between geomCG implementations on a $200 \times 200 \times 200$ random Gaussian tensor with multilinear rank 40 and subsampling factor = 0.1, both run for 20 iterations. Quantities are relative to the respective underlying solution tensor.

Since we take advantage of dense, multithreaded linear algebra routines, we find that our Matlab implementation is significantly faster than the mex code of [34] when $K \geq 20$ and $|\Omega|$ is a significant fraction of N^d , as is the case in the examples below.

In the examples considered below, we consider fixed multilinear ranks for each problem. There have been rank increasing strategies proposed for matrix completion in, among other places, [54, 39, 6] and as previously mentioned in [34], but we leave the question as to how to incorporate these heuristics in to the HT optimization framework for future research.

6.1. Seismic data. We briefly summarize the structure of seismic data in this section. Seismic data is typically collected via a boat equipped with an airgun and, for our purposes, a 2D array of receivers positioned on the ocean floor. The boat periodically fires a pressure wave in to the earth, which reflects off of subterranean discontinuities and produces a returning wave that is measured at the receiver array. The resulting data volume is five-dimensional, with two spatial source coordinates, denoted x_{src}, y_{src} , two receiver coordinates, denoted x_{rec}, y_{rec} , and time. For these experiments, we take a Fourier transform along the time axis and extract a single 4D volume by fixing a frequency and let \mathbf{D} denote the resulting *frequency slice* with dimensions $n_{src} \times n_{src} \times n_{rec} \times n_{rec}$.

From a practical point of view, the acquisition of seismic data from a physical system only allows us to subsample *receiver* coordinates, i.e., $\Omega = [n_{src}] \times [n_{src}] \times \mathcal{I}$ for some $\mathcal{I} \subset [n_{rec}] \times [n_{rec}]$ with $|\mathcal{I}| < n_{rec}^2$, rather than the standard tensor completion approach, which assumes that $\Omega \subset [n_{src}] \times [n_{src}] \times [n_{rec}] \times [n_{rec}]$ is random and unstructured. As a result, we use the dimension tree



for completing seismic data. With this choice, the fully sampled data \mathbf{D} has *quickly* decaying singular values in each matricization $D^{(t)}$ and is therefore represented well in the HT format. Additionally, the subsampled data $P_\Omega \mathbf{D}$ has increased singular values in all matricizations, and is poorly represented as a HT tensor with fixed ranks \mathbf{k} as a result. We examine this effect empirically in [12] and note that this data organization is used in [15] to promote low rank of the solution operators of the wave equation with respect to source and receiver coordinates. In a noise-free environment, seismic data volumes are modelled as the restriction of Green’s functions of the wave equation to the acquisition surface, which explains why this particular coordinate grouping decreases the rank of seismic data. Moreover, since we restrict our attention to relatively low frequency data, these volumes are relatively smooth, resulting in quickly decaying singular values [46].

Although this approach is limited to considerations of seismic data, for larger dimensions/different domains, potentially the method of [4] can choose an appropriate dimension tree automatically. In the next section, we also include the case when $\Omega \subset [n_{src}] \times [n_{src}] \times [n_{rec}] \times [n_{rec}]$, i.e. the “missing points” scenario, to demonstrate the added difficulty of the “missing receivers” case described above.

6.2. Single reflector data. For this data set, we generate data from a very simple seismic model consisting of two horizontal layers with a moderate difference in wavespeed and density between them. We generate this data with $n_{src} = n_{rec} = 50$ and extract a frequency slice at 4.21Hz, rescaled to have unit norm.

We consider the two sampling scenarios discussed in the previous section: we remove random points from the tensor, with results shown in Figure 3, and we remove random receivers from the tensor, with results shown in Figure 4. Here $\text{geomCG}(r_{\text{leaf}}) - w$ denote the Tucker interpolation algorithm with rank r_{leaf} in each mode and w rank continuation steps, i.e., the approach proposed in [34]. We also let $\text{HT}(r_{\text{leaf}}, r_{x_{src}x_{rec}})$ denote the HT interpolation method with rank r_{leaf} as in the Tucker interpolation and rank $r_{x_{src}x_{rec}}$ as the internal rank for the dimension tree. As is customary in the seismic literature, we measure recovery quality in terms of SNR, namely

$$\text{SNR}(\mathbf{X}, \mathbf{D}) = -20 \log_{10} \left(\frac{\|\mathbf{X}_{\Omega^c} - \mathbf{D}_{\Omega^c}\|}{\|\mathbf{D}_{\Omega^c}\|} \right) \text{ dB},$$

where \mathbf{X} is our interpolated signal, \mathbf{D} is our reference solution, and $\Omega^c = [n_{src}] \times [n_{src}] \times [n_{rec}] \times [n_{rec}] \setminus \Omega$. As we can see in Figure 3, the HT formulation is able to take advantage of low-rank separability of the seismic volume to produce a much higher quality solution than that of the Tucker tensor completion. The rank continuation scheme does not seem to be improving the recovery quality of the Tucker solution to the same degree as in [34], although it does seem to mitigate some of the overfitting errors for $\text{geomCG}(30)$. We display slices for fixed source coordinates and varying receiver coordinates in Figure 3 for randomly missing points and Figure 4 for randomly missing receivers. We summarize our results in Tables 2 and 3. By exploiting the low-rank structure of the HT format compared to the Tucker format, we are able to achieve much better results than Tucker tensor completion, especially for the realistic case of missing receiver samples.

