

**Max-Planck-Institut
für Mathematik
in den Naturwissenschaften
Leipzig**

**Convergence analysis of Riemannian
GaussNewton methods and its
connection with the geometric
condition number**

by

Paul Breiding and Nick Vannieuwenhoven

Preprint no.: 69

2017



Convergence analysis of Riemannian Gauss–Newton methods and its connection with the geometric condition number

Paul Breiding¹

Max-Planck-Institute for Mathematics in the Sciences, Leipzig, Germany.

Nick Vannieuwenhoven²

KU Leuven, Department of Computer Science, Leuven, Belgium.

Abstract

We obtain estimates of the multiplicative constants appearing in local convergence results of the Riemannian Gauss–Newton method for least squares problems on manifolds and relate them to the geometric condition number of [P. Bürgisser and F. Cucker, *Condition: The Geometry of Numerical Algorithms*, 2013].

Keywords: Riemannian Gauss–Newton method, convergence analysis, geometric condition number, CPD

1. Introduction

Many problems in science and engineering are *parameter identification problems* (PIPs). Herein, there is a parameter domain $\mathcal{M} \subset \mathbb{R}^M$ and a function $\Phi : \mathcal{M} \rightarrow \mathbb{R}^N$. Given a point \mathbf{y} in the image of Φ , the PIP asks to identify parameters $\mathbf{x} \in \mathcal{M}$ such that $\mathbf{y} = \Phi(\mathbf{x})$; note that there could be several such parameters. For example, computing *QR*, *LU*, Cholesky, polar, singular value and eigendecompositions of a given matrix $A \in \mathbb{R}^{m \times n} \simeq \mathbb{R}^{mn}$ are examples of this. In other cases we have a tensor $\mathfrak{A} \in \mathbb{R}^{m_1 \times \dots \times m_d}$ and need to compute CP, Tucker, block term, hierarchical Tucker, or tensor trains decompositions [7].

If the object $\tilde{\mathbf{y}} \in \mathbb{R}^N$ whose parameters should be identified originates from applications, then usually $\tilde{\mathbf{y}} \notin \Phi(\mathcal{M})$. Nevertheless, in this setting one seeks parameters $\mathbf{x} \in \mathcal{M}$ such that $\mathbf{y} := \Phi(\mathbf{x})$ is as close as possible to $\tilde{\mathbf{y}}$, e.g., in the Euclidean norm. This can be formulated as a nonlinear least squares problem:

$$\tilde{\mathbf{y}} \mapsto \arg \min_{\mathbf{x} \in \mathcal{M}} \frac{1}{2} \|\Phi(\mathbf{x}) - \tilde{\mathbf{y}}\|^2. \quad (1)$$

Here, we deal with functions Φ that offer differentiability guarantees, so that continuous optimization methods can be employed for solving (1). Specifically, we assume that \mathcal{M} is a smooth embedded submanifold³ of \mathbb{R}^M and that Φ is a smooth function on \mathcal{M} [11, Chapters 1 and 2]. Hence, (1) is a *Riemannian optimization problem* that can be solved using, e.g., Riemannian Gauss–Newton (RGN) methods [2]; see Section 2.

The sensitivity of $\mathbf{x} \in \mathcal{M}$ with respect to perturbations of $\mathbf{y} = \Phi(\mathbf{x})$ might impact the performance of these RGN methods. Let $\Psi : \mathcal{X} \rightarrow \mathcal{Y}$ be a smooth map between manifolds \mathcal{X} and \mathcal{Y} , and let $T_{\mathbf{x}}\mathcal{X}$

Email addresses: breiding@mis.mpg.de (Paul Breiding), nick.vannieuwenhoven@cs.kuleuven.be (Nick Vannieuwenhoven)

¹Funding: The author was partially supported by DFG research grant BU 1371/2-2.

²Funding: The author was supported by a Postdoctoral Fellowship of the Research Foundation—Flanders (FWO).

³Both the optimization problem (1) and the condition number of maps between manifolds can be defined for abstract manifolds. Nevertheless, we consider embedded manifolds because it greatly simplifies the proof of the main theorem, allowing us to compare tangent spaces in the ambient space using Wedin’s theorem [13, Chapter III, Theorem 3.9]. This is no longer possible for abstract manifolds, which would make the letter much more difficult to understand. In practice, many manifolds are naturally embedded.

denote the tangent space to the manifold \mathcal{X} at $\mathbf{x} \in \mathcal{X}$. We recall from [5, Section 14.3] that the geometric condition number $\kappa(\mathbf{x})$ characterizes to first-order the sensitivity of the output $\mathbf{y} = \Psi(\mathbf{x})$ to input perturbations as the spectral norm of the derivative operator $d_{\mathbf{x}}\Psi : T_{\mathbf{x}}\mathcal{X} \rightarrow T_{\Psi(\mathbf{x})}\mathcal{Y}$; that is, $\kappa(\mathbf{x}) := \|d_{\mathbf{x}}\Psi\| := \max_{\mathbf{t} \in T_{\mathbf{x}}\mathcal{X}} \|d_{\mathbf{x}}\Psi(\mathbf{t})/\|\mathbf{t}\|\|$. In the case of PIPs, the geometric condition number is derived as follows. Assume that there exists an open neighborhood \mathcal{N} of $\mathbf{x} \in \mathcal{M}$ such that $\mathcal{M} = \Phi(\mathcal{N})$ is a smooth manifold with $m = \dim \mathcal{M} = \dim \mathcal{N}$. Since $\Phi|_{\mathcal{N}} : \mathcal{N} \rightarrow \mathcal{M}$ is a smooth map between manifolds, the inverse function theorem for manifolds [11, Theorem 4.5] entails that there exists a unique inverse function $\Phi_{\mathbf{x}}^{-1}$ whose derivative satisfies $d_{\Phi(\mathbf{x})}\Phi_{\mathbf{x}}^{-1} = (d_{\mathbf{x}}\Phi)^{-1}$, provided that $d_{\mathbf{x}}\Phi$ is injective. Hence, the geometric condition number of the *parameters*⁴ \mathbf{x} is

$$\kappa(\mathbf{x}) := \|d_{\Phi(\mathbf{x})}\Phi_{\mathbf{x}}^{-1}\| = \|(d_{\mathbf{x}}\Phi)^{-1}\| = \frac{1}{\varsigma_m(d_{\mathbf{x}}\Phi)}, \quad (2)$$

where $\varsigma_m(A)$ is the m th largest singular value of the linear operator A . If the derivative is not injective, then the condition number is defined to be ∞ .

