# A Riemannian trust-region method for low-rank tensor completion

Gennadij Heidel[*] and Volker Schulz

*Fachbereich IV - Mathematik, Universität Trier, 54286 Trier, Germany*

SUMMARY

The goal of tensor completion is to fill in missing entries of a partially known tensor (possibly including some noise) under a low-rank constraint. This may be formulated as a least-squares problem. The set of tensors of a given multilinear rank is known to admit a Riemannian manifold structure, thus methods of Riemannian optimization are applicable.
In our work, we derive the Riemannian Hessian of an objective function on the low-rank tensor manifolds using the Weingarten map, a concept from differential geometry. We discuss the convergence properties of Riemannian trust-region methods based on the exact Hessian and standard approximations, both theoretically and numerically. We compare our approach to Riemannian tensor completion methods from recent literature, both in terms of convergence behaviour and computational complexity. Our examples include the completion of randomly generated data with and without noise and recovery of multilinear data from survey statistics. Copyright © 2017 John Wiley & Sons, Ltd.

## 1. INTRODUCTION

In this paper, we discuss optimization techniques on the manifold of tensors of a given rank. We consider least-squares problems of the form

$$\min_{\mathbf{X}} f(\mathbf{X}) = \frac{1}{2} \left\| \mathrm{P}_{\Omega}\, \mathbf{X} - \mathrm{P}_{\Omega}\, \mathbf{A} \right\|^2 \tag{1.1}$$
$$\text{s.\,t. } \mathbf{X} \in \mathcal{M}_{\mathbf{r}} \coloneqq \left\{ \mathbf{X} \in \mathbb{R}^{n_1 \times \cdots \times n_d} \mid \mathrm{rank}(\mathbf{X}) = \mathbf{r} \right\},$$

where $\mathrm{rank}(\mathbf{X}) \in \mathbb{R}^d$ denotes the multilinear rank of a tensor $\mathbf{X}$, and $\mathrm{P}_{\Omega} : \mathbb{R}^{n_1 \times \cdots \times n_d} \to \mathbb{R}^{n_1 \times \cdots \times n_d}$ is a linear operator. A typical choice found in the literature is

$$[\mathrm{P}_{\Omega}\, \mathbf{X}]_{i_1 \dots i_d} \coloneqq \begin{cases} x_{i_1 \dots i_d} & \text{if } (i_1, \dots, i_d) \in \Omega, \\ 0 & \text{otherwise,} \end{cases}$$

where $\Omega \subset \{1, \dots, n_1\} \times \cdots \times \{1, \dots, n_d\}$ denotes the sampling set, i.e. we assume that $\mathbf{A} \in \mathbb{R}^{n_1 \times \cdots \times n_d}$ is a tensor whose entries with indices in $\Omega$ are known.

The tensor completion problem is a generalization of the matrix completion problem, see the page by Ma et al. [1] for an overview of methods and applications in the context of convex optimization. Early work on tensor completion has been done by Liu et al. [2], who consider the problem

$$\min_{\mathbf{X}} \|\mathbf{X}\|_* \quad \text{s.\,t. } \mathrm{P}_{\Omega}\, \mathbf{X} = \mathrm{P}_{\Omega}\, \mathbf{A} \tag{1.2}$$

---

[*]Correspondence to: heidel@uni-trier.de

in the context of image data recovery, where $\|\cdot\|_*$ is a generalized nuclear norm. Note that (1.2) can be viewed as the dual of (1.1). It ensures convexity for the tensor completion problem at the cost of losing the underlying manifold structure of low-rank tensors. Specifically, it does not give a low-rank solution in the presence of noise, i. e. if $\mathbf{A} \notin \mathcal{M}_{\mathbf{r}}$; in this case, an additional routine may be needed to truncate the result to low rank. Signoretto et al. [3] and Gandy et al. [4] choose a Tikhonov-like approach by minimizing the a penalized unconstrained function

$$\min_{\mathbf{X}} \frac{1}{2} \big\| \operatorname{P}_\Omega \mathbf{X} - \operatorname{P}_\Omega \mathbf{A} \big\|^2 + \frac{\mu}{2} \|\mathbf{X}\|_*.$$

A Riemannian CG method for (1.1) has been proposed by Kressner et al. [5], which is an extension of Vandereycken's earlier work [6] for the matrix completion problem. The authors show rapid linear convergence of their method with satisfactory reconstruction of missing data for a range of applications. Other Riemannian approaches for matrix completion include the work by Ngo/Saad [7] and Mishra et al. [8], who use a product Graßmann quotient manifold structure.

In recent research, second-order methods in Riemannian optimization have generated considerable interest in order to find superlinearly converging methods, see the overview by Absil et al. [9, Chapters 6–8] and the references therein. Boumal/Absil [10] apply these techniques to matrix completion in the Graßmannian framework. Vandereycken [6, Subsection 2.3] derives the Hessian for Riemannian matrix completion with an explicit expression of the singular values. In the higher-order tensor case, Eldén/Savas [11] propose a Newton method for computing a rank-$\mathbf{r}$ tensor approximation, using a Graßmannian approach. Ishteva et al. [12] extend these ideas to construct a Riemannian trust-region scheme.

In this paper, we propose a Riemannian trust-region scheme for (1.1) using explicit Tucker decompositions and compare it to a state-of-the-art Riemannian CG as used in [5]. We derive the exact expression of the Riemannian Hessian on $\mathcal{M}_{\mathbf{r}}$ for this manifold geometry by using the Weingarten map proposed by Absil et al. [13]. Our work focuses on the application case of tensor completion and contains tensor approximation as the special case of full sampling, i. e. $|\Omega| = \prod_i n_i$.

The rest of the paper is organized as follows. In Section 2, we cite some basic results about tensor arithmetic and the manifold of low-rank tensors. In Section 3, we present a brief overview of Riemannian optimization and prove our main result, the Riemannian Hessian on $\mathcal{M}_{\mathbf{r}}$. In Section 4, we explain the Riemannian trust-region methods based on exact and approximate Hessian evaluations. In Section 5, we present the some numerical experiments for our method on synthetic data and a standard test data set from multilinear statistics.

## 2. LOW-RANK TENSORS

In this section, we collect some basic concepts and results on the Tucker decomposition and multilinear rank of tensors needed for our work. First, we define notations and results of general tensor arithmetic, as laid out in the survey paper [14]. Then, we introduce the manifold geometry of $\mathcal{M}_{\mathbf{r}}$, see [15, 16, 5].

### 2.1. Multilinear rank and Tucker decomposition

For a tensor $\mathbf{A} \in \mathbb{R}^{n_1 \times \cdots \times n_d}$, the matrix

$$A_{(i)} \in \mathbb{R}^{n_i \times \prod_{j \neq i} n_j},$$

such that the row index of $A_{(i)}$ is the $i$th modes of $\mathbf{A}$ and the column index is a multi-index of the remaining $d-1$ modes, in lexicographic order, is called the *mode-$i$ matricization* of $\mathbf{A}$. It may be viewed as a $d$-order generalization of the matrix transpose, since, for $d=2$, it holds that $A_{(1)} = A$ and $A_{(2)} = A^{\mathrm{T}}$. We denote the re-tensorization of a matricized tensor by a superscript index, i. e. $(A_{(i)})^{(i)} = \mathbf{A}$.

The *multilinear rank* of a tensor $\mathbf{A}$ is the $d$-tuple

$$\operatorname{rank}(\mathbf{A}) = \big( \operatorname{rank}(A_{(1)}), \ldots, \operatorname{rank}(A_{(d)}) \big),$$

with $\mathrm{rank}(\cdot)$ on the right-hand side of the equation denoting the matrix rank. In contrast to the matrix case, the ranks of different matricizations of a tensor may be different, e. g. consider $\mathbf{A} \in \mathbb{R}^{2\times2\times2}$, given by its mode-1 matricization

$$A_{(1)} = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \end{bmatrix}.$$

Then, the other matricizations are

$$A_{(2)} = A_{(1)}, \; A_{(3)} = \begin{bmatrix} 1 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 \end{bmatrix},$$

so clearly $\mathrm{rank}(\mathbf{X}) = (2,2,1)$.

The *i-mode product* of $\mathbf{A}$ with a matrix $M \in \mathbb{R}^{m \times n_i}$ is defined as

$$\mathbf{B} = \mathbf{A} \times_i M \iff B_{(i)} = M A_{(i)}, \; \mathbf{B} \in \mathbb{R}^{n_1 \times \cdots \times n_{i-1} \times m \times n_{i+1} \cdots \times n_d}.$$

It is worth noting that, for different modes, the order of multiplications is irrelevant, i. e.

$$\mathbf{A} \times_i M \times_j N = \mathbf{A} \times_j N \times_i M \quad \text{if } i \neq j. \tag{2.1}$$

If the modes are equal, then
$$\mathbf{A} \times_i M \times_i N = \mathbf{A} \times_i (NM). \tag{2.2}$$

A Frobenius inner product on $\mathbb{R}^{n_1 \times \cdots \times n_d}$ is given by

$$\langle \mathbf{A}, \mathbf{B} \rangle := \mathrm{tr}\left( A_{(1)}^{\mathrm{T}} B_{(1)} \right) = \cdots = \mathrm{tr}\left( A_{(d)}^{\mathrm{T}} B_{(d)} \right) = \sum_{i_1=1}^{n_1} \cdots \sum_{i_d=1}^{n_d} a_{i_1 \ldots i_d} b_{i_1 \ldots i_d},$$

with the induced norm $\|\mathbf{A}\| := \sqrt{\langle \mathbf{A}, \mathbf{A} \rangle}$.