In all instances for these experiments, the HT tensor completion outperforms the conventional Tucker approach both in terms of recovery quality and recovery speed. We note that geomCG does not scale as well computationally as our HT algorithm for $d > 3$, as the complexity analysis in [34] predicts. As such, we only consider the HT interpolation for the next sections, where we will solve the tensor completion problem for much larger data volumes.

6.3. Convergence speed. We demonstrate the superior convergence of the Gauss-Newton method compared to the Steepest Descent and Conjugate Gradient methods when used to complete the simple, single reflector data volume with 50% missing receivers and all ranks set to 20. We start each method with the same initial point and allow each to run for at most 1000 iterations. As we can see in Figure 5, the Gauss-Newton method converges much faster than the other two while simultaneously having per-iteration computational costs that are comparable to the other methods.

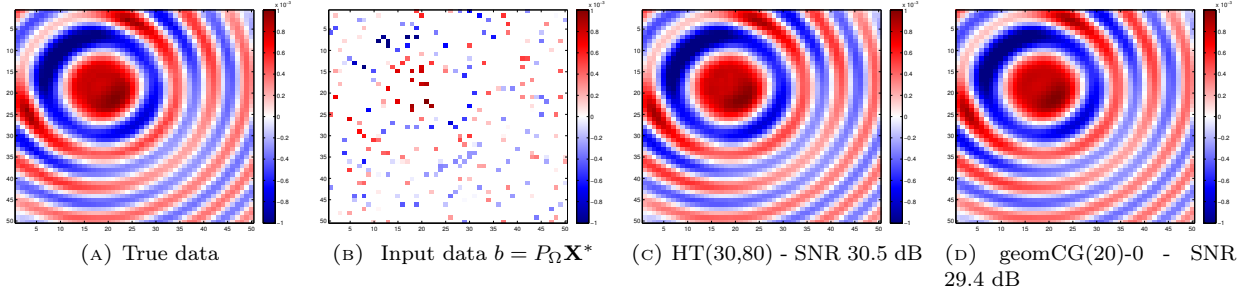


FIGURE 3. Reconstruction results for 90% missing points, best results for geomCG and HTOpt.

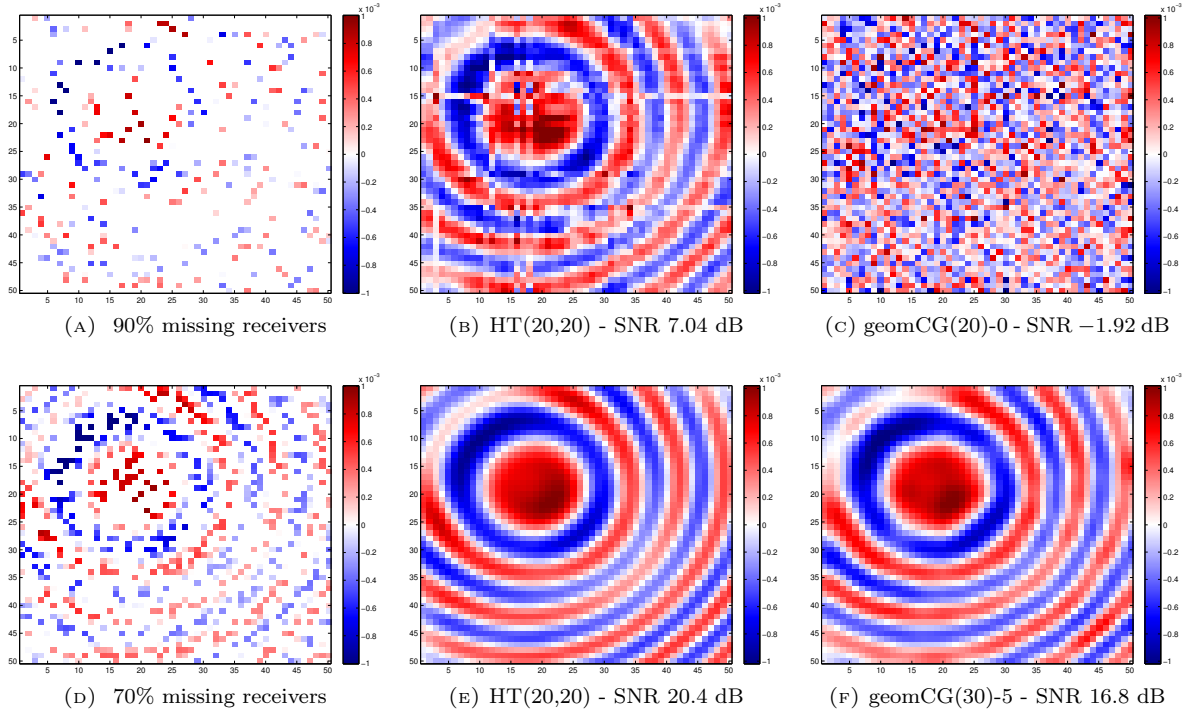


FIGURE 4. Reconstruction results for sampled receiver coordinates, best results for geomCG and HTOpt. (a-c) 90% missing receivers. (d-f): 70% missing receivers.