In this letter, we show that the condition number of the parameters \mathbf{x} in (2) appears naturally in the multiplicative constants in convergence estimates of RGN methods. Our main contribution is Theorem 1.

2. The Riemannian Gauss–Newton method

Recall that a Riemannian manifold $(\mathcal{M}, \langle \cdot, \cdot \rangle)$ is a smooth manifold \mathcal{M} , where for each $p \in \mathcal{M}$ the tangent space $T_p\mathcal{M}$ is equipped with an inner product $\langle \cdot, \cdot \rangle_p$ that varies smoothly with p ; see [11, Chapter 13]. The zero element of $T_p\mathcal{M}$ is denoted by 0_p . Since we deal exclusively with embedded submanifolds $\mathcal{M} \subset \mathbb{R}^M$, we take $\langle \mathbf{a}, \mathbf{b} \rangle_p := \mathbf{a}^T \mathbf{b}$ equal to the standard inner product on \mathbb{R}^M . In the following we drop the subscript “ p .” The induced norm is $\|\mathbf{v}\| = \sqrt{\langle \mathbf{v}, \mathbf{v} \rangle}$. The *tangent bundle* of a manifold \mathcal{M} is the smooth vector bundle $T\mathcal{M} := \{(p, \mathbf{v}) \mid p \in \mathcal{M}, \mathbf{v} \in T_p\mathcal{M}\}$.

In the remainder of this letter, we let $\mathcal{M} \subset \mathbb{R}^M$ be an embedded submanifold with $m = \dim \mathcal{M} \leq M$ equipped with the standard Riemannian metric inherited from \mathbb{R}^M . Riemannian optimization methods can be applied to the minimization of a least-squares cost function

$$f : \mathcal{M} \rightarrow \mathbb{R}, p \mapsto \frac{1}{2}\|F(p)\|^2 \quad \text{with } F : \mathcal{M} \rightarrow \mathbb{R}^N. \quad (3)$$

Recall that Newton’s method for minimizing f consists of choosing a $x_0 \in \mathcal{M}$ and then generating a sequence of iterates x_1, x_2, \dots in \mathcal{M} according to the following process:

$$x_{k+1} \leftarrow R_{x_k}(\eta_k) \quad \text{with } (\nabla_{x_k}^2 f)\eta_k = -\nabla_{x_k} f; \quad (4)$$

herein, $\nabla_{x_k} f : T_{x_k}\mathcal{M} \rightarrow \mathbb{R}$ is the *Riemannian gradient*, and $\nabla_{x_k}^2 f : T_{x_k}\mathcal{M} \rightarrow T_{x_k}\mathcal{M}$ is the *Riemannian Hessian*; for details see [2, Chapter 6]. The map $R_{x_k} : T_{x_k}\mathcal{M} \rightarrow \mathcal{M}$ is a *retraction operator*.

Definition 1 (Retraction [2, 9]). *A retraction R is a map from an open subset $T\mathcal{M} \supset \mathcal{U} \rightarrow \mathcal{M}$ that satisfies all of the following properties for every $p \in \mathcal{M}$:*

1. $R(p, 0_p) = p$;
2. \mathcal{U} contains a neighborhood \mathcal{N} of $(p, 0_p)$ such that the restriction $R|_{\mathcal{N}}$ is smooth;
3. R satisfies the local rigidity condition $d_{0_x}R(x, \cdot) = \text{id}_{T_x\mathcal{M}}$ for all $(x, 0_x) \in \mathcal{N}$.

We let $R_p(\cdot) := R(p, \cdot)$ be the retraction R with foot at p .

A retraction is a first-order approximation of the exponential map [2]; the following result is well-known.

⁴Note that this is the geometric condition number at the output rather than the input of $\Phi_{\mathbf{x}}^{-1}$. The reason is that the PIP can have several $\mathbf{x}_i \in \mathcal{M}$ as solutions. Since the RGNs will only output one of these solutions, say \mathbf{x}_1 , the natural question is whether this computed solution \mathbf{x}_1 is stable to perturbations of $\Phi(\mathbf{x}_1)$.

Lemma 1. *Let R be a retraction. Then for all $x \in \mathcal{M}$ there exists some $\delta_x > 0$ such that for all $\eta \in \mathbb{T}_x \mathcal{M}$ with $\|\eta\| < \delta_x$ one has $R_x(\eta) = x + \eta + \mathcal{O}(\|\eta\|^2)$.*

As stated in [2, Section 8.4.1], the RGN method for minimizing f is obtained by replacing the $\nabla_{x_k}^2 f$ in the Newton process (4) by the Gauss–Newton approximation $(d_{x_k} F)^* \circ (d_{x_k} F)$; herein A^* denotes the adjoint of the bounded linear operator A with respect to the inner product $\langle \cdot, \cdot \rangle$. Note that an explicit expression for the update direction η_k can be obtained. The Riemannian gradient is

$$\nabla_{x_k} f = \nabla_{x_k} \frac{1}{2} \langle F(x), F(x) \rangle = (d_{x_k} F)^*(F(x_k)); \quad (5)$$

see [2, Section 8.4.1]. If $d_{x_k} F$ is injective, then the solution of the system in (4) with the Riemannian Hessian replaced by the Gauss–Newton approximation is given explicitly by

$$\eta_k = -((d_{x_k} F)^* \circ (d_{x_k} F))^{-1} (d_{x_k} F)^*(F(x_k)) =: -(d_{x_k} F)^\dagger(F(x_k)).$$

3. Main result: Convergence analysis of the RGN method

We prove in this section that both the convergence rate and radius of the RGN method are influenced by the condition number of the PIP at the local minimizer. In the case of PIPs, we have $F(\mathbf{x}) := \Phi(\mathbf{x}) - \tilde{\mathbf{y}}$ for some fixed $\tilde{\mathbf{y}} \in \mathbb{R}^N$. Hence, $d_{\mathbf{x}} F = d_{\mathbf{x}} \Phi$, so that the next theorem relates the geometric condition number (2) to the convergence properties of the RGN method for solving the least-squares problem (3).

Remark 1. *The RGN method is only locally convergent. Practical methods are obtained by adding a globalization strategy [2, 12] such as a line search or trust region scheme. The goal of these strategies is guaranteeing sufficient descent for global convergence, while preserving the local rate of convergence. In the main theorem, we present the analysis without globalization strategy, so as to focus on the main idea of the proof. In case of a trust region scheme, the usual approach for extending the proof consists of showing that close to a local minimizer, the unconstrained Newton step is always contained in the trust region and hence selected. This will be true if the starting point is sufficiently close to the local minimizer. Then, the local rate of convergence will be the same as when no trust region scheme is employed.*

In the remainder of this section, let $B_\tau(\mathbf{x})$ denote the ball of radius τ centered at $\mathbf{x} \in \mathbb{R}^M$. The following is the main theorem of this letter.