A tensor $\mathbf{X}$ with $\mathrm{rank}(\mathbf{X}) = \mathbf{r} = (r_1, \ldots, r_d)$ can be represented in the *Tucker decomposition* [17]

$$\mathbf{X} = \mathbf{C} \times_1 U_1 \cdots \times_d U_d = \mathbf{C} \bigtimes_{i=1}^{d} U_i, \tag{2.3}$$

with a *core tensor* $\mathbf{C} \in \mathbb{R}^{r_1 \times \cdots \times r_d}$ with $\mathrm{rank}(\mathbf{C}) = \mathbf{r}$ and *basis matrices* $U_i \in \mathbb{R}^{n_i \times r_i}$ with linearly independent columns. Without loss of generality, it can be assumed that the basis matrices have orthonormal columns, i. e. $U_i^{\mathrm{T}} U_i = I$. If for some $i$ this is not the case, a QR factorization $U_i = \widetilde{U}_i R$, with $\widetilde{U}_i$ orthonormal and $R$ regular and $\widetilde{\mathbf{C}} = (RC_{(i)})^{(i)}$ gives the required property.

A rank-$\mathbf{r}$ approximation to a tensor $\mathbf{A}$ can be computed by the *truncated higher-order SVD (HOSVD)* [18]: Let $\mathrm{P}_{r_i}^i$ be a the best rank-$r_i$ approximation operator in the $i$th mode, i. e. $\mathrm{P}_{r_i}^i \mathbf{A} = (U_i U_i^{\mathrm{T}} A_{(i)})^{(i)}$, where $U_i$ denotes the matrix of the $r_i$ dominant left singular vectors of $A_{(i)}$. Then the rank-$\mathbf{r}$ truncated HOSVD operator $\mathrm{P}_{\mathbf{r}}^{\mathrm{HO}}$ is given by

$$\mathrm{P}_{\mathbf{r}}^{\mathrm{HO}} \mathbf{A} := \mathrm{P}_{r_1}^1 \cdots \mathrm{P}_{r_d}^d \mathbf{A}. \tag{2.4}$$

In contrast to the matrix case, the HOSVD does not yield a best rank-$\mathbf{r}$ approximation, but only a quasi-best-approximation [18, Property 10] with a constant which deteriorates with respect to the number of modes:

$$\left\| \mathbf{A} - \mathrm{P}_{\mathbf{r}}^{\mathrm{HO}} \mathbf{A} \right\| \leq \sqrt{d} \min_{\mathbf{X} \in \mathcal{M}_{\mathbf{r}}} \| \mathbf{A} - \mathbf{X} \|. \tag{2.5}$$

*2.2. Riemannian manifold structure*

In [15], the authors show that the set $\mathcal{M}_{\mathbf{r}}$ of tensors of fixed multilinear rank $\mathbf{r} = (r_1, \ldots, r_d)$ forms a smooth embedded submanifold of $\mathbb{R}^{n_1 \times \cdots \times n_d}$. By counting the degrees of freedom in (2.3), it

follows that

$$\dim(\mathcal{M}_{\mathbf{r}}) = \prod_{i=1}^{d} r_i + \sum_{i=1}^{d} r_i n_i - r_i^2,$$

where the last term accounts for the fact that the Tucker decomposition is invariant to simultaneous transformation of the basis matrix with an invertible matrix and the core tensor with its inverse; as described in the previous subsection. Being a submanifold of the Euclidean space $(\mathbb{R}^{n_1 \times \cdots \times n_d}, \langle \cdot, \cdot \rangle)$, the manifold $\mathcal{M}_{\mathbf{r}}$ can be endowed with a Riemannian structure in a natural way with the Frobenius inner product $\langle \cdot, \cdot \rangle$ as the Riemannian metric.

As is proven in [16, Subsection 2.3], the tangent space of $\mathcal{M}_{\mathbf{r}}$ at $\mathbf{X} = \mathbf{C} \bigtimes_{i=1}^{d} U_i$ is parametrized as

$$T_{\mathbf{X}} \mathcal{M}_{\mathbf{r}} = \left\{ \dot{\mathbf{C}} \bigtimes_{i=1}^{d} U_i + \sum_{i=1}^{d} \mathbf{C} \times_i \dot{U}_i \bigtimes_{j \neq i} U_j \ \middle| \ \dot{\mathbf{C}} \in \mathbb{R}^{r_1 \times \cdots \times r_d}, \ \dot{U}_i \in \mathbb{R}^{n_i \times r_i} \text{ with } \dot{U}_i^{\mathrm{T}} U_i = O \right\},$$
(2.6)

and the orthogonal projection $\mathrm{P}_{\mathbf{X}} : \mathbb{R}^{n_1 \times \cdots \times n_d} \to T_{\mathbf{X}} \mathcal{M}_{\mathbf{r}}$ is given by

$$\mathbf{A} \mapsto \left( \mathbf{A} \bigtimes_{j=1}^{d} U_j^{\mathrm{T}} \right) \bigtimes_{i=1}^{d} U_i + \sum_{i=1}^{d} \mathbf{A} \times_i \left( \mathrm{P}_{U_i}^{\perp} \left[ \mathbf{A} \bigtimes_{j \neq i} U_j^{\mathrm{T}} \right]_{(i)} C_{(i)}^+ \right) \bigtimes_{k \neq i} U_k, \qquad (2.7)$$

where $C_{(i)}^+$ denotes the Moore-Penrose pseudoinverse of $C_{(i)}$. Note that $C_{(i)}$ has full row rank, i.e. $C_{(i)}^+ = C_{(i)}^{\mathrm{T}} (C_{(i)} C_{(i)}^{\mathrm{T}})^{-1}$. We use $\mathrm{P}_{U_i}^{\perp} = I_{n_i} - U_i U_i^{\mathrm{T}}$ to denote the orthogonal projection onto $\mathrm{span}(U_i)^{\perp}$.

Furthermore, it can be shown that the HOSVD (2.4) is locally a $C^\infty$ function in the manifold topology of $\mathcal{M}_{\mathbf{r}}$, see [5, Proposition 2.1] for further details. This allows us its use in continuous optimization, as we will se in the next section.

## 3. THE GEOMETRY OF $\mathcal{M}_{\mathbf{r}}$ AND RIEMANNIAN OPTIMIZATION
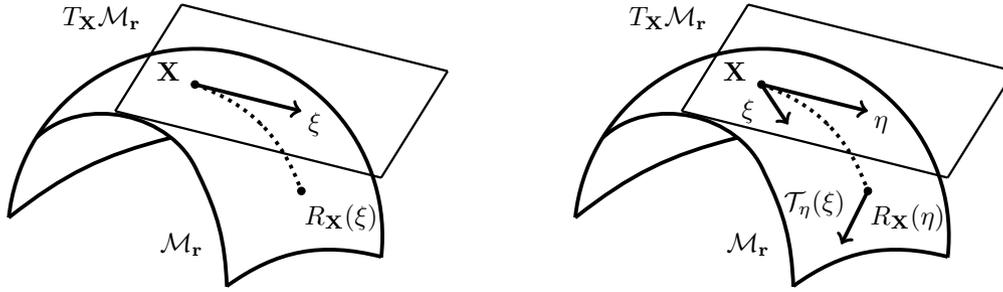
To construct optimization methods on $\mathcal{M}_{\mathbf{r}}$, we collect some basic concepts from the theory of optimization on manifolds. Our exposition follows the overview book [9]. Furthermore, we need to define and calculate the first and second derivatives of functions on $\mathcal{M}_{\mathbf{r}}$. In Corollary 3.7, we prove our main result, an explicit expression for the Riemannian Hessian on $\mathcal{M}_{\mathbf{r}}$. In the following, we will denote a Riemannian manifold by $\mathcal{M}$ and its elements by $x, y, \ldots \in \mathcal{M}$, when citing general results, and the manifold of tensors of fixed multilinear rank by $\mathcal{M}_{\mathbf{r}}$ and its elements by $\mathbf{X}, \mathbf{Y}, \ldots \in \mathcal{M}_{\mathbf{r}}$.

### 3.1. Retraction and vector transport

Since a manifold is in general not a linear space, the calculations required for a continuous optimization method need to be performed in a tangent space. Therefore, in each step, the need arises to map points from a tangent space to the manifold in order to generate the new iterate. The theoretically superior choice of such a mapping is the *exponential map*, which moves a point $x$ on the manifold along the geodesic locally defined by a vector in the tangent space $T_x \mathcal{M}$. However, computing the exponential map is prohibitively expensive in most situations, and it is shown in [9] that a first-order approximation, as specified in the following definition, is sufficient for many convergence results.

**Definition 3.1.** A *retraction* on a manifold $\mathcal{M}$ is a smooth mapping $R$ from the tangent bundle $T\mathcal{M}$ onto $\mathcal{M}$ with the following properties. Let $R_x$ denote the restriction of $R$ to $T_x \mathcal{M}$.

(i) $R_x(0_x) = x$, where $0_x$ denotes the zero element of $T_x \mathcal{M}$.

Figure 1. Retraction (left) and vector transport (right) on $\mathcal{M}_\mathbf{r}$.

(ii) With the canonical identification $T_{0_x} T_x \mathcal{M} \simeq T_x \mathcal{M}$, the mapping $R_x$ satisfies the *rigidity condition*

$$\mathrm{D} R_x(0_x) = \mathrm{id}_{T_x \mathcal{M}},$$

where $\mathrm{id}_{T_x \mathcal{M}}$ denotes the identity mapping on $T_x \mathcal{M}$.

Furthermore, "comparing" tangent vectors at distinct points on the manifold will be useful. The following definition gives us a way to transport a tangent vector $\xi \in T_x \mathcal{M}$ to the tangent space $T_{R_x(\eta)} \mathcal{M}$ for some $\eta \in T_x \mathcal{M}$ and some retraction $R$.

