

Adaptive Localized Cayley Parametrization Technique for Smooth Optimization over the Stiefel Manifold

Keita Kume, Isao Yamada

Dept. of Information and Communications Engineering, Tokyo Institute of Technology

Email: {kume,isao}@sp.ce.titech.ac.jp

Abstract—We propose a novel computational strategy, named the adaptive localized Cayley parametrization technique for acceleration of optimization over the Stiefel manifold. The proposed optimization algorithm is designed as a gradient descent type scheme for the composite of the original cost function and the inverse of the localized Cayley transform defined on the vector space of all skew-symmetric matrices. Thanks to the adaptive localized Cayley transform which is a computable diffeomorphism between the orthogonal group and the vector space of the skew-symmetric matrices, the proposed algorithm (i) is free from the singularity issue, which can cause performance degradation, observed in the dual Cayley parametrization technique [Yamada-Ezaki'03] as well as (ii) can enjoy powerful arts for acceleration on the vector space without suffering from the nonlinear nature of the Stiefel manifold. We also present a convergence analysis, for the prototype algorithm employing the Armijo's rule, that shows the gradient of the composite function at zero in the range space of the localized Cayley transform is guaranteed to converge to zero. Numerical experiments show excellent performance compared with major optimization algorithms designed essentially with retractions on the tangent space of the Stiefel manifold [Absil-Mahony-Sepulcher'08, Wen-Yin'13].

Index Terms—Stiefel manifold optimization, orthogonal group optimization, Riemannian manifold optimization, Cayley transform, Anderson acceleration

I. INTRODUCTION

In this paper, we consider an optimization problem over the Stiefel manifold $St(p, N) = \{\mathbf{W} \in \mathbb{R}^{N \times p} \mid \mathbf{W}^T \mathbf{W} = \mathbf{I}\}$:

$$\text{minimize } f(\mathbf{U}) \text{ subject to } \mathbf{U} \in St(p, N), \quad (1)$$

where $f : \mathbb{R}^{N \times p} \rightarrow \mathbb{R}$ is differential. This problem has rich applications in data science including signal processing and machine learning, such as nearest low-rank correlation matrix problem [1], nonlinear eigenvalue problem [2], sparse principal component analysis [3], 1-bit compressed sensing [4], dimension reduction for ICA [5], face descriptor for recognition [6] and improvement of generalization for deep neural network [7]. However, the problem (1) is not simple at all due to the non-convexity of $St(p, N)$, which bears many computational ideas to deal with $St(p, N)$ for (1).

A major strategies [8]–[15] deal with $St(p, N)$ as a Riemannian manifold endowed with a Riemannian metric to build optimization schemes in its tangent bundle $TSt(p, N) := \bigcup_{\mathbf{U} \in St(p, N)} (\{\mathbf{U}\} \times T_{\mathbf{U}}St(p, N))$, where $T_{\mathbf{U}}St(p, N)$ stands for $Np - p(p+1)/2$ -dimensional tangent vector space whose zero corresponds to \mathbf{U} . Along these strategies, classical optimization algorithms, e.g., the steepest descent method, the Newton's method and the conjugate gradient method,

designed originally on a vector space, have been extended to those on the tangent spaces of $St(p, N)$. At each step, the extended algorithms try to update the current estimate $\mathbf{U} \in St(p, N)$ of the optimal point ideally along the geodesic on $St(p, N)$ in the direction \mathbf{D} determined on the tangent space. However the computation of the geodesic requires infinite calculations, its approximated map called *retraction* $R : TSt(p, N) \rightarrow St(p, N) : (\mathbf{U}, \mathbf{D}) \mapsto R_{\mathbf{U}}(\mathbf{D})$, which is inherited two natures from the geodesic: (i) $R_{\mathbf{U}}(\mathbf{0}) = \mathbf{U}$ and (ii) $\left. \frac{d}{dt} R_{\mathbf{U}}(t\mathbf{D}) \right|_{t=0} = \mathbf{D}$, is proposed, e.g., with QR decomposition/polar decomposition, the Euclidean projection [8], [13], [14], and the Cayley transform [9], [13], [16]. These strategies are realized by the following procedure: (i) Determine a search direction $\mathbf{D} \in T_{\mathbf{U}}St(p, N)$; (ii) Apply a retraction which assigns to the current estimate $\mathbf{U} \in St(p, N)$ a new estimate on $St(p, N)$ along the approximated geodesic in the search direction \mathbf{D} determined in (i). Although these strategies can also be extended to optimization problems over general manifolds, not necessary over the Stiefel manifold, they do not seem to exploit inherently algebraic properties of the Stiefel manifold, e.g., $St(p, N)$ can be seen as a canonical projection of the orthogonal group $O(N) := St(N, N)$. This situation suggests a possibility toward alternative powerful strategy for optimization over the Stiefel manifold if we find some ideas to exploit algebraic property of $O(N)$.