Theorem 1. *Assume that $\mathbf{x}_* \in \mathcal{M}$ is a local minimum of the objective function f from (3), where $d_{\mathbf{x}_*} F$ is injective. Let $\kappa := (\varsigma_m(d_{\mathbf{x}_*} F))^{-1} > 0$. Then, there exists $\epsilon' > 0$ such that for all $0 < \alpha < 1$ there exists a universal constant $c > 0$ depending on ϵ' , F , \mathbf{x}_* , \mathcal{M} , and R so that the following holds.*

1. (Linear convergence): *If $\frac{c\kappa^2 \|F(\mathbf{x}_*)\|}{\alpha} < 1$, then for all $\mathbf{x}_0 \in B_\epsilon(\mathbf{x}_*) \cap \mathcal{M}$ with*

$$\epsilon := \min \left\{ \frac{1 - \alpha}{c\kappa}, \frac{\alpha\epsilon'}{1 + \alpha + c\kappa^2 \|F(\mathbf{x}_*)\|} \right\},$$

the RGN method generates a sequence $\mathbf{x}_0, \mathbf{x}_1, \dots$ that converges linearly to \mathbf{x}_ . In fact,*

$$\|\mathbf{x}_* - \mathbf{x}_{k+1}\| \leq \frac{c\kappa^2 \|F(\mathbf{x}_*)\|}{\alpha} \|\mathbf{x}_* - \mathbf{x}_k\| + \mathcal{O}(\|\mathbf{x}_* - \mathbf{x}_k\|^2).$$

2. (Quadratic convergence): *If \mathbf{x}_* is a zero of the objective function f , then for all $\mathbf{x}_0 \in B_\epsilon(\mathbf{x}_*) \cap \mathcal{M}$ with*

$$\epsilon := \min \left\{ \frac{1 - \alpha}{c\kappa}, \frac{\alpha\epsilon'}{1 + \alpha} \right\},$$

the RGN method generates a sequence $\mathbf{x}_0, \mathbf{x}_1, \dots$ that converges quadratically to \mathbf{x}_ . In fact,*

$$\|\mathbf{x}_* - \mathbf{x}_{k+1}\| \leq \frac{c(\kappa + 1)}{\alpha} \|\mathbf{x}_* - \mathbf{x}_k\|^2 + \mathcal{O}(\|\mathbf{x}_* - \mathbf{x}_k\|^3).$$

Remark 2. The order of convergence may also be established from [2, Theorem 8.2.1]. However, intrinsic multiplicative constants are not derived there, as their analysis is founded on coordinate expressions that depend on the chosen chart; they thus only derive chart-dependent multiplicative constants.

In the following let P_A denote the orthogonal projection onto the linear subspace $A \subset \mathbb{R}^N$. Recall the next lemma from [3, Section 2], which we need in the proof of Theorem 1.

Lemma 2. Let $F : \mathcal{M} \rightarrow \mathbb{R}^N$ be a smooth function and $\mathbf{x} \in \mathcal{M}$. Then, there exist constants $r_F > 0$ and $\gamma_F \geq 0$ such that for all $\mathbf{y} \in B_{r_F}(\mathbf{x}) \cap \mathcal{M}$ we have $F(\mathbf{y}) = F(\mathbf{x}) + (d_{\mathbf{x}}F)P_{T_{\mathbf{x}}\mathcal{M}}\Delta + \mathbf{v}_{\mathbf{x},\mathbf{y}}$, where $\Delta = (\mathbf{y} - \mathbf{x}) \in \mathbb{R}^N$ and $\|\mathbf{v}_{\mathbf{x},\mathbf{y}}\| \leq \gamma_F\|\Delta\|^2$.

We can now prove Theorem 1.

Proof of Theorem 1. We begin with some general considerations: In Lemma 1 we choose δ small enough such that it applies to all $\mathbf{x} \in B_\delta(\mathbf{x}_*) \cap \mathcal{M}$. Let $0 < \epsilon' \leq \delta$. Then, there exists a constant $\gamma_R > 0$ depending on the retraction operator R , such that for all $\mathbf{x} \in B_{\epsilon'}(\mathbf{x}_*) \cap \mathcal{M}$ we have

$$\|R_{\mathbf{x}}(\eta) - (\mathbf{x} + \eta)\| \leq \gamma_R\|\eta\|^2 \text{ for every } \eta \in B_{\epsilon'}(0) \cap T_{\mathbf{x}}\mathcal{M}. \quad (6)$$

By applying Lemma 2 to the smooth functions F and $\text{id}_{\mathcal{M}}$ respectively and using the smoothness of (the derivative of) F , we see that there exist constants $\gamma_F, \gamma_I > 0$ so that for all $\mathbf{x} \in B_{\epsilon'}(\mathbf{x}_*) \cap \mathcal{M}$ we have

$$F(\mathbf{x}_*) - F(\mathbf{x}) - (d_{\mathbf{x}}F)P_{T_{\mathbf{x}}\mathcal{M}}(\mathbf{x}_* - \mathbf{x}) = \mathbf{v} \text{ with } \|\mathbf{v}\| \leq \gamma_F\|\mathbf{x}_* - \mathbf{x}\|^2, \text{ and} \quad (7)$$

$$\mathbf{x}_* - \mathbf{x} - P_{T_{\mathbf{x}}\mathcal{M}}(\mathbf{x}_* - \mathbf{x}) = \mathbf{w} \text{ with } \|\mathbf{w}\| \leq \gamma_I\|\mathbf{x}_* - \mathbf{x}\|^2. \quad (8)$$

Moreover, we define the Lipschitz constant

$$C := \max_{\mathbf{x} \in B_{\epsilon'}(\mathbf{x}_*) \cap \mathcal{M}} \frac{\|d_{\mathbf{x}_*}F \circ P_{T_{\mathbf{x}_*}\mathcal{M}} - d_{\mathbf{x}}F \circ P_{T_{\mathbf{x}}\mathcal{M}}\|}{\|\mathbf{x}_* - \mathbf{x}\|}. \quad (9)$$

We choose a constant c , depending on R, F and ϵ' , that satisfies

$$\epsilon \leq \min\left\{\frac{1}{\kappa\gamma_F}, \frac{1-\alpha}{C\kappa}, \left(1 + \frac{1 + \frac{1}{2}(1 + \sqrt{5})C\kappa^2\|F(\mathbf{x}_*)\|}{\alpha}\right)^{-1} \epsilon'\right\} \text{ and } c \geq \max\left\{\frac{1}{2}(1 + \sqrt{5})C, \gamma_F, \gamma_R, \gamma_I\right\}. \quad (10)$$