	Single reflector data - sampling percentage (missing points)					
	10%		30%		50%	
	SNR [dB]	time [s]	SNR [dB]	time [s]	SNR [dB]	time [s]
geomCG(20) - 0	28.5	1023	30.5	397	30.7	340
geomCG(30) - 0	-6.7	1848	21.8	3621	31.5	2321
geomCG(30) - 5	16.1	492	13.8	397	15.5	269
HTOpt(20,60)	30.1	83	30.4	59	30.4	57
HTOpt(20,80)	30.3	121	30.8	75	30.8	53
HTOpt(30,80)	31.6	196	32.9	133	33.1	114

TABLE 2. Reconstruction results for single reflector data - missing points - mean SNR over 5 random training sets

	Single reflector data - sampling percentage (missing receivers)					
	10%		30%		50%	
	SNR [dB]	time [s]	SNR [dB]	time [s]	SNR [dB]	time [s]
geomCG(20) - 0	-5.1	899	9.9	898	18.5	891
geomCG(30) - 0	-3.6	1796	-4.7	1834	6.1	1802
geomCG(30) - 5	-6.4	727	11.1	670	14.2	356
HTOpt(20,20)	6.1	111	19.8	101	20.1	66
HTOpt(30,20)	2.8	117	18.1	109	19.8	94
HTOpt(30,40)	0.0	130	13.4	126	21.6	108

TABLE 3. Reconstruction results for single reflector data - missing receivers - mean SNR over 5 random training test sets

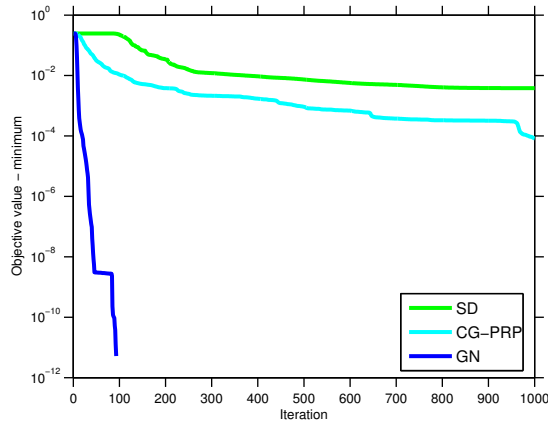


FIGURE 5. Convergence speed of various optimization methods

6.4. Performance. We investigate the empirical performance scaling of our approach as N, d, K , and $|\Omega|$ increase, as well as the number of processors for the parallel case, in Figure 6. Here we denote the use of Algorithm 1 as the “dense” case and Figure 2 as the “sparse” case. We run our optimization code in Steepest Descent mode with a single iteration for the line search, and average the running time over 10 iterations and 5 random problem instances. Our empirical performance results agree very closely with the theoretical complexity estimates, which are $O(N^d K)$ for the dense case and $O(|\Omega| d K^3)$ for the sparse case. Our parallel implementation for the sparse case scales very close to the theoretical time $O(1/\# \text{ processors})$.

6.5. Synthetic BG Compass data. This data set was provided to us by the BG Group company and consists of 5D data generated from an unknown synthetic model. Here $n_{src} = 68$ and $n_{rec} = 401$ and we extract frequency slices at 4.86 Hz, 7.34 Hz, and 12.3 Hz. On physical grounds, we expect a slower decay of the singular values at higher frequencies and thus the problem is much more difficult at 12.3 Hz compared to 4.86 Hz.

At these frequencies, the data has relatively low spatial frequency content in the receiver coordinates, and thus we subsample the receivers by a factor of 2 to $n_{rec} = 201$, for the purposes of speeding up the overall computation and ensuring that the intermediate vectors in the optimization are able to fit in memory. Our overall data volume has dimensions $\mathbf{D} \in \mathbb{R}^{68 \times 68 \times 201 \times 201}$.

We randomly remove varying amounts of receivers from this reduced data volume and interpolate using 50 iterations of the GN method discussed earlier. We display several recovered slices for fixed source coordinates and varying receiver coordinates (so-called *common source gathers* in seismic terminology) in Figure 9.

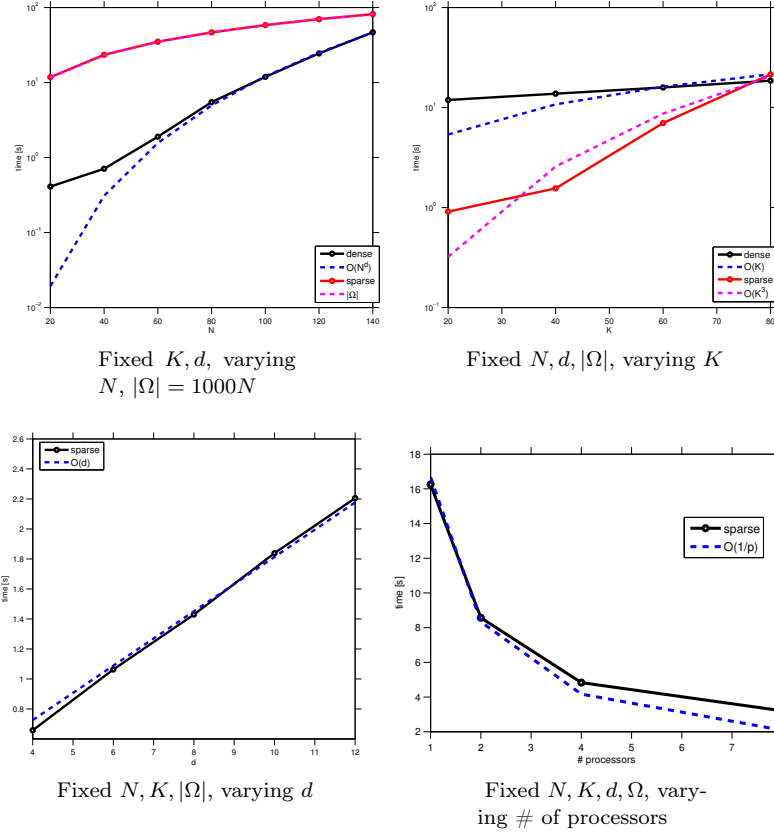


FIGURE 6. Dense & sparse objective, gradient performance.