A dense subset of $O(N)$ can be parameterized in terms of skew-symmetric matrices by using a diffeomorphism known as the *Cayley transform* [17], [18] $\Phi : O(N) \setminus E \rightarrow Q(N) : \mathbf{U} \mapsto (\mathbf{I} - \mathbf{U})(\mathbf{I} + \mathbf{U})^{-1}$, and its inverse $\Phi^{-1} : Q(N) \rightarrow O(N) \setminus E : \mathbf{V} \mapsto (\mathbf{I} - \mathbf{V})(\mathbf{I} + \mathbf{V})^{-1}$, where $Q(N) := \{\mathbf{V} \in \mathbb{R}^{N \times N} \mid \mathbf{V}^T = -\mathbf{V}\}$ is the vector space of all skew-symmetric matrices and $E := \{\mathbf{U} \in O(N) \mid \det(\mathbf{I} + \mathbf{U}) = 0\}$ is the set of all singular points of Φ . Based on the fact that $Q(N)$ constitutes a real vector space, gradient descent type schemes including the steepest descent method and the Newton's method for the composite of f and Φ^{-1} defined on $Q(N)$, called the *dual Cayley parametrization technique*, were introduced in [17] as algorithms specialized for (1) in the case of $p = N$. The update of the schemes can be expressed by $\mathbf{U}_{n+1} = \Phi^{-1}(\mathcal{A}(\Phi(\mathbf{U}_n)))$, where \mathcal{A} denotes the update rule from $\mathbf{V}_n := \Phi(\mathbf{U}_n) \in Q(N)$ to $\mathbf{V}_{n+1} := \mathcal{A}(\mathbf{V}_n) := \Phi(\mathbf{U}_{n+1}) \in Q(N)$. Any computational technique developed for minimization of $f \circ \Phi^{-1}$ over vector spaces, e.g., the gradient descent type method and the Newton's method, can be applied as building blocks of the update rule \mathcal{A} . However their

numerical performances tend to slow down severely around E due to the appearance of undesirable plateau [17].

In this paper, to overcome completely the weakness caused by the set E , while preserving the advantage in [17], we newly propose the *adaptive localized Cayley parametrization technique* as an adaptive extension and acceleration of the dual Cayley parametrization technique [17]. The proposed adaptive localized Cayley parametrization is a time varying modification of [17] for adjustment of the Cayley type parametrization itself in order for the point sequence not to approach the modified singular set. Similar modifications of the parametrization are found, e.g., in [19], where a measure of proximity to a singular point of the *standard Cayley transform* was introduced to judge whether a certain modification, called "pivoting", of Φ is to be applied or not, and in [9], where Φ^{-1} is used as a retraction but their approach can also be interpreted as a time varying modification, at every update of the orthogonal matrix, of the parametrization in [17]. However so far the existing strategies do not seem to have succeeded yet in exploiting maximally the potential advantage hidden in the translation of (1), achieved essentially by the dual Cayley parametrization, into the simplified task of optimization over the vector space. Motivated by these situations, in this paper, we propose to combine a time varying modification of the Cayley parametrization and powerful acceleration techniques [20] including *Anderson acceleration* [21] (See II.C) which was developed originally for fixed-point approximation and is applicable to optimization tasks but only over vector spaces. We also present a convergence analysis, for the adaptive localized Cayley parametrization for a gradient descent type algorithm with the *Armijo's rule* [22], that shows the gradient sequence of the composite function at zero in the range of space of the localized Cayley transform is guaranteed to converge to zero. Numerical experiments in the scenarios of a joint diagonalization for symmetric matrices and an eigenbasis extraction demonstrate that the proposed adaptive extension with acceleration of [17] successfully avoids undesirable plateau and enjoy dramatical improvement in the speed of convergence, outperforming major optimization algorithms designed with retractions on the tangent space of the Stiefel manifold.

II. PRELIMINARIES

A. Notation

Let $\mathbf{I}_{N \times p} \in \mathbb{R}^{N \times p}$ stand for the first p -th columns of the identity matrix $\mathbf{I} \in \mathbb{R}^{N \times N}$. Let $SO(N) := \{\mathbf{U} \in O(N) \mid \det(\mathbf{U}) = 1\}$. $\sigma_{\min}(\mathbf{X})$ stands for the smallest singular value of $\mathbf{X} \in \mathbb{R}^{N \times N}$ and $o(\epsilon)$ is a matrix-valued function which satisfies $\lim_{\epsilon \rightarrow 0} \left\| \frac{o(\epsilon)}{\epsilon} \right\|_2 = 0$, where $\|\cdot\|_2$ stands for the spectral norm. We also use $\|\cdot\|$ to denote the standard Euclidean norm for vectors.

B. Dual Cayley parametrization techniques

The dual Cayley parametrization technique [17] was proposed to solve (1) in the case of $p = N$:

$$\text{minimize } f(\mathbf{U}) \text{ subject to } \mathbf{U} \in O(N). \quad (2)$$

The idea of the technique is to regard the optimization problem as optimization over the vector space $Q(N)$ by parametrization of the subset of $O(N)$ with $Q(N)$ using the Cayley transform Φ . Although Φ can parametrize $O(N) \setminus E \subsetneq SO(N)$, can not parametrize any subset of $O(N) \setminus SO(N) = \{\mathbf{U} \in O(N) \mid \det(\mathbf{U}) = -1\} = \{\mathbf{U}\mathbf{T} \mid \mathbf{U} \in SO(N)\}$, where $\mathbf{T} \in O(N) \setminus SO(N)$. To parametrize also a subset of $O(N) \setminus SO(N)$, the dual Cayley parametrization [17] was proposed.