The rest of the proof is by induction. Suppose that the RGN method applied to starting point $\mathbf{x}_0 \in \mathcal{M}$ generated the sequence of points $\mathbf{x}_0, \dots, \mathbf{x}_k \in \mathcal{M}$. First, we show that $d_{\mathbf{x}_k}F$ is injective, so that the update direction $\eta = -(d_{\mathbf{x}_k}F)^\dagger F(\mathbf{x}_k)$, and, hence, \mathbf{x}_{k+1} is defined. Thereafter, we prove the asserted bounds on $\|\mathbf{x}_k - \mathbf{x}_{k+1}\|$. For avoiding subscripts, let $\mathbf{x} := \mathbf{x}_k \in \mathcal{M}$ and $\mathbf{y} := \mathbf{x}_{k+1} = R_{\mathbf{x}}(-(d_{\mathbf{x}}F)^\dagger F(\mathbf{x})) \in \mathcal{M}$.

By induction, we can assume $\|\mathbf{x}_k - \mathbf{x}_*\| \leq \|\mathbf{x}_0 - \mathbf{x}_*\|$; indeed the base case $k = 0$ is trivially true, and we will prove that it also holds for $k + 1$ at the end of this proof, completing the induction. Hence, $\mathbf{x} \in B_{\epsilon'}(\mathbf{x}_*) \cap \mathcal{M}$. Let $J \in \mathbb{R}^{N \times M}$ denote the matrix of $d_{\mathbf{x}}F$ with respect to the standard bases on \mathbb{R}^M and \mathbb{R}^N . Let $J = U\Sigma V^T$ be its (compact) singular value decomposition (SVD), where $U \in \mathbb{R}^{N \times m}$ and $V \in \mathbb{R}^{M \times m}$ have orthonormal columns, the columns of V span $T_{\mathbf{x}}\mathcal{M}$ and $\Sigma \in \mathbb{R}^{m \times m}$ is diagonal matrix containing the singular values. Then, the matrix of $(d_{\mathbf{x}}F)^\dagger$ with respect to the standard bases is J^\dagger , i.e., the Moore–Penrose pseudoinverse of J , and $\kappa(\mathbf{x}) = \|J^\dagger\|$. Similarly, let $J_* \in \mathbb{R}^{N \times M}$ denote the matrix of $d_{\mathbf{x}_*}F$, and let $U_*\Sigma_*V_*^T$ be its SVD.

By assumption, $d_{\mathbf{x}_*}F$ is injective and thus we have $\kappa^{-1} = \varsigma_{\min}(d_{\mathbf{x}_*}F) = \varsigma_{\min}(J_*) > 0$. The matrix of $P_{T_{\mathbf{x}}\mathcal{M}}$ is VV^T , and similarly for $P_{T_{\mathbf{x}_*}\mathcal{M}}$. Then, by the definition of C in (9), we have

$$\|J_* - J\| = \|J_*(V_*V_*^T) - J(VV^T)\| \leq C\|\mathbf{x}_* - \mathbf{x}\|, \quad (11)$$

and hence $\|J_* - J\| \leq C\epsilon \leq C\frac{1-\alpha}{c\kappa} \leq (1-\alpha)\varsigma_{\min}(J_*)$, because $\mathbf{x} \in B_\epsilon(\mathbf{x}_*)$ and the definition of c . From Weyl's perturbation lemma it follows that $|\varsigma_{\min}(J_*) - \varsigma_{\min}(J)| \leq \|J_* - J\| \leq (1-\alpha)\varsigma_{\min}(J_*)$. We obtain $\varsigma_{\min}(J) > \alpha\varsigma_{\min}(J_*) > 0$, where the last inequality is by the assumption $\alpha > 0$. It follows that

$$\|J^\dagger\| = (\varsigma_{\min}(J))^{-1} < \kappa\alpha^{-1} < \infty, \quad (12)$$

so that $d_{\mathbf{x}}F$ is indeed injective. This shows that the RGN update direction η is well defined.

It remains to prove the bound on $\|\mathbf{x}_* - \mathbf{y}\|$. First we show that $\|\eta\| = \|-J^\dagger F(\mathbf{x})\| \leq \epsilon' < \delta$, so that the retraction would satisfy (6). By assumption \mathbf{x}_* is a local minimum of (3), so that from (5) we obtain $0 = \nabla_{\mathbf{x}_*} f = J_*^T F(\mathbf{x}_*) = V_* \Sigma_* U_*^T F(\mathbf{x}_*)$. By [13, Chapter III, Theorem 1.2 (9)] and the assumption that $d_{\mathbf{x}_*} F$ is injective, we have $J_*^\dagger = V_* \Sigma_*^{-1} U_*^T$ from which we conclude $J_*^\dagger F(\mathbf{x}_*) = 0$. Let $P = P_{T_{\mathbf{x}} \mathcal{M}}$. From (7),

$$J^\dagger F(\mathbf{x}) = J^\dagger F(\mathbf{x}_*) - P(\mathbf{x}_* - \mathbf{x}) - J^\dagger \mathbf{v}, \quad (13)$$

so that

$$\|\eta\| = \|- (J^\dagger - J_*^\dagger) F(\mathbf{x}_*) + P(\mathbf{x}_* - \mathbf{x}) + J^\dagger \mathbf{v}\| \leq \|J^\dagger - J_*^\dagger\| \|F(\mathbf{x}_*)\| + \|P\| \|\mathbf{x}_* - \mathbf{x}\| + \|J^\dagger\| \|\mathbf{v}\|. \quad (14)$$

From Wedin's theorem [13, Chapter III, Theorem 3.9] we obtain

$$\|J^\dagger - J_*^\dagger\| \leq \frac{1 + \sqrt{5}}{2} \|J^\dagger\| \|J_*^\dagger\| \|J - J_*\| \leq \frac{(1 + \sqrt{5}) C \kappa^2}{2\alpha} \|\mathbf{x}_* - \mathbf{x}\|, \quad (15)$$

where the last step is because of (11) and (12). Using $\|P\| = 1$ for orthogonal projectors, the assumption $\|\mathbf{x}_* - \mathbf{x}\| \leq (\kappa \gamma_F)^{-1}$, (12), (15), and the bound on $\|\mathbf{v}\|$ in (7), it follows from (14) that

$$\|\eta\| \leq \left(1 + \frac{\kappa \gamma_F \|\mathbf{x}_* - \mathbf{x}\| + \frac{1}{2}(1 + \sqrt{5}) C \kappa^2 \|F(\mathbf{x}_*)\|}{\alpha}\right) \|\mathbf{x}_* - \mathbf{x}\|. \quad (16)$$