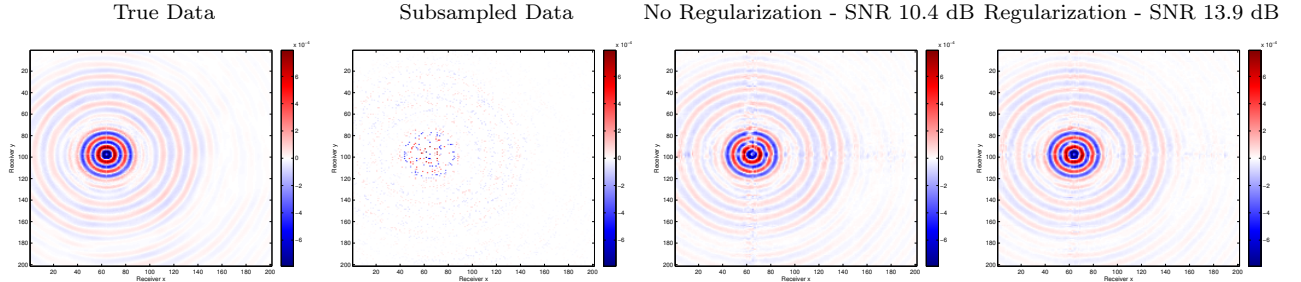


FIGURE 7. Regularization reduces some of the spurious artifacts and reduces overfitting in the case where there is very little data. 4.86 Hz data, 90% missing receivers.

We summarize our recovery results for tensor completion on these data sets from missing receivers in Table 4 and the various recovery parameters we use in Table 5. When the subsampling rate is extremely high (90% missing receivers in these examples), the recovery can suffer from overfitting issues, which leads to spurious artifacts in the recovered volume and lower SNRs overall. Using the Gramian-based regularization method discussed earlier, we can mitigate some of those artifacts and boost recovered SNRs, as seen in Figure 7.

6.6. Effect of varying regularizer strength. In this section, we examine the effect of varying the λ parameter based on the scenario described in Figure 7. That is to say, we use the same data volume \mathbf{D}

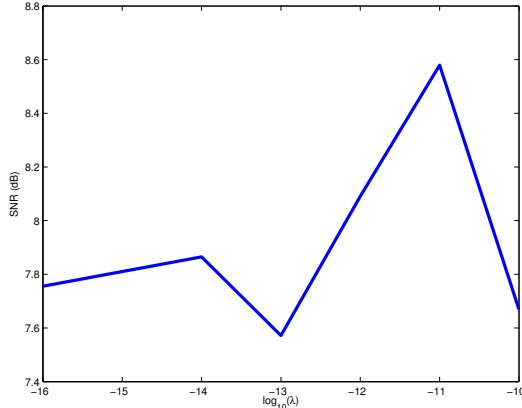


FIGURE 8. Recovery SNR versus $\log_{10}(\lambda)$ for 4.68Hz data., 90% missing receivers.

at as in the previous section with 90% of the receiver points removed and recover this volume using 50 iterations of the Gauss-Newton method for $\lambda = 0, 10^{-15}, \dots, 10^{-10}$. We plot the results of the test SNRs in Figure 8 and indicate 0 on the log-scale graph by -16 , since 10^{-16} is effectively at machine precision. For this problem with very little data, choosing an appropriate λ parameter can mitigate some of the overfitting errors and in our experiments, we found a parameter range from $10^{-12} - 10^{-10}$ improved our recovery SNRs, i.e., small but nonzero values. Higher values of λ than those shown here resulted in underfitting of the data and significantly worse SNRs and were omitted. In practice, one would employ a cross-validation scheme to choose an appropriate λ .

7. CONCLUSIONS AND DISCUSSION

In this work we have developed the algorithmic components to solve optimization problems on the manifold of fixed-rank Hierarchical Tucker tensors. By exploiting this manifold structure, we solve the tensor completion problem where the tensors of interest exhibit low-rank behavior. Our algorithm is computationally efficient because we mostly rely on operations on the small HT parameter space. The manifold optimization itself guarantees that we do not run into convergence issues, which arise when we ignore the quotient structure of the HT format. Our application of this framework to seismic examples confirms the validity of our new approach and outperforms existing Tucker-based approaches for large data volumes. To stabilize the recovery for high subsampling ratios, we introduced an additional regularization term that exploits properties of the Gramian matrices without the need to compute SVDs in the ambient space.

While the method clearly performs well on large-scale problems, there are still a number of theoretical questions regarding the performance of this approach. In particular, the generalization of matrix completion recovery guarantees to the HT format remains an open problem. As in many alternative approaches to matrix/tensor completion, the selection of the rank parameters and regularization parameters remain challenging both theoretically and from a practical point of view. However, the paper clearly illustrates that the HT format is a viable option to represent and complete high-dimensional data volumes in a computationally feasible manner.

8. ACKNOWLEDGEMENTS

We would like to thank the sponsors of the SINBAD consortium for their continued support. This work was supported by the CRD grant DNOISE II (CRDPJ 375142-08) from the National Sciences and Engineering Research Council of Canada (NSERC). We would also like to thank the BG Group company for providing us with the Compass data set.