Definition 1 (Dual Cayley parametrization).

$$\begin{cases} \Phi^{-1} & : Q(N) \rightarrow O(N) \setminus E : \mathbf{V} \mapsto \Phi^{-1}(\mathbf{V}) \\ \mathbf{T}\Phi^{-1} & : Q(N) \rightarrow \mathbf{T}(O(N) \setminus E) : \mathbf{V} \mapsto \Phi^{-1}(\mathbf{V})\mathbf{T}, \end{cases} \quad (3)$$

where $\mathbf{T} \in O(N) \setminus SO(N)$ and $\mathbf{T}(O(N) \setminus E) := \{\mathbf{U}\mathbf{T} \mid \mathbf{U} \in O(N) \setminus E\} \subsetneq O(N) \setminus SO(N)$.

Since $(O(N) \setminus E) \cup \mathbf{T}(O(N) \setminus E)$ is dense in $O(N)$, [17] formulated an optimization problem over $(O(N) \setminus E) \cup \mathbf{T}(O(N) \setminus E)$ corresponding to (2) by translating the optimization further into

$$\begin{cases} \text{find } \mathbf{U}^{*(1)} \in \Phi^{-1} \left(\arg \min_{\mathbf{V} \in Q(N)} \hat{f}_1(\mathbf{V}) \right) \\ \text{find } \mathbf{U}^{*(2)} \in \mathbf{T} \left(\Phi^{-1} \left(\arg \min_{\mathbf{V} \in Q(N)} \hat{f}_2(\mathbf{V}) \right) \right), \end{cases} \quad (4)$$

where $\hat{f}_1(\mathbf{V}) := f(\Phi^{-1}(\mathbf{V}))$ and $\hat{f}_2(\mathbf{V}) := f(\Phi^{-1}(\mathbf{V})\mathbf{T})$ for $\mathbf{V} \in Q(N)$. Since $Q(N)$ is a vector space, standard optimization techniques, e.g., the steepest descent method, the Newton's method, the conjugate gradient method, can be applied to the optimization problem (4). Unfortunately, the numerical performance of the technique tends to slow down severely due to the appearance of undesirable plateau around the singular points E .

C. Anderson acceleration

The Anderson acceleration [20], [21], [23] is a well-established technique to accelerate Picard type iterative schemes: $\mathbf{x}_{n+1} = g(\mathbf{x}_n) \in \mathbb{R}^M$ for finding a fixed-point $\mathbf{x}^* \in \mathbb{R}^M$ of wide range of nonlinear operators $g : \mathbb{R}^M \rightarrow \mathbb{R}^M$. This remarkably simple acceleration is known to be very effective [20], [23] for many algorithms using $g : \text{Id} - \gamma \nabla J$ with $\gamma > 0$ to minimize of J . Anderson acceleration is a restarting technique of the Picard type update whose new initial point is computed as $\mathbf{x}_{\text{ext}} := \sum_{k=0}^{K-1} c_k^* \mathbf{x}_k$, where $(\mathbf{x}_k)_{k=0}^{K-1}$ are the K latest estimates and \mathbf{c}^* is a minimizer of $\| \sum_{k=0}^{K-1} c_k (g(\mathbf{x}_k) - \mathbf{x}_k) \|$, as an approximation of $\| g(\sum_{k=0}^{K-1} c_k \mathbf{x}_k) - \sum_{k=0}^{K-1} c_k \mathbf{x}_k \|$, subject to $\sum_{k=0}^{K-1} c_k = 1$. Therefore we have to solve the approximated problem

$$\text{find } \mathbf{c}^* \in \arg \min_{\mathbf{c}^T \mathbf{1} = 1} \| \mathbf{R}\mathbf{c} \|, \quad (5)$$

where $\mathbf{R} = [\nabla J(\mathbf{x}_0) \quad \nabla J(\mathbf{x}_1) \quad \dots \quad \nabla J(\mathbf{x}_{K-1})]$. However, since \mathbf{R} tends to become singular or nearly-singular even if K is small, a Tikhonov type regularization technique has also been discussed as *regularized nonlinear acceleration* (RNA)

Algorithm 1 Regularized Nonlinear Acceleration [24] as an Anderson type acceleration [21]