By the definition of ϵ and the assumption $\|\mathbf{x}_* - \mathbf{x}\| < \epsilon$, we have $\|\mathbf{x}_* - \mathbf{x}\| < \epsilon < \frac{1}{\kappa \gamma_F}$, so that by (16),

$$\|\eta\| \leq \left(1 + \frac{1 + \frac{1}{2}(1 + \sqrt{5}) C \kappa^2 \|F(\mathbf{x}_*)\|}{\alpha}\right) \|\mathbf{x}_* - \mathbf{x}\|. \quad (17)$$

Using the third bound on ϵ in (10), we have

$$\|\mathbf{x}_* - \mathbf{x}\| < \epsilon < \left(1 + \frac{1 + \frac{1}{2}(1 + \sqrt{5}) C \kappa^2 \|F(\mathbf{x}_*)\|}{\alpha}\right)^{-1} \epsilon',$$

which when plugged into (17) yields $\|\eta\| < \epsilon'$.

From the foregoing discussion, we conclude that (6) applies to $R_{\mathbf{x}}(\eta) = R_{\mathbf{x}}(-J^\dagger F(\mathbf{x}))$, so that

$$\|\mathbf{y} - \mathbf{x}_*\| = \|R_{\mathbf{x}}(-J^\dagger F(\mathbf{x})) - \mathbf{x}_*\| \leq \|\mathbf{x} - J^\dagger F(\mathbf{x}) - \mathbf{x}_*\| + \gamma_R \|\eta\|^2. \quad (18)$$

Let $\zeta := \|\mathbf{x} - J^\dagger F(\mathbf{x}) - \mathbf{x}_*\|$. We use $J_*^\dagger F(\mathbf{x}_*) = 0$ and the formula from (13) to derive that

$$\begin{aligned} \zeta &= \|\mathbf{x} - \mathbf{x}_* - (J^\dagger F(\mathbf{x}) - J_*^\dagger F(\mathbf{x}_*))\| = \|\mathbf{x} - \mathbf{x}_* - (J^\dagger - J_*^\dagger) F(\mathbf{x}_*) + P(\mathbf{x}_* - \mathbf{x}) + J^\dagger \mathbf{v}\| \\ &= \|- (J^\dagger - J_*^\dagger) F(\mathbf{x}_*) - \mathbf{w} + J^\dagger \mathbf{v}\| \\ &\leq \|J^\dagger - J_*^\dagger\| \|F(\mathbf{x}_*)\| + \gamma_I \|\mathbf{x}_* - \mathbf{x}\|^2 + \|J^\dagger\| \gamma_F \|\mathbf{x}_* - \mathbf{x}\|^2, \end{aligned}$$

where the second-to-last equality is due to (8), and in the last line we have used the triangle inequality and the bounds on $\|\mathbf{v}\|$ and $\|\mathbf{w}\|$ from (7) and (8). Combining this with (15) and (12) yields

$$\zeta \leq \frac{\frac{1}{2}(1 + \sqrt{5}) C \kappa^2 \|F(\mathbf{x}_*)\|}{\alpha} \|\mathbf{x}_* - \mathbf{x}\| + \left(\gamma_I + \frac{\gamma_F \kappa}{\alpha}\right) \|\mathbf{x}_* - \mathbf{x}\|^2, \quad (19)$$

Note that we have chosen the constant c large enough, so that $\frac{1}{2}(1 + \sqrt{5}) C < c$. Plugging (19) and (17) into (18) yields the first bound.

For the second assertion we have the additional assumption that \mathbf{x}_* is a zero of the objective function $f(\mathbf{x}) = \frac{1}{2} \|F(\mathbf{x})\|^2$. From (19) we obtain $\zeta \leq \left(\frac{\kappa \gamma_F}{\alpha} + \gamma_I\right) \|\mathbf{x}_* - \mathbf{x}\|^2$. From (14) we get $\|\eta\| = \|P(\mathbf{x}_* - \mathbf{x}) + J^\dagger \mathbf{v}\|$ so that we can bound $\|\eta\|^2$ by

$$\|P(\mathbf{x}_* - \mathbf{x})\|^2 + 2|\langle P(\mathbf{x}_* - \mathbf{x}), J^\dagger \mathbf{v} \rangle| + \|J^\dagger \mathbf{v}\|^2 \leq \|\mathbf{x}_* - \mathbf{x}\|^2 + 2\gamma_F \|J^\dagger\| \|\mathbf{x}_* - \mathbf{x}\|^3 + \gamma_F^2 \|J^\dagger\|^2 \|\mathbf{x}_* - \mathbf{x}\|^4,$$

where the inequality is by the Cauchy-Schwartz inequality and the fact that $\|P\| = 1$ for orthogonal projectors. As before, plugging these bounds for ζ and $\|\eta\|$ into (18) and exploiting that $c \geq \max\{\gamma_F, \gamma_I, \gamma_R\}$, the second bound is obtained. \square

A reviewer asked how critical the injectivity assumption on the derivative $d_{\mathbf{x}}F$ in the above theorem is. The brief answer is that it is usually a very weak assumption in practice. First, we need a lemma.

Lemma 3. *Let $\mathcal{M} \subset \mathbb{R}^M$ be an embedded manifold whose projectivization is a smooth projective variety, and let $\Phi : \mathcal{M} \rightarrow \mathbb{R}^N$ be a regular map. Let \mathcal{N} denote the \mathbb{R} -variety that is the Zariski closure of the image $\Phi(\mathcal{M})$. If the dimensions satisfy $\dim \mathcal{M} = \dim \mathcal{N}$, then the locus of points $\mathcal{G} := \{\mathbf{x} \in \mathcal{M} \mid d_{\mathbf{x}}\Phi \text{ is injective and } \Phi(\mathbf{x}) \text{ is a smooth point of } \mathcal{N}\}$ is a dense subset of \mathcal{M} in the Euclidean topology.*

Proof. This is essentially a restatement of [8, Theorem 11.12]. \square

The following proposition shows, under the assumptions of Lemma 3, that the local optimizer \mathbf{x}_* in Theorem 1 has an injective derivative $d_{\mathbf{x}_*}F = d_{\mathbf{x}_*}\Phi$ on a set of inputs (in \mathbb{R}^N) of positive measure.