**Definition 3.2.** A *vector transport* on a manifold $\mathcal{M}$ is a smooth mapping

$$T\mathcal{M} \oplus T\mathcal{M} \to T\mathcal{M} : (\eta, \xi) \mapsto \mathcal{T}_\eta(\xi),$$

satisfying the following properties for all $x \in \mathcal{M}$:

(i) (Associated retraction) There exists a retraction $R$, called the *retraction associated with $\mathcal{T}$*, such that, for all $\eta, \xi$, it holds that $\mathcal{T}_\eta \xi \in T_{R_x(\eta)} \mathcal{M}$.

(ii) (Consistency) $\mathcal{T}_{0_x} \xi = \xi$ for all $\xi \in T_x \mathcal{M}$.

(iii) (Linearity) The mapping $\mathcal{T}_\eta : T_x \mathcal{M} \to T_{R_x(\eta)} \mathcal{M}$, $\xi \mapsto \mathcal{T}_\eta \xi$ is linear.

For $\mathcal{M} = \mathcal{M}_\mathbf{r}$, a retraction is given by the HOSVD, i.e. $R_\mathbf{X}(\xi) = \mathrm{P}_\mathbf{r}^{\mathrm{HO}}(\mathbf{X} + \xi)$. This is a consequence of the smoothness of the HOSVD (cf. Subsection 2.2) and the quasi-best approximation property (2.5). Details may be found in [5, Proposition 3]. A vector transport associated with a retraction $R$ is given by the orthogonal projection onto the tangent space, i.e. $\mathcal{T}_\eta(\xi) = \mathrm{P}_{R_\mathbf{X}(\eta)}(\xi)$, see [9, Subsection 8.1.3]; in our case, this is the formula (2.7). The efficient implementation of these operations is discussed in [5, Subsections 3.3–3.4]. A geometrical interpretation is shown in Figure 1.

### 3.2. The Riemannian gradient

The low-rank Tucker manifold $\mathcal{M}_\mathbf{r}$ being a submanifold a Euclidean space, the gradient of a real-valued function defined on it can be easily calculated by projecting the Euclidean gradient onto the tangent space.

**Lemma 3.3.** *[9, Section 3.6.1] Let $\mathcal{M}$ be a Riemannian submanifold of a Euclidean space $E$. Let $\bar{f} : E \to \mathbb{R}$ be a function with Euclidean gradient $\mathrm{grad}\,\bar{f}(x)$ at point $x \in \mathcal{M}$. Then the Riemannian gradient of $f := \bar{f}|_\mathcal{M}$ is given by $\mathrm{grad}\,f(x) = \mathrm{P}_x \mathrm{grad}\,\bar{f}(x)$, where $\mathrm{P}_x$ denotes the orthogonal projection onto the tangent space $T_x \mathcal{M}$.*

Then, by Lemma 3.3, the Riemannian gradient of the tensor completion cost function is given by

$$\mathrm{grad}\,f(\mathbf{X}) = \mathrm{P}_\mathbf{X}(\mathrm{P}_\Omega \mathbf{X} - \mathrm{P}_\Omega \mathbf{A}). \tag{3.1}$$

Using the sparsity of $\mathrm{P}_\Omega \mathbf{X} - \mathrm{P}_\Omega \mathbf{A}$, a gradient evaluation requires $\mathcal{O}(r^d(|\Omega| + n) + r^{d+1})$ operations, cf. [5, Subsection 3.1], where we assume that the $r_i$ and $n_i$ are constant in each mode for simplicity of notation.

### 3.3. The Riemannian Hessian

By definition, the *Riemannian Hessian* of a real-valued function $f$ on a Riemannian manifold $\mathcal{M}$ is a linear mapping

$$\mathrm{Hess}\, f(x)[\xi] = \nabla_\xi \operatorname{grad} f(x), \tag{3.2}$$

where $\nabla$ denotes the *Riemannian connection* on $\mathcal{M}$, cf. [9, Definition 5.5.1]. A finite-difference approximation can be defined in different ways. An intuitive formula is given by

$$H^{\mathrm{FD}}[\xi] = \frac{\mathcal{T}_\xi \operatorname{grad} f(R_x(h\xi)) - \operatorname{grad} f(x)}{h}, \tag{3.3}$$

see, for example, [9, Subsection 8.2.1]. However, such a mapping will in general not be linear [19], and should be applied with care, as theoretical understanding is yet incomplete.

On a Riemannian submanifold of a Euclidean space, the Riemannian connection is just the orthogonal projection of the directional derivative, i. e.

$$\mathrm{Hess}\, f(x)[\xi_x] = \mathrm{P}_x \left( \mathrm{D}\operatorname{grad} f(x)[\xi_x] \right), \tag{3.4}$$

and using Lemma 3.3, we get the following result.

**Lemma 3.4.** *[9, Section 5.3.3] Let $\mathcal{M}$ be a Riemannian submanifold of a Euclidean space $E$. Let $\bar{f} : E \to \mathbb{R}$ be a function with Euclidean gradient $\operatorname{grad} \bar{f}(x)$ at point $x \in \mathcal{M}$. Then the Riemannian Hessian of $f := \bar{f}|_{\mathcal{M}}$ is given by*

$$\mathrm{Hess}\, f(x)[\xi_x] = \mathrm{P}_x \,\mathrm{D} \left( \mathrm{P}_x \operatorname{grad} \bar{f}(x) \right). \tag{3.5}$$

Using the chain rule, we can write (3.5) as

$$\begin{aligned}
\mathrm{Hess}\, f(x)[\xi] &= \mathrm{P}_x \,\mathrm{D} \left( \mathrm{P}_x \operatorname{grad} \bar{f}(x) \right) \\
&= \mathrm{P}_x \,\mathrm{Hess}\, \bar{f}(x)[\xi_x] + \mathrm{P}_x \,\mathrm{D}_\xi \,\mathrm{P}_x \operatorname{grad} \bar{f}(x),
\end{aligned} \tag{3.6}$$

where we view $x \mapsto \mathrm{P}_x$ as an operator-valued function and denote its directional derivative by $\mathrm{D}_\xi$. We observe that he first term in (3.6) is just the orthogonal projection of the Euclidean Hessian, while the second one depends on the curvature of the manifold $\mathcal{M}$. Indeed, the second term is equal to zero when $\mathcal{M}$ is flat, i. e. a linear subspace of the embedding Euclidean space, cf. [20, Subsection 4.1]. Clearly, the main challenge in calculating the Riemannian Hessian in (3.6) is the derivative of the projection operator. In [13, Section 3], the authors show the following result using the *Weingarten map*.

**Lemma 3.5.** *Let $\mathcal{M}$ be a Riemannian submanifold of a Euclidean space $\mathcal{E}$. For any $x \in \mathcal{M}$, let $\mathrm{P}_x$ denote the orthogonal projection onto the tangent space $T_x\mathcal{M}$, and $\mathrm{P}_x^\perp := \mathrm{id}_{\mathcal{E}} - \mathrm{P}_x$ the orthogonal projection on its orthogonal complement $(T_x\mathcal{M})^\perp$. We view $x \mapsto \mathrm{P}_x$ as an operator-valued function and denote its Gâteaux derivative at point $x$ in the direction of $\xi \in T_x\mathcal{M}$ by $\mathrm{D}_\xi \,\mathrm{P}_x$. Then*

$$\mathrm{P}_x \,\mathrm{D}_\xi \,\mathrm{P}_x \, u = \mathrm{P}_x \,\mathrm{D}_\xi \,\mathrm{P}_x \left( \mathrm{P}_x^\perp u \right), \tag{3.7}$$

*for all $x \in \mathcal{M}$, $\xi \in T_x\mathcal{M}$ and $u \in \mathcal{E}$.*

This result can be applied to the case of the low-rank Tucker manifold $\mathcal{M} = \mathcal{M}_{\mathbf{r}}$. First, we calculate the derivative $\mathrm{D}_\xi \,\mathrm{P}_{\mathbf{X}}$.

**Lemma 3.6.** *Let $\mathbf{X} \in \mathcal{M}_{\mathbf{r}}$ be a tensor on the low-rank manifold, given by the factorization $\mathbf{X} = \mathbf{C} \bigtimes_{i=1}^{d} U_i$, and let $\xi \in T_{\mathbf{X}}\mathcal{M}_{\mathbf{r}}$, given by the variations*

$$\xi = \dot{\mathbf{C}} \bigtimes_{i=1}^{d} U_i + \sum_{i=1}^{d} \mathbf{C} \times_i \dot{U}_i \bigtimes_{j \neq i} U_j.$$