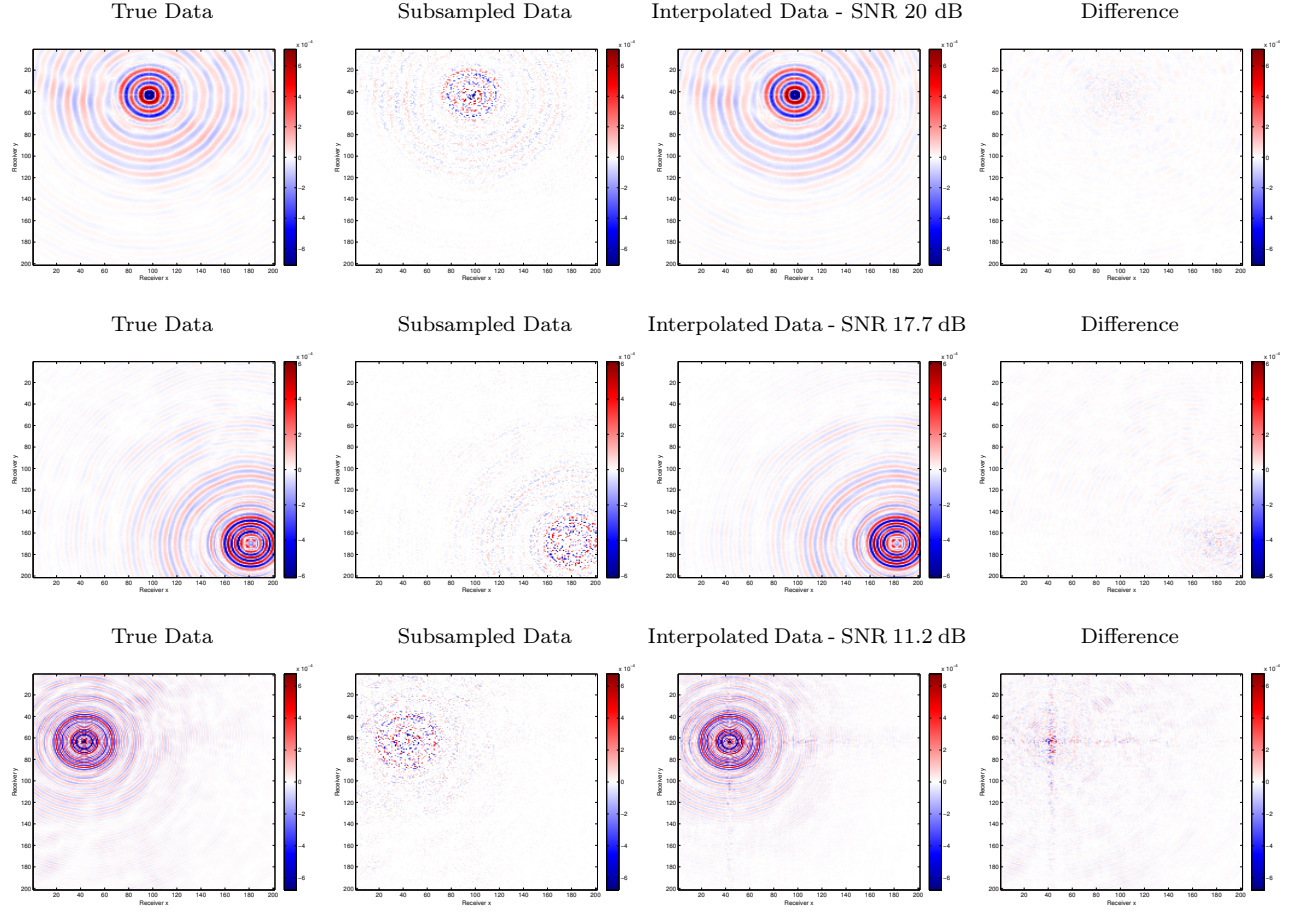


FIGURE 9. 75% missing receivers, fixed source coordinates. *Top*: 4.68 Hz, *Middle*: 7.34 Hz, *Bottom*: 12.3 Hz.

Frequency	% Missing	Train SNR (dB)	Test SNR (dB)	Runtime (s)
4.86 Hz	25%	21.2	21	4033
	50%	21.3	20.9	4169
	75%	21.5	19.9	4333
	90%	19.9	10.4	4679
	90%*	20.8*	13.0*	5043
7.34 Hz	25%	17.3	17.0	4875
	50%	17.4	16.9	4860
	75%	17.7	16.5	5422
	90%	16.6	9.82	4582
	90%*	16.6*	10.5*	4947
12.3 Hz	25%	14.9	14.2	5950
	50%	15.2	13.8	7083
	75%	15.8	9.9	7387
	90%	13.9	5.39	4578
	90%*	14*	6.5*	4966

TABLE 4. HT Recovery results - randomly missing receivers. Starred quantities are computed with regularization.

Frequency	$r_{x_{src}x_{rec}}$	$r_{x_{src}}$	$r_{x_{rec}}$	HT-SVD SNR (dB)
4.86 Hz	150	68	120	21.1
7.34 Hz	200	68	120	17.0
12.3 Hz	250	68	150	13.9

TABLE 5. HT parameters for each data set and the corresponding SNR of the HT-SVD approximation of each data set. The 12.3 Hz data is of much higher rank than the other two data sets and thus is much more difficult to recover.

APPENDIX A. PROOF OF PROPOSITION 3

We first explicitly define the orthogonal parameter space and its corresponding tangent space.