Proposition 2. *Let \mathcal{M} , \mathcal{G} , \mathcal{N} , and Φ be as in Lemma 3. Assume that we have the equality $\dim \mathcal{M} = \dim \mathcal{N}$. Let $\mathcal{B} := \{\mathbf{y} \in \mathbb{R}^N \mid \arg \min_{\mathbf{x} \in \mathcal{M}} \|\Phi(\mathbf{x}) - \mathbf{y}\| \subset \mathcal{G}\}$ be the set of points \mathbf{y} all of whose closest approximations on $\Phi(\mathcal{M})$, i.e., $C_{\mathbf{y}} := \arg \min_{\mathbf{x} \in \mathcal{M}} \|\Phi(\mathbf{x}) - \mathbf{y}\|$, lie in \mathcal{G} . Then, \mathcal{B} has positive Lebesgue measure, i.e., it is open in the Euclidean topology. Moreover, $\mathcal{N} \cap \mathcal{B}$ is Euclidean dense in $\Phi(\mathcal{M})$.*

Proof. Let $\mathcal{W}_1 \subset \mathcal{M}$ be the locus where the dimension of the fiber $\Phi^{-1}(\Phi(\mathbf{x}))$ is strictly positive, i.e., $\mathcal{W}_1 = \{\mathbf{x} \in \mathcal{M} \mid \dim \Phi^{-1}(\Phi(\mathbf{x})) > 0\}$. By [8, Theorem 11.12] \mathcal{W}_1 is a Zariski-closed set. The last claim of the proposition is also a corollary of this theorem and the assumption that the generic fiber is 0-dimensional.

It remains to show the first claim. Let $\mathcal{W}_2 \subset \mathcal{M}$ be the Zariski-closed subset of points $\mathbf{x} \in \mathcal{M}$ for which $\Phi(\mathbf{x})$ lies in the singular locus of \mathcal{N} . Set $\mathcal{W} := (\mathcal{W}_1 \cup \overline{\mathcal{W}_2}) \subset \mathcal{M}$, where the overline denotes the closure in the Zariski topology. Note that $\mathcal{M} \setminus \mathcal{W} \subset \mathcal{G}$.

Let $\mathbf{x} \notin \mathcal{W}$. Since the derivative $d_{\mathbf{x}}\Phi$ is injective, there exists a local diffeomorphism between an open neighborhood $\mathcal{M}_0 \subset \mathcal{M}$ of \mathbf{x} and an open neighborhood $\mathcal{N}_0 \subset \mathcal{N}$ of $\Phi(\mathbf{x})$. By restricting neighborhoods, we can assume that the Euclidean closure of \mathcal{N}_0 is contained in the smooth locus of \mathcal{N} and that \mathcal{M}_0 is contained in $\mathcal{M} \setminus \mathcal{W}$. Take a tubular neighborhood \mathcal{T} of $\mathcal{N}_0 \subset \mathbb{R}^N$ that does not intersect $\mathcal{N} \setminus \mathcal{N}_0$, and let h be its height; note that $h > 0$, because the closure of \mathcal{N}_0 does not contain singular points of \mathcal{N} . Then, there exists an open ball $B_\delta(\Phi(\mathbf{x}))$ in \mathbb{R}^N of positive radius $0 < \delta < h$, centered at $\Phi(\mathbf{x})$, whose intersection with \mathcal{N} is contained in \mathcal{N}_0 . By construction $B_\delta(\Phi(\mathbf{x})) \subset (\mathcal{N}_0 \cup \mathcal{T})$. It follows from the triangle inequality that the closest point on \mathcal{N} to any point of $\mathcal{B}_{\mathbf{x}} := B_{\delta/2}(\Phi(\mathbf{x}))$ is contained in $\mathcal{N}_0 \subset \Phi(\mathcal{M}) \subset \mathcal{N}$. Since $\mathcal{N}_0 = \Phi(\mathcal{M}_0)$ and because $\mathcal{M}_0 \subset \mathcal{G}$, it follows that $\mathcal{B}_{\mathbf{x}} \subset \mathcal{B}$ for all $\mathbf{x} \in \mathcal{M} \setminus \mathcal{W}$. \square

4. Numerical experiments

Here we experimentally verify the dependence of the multiplicative constant on the geometric condition number for a special case of PIP (1), namely the tensor rank decomposition (TRD) problem. The model is

$$\Phi : \mathcal{S} \times \cdots \times \mathcal{S} \rightarrow \mathbb{R}^N, (\mathbf{a}_1^1 \otimes \cdots \otimes \mathbf{a}_1^d, \dots, \mathbf{a}_r^1 \otimes \cdots \otimes \mathbf{a}_r^d) \mapsto \sum_{i=1}^r \mathbf{a}_i^1 \otimes \cdots \otimes \mathbf{a}_i^d,$$

where $d \geq 3$, $N = m_1 \cdots m_d$, and $\mathcal{S} \subset \mathbb{R}^N$ is the manifold of $m_1 \times \cdots \times m_d$ rank-1 tensors [8]. The image of Φ is called a *join set* and the PIP is a special case of the *join decomposition problem* [3]. To put emphasis on the join structure of the image of Φ , we denote $\mathcal{J} := \Phi(\mathcal{S} \times \cdots \times \mathcal{S})$.

In the numerical experiments of this section we apply a RGN method to $\min_{\mathbf{x} \in \mathcal{M}} \frac{1}{2} \|\Phi(\mathbf{x}) - \mathfrak{A}\|^2$, where $\mathcal{M} := \mathcal{S} \times \cdots \times \mathcal{S}$ is the r -fold product manifold of \mathcal{S} , and $\mathfrak{A} \in \mathbb{R}^N$ is the given tensor to approximate. We choose the retraction operator $R : \mathcal{T}\mathcal{M} \rightarrow \mathcal{M}$ from [4].

The projectivization of the manifold \mathcal{S} is called the Segre variety; it is a smooth, irreducible projective variety with affine dimension $\dim \mathcal{S} = 1 + \sum_{k=1}^d (m_k - 1)$. The problem of computing the dimension of the Zariski-closure $\overline{\mathcal{J}}$, which is called the r -secant variety of \mathcal{S} , has been classically studied; see [10, Section 5.5] for an overview. The results of [6] entail that the dimension equality $\dim \mathcal{M} = \dim \overline{\mathcal{J}}$ is satisfied for all $r \cdot \dim \mathcal{S} < N$ and $N \leq 15000$, subject to a few theoretically characterized exceptions. In the example below,

we take $r = 2$ for which the dimension equality is always satisfied [1]. Hence, Proposition 2 entails that the injectivity assumption in Theorem 1 is satisfied at least on a set of positive Lebesgue measure. Therefore, the convergence rate of the RGN method is influenced by the geometric condition number of the optimal parameters $\mathbf{x}_* \in \mathcal{M}$ that minimizes the objective function.