*We use the notations* $P_{U_i} = U_i U_i^{\mathrm{T}}$, $P_{U_i}^{\perp} = I_{n_i} - U_i U_i^{\mathrm{T}}$ *and* $\dot{P}_{U_i} = \dot{U}_i U_i^{\mathrm{T}} + U_i \dot{U}_i^{\mathrm{T}}$. *Then, for any* $\mathbf{E} \in \mathbb{R}^{n_1 \times \cdots \times n_d}$, *the derivative of* $P_{\mathbf{X}}$ *in the direction of* $\xi$ *is given by*

$$
\begin{aligned}
\mathrm{D}_\xi \, \mathrm{P}_{\mathbf{X}} \, \mathbf{E} = \sum_{i=1}^{d} \Bigg\{ & \mathbf{E} \times_i \dot{P}_{U_i} \bigtimes_{j \neq i} P_{U_j} \\
&+ \dot{\mathbf{C}} \times_i \left( P_{U_i}^{\perp} \left[ \mathbf{E} \bigtimes_{j \neq i} U_j^{\mathrm{T}} \right]_{(i)} C_{(i)} \right) \bigtimes_{k \neq i} U_k \\
&- \mathbf{C} \times_i \left( \dot{P}_{U_i} \left[ \mathbf{E} \bigtimes_{j \neq i} U_j^{\mathrm{T}} \right]_{(i)} C_{(i)} \right) \bigtimes_{k \neq i} U_k \\
&+ \sum_{l \neq i} \mathbf{C} \times_i \left( P_{U_i}^{\perp} \left[ \mathbf{E} \times_l \dot{U}_l^{\mathrm{T}} \bigtimes_{l \neq j \neq i} U_j^{\mathrm{T}} \right]_{(i)} C_{(i)} \right) \bigtimes_{k \neq i} U_k \\
&+ \mathbf{C} \times_i \left( P_{U_i}^{\perp} \left[ \mathbf{E} \bigtimes_{j \neq i} U_j^{\mathrm{T}} \right]_{(i)} \left[ \left( I - C_{(i)}^+ C_{(i)} \right) \dot{C}_{(i)}^{\mathrm{T}} C_{(i)}^{+\mathrm{T}} C_{(i)}^+ - C_{(i)}^+ \dot{C}_{(i)} C_{(i)}^+ \right] \right) \bigtimes_{k \neq i} U_k \\
&+ \sum_{l \neq i} \mathbf{C} \times_i \left( P_{U_i}^{\perp} \left[ \mathbf{E} \bigtimes_{j \neq i} U_j^{\mathrm{T}} \right]_{(i)} C_{(i)} \right) \times_l \dot{U}_l \bigtimes_{l \neq k \neq i} U_j \Bigg\},
\end{aligned}
$$

*where* $I = I_{\prod_{j \neq i} r_j}$ *is the identity matrix of the appropriate size.*

*Proof*

The formula can be obtained by identifying the tensor $\mathbf{X}$ with the factors in the Tucker decomposition and viewing the orthogonal projection defined in (2.7) as a function

$$
\mathrm{P}_{\cdot} \, \mathbf{E} : \mathbb{R}^{r_1 \times \cdots \times r_d} \times \mathbb{R}^{n_1 \times r_1} \times \cdots \times \mathbb{R}^{n_d \times r_d} \to \mathbb{R}^{n_1 \times \cdots \times n_d}, \ (\mathbf{C}, U_1, \ldots, U_d) \mapsto \mathrm{P}_{\mathbf{X}} \, \mathbf{E},
$$

for any $\mathbf{E} \in \mathbb{R}^{n_1 \times \cdots \times n_d}$. For calculating the derivative of the pseudoinverse, we use the formula given in [21, Theorem 4.3], i. e.

$$
\mathrm{D}_{\dot{C}} \left( C^+ \right) = \left( I - C^+ C \right) \dot{C}^{\mathrm{T}} C^{+\mathrm{T}} C^+ + C^+ C^{+\mathrm{T}} \dot{C}^{\mathrm{T}} \left( CC^+ - I \right) - C^+ \dot{C} C^+,
$$

and note that, here, the second term vanishes since $C = C_{(i)}$ has full row rank, and thus the pseudoinverse is a right inverse. $\qquad\square$

Using this result, we can immediately evaluate the curvature term in (3.6).

**Corollary 3.7.** *We use the setting of Lemma 3.6 and denote the orthogonal projection onto* $(T_{\mathbf{X}} \mathcal{M}_{\mathbf{r}})^{\perp}$ *by* $\mathrm{P}_{\mathbf{X}}^{\perp} := \mathrm{id} - \mathrm{P}_{\mathbf{X}}$. *Then*

$$
\mathrm{P}_{\mathbf{X}} \, \mathrm{D}_\xi \, \mathrm{P}_{\mathbf{X}} \, \mathrm{P}_{\mathbf{X}}^{\perp} \, \mathbf{E} = \widetilde{\mathbf{C}} \bigtimes_{i=1}^{d} U_i + \sum_{i=1}^{d} \mathbf{C} \times_i \widetilde{U}_i \bigtimes_{j \neq i} U_j \in T_{\mathbf{X}} \mathcal{M}_{\mathbf{r}},
$$

*with*

$$
\widetilde{\mathbf{C}} = \sum_{j=1}^{d} \left( \mathbf{E} \times_j \dot{U}_j^{\mathrm{T}} \bigtimes_{k \neq j} U_k^{\mathrm{T}} - \mathbf{C} \times_j \left( \dot{U}_j^{\mathrm{T}} \left[ \mathbf{E} \bigtimes_{k \neq j} U_j^{\mathrm{T}} \right]_{(j)} C_{(j)}^+ \right) \right),
$$

$$
\widetilde{U}_i = P_{U_i}^{\perp} \left( \left[ \mathbf{E} \bigtimes_{j \neq i} U_j^{\mathrm{T}} \right]_{(i)} \left( I - C_{(i)}^+ C_{(i)} \right) \dot{C}_{(i)}^{\mathrm{T}} C_{(i)}^{+\mathrm{T}} + \sum_{k \neq i} \left[ \mathbf{E} \times_k \dot{U}_k^{\mathrm{T}} \bigtimes_{k \neq j \neq i} U_j^{\mathrm{T}} \right]_{(i)} \right) C_{(i)}^+,
$$

*Proof*

The result follows by applying Lemma 3.6 to $P_{\mathbf{X}}^{\perp} \mathbf{E} \in \mathbb{R}^{n_1 \times \cdots \times n_d}$ after some lengthy but straightforward calculations, using the orthonormality relations $\dot{U}_i^{\mathrm{T}} U_i = O$, $U_i^{\mathrm{T}} U_i = I$ and the rules (2.1) and (2.2) for the matrix-tensor product. $\qquad\square$

Thus, the Riemannian Hessian of the function $f : \mathcal{M}_\mathbf{r} \to \mathbb{R}$,

$$f(\mathbf{X}) = \frac{1}{2} \big\| \, \mathrm{P}_\Omega \, \mathbf{X} - \mathrm{P}_\Omega \, \mathbf{A} \big\|^2,$$

can be written as

$$\mathrm{Hess}\, f(\mathbf{X})[\xi] = \mathrm{P}_\Omega(\xi) + \mathrm{P}_\mathbf{X} \, \mathrm{D}_\xi \, \mathrm{P}_\mathbf{X} \, \mathrm{P}_\mathbf{X}^\perp (\mathrm{P}_\Omega \, \mathbf{X} - \mathrm{P}_\Omega \, \mathbf{A}) \qquad (3.8)$$

and the second term can be evaluated with Corollary 3.7.

Note that for an efficient computation of the terms $\widetilde{U}_i$, it is advantageous to multiply out the term containing $I - C_{(i)}^+ C_{(i)}$. Then, the computation of $\mathrm{Hess}\, f(\mathbf{X})[\xi]$ for any given $\xi \in T_\mathbf{X}\mathcal{M}_\mathbf{r}$ has the same complexity as the computation of the gradient, i. e. $\mathcal{O}(r^d(|\Omega| + n) + r^{d+1})$.

*Remark 3.8.* For the matrix case $d = 2$, the Hessian expression (3.8) can be simplified to recover the expression shown in [6, 13],

$$\begin{aligned}
\mathrm{Hess}\, f(X)[\xi] = {} & \mathrm{P}_U \, \mathrm{P}_\Omega(\xi) \, \mathrm{P}_V + \mathrm{P}_U^\perp \big[ \, \mathrm{P}_\Omega(\xi) + \mathrm{P}_\Omega(X - A)\dot{V}\Sigma^{-1}V^\mathrm{T} \big] \, \mathrm{P}_V \\
& + \mathrm{P}_U \big[ \, \mathrm{P}_\Omega(\xi) + U\Sigma^{-1}\dot{U}^\mathrm{T} \, \mathrm{P}_\Omega(X - A) \big] \, \mathrm{P}_V^\perp,
\end{aligned}$$

where we identify the Tucker decomposition with the usual notation for the SVD, i. e. $U = U_1$, $V = U_2$ and $\Sigma = C$.

## 4. RIEMANNIAN MODELS AND TRUST-REGION METHODS

In principle, the results of the previous subsections can be used to conceive a Riemannian Newton method for the solution of problem (1.1). Such a method has been proposed in [22, pp. 279–283], where a convergence proof is given for strongly convex functions [22, Definition 1.1 in Chapter 7], using retraction by the exponential mapping (i. e. moving locally on a geodesic). [23, Theorem 4.4] proves quadratic convergence of the method to a critical point. [9, Theorem 6.3.2] provides a generalization for general retractions.

However, a plain Newton method has some well-known drawbacks:

1. The convergence radius may be small, i. e. if the initial guess is too far from a critical point the method may diverge.

2. Each step requires the solution of a linear system. This may be expensive and conceptually difficult if the Hessian operator is not even given explicitly but in terms of the action on a vector in the tangent space, as in (3.8).

There exists a number of strategies for remedying these problems. An intuitive method for globalizing the convergence of a Newton method is to modify the Hessian such that the solution $\xi$ of

$$\mathrm{Hess}\, f(x_k)[\xi] = -\operatorname{grad} f(x_k) \qquad (4.1)$$

defines a descent direction, see [24, Section 3.4] for an overview in the Euclidean case. In [9, Section 6.2] a generalization to the Riemannian case is proposed, replacing the Newton equation with

$$\big( \mathrm{Hess}\, f(x_k) + E_k \big)[\xi] = -\operatorname{grad} f(x_k),$$

where $E_k$ is a sequence of positive-definite linear operators on the tangent spaces $T_{x_k}\mathcal{M}$.