Below, let $W_{k_{t_l}, k_{t_r}, k_t}$ be the closed submanifold of $\mathbb{R}_*^{k_{t_l} \times k_{t_r} \times k_t}$, the set of 3-tensors with full multilinear rank, such that $\mathbf{W} \in W_{k_{t_l}, k_{t_r}, k_t}$ is orthonormal along modes 1 and 2, i.e., $(W^{(1,2)})^T W^{(1,2)} = I_{k_t}$ and let $\text{St}(n_t, k_t)$ be the $n_t \times k_t$ Stiefel manifold of $n \times k_t$ matrices with orthonormal columns. Our orthogonal parameter space \mathcal{M} is then

$$\mathcal{M} = \bigtimes_{t \in L} \text{St}(n_t, k_t) \times \bigtimes_{t \notin L \cup t_{\text{root}}} W_{k_{t_l}, k_{t_r}, k_t} \times \mathbb{R}_*^{k_{(t_{\text{root}})_l} \times k_{(t_{\text{root}})_r}}$$

with corresponding tangent space at $x = (U_t, \mathbf{B}_t) \in \mathcal{M}$

$$\mathcal{T}_x \mathcal{M} = \bigtimes_{t \in L} \mathcal{T}_{U_t} \text{St}(n_t, k_t) \times \bigtimes_{t \notin L \cup t_{\text{root}}} \mathcal{T}_{\mathbf{B}_t} W_{k_{t_l}, k_{t_r}, k_t} \times \mathbb{R}^{k_{(t_{\text{root}})_l} \times k_{(t_{\text{root}})_r}}.$$

Note that $\mathcal{T}_Y \text{St}(n, p) = \{Y\Omega + Y^\perp K : \Omega^T = -\Omega \in \mathbb{R}^{p \times p}, K \in \mathbb{R}^{(n-p) \times p}\}$. We omit an explicit description of $\mathcal{T}_{B_t} W_{k_{t_l}, k_{t_r}, k_t}$ for brevity.

Proof. It is easy to see that the first point in Definition 7 is satisfied, since for $X \in \text{St}(n, p)$, $\text{qf}(X) = X$

Let $x = (U_t, \mathbf{B}_t) \in \mathcal{M}$ and $\eta = (\delta U_t, \delta \mathbf{B}_t) \in \mathcal{T}_x \mathcal{M}$. To avoid notational overload, we use the slight abuse of notation that $B_t := B_t^{(1,2)}$ for $t \neq t_{\text{root}}$.

Let $s \in [0, t] \mapsto x(s)$ be a curve in the parameter space \mathcal{M} with $x(0) = x$ and $x'(0) = \eta$ and $x(s) = (U_t(s), B_t(s))$ and $x'(s) = (\delta U_t(s), \delta B_t(s))$.

Then we have that, in Kronecker form,

$$D\mathcal{R}_x(0_x)[\eta] = \begin{cases} \frac{d}{ds} \text{qf}(x(s)_t) \big|_{s=0} & \text{if } t \in L \\ \frac{d}{ds} \text{qf}((R_{t_r}(s) \otimes R_{t_l}(s))(x(s)_t)) \big|_{s=0} & \text{if } t \notin t_{\text{root}} \cup L \\ \frac{d}{ds} (R_{t_r}(s) \otimes R_{t_l}(s))(x(s)_t) \big|_{s=0} & \text{if } t = t_{\text{root}} \end{cases}$$

The fact that $D\mathcal{R}_x(0_x)[\eta]_t = \delta U_t$ for $t \in L$ follows from Example 8.1.5 in [1].

To compute $D\mathcal{R}_x(0_x)[\eta]_t$ for $t \notin L \cup t_{\text{root}}$, we first note the formula from [1]

$$(22) \quad D \text{qf}(Y)[U] = \text{qf}(Y) \rho_{\text{skew}}(\text{qf}(Y)^T U (\text{qf}(Y)^T Y)^{-1}) + (I - \text{qf}(Y) \text{qf}(Y)^T) U (\text{qf}(Y)^T Y)^{-1}$$

where $Y \in \mathbb{R}_*^{n \times k}$, $U \in T_Y \mathbb{R}_*^{n \times k} \simeq \mathbb{R}^{n \times k}$ and $\text{qf}(Y)$ is the Q-factor of the QR-decomposition of Y .

Therefore, if we set $Z(s) = (R_{t_r}(s) \otimes R_{t_l}(s))(x(s)_t)$, where $R_t(s)$ is the R -factor of the QR-decomposition of the matrix associated to node t , we have

$$Z'(0) = [(R'_{t_r}(0) \otimes I_{k_l}) + (I_{k_r} \otimes R'_{t_l}(0))] B_t + \delta B_t$$

As a result of the discussion in Example 8.1.5 in [1], since $R_t(0) = I_{k_t}$ we have that

$$R'_t(0) = \begin{cases} \rho_{UT}(U_t^T \delta U_t) & \text{for } t \in L \\ \rho_{UT}(B_t^T \delta B_t) & \text{for } t \notin L \cup t_{\text{root}} \end{cases}$$

where $\rho_{UT}(A)$ is the projection onto the upper triangular term of the unique decomposition of a matrix into the sum of a skew-symmetric term and an upper triangular term.

Since $U_t \in \text{St}(n_t, k_t)$ and $B_t \in \text{St}(k_{t_l} k_{t_r}, k_t)$, in light of the fact that for $X \in \text{St}(n, k)$,

$$T_X \text{St}(n, k) = \{X\Omega + X^\perp K : \Omega = -\Omega^T\},$$

then $X^T \delta X$ is skew symmetric, for any tangent vector δX , which implies that $\rho_{UT}(X^T \delta X)$ is zero.