Generate a pair of sequences $(\mathbf{x}_k)_{k=1}^K$, $(\mathbf{y}_k)_{k=0}^{K-1}$ from (6)
 $\mathbf{R} = [\mathbf{x}_1 - \mathbf{y}_0 \ \mathbf{x}_2 - \mathbf{y}_1 \ \dots \ \mathbf{x}_K - \mathbf{y}_{K-1}]$
 $\mathbf{c}^* = \frac{(\mathbf{R}^T \mathbf{R} + \lambda \mathbf{I})^{-1} \mathbf{1}}{\mathbf{1}^T (\mathbf{R}^T \mathbf{R} + \lambda \mathbf{I})^{-1} \mathbf{1}}$
 $\mathbf{x}_{\text{ext}} = \sum_{k=1}^K c_k^* \mathbf{x}_k$

in [23], [24] of which a version is summarized in Algorithm 1, where a pair of sequences $(\mathbf{x}_n)_{n=1}^K$, $(\mathbf{y}_n)_{n=0}^{K-1}$ generated by

$$\begin{cases} \mathbf{x}_n &= g(\mathbf{y}_{n-1}) \\ \mathbf{y}_n &= \frac{\gamma_n}{\gamma} \mathbf{x}_n + \left(1 - \frac{\gamma_n}{\gamma}\right) \mathbf{y}_{n-1}. \end{cases} \quad (6)$$

III. OPTIMIZATION ALGORITHMS VIA ADAPTIVE LOCALIZED CAYLEY PARAMETRIZATION

A. The mobility bound of the Cayley parametrization

To see the influence of the singular points of Φ in the gradient type scheme algorithm in [17], we evaluate the mobility at \mathbf{V} as

$$\|\Phi^{-1}(\mathbf{V} + \epsilon \mathbf{F}) - \Phi^{-1}(\mathbf{V})\|_2 \leq \frac{2\|\epsilon \mathbf{F}\|_2}{\sigma_{\min}^2(\mathbf{I} + \mathbf{V})} + \|o(\epsilon)\|_2, \quad (7)$$

where \mathbf{V} , $\mathbf{F} \in Q(N)$, $\epsilon \in \mathbb{R}$. From this inequality, the upper bound of the mobility is inversely proportional to the square of $\sigma_{\min}(\mathbf{I} + \mathbf{V})$ and its maximum is achieved by $\mathbf{V} = \mathbf{0}$, implying thus undesirable plateau hardly appears around $\mathbf{0}$. Unfortunately, the dual Cayley parametrization technique [17] can not obtain the estimate $\mathbf{V}_n \in Q(N)$ while keeping $\sigma_{\min}(\mathbf{I} + \mathbf{V}_n)$ small when \mathbf{V}^* , which corresponds to the optimal point \mathbf{U}^* , is located far away from $\mathbf{0}$.

B. Proposed algorithms

In the following, we present an extension of the Cayley transform in order to keep the mobility large in optimization algorithms for problem (2).

Definition 2 (Localized Cayley transform). For a given $\mathbf{S} \in O(N)$, the localized Cayley transform $\Phi_{\mathbf{S}} : O(N) \setminus E(\mathbf{S}) \rightarrow Q_{\mathbf{S}}(N)$ centered at \mathbf{S} is defined by

$$\Phi_{\mathbf{S}}(\mathbf{U}) := (\mathbf{S} - \mathbf{U})(\mathbf{S} + \mathbf{U})^{-1} \quad (\mathbf{U} \in O(N) \setminus E(\mathbf{S})), \quad (8)$$

where $E(\mathbf{S}) := \{\mathbf{U} \in O(N) \mid \det(\mathbf{S} + \mathbf{U}) = 0\}$ is the set of all singular points of $\Phi_{\mathbf{S}}$, and $Q_{\mathbf{S}}(N) := \{\Phi_{\mathbf{S}}(\mathbf{U}) \mid \mathbf{U} \in O(N) \setminus E(\mathbf{S})\}$ is also the set of all skew-symmetric matrices but for parametrization of $O(N) \setminus E(\mathbf{S})$.

Moreover, we give the inverse of $\Phi_{\mathbf{S}}$ is given in the following lemma.

Lemma 1 (Inverse of the localized Cayley transform). The localized Cayley transform $\Phi_{\mathbf{S}}$ centered at $\mathbf{S} \in O(N)$ is diffeomorphic and its inverse mapping $\Phi_{\mathbf{S}}^{-1} : Q_{\mathbf{S}}(N) \rightarrow O(N) \setminus E(\mathbf{S})$ is given by

$$\Phi_{\mathbf{S}}^{-1}(\mathbf{V}) := (\mathbf{I} - \mathbf{V})(\mathbf{I} + \mathbf{V})^{-1} \mathbf{S} \quad (\mathbf{V} \in Q_{\mathbf{S}}(N)). \quad (9)$$

Algorithm 2 Adaptive localized Cayley parametrization with a steepest descent method at every K -th iteration

Choose $\mathbf{U}_0^{(1)} \in SO(N)$, $\mathbf{U}_0^{(2)} \in T(SO(N))$
 $\mathbf{S}_0^{(t)} = \mathbf{U}_0^{(t)}$ ($t = 1, 2$)
for $n = 0, 1, 2, \dots$ **do**
 $\mathbf{V}_{n+1}^{(t)} = \mathbf{V}_n^{(t)} - \gamma_n^{(t)} \nabla (f \circ \Phi_{\mathbf{S}_n^{(t)}}^{-1})(\mathbf{V}_n^{(t)})$ ($\gamma_n^{(t)} \in [0, \infty)$)
 $\mathbf{U}_{n+1}^{(t)} = \Phi_{\mathbf{S}_n^{(t)}}^{-1}(\mathbf{V}_{n+1}^{(t)})$
if $n \bmod K = 0$ **then**
 $\mathbf{S}_{n+1}^{(t)} = \mathbf{U}_{n+1}^{(t)}$, $\mathbf{V}_{n+1}^{(t)} = \mathbf{0}$
else
 $\mathbf{S}_{n+1}^{(t)} = \mathbf{S}_n^{(t)}$
end if
end for ($t = 1, 2$)