We showed in [3, Section 5.1] that the condition number of the above PIP at $\mathbf{x} = (\mathbf{a}_i^1 \otimes \cdots \otimes \mathbf{a}_i^d)_{i=1}^r$ is $\kappa(\mathbf{x}) = (\varsigma_m(U))^{-1}$, where $m = r \cdot \dim \mathcal{S}$, and the matrix $U \in \mathbb{R}^{N \times m}$ is given by $U = [U_1 \cdots U_r]$ with

$$U_i := \left[\frac{\mathbf{a}_i^1}{\|\mathbf{a}_i^1\|} \otimes \cdots \otimes \frac{\mathbf{a}_i^d}{\|\mathbf{a}_i^d\|} \quad Q_{1,i} \otimes \frac{\mathbf{a}_i^2}{\|\mathbf{a}_i^2\|} \otimes \cdots \otimes \frac{\mathbf{a}_i^d}{\|\mathbf{a}_i^d\|} \quad \cdots \quad \frac{\mathbf{a}_i^1}{\|\mathbf{a}_i^1\|} \otimes \cdots \otimes \frac{\mathbf{a}_i^{d-1}}{\|\mathbf{a}_i^{d-1}\|} \otimes Q_{d,i} \right],$$

where $Q_{k,i} \in \mathbb{R}^{m_k \times (m_k - 1)}$ is a matrix containing an orthonormal basis of the orthogonal complement of \mathbf{a}_i^k in \mathbb{R}^{m_k} . These expressions allow us to compute the condition number at any given decomposition $\mathbf{x} \in \mathcal{M}$.

All of the following computations were performed in Matlab R2016b. For clearly illustrating the rates of convergence, we used variable precision arithmetic (vpa) with 400 digits of accuracy. Since performing experiments in vpa is very expensive, we consider only the tiny example of a rank-2 tensor of size $3 \times 3 \times 3$. We showed in [4] that an implementation of the RGN method with trust region globalization strategy applied to the above PIP formulation, can outperform state-of-the-art optimization methods for the tensor rank approximation problem on small-scale, dense problems with $r \sum_{k=1}^d m_k \lesssim 1000$.

4.1. Experiment 1: Random perturbations

Consider the following parametrized tensors in $\mathbb{R}^3 \otimes \mathbb{R}^3 \otimes \mathbb{R}^3$. For $s \geq 0$ we let $\mathbf{x}(s) = (x(s), \mathbf{e}_2^{\otimes 3}) \in \mathcal{S} \times \mathcal{S}$ where $x(0) := \mathbf{e}_1^{\otimes 3}$ and $x(s) := (\mathbf{e}_2 - 2^{-s} \mathbf{e}_1)^{\otimes 3}$ for $s > 0$ and $\mathbf{e}_k \in \mathbb{R}^3$ is the k th standard basis vector. Then, we define $\mathfrak{A}(s) := \Phi(\mathbf{x}(s)) = x(s) + \mathbf{e}_2^{\otimes 3}$.

For every $s = 0, 1, 3, 5$, we created a perturbed decomposition $\mathbf{x}'(s) = R(\mathbf{x}(s), 10^{-20} \cdot \frac{\mathfrak{X}}{\|\mathfrak{X}\|})$, where R is the aforementioned retraction and the entries of \mathfrak{X} are chosen from the standard normal distribution. We also sampled a perturbed tensor $\mathfrak{A}'(s) := \mathfrak{A}(s) + 10^{-10} \frac{\mathfrak{Z}}{\|\mathfrak{Z}\|}$, where the entries of \mathfrak{Z} are also standard normal.

For verifying the linear convergence, the RGN method was applied to $\mathfrak{A}'(s)$ while the quadratic convergence was checked by applying the RGN method to $\mathfrak{A}(s)$, both starting from $\mathbf{x}'(s)$. In all tested cases, the RGN method generated a sequence $\mathbf{x}_1(s), \mathbf{x}_2(s), \dots$ in \mathcal{M} converging to a local minimizer $\mathbf{x}_*(s) \in \mathcal{M}$. The residual $\|F(\mathbf{x}_*(s))\|$ was approximately $7 \cdot 10^{-11}$ in all cases.

The results are shown in Figures 1(a) and 1(c), illustrating respectively the predicted linear and quadratic convergence. The graphs confirm the prime message of this letter: *the convergence speed of the RGN method deteriorates when the geometric condition number increases*, as Theorem 1 predicts.

As the full lines in Figure 1(a) show, the multiplicative constants derived in Theorem 1 can be pessimistic, especially when the condition number is large. We attribute this to bound (15); while it is sharp [13, p. 152], it is very pessimistic in this experiment. A qualitatively better description of the convergence is shown as the dashed lines in Figure 1(a), where the constant in (15) was estimated heuristically as $E(s) := \frac{\|J^\dagger - J_*^\dagger\|}{\|\mathbf{x}_1(s) - \mathbf{x}_*(s)\|}$, where J is the matrix of $d_{\mathbf{x}_1(s)} F$ and J_* is the matrix of $d_{\mathbf{x}_*(s)} F$ as in the proof of Theorem 1.

4.2. Experiment 2: Adversarial perturbations.

For illustrating the sharpness of the bound (15), we performed an additional experiment with tensors in $\mathbb{R}^3 \otimes \mathbb{R}^3 \otimes \mathbb{R}^3 \cong \mathbb{R}^{27}$. This time we constructed an adversarially perturbed starting point $\mathbf{x}'(s)$ by generating a random tensor $\mathfrak{N} \in \mathbb{R}^{27}$ with entries sampled from the standard normal distribution, then computing numerically the gradient \mathbf{g} of the function $f(\mathbf{x}) = \frac{1}{2} \|((d_{\mathbf{x}(s)} \Phi)^\dagger - (d_{\mathbf{x}} \Phi)^\dagger) \mathfrak{N}\|^2$, and finally setting $\mathbf{x}'(s) = R(\mathbf{x}(s), 10^{-20} \mathbf{g})$. As adversarial perturbation of $\mathfrak{A}(s) = \Phi(\mathbf{x}(s))$, we chose $\mathfrak{Z} \in \mathbb{R}^{27}$ equal to the left singular vector $\mathbf{u}_{14} \in \mathbb{R}^{27}$ corresponding to the smallest nonzero singular value ς_{14} of $d_{\mathbf{x}'(s)} \Phi$; note that $\dim \mathcal{M} = r \cdot \dim \mathcal{S} = 2(1 + 2 + 2 + 2) = 14$. As before, we set $\mathfrak{A}'(s) = \mathfrak{A}(s) + 10^{-10} \frac{\mathfrak{Z}}{\|\mathfrak{Z}\|}$.