However, such perturbed Newton methods rely on heuristics, and their general convergence properties are not well understood. Moreover, they still require the solution of a linear system in each iteration. A way to circumvent this are trust-region methods [25], which find a critical point of the function $f$ by minimizing a sequence of constraint quadratic models $m_{x_k}$. Our exposition follows the generalization to Riemannian optimization as given by Absil et al. [26].

### 4.1. Models on a Riemannian manifold $\mathcal{M}$

For a real-valued function $f$ on a Riemannian manifold $\mathcal{M}$, a function $m_x$ is called an *order-q model*, $q > 0$, of $\mathcal{M}$ in $x \in \mathcal{M}$ if there exists a neighbourhood $\mathcal{U}$ of $x$ in $\mathcal{M}$ and a constant $c > 0$ such that

$$\big| f(y) - m_x(y) \big| \le c \big( \operatorname{dist}(x,y) \big)^{q+1}, \quad \text{for all } y \in \mathcal{U},$$

where $\operatorname{dist}$ denotes the Riemannian (geodesic) distance on $\mathcal{M}$. It can be shown [9, Proposition 7.1.3] that a model $m_x$ is order-$q$ if and only if there exists a neighbourhood $\mathcal{U}'$ of $x$ in $\mathcal{M}$ and a constant $c' > 0$ such that

$$\big| f(y) - m_x(y) \big| \le c \big\| R_x^{-1}(y) \big\|^{q+1}, \quad \text{for all } y \in \mathcal{U}.$$

i. e. the order of a model can be assessed using any retraction and we can avoid working with the exact geodesic.

Given a retraction $R$, this result allows to build a model for $f$ by simply taking a truncated Taylor expansion of

$$\widehat{f}_x := f \circ R_x,$$

for any $x \in \mathcal{M}$. The definition of $\widehat{f}_x : T_x \mathcal{M} \to \mathbb{R}$ as a real-valued function on a Euclidean space allows us to use standard results from multivariate analysis. A simple first-order model is then given by

$$\widehat{m}_x = \widehat{f}_x(0_x) + \mathrm{D}\,\widehat{f}_x(0_x)[\xi] = f(x) + \langle \operatorname{grad} f(x), \xi \rangle,$$

where the second equality follows form the rigidity condition of the retraction. A generic second-order model is given by

$$\begin{aligned} \widehat{m}_x &= \widehat{f}_x(0_x) + \mathrm{D}\,\widehat{f}_x(0_x)[\xi] + \tfrac{1}{2}\,\mathrm{D}^2\,\widehat{f}_x(0_x)[\xi,\xi] \\ &= f(x) + \langle \operatorname{grad} f(x), \xi \rangle + \tfrac{1}{2} \big\langle \operatorname{Hess} \widehat{f}(x)[\xi], \xi \big\rangle. \end{aligned}$$

A straightforward and useful modification is obtained by replacing the Euclidean Hessian on the tangent space $\operatorname{Hess} \widehat{f}(x)$ by the Riemannian expression $\operatorname{Hess} f(x)$. The following lemma shows that this can be done in a critical point of $f$ without any loss of information.

**Lemma 4.1.** *[9, Proposition 5.5.6] Let $R$ be a retraction and let $x^*$ be a critical point of a real-valued function $f$ on $\mathcal{M}$, i. e. $\operatorname{grad} f(x^*) = 0_{x^*}$. Then*

$$\operatorname{Hess} f(x^*) = \operatorname{Hess} \widehat{f}(0_{x^*}).$$

Thus, we can define a model

$$m_x = f(x) + \langle \operatorname{grad} f(x), \xi \rangle + \tfrac{1}{2} \big\langle \operatorname{Hess} f(x)[\xi], \xi \big\rangle,$$

which does not make any use of a retraction. However, Lemma 4.1 only guarantees that $m_x$ matches $f$ up to second order if $x$ is a critical point. In general, we can only prove that it will only give us a first-order model. The model $m_x$ can be shown to be of second order for general $x$ if the retraction $R$ is of second order, i. e. if it preserves second-order information of the exponential map, cf. [9, Proposition 5.5.5]. However, numerical results presented later in this section suggest that in our case the result also holds for general points on the manifold.

### 4.2. Models of different orders on $\mathcal{M}_{\mathbf{r}}$

We consider the manifold $\mathcal{M}_{\mathbf{r}}$ of fixed-rank tensors and would like to assess the quality of different model functions. In accordance with the previous subsection, we consider a first-order model

$$m_{\mathbf{X}}^{\mathrm{SD}}(\xi) := f(\mathbf{X}) + \langle \operatorname{grad} f(\mathbf{X}), \xi \rangle, \tag{4.2}$$

where the superscript indicates that this model corresponds to a steepest-descent method, and a second-order model

$$m_{\mathbf{X}}^{\mathrm{N}}(\xi) := f(\mathbf{X}) + \langle \operatorname{grad} f(\mathbf{X}), \xi \rangle + \tfrac{1}{2} \big\langle \operatorname{Hess} f(\mathbf{X})[\xi], \xi \big\rangle,$$
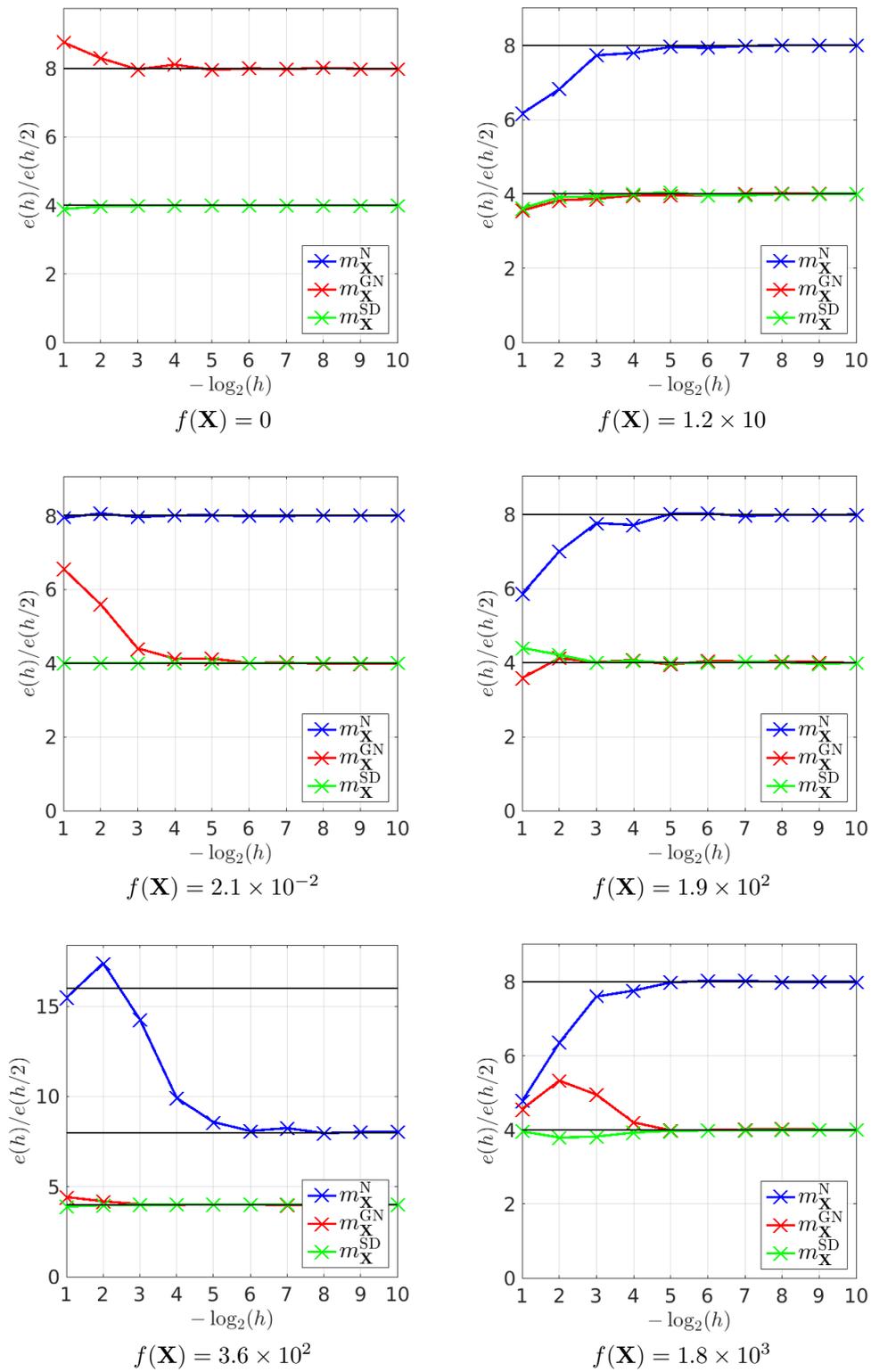
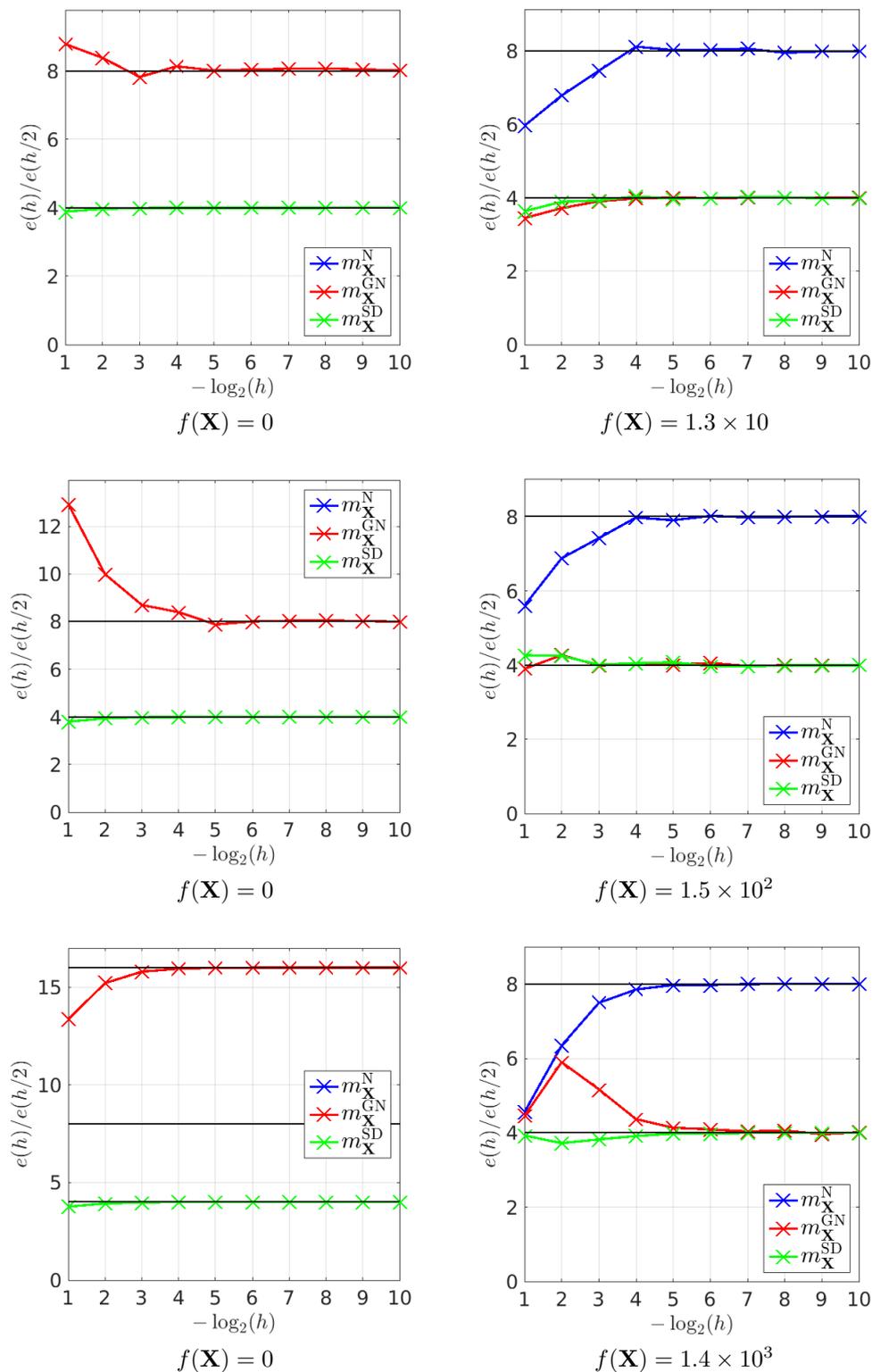Figure 2. The unknown tensor $\mathbf{A}$ has full rank, i. e. $\mathbf{A} \notin \mathcal{M}_{\mathbf{r}}$.