In a way similar to (7), we obtain

$$\|\Phi_{\mathbf{S}}^{-1}(\Phi_{\mathbf{S}}(\mathbf{U}) + \epsilon \mathbf{F}) - \Phi_{\mathbf{S}}^{-1}(\Phi_{\mathbf{S}}(\mathbf{U}))\|_2 \leq \frac{2\|\epsilon \mathbf{F}\|_2 + \|o(\epsilon)\|_2}{\sigma_{\min}^2(\mathbf{I} + \Phi_{\mathbf{S}}(\mathbf{U}))}$$

for any $\mathbf{U} \in O(N) \setminus E(\mathbf{S})$, which suggests that the upper bound of the mobility $\|\Phi_{\mathbf{S}}^{-1}(\Phi_{\mathbf{S}}(\mathbf{U}) + \epsilon \mathbf{F}) - \Phi_{\mathbf{S}}^{-1}(\Phi_{\mathbf{S}}(\mathbf{U}))\|_2$ at $\Phi_{\mathbf{S}}(\mathbf{U})$ is maximized by setting $\mathbf{S} := \mathbf{U}$ to ensure $\Phi_{\mathbf{S}}(\mathbf{U}) = \mathbf{0}$.

Based on this fact, we propose the *adaptive localized Cayley parametrization technique*: $\mathbf{U}_{n+1} := \Phi_{\mathbf{S}_n}^{-1}(\mathcal{A}(\Phi_{\mathbf{S}_n}(\mathbf{U}_n)))$ for (2), where \mathcal{A} is the update rule from $\mathbf{V}_n := \Phi_{\mathbf{S}_n}(\mathbf{U}_n) \in Q_{\mathbf{S}_n}(N)$ to $\mathbf{V}_{n+1} := \mathcal{A}(\mathbf{V}_n) := \Phi_{\mathbf{S}_n}(\mathbf{U}_{n+1}) \in Q_{\mathbf{S}_n}(N)$ for finding $\mathbf{V}_{n+1} \in Q_{\mathbf{S}_n}(N)$ and $\mathbf{U}_{n+1} := \Phi_{\mathbf{S}_n}^{-1}(\mathbf{V}_{n+1})$ to suppress $f \circ \Phi_{\mathbf{S}_n}^{-1}$ over $Q_{\mathbf{S}_n}(N)$, equivalently f over $O(N) \setminus E(\mathbf{S}_n)$, and \mathbf{S}_n is updated to \mathbf{U}_n at every K -th iteration. Actually, it is ideal to keep the upper bound of the mobility the maximum, i.e., \mathbf{S}_n is updated to \mathbf{U}_n at every iteration, but it is expected that the upper bound remains near the maximum if K is not large. Thanks to the adaptive localized Cayley parametrization technique, any optimization technique over a vector space can be employed as the update rule \mathcal{A} , and as a simplest example, we present a steepest descent method for (2) in Algorithm 2, where (i) $\mathbf{U}_n^{(1)} \in SO(N)$ and $\mathbf{U}_n^{(2)} \in T(SO(N))$ are generated simultaneously by using suitable stepsizes $\gamma_n^{(t)}$ ($t = 1, 2$), e.g., the Armijo's rule (See Theorem 1) to achieve monotone decreasing of $f(\mathbf{U}_n^{(t)})$ ($n = 1, 2, \dots$) for $t = 1, 2$, and (ii) the gradient is given by $\nabla (f \circ \Phi_{\mathbf{S}}^{-1})(\mathbf{V}) = 2(\mathbf{W} - \mathbf{W}^T)$ with $\mathbf{W} = (\mathbf{I} + \mathbf{V})^{-1} \mathbf{S} \nabla f(\Phi_{\mathbf{S}}^{-1}(\mathbf{V}))^T (\mathbf{I} + \mathbf{V})^{-1}$.

Remark 1. A specialization of Algorithm 2 with $K = 1$ reproduces a Riemannian optimization technique for (1) with the Cayley transform as a retraction [9]. We emphasize that the strategies with $K > 1$ [17], [19] and with $K = 1$ [9] have distinct difference in the sense that the former can enjoy further special arts developed for optimization tasks over a common vector space during K iterations (See, e.g., Algorithm 3) but the latter can not.