It follows that $R'_t(0) = 0$ for all $t \in T \setminus \mathfrak{t}_{\text{root}}$, and therefore

$$Z'(0) = \delta B_t$$

from which we immediately obtain

$$D\mathcal{R}_x(0_x)[\eta]_t = \delta B_t \quad \text{for } t \notin L \cup \mathfrak{t}_{\text{root}}$$

A similar approach holds when $t = \mathfrak{t}_{\text{root}}$, and therefore, $\mathcal{R}_x(\eta)$ is a retraction on \mathcal{M} . \square

REFERENCES

- [1] P.A. Absil, R. Mahony, and R. Sepulchre. *Optimization algorithms on matrix manifolds*. Princeton Univ Press, 2008.
- [2] Evrim Acar, Daniel M Dunlavy, and Tamara G Kolda. A scalable optimization approach for fitting canonical tensor decompositions. *Journal of Chemometrics*, 25(2):67–86, 2011.
- [3] B. W. Bader, T. G. Kolda, et al. Matlab tensor toolbox version 2.5. <http://www.sandia.gov/~tgkolda/TensorToolbox/>, January 2012.
- [4] Jonas Ballani and Lars Grasedyck. Tree adaptive approximation in the hierarchical tensor format. *Preprint*, 141, 2013.
- [5] Jonas Ballani, Lars Grasedyck, and Melanie Kluge. Black box approximation of tensors in hierarchical tucker format. *Linear Algebra and its Applications*, 438(2):639 – 657, 2013. Tensors and Multilinear Algebra.
- [6] J.D. Blanchard, J. Tanner, and Ke Wei. Conjugate gradient iterative hard thresholding: Observed noise stability for compressed sensing. *Signal Processing, IEEE Transactions on*, 63(2):528–537, Jan 2015.
- [7] Francesca Cagliari, Barbara Di Fabio, and Claudia Landi. The natural pseudo-distance as a quotient pseudo-metric, and applications. *AMS Acta, Universita di Bologna*, 3499, 2012.
- [8] Jian-Feng Cai, Emmanuel J Candès, and Zuowei Shen. A singular value thresholding algorithm for matrix completion. *SIAM Journal on Optimization*, 20(4):1956–1982, 2010.
- [9] E.J. Candès and B. Recht. Exact matrix completion via convex optimization. *Foundations of Computational mathematics*, 9(6):717–772, 2009.
- [10] Emmanuel J Candès and Terence Tao. The power of convex relaxation: Near-optimal matrix completion. *Information Theory, IEEE Transactions on*, 56(5):2053–2080, 2010.
- [11] J. Carroll and J. Chang. Analysis of individual differences in multidimensional scaling via an n-way generalization of “eckart-young” decomposition. *Psychometrika*, 35:283–319, 1970. 10.1007/BF02310791.
- [12] Curt Da Silva and Felix J. Herrmann. Hierarchical tucker tensor optimization - applications to tensor completion. In *10th international conference on Sampling Theory and Applications (SampTA 2013)*, pages 384–387, Bremen, Germany, July 2013.
- [13] L. De Lathauwer, B. De Moor, and J. Vandewalle. A multilinear singular value decomposition. *SIAM Journal on Matrix Analysis and Applications*, 21(4):1253–1278, 2000.
- [14] V. De Silva and L.H. Lim. Tensor rank and the ill-posedness of the best low-rank approximation problem. *SIAM Journal on Matrix Analysis and Applications*, 30(3):1084–1127, 2008.
- [15] Laurent Demanet. *Curvelets, Wave Atoms, and Wave Equations*. PhD thesis, California Institute of Technology, 2006.
- [16] M. P. Friedlander E. van den Berg. Spot – a linear-operator toolbox. <http://www.cs.ubc.ca/labs/scl/spot/>.
- [17] Lars Eldén and Berkant Savas. A newton-grassmann method for computing the best multilinear rank- (r_1, r_2, r_3) approximation of a tensor. *SIAM Journal on Matrix Analysis and applications*, 31(2):248–271, 2009.
- [18] Antonio Falcó, Wolfgang Hackbusch, Anthony Nouy, et al. Geometric structures in tensor representations (release 2). 2014.
- [19] S. Gandy, B. Recht, and I. Yamada. Tensor completion and low-n-rank tensor recovery via convex optimization. *Inverse Problems*, 27(2):025010, 2011.
- [20] L. Grasedyck. Hierarchical singular value decomposition of tensors. *SIAM Journal on Matrix Analysis and Applications*, 31(4):2029–2054, 2010.