Figure 3. The unknown tensor $\mathbf{A}$ has low rank, i. e. $\mathbf{A} \in \mathcal{M}_{\mathbf{r}}$.

where the superscript indicates that this model corresponds to a Newton method. Furthermore, we would like to assess the quality of a Hessian approximation which drops the curvature term in Corollary 3.7 and thus ignores the second-order geometry of $\mathcal{M}_\mathbf{r}$. This is given by omitting the second term in (3.8), and just considering the projection of the Euclidean Hessian i.e.

$$\widetilde{\mathrm{Hess}}f(\mathbf{X})[\xi] = \mathrm{P}_\Omega\,\xi. \tag{4.3}$$

Omitting the curvature term of the Hessian corresponds to a Riemannian Gauß–Newton method, as described in [9, Subsection 8.4.1]. Thus, we consider the model function

$$m_\mathbf{X}^{\mathrm{GN}}(\xi) := f(\mathbf{X}) + \langle \mathrm{grad}\,f(\mathbf{X}), \xi \rangle + \tfrac{1}{2}\langle \widetilde{\mathrm{Hess}}f(\mathbf{X})[\xi], \xi \rangle. \tag{4.4}$$

As usual, we can expect a Gauß–Newton method to converge superlinearly (as the corresponding model to be of order higher than 1) if the residual of the least-squares problem is low. This can be seen in terms of (3.8), where the curvature term is given as

$$\big(\mathrm{Hess}\,f(\mathbf{X}) - \widetilde{\mathrm{Hess}}f(\mathbf{X})\big)[\xi] = \mathrm{P}_\mathbf{X}\,\mathrm{D}_\xi\,\mathrm{P}_\mathbf{X}\,\mathrm{P}_\mathbf{X}^\perp(\mathrm{P}_\Omega\,\mathbf{X} - \mathrm{P}_\Omega\,\mathbf{A}),$$

which is clearly equal to zero if $\mathrm{P}_\Omega\,\mathbf{X} = \mathrm{P}_\Omega\,\mathbf{A}$ and hence $f(\mathbf{X}) = 0$.

To assess the order of a model, we define for a given $m_\mathbf{X}$ the *model error*

$$e(\xi, h) := \big|\widehat{f}_\mathbf{X}(h\xi) - m_\mathbf{X}(h\xi)\big|,$$

for $\xi \in T_\mathbf{X}\mathcal{M}_\mathbf{r}$ and $h \geq 0$. Then $m_\mathbf{X}$ is an order-$q$ model in $\mathbf{X}$ if and only if

$$e(\xi, h) = \mathcal{O}\big(h^{q+1}\big), \quad \text{for all } \xi \in T_\mathbf{X}\mathcal{M}_\mathbf{r}.$$

In Figures 2 and 3, we test the model orders of (4.2)–(4.4). We generate random tensors $\mathbf{B}_1, \ldots, \mathbf{B}_{1000} \in \mathbb{R}^{10 \times 10 \times 10}$ with normally distributed entries and project them onto a given tangent space of $\mathcal{M}_{(3,3,3)}$ to get $\xi_i = \mathrm{P}_\mathbf{X}(\mathbf{B}_i)$. We normalize the resulting vectors to get $\|\xi_i\| = 1$. We compute the errors $e(\xi_i, 2^{-j})$ for $j = 0, \cdots, 10$, and plot the geometric mean of the factors $(\xi_i, 2^{-(j+1)})/(\xi_i, 2^{-j})$ over all $i$. The first columns contain the results for a stationary point of $f$, i.e. $\|\mathrm{grad}\,f(\mathbf{X}^*)\| = 0$, the second columns contain the results for an arbitrary point on the manifold with $\|\mathrm{grad}\,f(\mathbf{X})\| \neq 0$. The first, second and third rows contain results for different sampling sizes, with $|\Omega| = 10, 100, 1000$, respectively. Note that $|\Omega| = 1000$ represents full sampling, i.e. vector approximation. We write $f(\mathbf{X}) = 0$ whenever the function value computed is smaller that the machine precision of $10^{-16}$.

We observe that the model function $m_\mathbf{X}^{\mathrm{SD}}$, indeed, provides results of first order in all cases. The model function $m_\mathbf{X}^{\mathrm{N}}$ provides results of second order not only in critical points, as has been proved by theory, but also in general points on the manifold. This can be seen as an indication that the retraction by HOSVD preserves second-order information although we cannot prove this. We also observe that the Gauß–Newton type model function $m_\mathbf{X}^{\mathrm{GN}}$ gives second-order results whenever the curvature term is small enough, otherwise it is only a first-order model; this matches the theoretical predictions we made earlier. It is especially worth noting that, for $\mathbf{A} \in \mathcal{M}_\mathbf{r}$, a Gauß–Newton model is sufficient; however, this result is not robust if we add some noise. Note that in the cases where the blue curve cannot be seen in the plot, the models $m_\mathbf{X}^{\mathrm{GN}}$ and $m_\mathbf{X}^{\mathrm{N}}$ match almost exactly.

We also remark that in the case of exact tensor reconstruction, i.e. $\mathbf{A} \in \mathcal{M}_\mathbf{r}$ and $|\Omega| = \prod_i n_i$ (the lower-left plot in Figure 3), both $m_\mathbf{X}^{\mathrm{N}}$ and $m_\mathbf{X}^{\mathrm{GN}}$ seem to be models of order 3, which means that the third-order term in the Taylor expansion of $f$ vanishes. This may be attributed to a possible symmetry of $f$ around the local minimizer $\mathbf{X}^* = \mathbf{A}$ in this case, i.e. $f(\mathrm{Exp}_{\mathbf{X}^*}(\xi)) = f(\mathrm{Exp}_{\mathbf{X}^*}(-\xi))$, where $\mathrm{Exp}$ denotes the exponential map. This means that the odd-exponent terms in the Taylor expansion are equal to zero. However, we cannot verify this theoretically as we do not have a closed-form expression for the exponential map on $\mathcal{M}_\mathbf{r}$.