The result of applying the RGN method to $\mathfrak{A}'(s)$ from starting point $\mathbf{x}'(s)$ is shown in Figure 1(b). In all cases, the method converged. The condition numbers at the local minimizers $\mathbf{x}_*(s)$ are about the same as in the previous experiment: the respective relative differences were less than 10^{-2} . The final residuals $\|F(\mathbf{x}_*(s))\|$ depended on s , however; they were $6.90 \cdot 10^{-46}$, $8.60 \cdot 10^{-34}$, $3.58 \cdot 10^{-31}$ and $1.10 \cdot 10^{-26}$ for

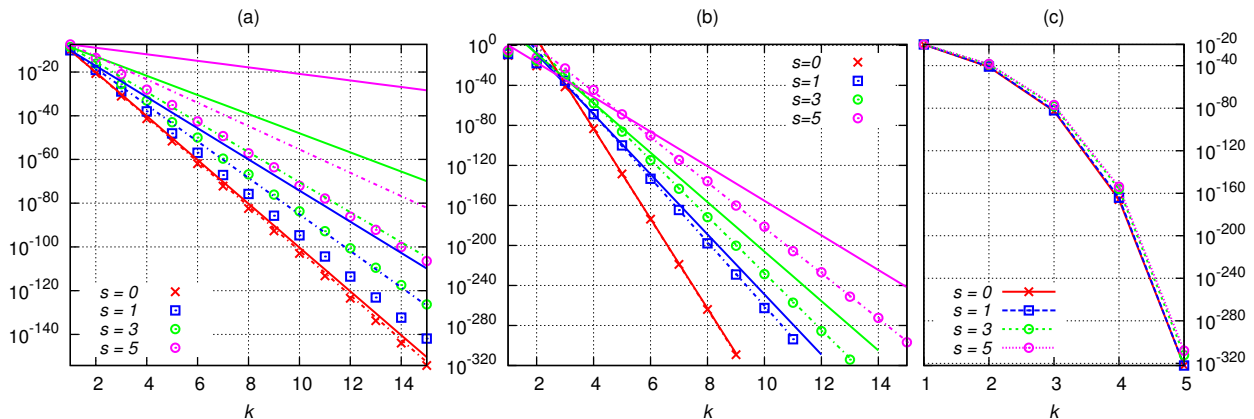


Figure 1: The data points show the distance $\|\mathbf{x}_k(s) - \mathbf{x}_*(s)\|$ for the sequence of points $\mathbf{x}_k(s)$, $k = 1, 2, \dots$, computed by the RGN method in function of k for $s = 0, 1, 3, 5$. The condition numbers of the local optimizers, rounded to two significant digits, were $\kappa = 1.0 \cdot 10^0$, $2.7 \cdot 10^1$, $1.5 \cdot 10^3$, and $9.3 \cdot 10^4$ for respectively $s = 0, 1, 3$, and 5 . In (a) and (b) linear convergence rates are illustrated when $\|F(\mathbf{x}_*(s))\| \neq 0$, and (c) shows the quadratic convergence rate when the residual $\|F(\mathbf{x}_*)\|$ vanishes. Figures (a) and (c) show instances of random perturbations, while (b) employed adversarial perturbations. In figure (b), the linear rate of convergence is not immediately observed; therefore the theoretical estimates were applied starting from the least k where $\|\mathbf{x}_k(s) - \mathbf{x}_*(s)\| \leq 10^{-50}$. In figures (a) and (b), the full lines (—, —, —, —) indicate the theoretical upper bounds from Theorem 1, i.e., $\frac{(1+\sqrt{5})C\kappa^2}{2\alpha}$ from (15). The dashed lines (---, ---, ---, ---) indicate the upper bounds obtained from Theorem 1, where the constants on the right-hand sides of (15) are *estimated heuristically* as $E(s)$.

respectively $s = 0, 1, 3$, and 5 . This is why the convergence may appear at first sight to be better than in the case of random perturbations. Nevertheless, it is observed that the theoretical estimate in Theorem 1 is indeed much closer to the observed convergence. In fact, the bounds involving the heuristic estimate $E(s)$ are visually indistinguishable from the actual data. Note in particular for $s = 0$, where $\kappa = 1$, that also the theoretical convergence rate from Theorem 1 is visually indistinguishable from the data, illustrating the sharpness of the bound in (15).

Acknowledgements

We thank two anonymous reviewers for their insightful and critical remarks that improved this letter.

References

- [1] Abo, H., Ottaviani, G., Peterson, C., 2009. Induction for secant varieties of Segre varieties. *Trans. Amer. Math. Soc.* 361, 767–792.
- [2] Absil, P.-A., Mahony, R., Sepulchre, R., 2008. *Optimization Algorithms on Matrix Manifolds*. Princeton University Press.
- [3] Breiding, P., Vannieuwenhoven, N., 2017. The condition number of join decompositions. arXiv:1611.08117. Submitted.
- [4] Breiding, P., Vannieuwenhoven, N., 2017. A Riemannian trust region method for the canonical tensor rank approximation problem. arXiv:1709.00033. Submitted.
- [5] Bürgisser, P., Cucker, F., 2013. *Condition: The Geometry of Numerical Algorithms*. Springer, Heidelberg.
- [6] Chiantini, L., Ottaviani, G., Vannieuwenhoven, N., 2014. An algorithm for generic and low-rank specific identifiability of complex tensors. *SIAM J. Matrix Anal. Appl.* 35 (4), 1265–1287.
- [7] Grasedyck, L., Kressner, D., Tobler, C., 2013. A literature survey of low-rank tensor approximation techniques. *GAMM Mitteilungen* 36 (1), 53–78.
- [8] Harris, J., 1992. *Algebraic Geometry, A First Course*. Vol. 133 of Graduate Text in Mathematics. Springer-Verlag.
- [9] Kressner, D., Steinlechner, M., Vandereycken, B., 2014. Low-rank tensor completion by Riemannian optimization. *BIT Numer. Math.* 54 (2), 447–468.
- [10] Landsberg, J. M., 2012. *Tensors: Geometry and Applications*. Vol. 128 of Graduate Studies in Mathematics. AMS, Providence, Rhode Island.
- [11] Lee, J. M., 2013. *Introduction to Smooth Manifolds*, 2nd Edition. Springer, New York, USA.
- [12] Nocedal, J., Wright, S. J., 2006. *Numerical Optimization*, 2nd Edition. Springer Series in Operation Research and Financial Engineering. Springer.
- [13] Stewart, G. W., Sun, J.-G., 1990. *Matrix Perturbation Theory*. Academic Press.