Moreover, to make the best of the advantage of the proposed

Algorithm 3 Adaptive localized Cayley parametrization with an Anderson type acceleration at every K -th iteration

Choose $\mathbf{U}_0^{(1)} \in SO(N)$, $\mathbf{U}_0^{(2)} \in T(SO(N))$, $\gamma \in [0, \infty)$
 $\mathbf{S}_0^{(t)} = \mathbf{U}_0^{(t)}$, $\mathbf{Y}_0^{(t)} = \mathbf{0}$ ($t = 1, 2$)
for $n = 0, 1, 2, \dots$ **do**
 $\mathbf{X}_{n+1}^{(t)} = \mathbf{Y}_n^{(t)} - \gamma \nabla \left(f \circ \Phi_{\mathbf{S}_n^{(t)}}^{-1} \right) (\mathbf{Y}_n^{(t)})$
if $n \bmod K = 0$ **then**
 $\mathbf{X}_{\text{ext}}^{(t)} = \text{RNA} \left((\mathbf{X})_{k=n-K+1}^K, (\mathbf{Y})_{k=n-K}^K, \lambda \right)$
 $\mathbf{U}_{n+1}^{(t)} = \Phi_{\mathbf{S}_n^{(t)}}^{-1} (\mathbf{X}_{\text{ext}}^{(t)})$
 $\mathbf{S}_{n+1}^{(t)} = \mathbf{U}_{n+1}^{(t)}$, $\mathbf{Y}_{n+1}^{(t)} = \mathbf{0}$
else
 $\mathbf{U}_{n+1}^{(t)} = \Phi_{\mathbf{S}_n^{(t)}}^{-1} (\mathbf{X}_{n+1}^{(t)})$
 $\mathbf{S}_{n+1}^{(t)} = \mathbf{S}_n^{(t)}$
 $\mathbf{Y}_{n+1}^{(t)} = \frac{\gamma_n^{(t)}}{\gamma} \mathbf{X}_{n+1}^{(t)} + \left(1 - \frac{\gamma_n^{(t)}}{\gamma} \right) \mathbf{Y}_n^{(t)}$ ($\gamma_n \in [0, \infty)$)
end if
end for ($t = 1, 2$)

parametrization with $K > 1$ in Remark 1, we also propose an enhancement as Algorithm 3 of the parametrization technique with an Anderson type acceleration, RNA [23], which can be applied to optimization tasks but only to over a vector spaces. Since both of these two computational ideas can be implemented at every K -th iteration, the combined algorithm in Algorithm 3 can enjoy excellent synergy effects (See sec. IV).

Remark 2. The singular point set $E(\mathbf{S})$ of $\Phi_{\mathbf{S}}$ is different from $E(\mathbf{S}')$ of $\Phi_{\mathbf{S}'}$, for all $\mathbf{S}' \in O(N) \setminus E(\mathbf{S})$. Moreover, since the maximum of the determinant $|\det(\mathbf{S} + \Phi_{\mathbf{S}}^{-1}(\mathbf{V}))| = \frac{2^N}{\det(\mathbf{I} + \mathbf{V})}$ is achieved at $\mathbf{V} = \mathbf{0}$, \mathbf{S} is ensured to locate away from $E(\mathbf{S})$.

Next, we present an extension of the adaptive localized Cayley parametrization techniques for (1). To achieve this goal, we use the canonical projection: $\Xi : O(N) \rightarrow St(p, N) : \mathbf{S} \mapsto \mathbf{S}\mathbf{I}_{N \times p}$, for translating cost function f into $h := f \circ \Xi$. The problem (1) with $p < N$ can be translated into minimization of h over $SO(N)$ because $St(p, N)$ for $p < N$ is a connected manifold unlike $O(N)$. Thus, the proposed parametrization techniques can be applied to problem (1) for general case by passing through Ξ .

C. Convergence Analysis

To establish a convergence analysis for Algorithm 2, we employ the *backtracking algorithm* for stepsizes $\gamma_n^{(t)}$ to satisfy the *Armijo's rule* (See [22]).

Theorem 1 (Convergence Analysis of Algorithm 2). Suppose that $f : \mathbb{R}^{N \times N} \rightarrow \mathbb{R}$ is differentiable and $(\mathbf{U}_n^{(t)})_{n=0}^{\infty}$ ($t = 1, 2$) are sequences generated by Algorithm 2 using the backtracking algorithm [22] for stepsizes. Then we have

$$\lim_{n \rightarrow \infty} \left\| \nabla \left(f \circ \Phi_{\mathbf{U}_n^{(t)}}^{-1} \right) (\mathbf{0}) \right\| = 0. \quad (10)$$

IV. NUMERICAL EXPERIMENTS

To demonstrate the performances of Algorithm 2 and Algorithm 3, we consider two optimization problems with many applications in data sciences. We compare Algorithm 2 and Algorithm 3 with (a) the original dual Cayley parametrization technique [17], which is equivalent to Algorithm 2 with $K = \infty$, (b) the steepest descent methods for optimization over a Riemannian manifold employing, as retractions, (b-1) the Cayley transform [9], which is equivalent to Algorithm 2 with $K = 1$ (See Remark 2), and (b-2) QR decomposition [8], called Algorithm 4 in this paper. Each stepsize of all algorithms are determined by the backtracking algorithm [22]. A parameter λ of Algorithm 3 for regularization is set as 10^{-17} .