### 4.3. Riemannian trust-region method

The main idea of trust-region methods is solving a model problem

$$
\min_{\eta \in T_{\mathbf{X}_k}\mathcal{M}_{\mathbf{r}}} m_{\mathbf{X}_k}(\eta)
$$
$$
\text{s.\,t. } \|\eta\| \leq \Delta_k,
\tag{4.5}
$$

for some $\Delta_k \geq 0$ in each iteration $k$ to obtain a search direction $\eta_k$. To get meaningful results it is crucial to check how well the model $m_{\mathbf{X}_k}$ approximates $\widehat{f}$ in $T_{\mathbf{X}_k}\mathcal{M}_{\mathbf{r}}$ in the neighbourhood of $0_{\mathbf{X}_k}$. This can be expressed in the form of the quotient

$$
\rho_k := \frac{\widehat{f}(0_{\mathbf{X}_k}) - \widehat{f}(\eta_k)}{m_{\mathbf{X}_k}(0_{\mathbf{X}_k}) - m_{\mathbf{X}_k}(\eta_k)}.
\tag{4.6}
$$

If $\rho_k$ is small (convergence theory suggests that $\rho' < \frac{1}{4}$ is an appropriate threshold), then the model is very inaccurate: the step must be rejected, and the trust-region radius $\Delta_k$ must be reduced. If $\rho_k$ is small but less dramatically so, then the step is accepted but the trust-region radius is reduced. If $\rho_k$ is close to 1, then there is a good agreement between the model and the function over the step, and the trust-region radius can be expanded. If $\rho_k \gg 1$, then the model is inaccurate, but the overall optimization iteration is producing a significant decrease in the cost. If this is the case and the restriction in (4.5) is active, we can try to expand the trust-region radius as long as we stay below a predefined bound $\bar{\Delta} > 0$. This method is summarized in Algorithm 4.2, cf. [26, Algorithm 1].

---

**Algorithm 4.2** Riemannian trust-region method for $\mathcal{M}_{\mathbf{r}}$

---

**Input:** Initial iterate $x_0 \in \mathcal{M}$; parameters $\bar{\Delta} > 0$, $\Delta_0 \in (0, \bar{\Delta})$, $\rho' \in (0, \frac{1}{4})$.
1: **for** $k = 0$ **until** convergence **do**
2:     Obtain $\eta_k$ by approximately solving (4.5)
3:     <Test for convergence>
4:     Evaluate $\rho_k$ from (4.6)
5:     **if** $\rho_k < \frac{1}{4}$ **then**
6:         $\Delta_{k+1} = \frac{1}{4}\Delta_k$
7:     **else if** $\rho_k > \frac{3}{4}$ and $\|\eta_k\| = \Delta_k$ **then**
8:         $\Delta_{k+1} = \min(2\Delta_k, \bar{\Delta})$
9:     **else**
10:        $\Delta_{k+1} = \Delta_k$
11:     **end if**
12:     **if** $\rho_k > \rho'$ **then**
13:        $\mathbf{X}_{k+1} = R_{\mathbf{X}_k}(\eta_k)$
14:     **else**
15:        $\mathbf{X}_{k+1} = \mathbf{X}_k$
16:     **end if**
17: **end for**

---

There exist different strategies for (approximately) solving the trust-region subproblems (4.5). We apply a truncated CG method [26, Algorithm 2], which is a straightforward adaptation of Steighaug's method [27] for problems in $\mathbb{R}^n$. It ensures that the CG method is stopped after a fixed maximal number of iterations $K_{\max}$. Since a CG iteration just requires a fixed number of matrix-vector products, the total cost of the trust-region method with exact Hessian evaluation is given by

$$
\mathcal{O}\big(K_{\max}(r^d(|\Omega| + n) + r^{d+1})\big)
$$

The convergence theory follows standard techniques from Euclidean optimization [25]. Under some technical assumptions, it can be shown that Algorithm 4.2 converges globally to a stationary point [26, Theorem 4.4] of $f$. Locally superlinear convergence to a nondegenerate local minimum

can be shown [26, Theorem 4.12] as long as the quadratic term in $m_{\mathbf{X}_k}$ is a sufficiently good Hessian approximation of $f$.

In general, we cannot rule out Algorithm 4.2 converging to a nonregular minimum if $|\Omega| < \dim(\mathcal{M}_{\mathbf{r}})$. If this causes problems, we can enforce positive-definiteness of the Hessian by considering a cost function regularized with an identity term

$$f_\mu(\mathbf{X}) = \frac{1}{2}\big\| \mathrm{P}_\Omega\, \mathbf{X} - \mathrm{P}_\Omega\, \mathbf{A} \big\|^2 + \frac{\mu}{2}\|\mathbf{X}\|^2,$$

for some $\mu > 0$. However, such a problem may not be well-posed since there is not enough information provided to recover $\mathbf{X}$ in a meaningful way. Moreover, in our practical experiments we did not have a need to use this regularization.

## 5. NUMERICAL EXPERIMENTS

We implemented our method in MATLAB version 2015b using the Tensor Toolbox version 2.6 [28, 29] for the basic tensor arithmetic and Manopt version 3.0 [30] for handling the Riemannian trust-region scheme. All tests were performed on a quad-core Intel i7-2600 CPU with 8 GB of RAM running 64-Bit Ubuntu 16.04 Linux.

In Algorithm 4.2, we choose the standard parameters $\bar\Delta = \dim(\mathcal{M}_{\mathbf{r}})$, $\Delta_0 = \bar\Delta/8$, $\rho' = 0.1$. The initial guess $\mathbf{X}_0$ is generated randomly by a uniform distribution on $(0,1)$ for each entry in the factors in the Tucker decomposition. We apply a QR factorization in each mode to ensure that the basis matrices are orthogonal. The sampling set $\Omega$ is chosen from a uniform distribution on the index set.

### 5.1. Uniformly distributed random data

We test the convergence behaviour of Algorithm 4.2 for the recovery of a partially known tensor $\mathbf{A}$ with uniformly distributed entries. We observe that the trust-region method with exact Hessian computation yields superlinear convergence after a small number of iterations in all cases observed here. The finite difference Hessian approximation shows similar behaviour, however, the convergence is slower and becomes less reliable for a large gradient norm reduction. The Gauß–Newton Hessian approximation shows shows superlinear convergence behaviour if $\mathbf{A} \in \mathcal{M}_{\mathbf{r}}$, but not in the case $\mathbf{A} \notin \mathcal{M}_{\mathbf{r}}$, as predicted in the previous sections. The state-of-the-art Riemannian method, nonlinear CG [5], shows linear convergence with convergence rate superior to steepest descent, but the convergence rate may slow, especially in the case of noise.

### 5.2. Survey data

In survey statistics, data in the form of order-3 tensors arises in a natural way: for $n_1$ of individuals, $n_2$ properties are collected over $n_3$ time points; see, for example, [32]. We choose a standard data set [31], containing reading proficiency test measures of schoolchildren over a period of time. A typical problem in such data sets in practice is missing entries, resulting from nonresponse or failure to enter some of the data points correctly; see [33]. A typical application case is a sampling set greater or equal to haf the total tensor size. As Figure 5 shows, data of this type shows rapidly decaying singular values, especially in the time mode ($i = 3$) and our trust-region method can be used to retrieve deleted data in a low-rank framework. The trust-region method also converges superlinearly in this case. The simplified Gauß–Newton trust-region scheme does not show superlinear convergence since noise is present in this application case. The trust-region methods also compares favorably with nonlinear CG in this case. Our results can be seen as an indication that Riemannian trust-region methods can be used for statistical data recovery.
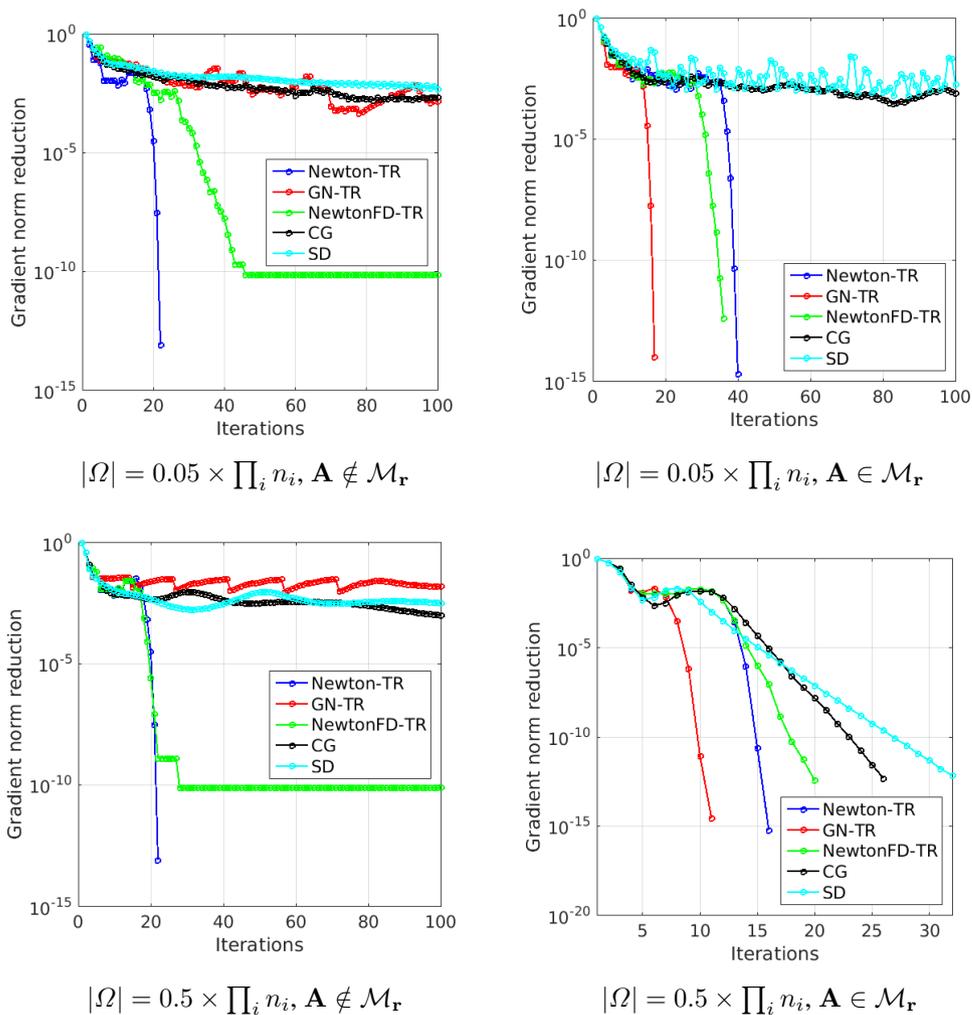
Figure 4. Convergence of Riemannian methods for (1.1) with $n_i \equiv 20$ and $r_i \equiv 2$.
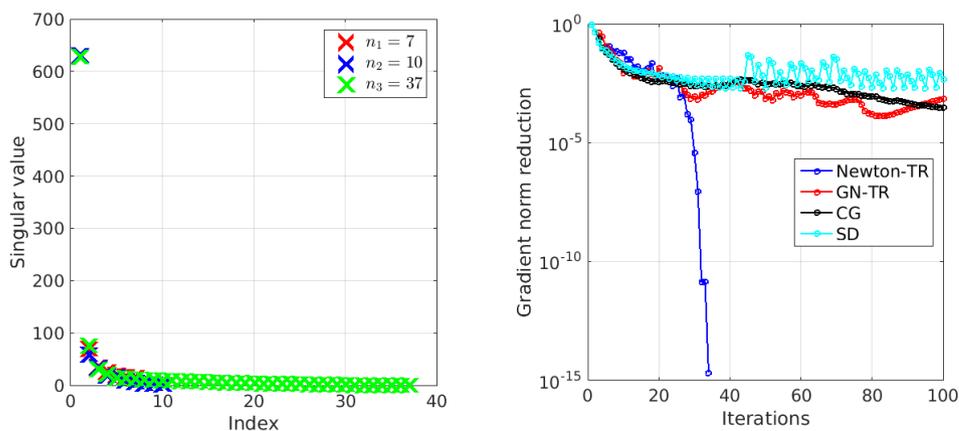