- [21] Lars Grasedyck, Melanie Kluge, and Sebastian Krämer. Alternating directions fitting (adf) of hierarchical low rank tensors. *Preprint*, 149, 2013.
- [22] Lars Grasedyck, Daniel Kressner, and Christine Tobler. A literature survey of low-rank tensor approximation techniques. *GAMM-Mitteilungen*, 36(1):53–78, 2013.
- [23] W. Hackbusch and S. Kühn. A new scheme for the tensor representation. *Journal of Fourier Analysis and Applications*, 15(5):706–722, 2009.
- [24] Wolfgang Hackbusch. *Tensor spaces and numerical tensor calculus*, volume 42. Springer, 2012.
- [25] Jutho Haegeman, Tobias J Osborne, and Frank Verstraete. Post-matrix product state methods: To tangent space and beyond. *Physical Review B*, 88(7):075133, 2013.
- [26] R.A. Harshman. Foundations of the parafac procedure: models and conditions for an "explanatory" multimodal factor analysis. *UCLA Working Papers in Phonetics*, 16:1–84, 1970.
- [27] Sebastian Holtz, Thorsten Rohwedder, and Reinhold Schneider. The alternating linear scheme for tensor optimization in the tensor train format. *SIAM Journal on Scientific Computing*, 34(2):A683–A713, 2012.
- [28] Sebastian Holtz, Thorsten Rohwedder, and Reinhold Schneider. On manifolds of tensors of fixed tt-rank. *Numerische Mathematik*, 120(4):701–731, 2012.
- [29] Bo Huang, Cun Mu, Donald Goldfarb, and John Wright. Provable low-rank tensor recovery. 2014.
- [30] Boris N. Khoromskij. Tensors-structured numerical methods in scientific computing: Survey on recent advances. *Chemometrics and Intelligent Laboratory Systems*, 110(1):1 – 19, 2012.
- [31] Tamara G Kolda and Brett W Bader. Tensor decompositions and applications. *SIAM review*, 51(3):455–500, 2009.
- [32] N. Kreimer and M.D. Sacchi. Tensor completion via nuclear norm minimization for 5d seismic data reconstruction. In *SEG Technical Program Expanded Abstracts 2012*, pages 1–5. Society of Exploration Geophysicists, 2012.
- [33] Nadia Kreimer and Mauricio D Sacchi. A tensor higher-order singular value decomposition for prestack seismic data noise reduction and interpolation. *Geophysics*, 77(3):V113–V122, 2012.
- [34] D. Kressner, M. Steinlechner, and B. Vandereycken. Low-rank tensor completion by riemannian optimization. Technical Report 2, 2014.
- [35] Daniel Kressner and Christine Tobler. Algorithm 941: Htucker—a matlab toolbox for tensors in hierarchical tucker format. *ACM Trans. Math. Softw.*, 40(3):22:1–22:22, April 2014.
- [36] Christian Lubich. *From quantum to classical molecular dynamics: reduced models and numerical analysis*. European Mathematical Society, 2008.
- [37] Christian Lubich, Thorsten Rohwedder, Reinhold Schneider, and Bart Vandereycken. Dynamical Approximation by Hierarchical Tucker and Tensor-Train Tensors. *SIAM Journal on Matrix Analysis and Applications*, 34(2):470–494, April 2013.
- [38] B. Mishra and R. Sepulchre. R3mc: A riemannian three-factor algorithm for low-rank matrix completion. In *53rd IEEE Conference on Decision and Control*, 2014.
- [39] Bamdev Mishra, Gilles Meyer, Francis Bach, and Rodolphe Sepulchre. Low-rank optimization with trace norm penalty. *SIAM Journal on Optimization*, 23(4):2124–2149, 2013.
- [40] Cun Mu, Bo Huang, John Wright, and Donald Goldfarb. Square deal: Lower bounds and improved relaxations for tensor recovery. *arXiv preprint arXiv:1307.5870*, 2013.
- [41] Ivan V Oseledets. Tensor-train decomposition. *SIAM Journal on Scientific Computing*, 33(5):2295–2317, 2011.
- [42] Ivan V Oseledets and SV Dolgov. Solution of linear systems and matrix inversion in the tt-format. *SIAM Journal on Scientific Computing*, 34(5):A2718–A2739, 2012.
- [43] Ivan V Oseledets and Eugene E Tyrtyshnikov. Breaking the curse of dimensionality, or how to use svd in many dimensions. *SIAM Journal on Scientific Computing*, 31(5):3744–3759, 2009.
- [44] Samet Oymak, Amin Jalali, Maryam Fazel, Yonina C Eldar, and Babak Hassibi. Simultaneously structured models with application to sparse and low-rank matrices. *arXiv preprint arXiv:1212.3753*, 2012.
- [45] Holger Rauhut, Reinhold Schneider, and Zeljka Stojanac. Tensor completion in hierarchical tensor representations. *arXiv preprint arXiv:1404.3905*, 2014.
- [46] Reinhold Schneider and André Uschmajew. Approximation rates for the hierarchical tensor format in periodic sobolev spaces. *Journal of Complexity*, 30(2):56–71, 2014.

- [47] Reinhold Schneider and André Uschmajew. Convergence results for projected line-search methods on varieties of low-rank matrices via $\{\backslash L\}$ ojasiewicz inequality. *arXiv preprint arXiv:1402.5284*, 2014.
- [48] Z J Shi and J Shen. New Inexact Line Search Method for Unconstrained Optimization. *Journal of Optimization Theory and Applications*, 127(2):425–446, November 2005.
- [49] Zhen-Jun Shi. Convergence of line search methods for unconstrained optimization. *Applied Mathematics and Computation*, 157(2):393–405, 2004.
- [50] M. Signoretto, R. Van de Plas, B. De Moor, and J. AK Suykens. Tensor versus matrix completion: a comparison with application to spectral data. *Signal Processing Letters, IEEE*, 18(7):403–406, 2011.
- [51] C. Tobler. *Low Rank Tensor Methods for Linear Systems and Eigenvalue Problems*. PhD thesis, ETH Zürich, 2012.
- [52] A. Uschmajew. Local convergence of the alternating least squares algorithm for canonical tensor approximation. *SIAM Journal on Matrix Analysis and Applications*, 33(2):639–652, 2012.
- [53] A. Uschmajew and B. Vandereycken. The geometry of algorithms using hierarchical tensors. *Linear Algebra and its Applications*, 439(1):133–166, July 2013.
- [54] A. Uschmajew and B. Vandereycken. Line-search methods and rank increase on low-rank matrix varieties. In *Proceedings of the 2014 International Symposium on Nonlinear Theory and its Applications*, 2014.
- [55] Yangyang Xu and Wotao Yin. A block coordinate descent method for regularized multiconvex optimization with applications to nonnegative tensor factorization and completion. *SIAM Journal on Imaging Sciences*, 6(3):1758–1789, 2013.