A. Joint diagonalization. We consider

$$\text{minimize } f(\mathbf{U}) := \sum_{m=1}^{10} \text{off}(\mathbf{U}^{-1} \mathbf{A}_m \mathbf{U}) \text{ subject to } \mathbf{U} \in O(10),$$

where $\text{off}(\mathbf{X}) := \sum_{i \neq j} x_{ij}^2$, x_{ij} denotes the (i, j) -th entry of $\mathbf{X} \in \mathbb{R}^{10 \times 10}$ and $\mathbf{A}_m \in \mathbb{R}^{10 \times 10}$ ($m = 1, 2, \dots, 10$) are obtained by $\mathbf{U}^* \mathbf{\Lambda}_m \mathbf{U}^{*-1}$ with $\mathbf{U}^* \in O(10)$ and a randomly chosen diagonal matrix $\mathbf{\Lambda}_m \in \mathbb{R}^{10 \times 10}$. A fixed stepsize γ of Algorithm 3 is set as 0.0005. To check the avoidance of undesirable plateau and the performance for the proposed algorithms, we consider the following two settings: (i) the optimal point and the initial point are around singular set E of Φ :

$$\mathbf{U}^* := \begin{bmatrix} -\mathbf{I}_2 & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_8 \end{bmatrix}, \mathbf{U}_0 := \begin{bmatrix} \mathbf{R}(15\pi/16) & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_8 \end{bmatrix}$$

and (ii) the optimal point and the initial point are far from E :

$$\mathbf{U}^* := \begin{bmatrix} \mathbf{R}(\pi/6) & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_8 \end{bmatrix}, \mathbf{U}_0 := \mathbf{I},$$

where $\mathbf{R}(\theta) \in \mathbb{R}^{2 \times 2}$ stands for a rotation matrix of a given angle θ .

B. Eigenbasis extraction. We consider

$$\text{minimize } f(\mathbf{U}) := -\text{Tr}(\mathbf{U}^T \mathbf{A} \mathbf{U}) \text{ subject to } \mathbf{U} \in St(7, 500),$$

where $\mathbf{A} \in \mathbb{R}^{500 \times 500}$ is randomly chosen such that \mathbf{A} is symmetric (Note: Any solution of this problem is known to be an orthonormal eigenbasis associated with the 7 largest eigenvalues of \mathbf{A} [8]). We also choose an initial point $\mathbf{U}_0 \in St(7, 500)$ randomly. A fixed stepsize γ of Algorithm 3 is set as 0.000001.

Results of our experiments of each algorithm for the *joint diagonalization* with the above settings (i), (ii) and for the *eigenbasis extraction* are illustrated respectively in Fig. 1, Fig. 2 and Fig. 3. Fig. 1 shows that the proposed algorithms (Algorithm 2 with $K = 5$ and Algorithm 3 with $K = 5$) succeed in avoiding the undesirable plateau which is observed in the original dual Cayley parametrization technique (Algorithm 2 with $K = \infty$). Compared with Algorithm 2 with $K = 1$ and Algorithm 4, Algorithm 2 with $K = 5$ in Fig. 1 and Fig. 2 shows better performance in the first scenario. Fig. 3 shows that Algorithm 2 ($K = 5$) has competitive

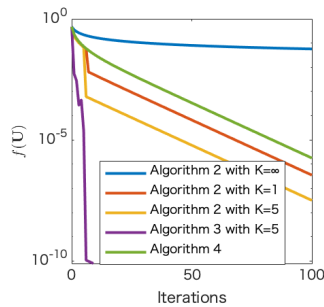


Fig. 1. The value of the cost function versus the number of iteration for different algorithms for the joint diagonalization with the setting (i) the initial point and the optimal point are around E .

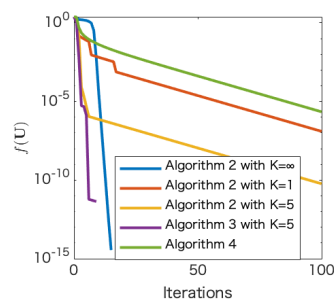


Fig. 2. The value of the cost function versus the number of iteration for different algorithms for the joint diagonalization with the setting (ii) the initial point and the optimal point are far from E .

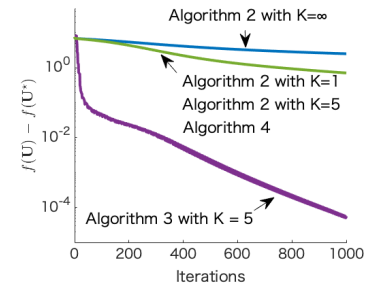


Fig. 3. The subtraction of the value of the cost function for the optimal point from one for each estimate versus the number of iteration for different algorithms for eigenbasis extraction.

performance with Algorithm 2 ($K = 1$) and Algorithm 4 in the second scenario. These results show that the prototype algorithm, Algorithm 2, has at least competitive performance with standard optimization techniques over Riemannian manifolds, which suggests great potential for improvement of the speed of convergence by introducing further advanced techniques applicable over vector spaces. Indeed, we can observe that Anderson type acceleration (Algorithm 3 with $K = 5$) succeeds in accelerating dramatically Algorithm 2 with $K = 5$, and Algorithm 3 with $K = 5$ overwhelms the others in all scenarios.