Figure 5. Left: Singular values of the data set [31]; right: convergence of Riemannian methods for tensor completion with $|\Omega| = 0.5 \times \prod_i n_i$ and $\mathbf{r} = (3, 5, 5)$.

## 6. CONCLUSIONS AND DISCUSSION

We have derived the Riemannian Hessian for functions on the manifold of tensors of fixed multilinear rank in Tucker format. We have shown that it can be used to construct a rapidly and robustly converging trust-region scheme for tensor completion. Furthermore, this is the first theoretical result on the second-order properties of the given manifold; we believe this to be useful for an improved understanding of the underlying geometry. Our numerical results also indicate that Riemannian optimization is a suitable technique for the recovery of missing entries from multilinear survey data with low-rank structure. We believe that this aspect merits further exploration; a comparison of Riemannian techniques with standard imputation methods from statistics [33] may reveal opportunities and limitations of this approach. For this, a better understanding of the sensitivity of the Tucker decomposition to perturbations is required.

Another well-known way to obtain superlinear convergence is a Riemannian BFGS method. In recent research, several schemes have been proposed, generalizing this standard method from Euclidean optimization to the Riemannian case; see [34, Subsection 5.2] for an application to the manifold of matrices of fixed rank. Extending this idea to tensors merits some examination. For high-dimensional applications with $d \gg 3$, hierarchical tensor formats [15, 35] are crucial; see [36] for a Riemannian optimization approach.

## ACKNOWLEDGEMENTS

## References

1. Ma Y, Min K, *et al.*. Low-rank matrix recovery and completion via convex optimization. http://perception.csl.illinois.edu/matrix-rank/. Accessed: 24 March 2017.
2. Liu J, Musialski P, Wonka P, Ye J. Tensor completion for estimating missing values in visual data. *IEEE Trans. Pattern Anal. Mach. Intell.* 2013; **35**(1):208–220.
3. Signoretto M, de Lathauwer L, Suykens JAK. Nuclear norms for tensors and their use for convex multilinear estimation. *Technical Report* 2010.
4. Gandy S, Recht B, Yamada I. Tensor completion and low-n-rank tensor recovery via convex optimization. *Inverse Probl.* 2011; **27**(2):025 010.
5. Kressner D, Steinlechner M, Vandereycken B. Low-rank tensor completion by Riemannian optimization. *BIT* 2014; **54**(2):447–468.
6. Vandereycken B. Low-rank matrix completion by Riemannian optimization. *SIAM J. Optim.* 2013; **23**(2):1214–1236.
7. Ngo TT, Saad Y. Scaled gradients on Grassmann manifolds for matrix completion. *Advances in Neural Information Processing Systems*, Pereira F, Burges C, Bottou L, Weinberger K (eds.), 25, 2012; 1412–1420.
8. Mishra B, Meyer G, Bonnabel S, Sepulchre R. Fixed-rank matrix factorizations and Riemannian low-rank optimization. *Comput. Stat.* 2014; **29**(3):591–621.
9. Absil PA, Mahoney R, Sepulchre R. *Optimization Algorithms on Matrix Manifolds*. Princeton University Press: Princeton, 2008.
10. Boumal N, Absil PA. RTRMC: a Riemannian trust-region method for low-rank matrix completion. *Advances in Neural Information Processing Systems*, Shawe-Taylor J, Zemel R, Bartlett P, Pereira F, Weinberger K (eds.), 24, 2011; 406–414.
11. Eldén L, Savas B. A Newton–Grassmann method for computing the best multilinear rank-$(r_1, r_2, r_3)$ approximation of a tensor. *SIAM. J. Matrix Anal. Appl.* 2009; **31**(2):248–271.
12. Ishteva M, Absil PA, van Huffel S, de Lathauwer L. Best low multilinear rank approximation of higher-order tensors, based on the Riemannian trust-region scheme. *SIAM J. Matrix Anal. Appl.* 2011; **32**(1):115–135.
13. Absil PA, Mahony R, Trumpf J. Optimization techniques on Riemannian manifolds. *Geometric Science of Information*, Nielsen F, Barbaresco F (eds.), 1, 2013; 361–368.
14. Kolda TG, Bader BW. Tensor decompositions and applications. *SIAM Rev.* 2009; **51**(3):131–173.
15. Uschmajew A, Vandereycken B. The geometry of algorithms using hierarchical tensors. *Linear Algebra Appl.* 2013; **439**(1):133–166.
16. Koch O, Lubich C. Dynamical tensor approximation. *SIAM J. Matrix Anal. Appl.* 2007; **31**(5):2360–2375.
17. Tucker LR. Some mathematical notes on three-mode factor analysis. *Psychometrika* 1966; **31**(3):279–311.

18. De Lathauwer L, de Moor B, Vandewalle J. A multilinear singular value decomposition. *SIAM J. Matrix Anal. Appl.* 2000; **21**(4):1253–1278.
19. Boumal N. Riemannian trust regions with finite-difference Hessian approximations are globally convergent. *Geometric Science of Information*, Nielsen F, Barbaresco F (eds.), 2, 2015; 467–475.
20. Kressner D, Steinlechner M, Vandereycken B. Preconditioned low-rank Riemannian optimization for linear systems with tensor product structure. *SIAM J. Sci. Comp.* 2016; **38**(4):A2018–A2044.
21. Golub GH, Pereyra V. The differentiation of pseudo-inverses and nonlinear least squares problems whose variables separate. *SIAM J. Numer. Anal.* 1973; **10**(2):413–432.
22. Udriște C. *Convex Functions and Optimization Methods on Riemannian Manifolds*. Kluwer Academic Publishers: Dordrecht, 1994.
23. Smith ST. Optimization techniques on Riemannian manifolds. *Hamiltonian and Gradient Flows, Algorithms and Control*, Bloch A (ed.). AMS: Providence, 1994; 397–434.
24. Nocedal J, Wright SJ. *Numerical Optimization*. Springer: New York, 2006.
25. Conn AR, Gould NIM, Toint PL. *Trust Region Methods*. SIAM: Philadelphia, 2000.
26. Absil PA, Baker CG, Gallivan KA. Trust-region methods on Riemannian manifolds. *Found. Comput. Math.* 2007; **7**(3):303–330.
27. Steihaug T. The conjugate gradient method and trust regions in large scale optimization. *SIAM J. Numer. Anal* 1983; **20**(3):626–637.
28. Bader BW, Kolda TG. Algorithm 862: MATLAB tensor classes for fast algorithm prototyping. *ACM Trans. Math. Softw.* 2011; **32**(4):635–653.
29. Bader BW, Kolda TG, *et al.*. MATLAB Tensor Toolbox Version 2.6. Available online 2015.
30. Boumal N, Mishra B, Absil PA, Sepulchre R, *et al.*. Manopt, a Matlab toolbox for optimization on manifolds. *J. Mach. Learn. Res.* 2014; **15**(1):1455–1459.
31. Kroonenberg PM. Information on Bus' learning-to-read data. http://www.leidenuniv.nl/fsw/three-mode/data/businfo.htm. Accessed: 24 March 2017.
32. Kroonenberg PM. Three-Mode Principal Component Analysis: Theory and Applications. PhD Thesis, Universiteit Leiden 1983.
33. Little RJA, Rubin DB. *Statistical Analysis with Missing Data*. Wiley: New York, 2014.
34. Huang W, Absil PA, Gallivan KA. Intrinsic representation of tangent vectors and vector transports on matrix manifolds. *Numer. Math.* 2016; .
35. Grasedyck L, Kressner D, Tobler C. A literature survey of low-rank tensor approximation techniques. *GAMM-Mitt.* 2013; **36**(1):53–78.
36. Da Silva C, Herrmann FJ. Optimization on the hierarchical Tucker manifold – applications to tensor completion. *Linear Algebra Appl.* 2015; **481**:131–173.