V. CONCLUSION

We presented a novel computational strategy, named the adaptive localized Cayley parametrization technique for acceleration of optimization over the Stiefel manifold. The proposed algorithm (i) is free from the singularity issue, which can cause performance degradation, observed in the dual Cayley parametrization technique [Yamada- Ezaki '03] as well as (ii) can enjoy powerful arts for acceleration on the vector space without suffering from the nonlinear nature of the Stiefel manifold. Numerical experiments show that the proposed algorithm outperforms major optimization algorithms designed with retractions on the tangent space of the Stiefel manifold.

REFERENCES

- [1] R. Pietersz and P. J. F. Groenen, "Rank reduction of correlation matrices by majorization," *Quantitative Finance*, vol. 4, no. 6, pp. 649–662, 2004.
- [2] Z. Zhao, Z. Bai, and X. Jin, "A Riemannian Newton algorithm for nonlinear eigenvalue problems," *SIAM Journal on Matrix Analysis and Applications*, vol. 36, no. 2, pp. 752–774, 2015.
- [3] Z. Lu and Y. Zhang, "An augmented Lagrangian approach for sparse principal component analysis," *Mathematical Programming*, vol. 135, no. 1, pp. 149–193, 2012.
- [4] P. T. Boufounos and R. G. Baraniuk, "1-bit compressive sensing," in *CISS*, Princeton, N.J., 2008, pp. 16–21.
- [5] F. J. Theis, T. P. Cason, and P. A. Absil, "Soft dimension reduction for ICA by joint diagonalization on the Stiefel manifold," in *ICA*. Springer Berlin Heidelberg, 2009, pp. 354–361.
- [6] J. Lu, V. E. Liang, X. Zhou, and J. Zhou, "Learning compact binary face descriptor for face recognition," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 37, no. 10, pp. 2041–2056, 2015.
- [7] L. Huang, X. Liu, B. Lang, A. W. Yu, Y. Wang, and B. Li, "Orthogonal weight normalization: Solution to optimization over multiple dependent Stiefel manifolds in deep neural networks," in *AAAI*, 2018.
- [8] P.-A. Absil, R. Mahony, and R. Sepulchre, *Optimization Algorithms on Matrix Manifolds*. Princeton University Press, 2008.
- [9] Z. Wen and W. Yin, "A feasible method for optimization with orthogonality constraints," *Mathematical Programming*, vol. 142, no. 1, pp. 397–434, 2013.
- [10] X. Zhu, "A Riemannian conjugate gradient method for optimization on the Stiefel manifold," *Computational Optimization and Applications*, vol. 67, no. 1, pp. 73–110, 2017.
- [11] C. Fraikin, K. Hüper, and P. V. Dooren, "Optimization over the Stiefel manifold," *PAMM*, vol. 7, no. 1, pp. 1 062 205–1 062 206, 2007.
- [12] B. Jiang and Y.-H. Dai, "A framework of constraint preserving update schemes for optimization on Stiefel manifold," *Mathematical Programming*, vol. 153, no. 2, pp. 535–575, 2015.
- [13] T. E. Abrudan, J. Eriksson, and V. Koivunen, "Steepest descent algorithms for optimization under unitary matrix constraint," *IEEE Transactions on Signal Processing*, vol. 56, no. 3, pp. 1134–1147, 2008.
- [14] P.-A. Absil and J. Malick, "Projection-like retractions on matrix manifolds," *SIAM Journal on Optimization*, vol. 22, no. 1, pp. 135–158, 2012.
- [15] M. Nikpour, J. H. Manton, and G. Hori, "Algorithms on the Stiefel manifold for joint diagonalisation," in *ICASSP*, vol. 2. IEEE, 2002, pp. II-1481.
- [16] S. Fiori, T. Kaneko, and T. Tanaka, "Learning on the compact Stiefel manifold by a cayley-transform-based pseudo-retraction map," in *IJCNN*, 2012, pp. 1–8.
- [17] I. Yamada and T. Ezaki, "An orthogonal matrix optimization by dual Cayley parametrization technique," in *ICA*, Nara, Japan, 2003, pp. 35–40.
- [18] I. Satake, *Linear algebra*. NY:Marcel Dekker, Inc., 1975.
- [19] G. Hori and T. Tanaka, "Pivoting in Cayley transform-based optimization on orthogonal groups," in *APSIPA*. Biopolis Singapore, 2010, pp. 181–184.
- [20] V. Eyert, "A comparative study on methods for convergence acceleration of iterative vector sequences," *Journal of Computational Physics*, vol. 124, no. 2, pp. 271 – 285, 1996.
- [21] D. G. Anderson, "Iterative procedures for nonlinear integral equations," *J. ACM*, vol. 12, no. 4, pp. 547–560, Oct. 1965.
- [22] A. Wächter and L. T. Biegler, "On the implementation of an interior-point filter line-search algorithm for large-scale nonlinear programming," *Mathematical Programming*, vol. 106, no. 1, pp. 25–57, 2006.
- [23] D. Scieur, A. d'Aspremont, and F. Bach, "Regularized nonlinear acceleration," in *NIPS*, 2016, pp. 712–720.
- [24] D. Scieur, "Acceleration in optimization," Ph.D. Thesis, PSL Research University, Sep. 2018.