

A Nesterov-type Acceleration with Adaptive Localized Cayley Parametrization for 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—Despite certain singular-point issues, the Cayley parametrization (CP) has great potential to serve as a key to import many powerful strategies, developed originally for optimization over a vector space, into the task for optimization over the Stiefel manifold. In this paper, we newly present (i) a computationally efficient CP that can circumvent the singular-point issues and (ii) a Nesterov type accelerated gradient method, based on the proposed CP, with its convergence analysis. To guarantee the convergence, we also evaluate a Lipschitz constant of the gradient of the cost function in the CP domain. Numerical experiments show excellent performance of the proposed accelerated algorithm compared with the standard algorithms, e.g., the Barzilai-Borwein method and L-BFGS method, combined with a vector transport for optimization over the Stiefel manifold as a special instance of the Riemannian manifold.

Index Terms—Stiefel manifold, orthogonal group, Cayley parametrization, non-convex optimization, Nesterov acceleration

I. INTRODUCTION

We consider orthogonal constrained optimization problem^{a)}:

Problem I.1. Let $\text{St}(p, N) := \{\mathbf{X} \in \mathbb{R}^{N \times p} \mid \mathbf{X}^T \mathbf{X} = \mathbf{I}_p \in \mathbb{R}^{p \times p}\}$ with $p \leq N$, where \mathbf{I}_p is the identity matrix. For a given continuously differentiable function $f: \mathbb{R}^{N \times p} \rightarrow \mathbb{R}$,

$$\text{Find } \mathbf{X}^* \in \arg \min f(\mathbf{X}) \text{ s.t. } \mathbf{X} \in \text{St}(p, N) \quad (1)$$

The feasible region $\text{St}(p, N)$ of (1) is called the Stiefel manifold, and particularly $\text{O}(N) := \text{St}(N, N)$ is called the orthogonal group. This optimization problem has rich applications in data sciences including signal processing and machine learning, such as joint diagonalization problem [1], orthogonal Procrustes problem [2], enhancement of the generalization performance in deep neural network [3]. However, development of numerically efficient algorithms for (1) has been challenging due to the severe non-linearity of $\text{St}(p, N)$.

The standard optimization strategy [4] for the Problem I.1 are realized by two steps at n th iteration: (i) find a search direction \mathbf{D}_n in the tangent space $T_{\mathbf{X}_n} \text{St}(p, N)$ to $\text{St}(p, N)$ at the latest estimate $\mathbf{X}_n \in \text{St}(p, N)$; (ii) assign $R_{\mathbf{X}_n}(\mathbf{D}_n) \in \text{St}(p, N)$ to a new estimate \mathbf{X}_{n+1} with a retraction $R_{\mathbf{X}_n}: T_{\mathbf{X}_n} \text{St}(p, N) \rightarrow \text{St}(p, N)$ as a certain approximation of the exponential map^{b)}. Many examples of retractions for $\text{St}(p, N)$ are known, e.g., QR decomposition,

^{a)}The existence of a minimizer in (1) is automatically guaranteed by the compactness of $\text{St}(p, N)$ and the continuity of f .

^{b)}The exponential map $\text{Exp}_{\mathbf{X}}: T_{\mathbf{X}} \text{St}(p, N) \rightarrow \text{St}(p, N)$ at $\mathbf{X} \in \text{St}(p, N)$ assigns a given direction $\mathbf{D} \in T_{\mathbf{X}} \text{St}(p, N)$ to a point on the geodesic of $\text{St}(p, N)$ with the initial velocity \mathbf{D} .

polar decomposition and the Cayley transform [4]–[7]. Along this strategy, optimization algorithms designed originally on a vector space have been extended for (1) combined with a vector transport [4], which translates a tangent vector into one in the other tangent space. Such extensions include the steepest descent method, the conjugated gradient method, Newton’s method, quasi-Newton’s method, the Barzilai–Borwein method and the trust-region method [1], [5], [7]–[12].

Recently, to gain the speed of the convergence to a solution, the way for exploiting the ideas in the *Nesterov Accelerated Gradient (NAG)* [13] for its applications to (1) has been studied (see, e.g., [14]–[16]). The NAG [13] over a vector space, say \mathcal{X} , generates the sequence $(\mathbf{Y}_n)_{n=0}^{\infty} \subset \mathcal{X}$ as

$$\left. \begin{aligned} \mathbf{Y}_{n+1} &= \mathbf{Z}_n - \gamma_n \nabla f(\mathbf{Z}_n) \\ \mathbf{Z}_{n+1} &= \mathbf{Y}_{n+1} + \beta_n (\mathbf{Y}_{n+1} - \mathbf{Y}_n) \end{aligned} \right\} \quad (2)$$

with initial points $\mathbf{Y}_0 = \mathbf{Z}_0 \in \mathcal{X}$. This remarkably simple method accelerates the convergence of the gradient descent method although NAG does not guarantee a monotonically decrease of the value of f . In particular for convex problems, NAG is known to achieve the optimal convergence rate [17], in the first-order methods, with suitable stepsizes $\gamma_n > 0$ and momentum parameters $\beta_n \in \mathbb{R}$. In [18], an extension of NAG has been made for unconstrained smooth non-convex problems and its efficacy has also been examined.

Since the second step in (2) of NAG relies on the linear structure of \mathcal{X} , its extensions to be applicable to $\text{St}(p, N) \subsetneq \mathcal{X} := \mathbb{R}^{N \times p}$ has to face inherent difficulty caused by the loss of linearity. For example, [14], [15] proposed to overcome this issue by translating $\mathbf{Y}_n \in \text{St}(p, N)$ into a point on the tangent space $T_{\mathbf{Y}_n} \text{St}(p, N)$ with the inversion of $\text{Exp}_{\mathbf{Y}_n}$, but the computation of the inversion map itself requires iterative algorithm [19] and therefore their methods are hard to be implemented for Problem I.1 as remarked in [16]. The author of [16] proposed to extend NAG for Problem I.1 by replacing the exponential map with the Cayley retraction in [7]. To ensure a monotonically decrease of the value of f , [16] employs an adaptive restart scheme, which starts the algorithm again with the latest estimate as an initial point, and numerical performance of [16] for Problem I.1 was verified experimentally. Although NAG in [16] guarantees the sequence of the gradients converges to zero for non-convex optimization on \mathcal{X} [16, Theorem 4.1], the condition corresponding to [16, Eq. (4.6)] has not been justified for the algorithm [16, Alg. 4.1] because different tangent spaces

are employed at every iteration and therefore [16, Alg. 4.1], designed for Problem I.1, has no guarantee of convergence.

From a totally different view in utilizing the Cayley transform for Problem I.1 with $p = N$, we recently proposed an Adaptive Localized Cayley Parametrization technique (ALCP) in [20]. In a Cayley parametrization technique (CP) [3], [20], [21], we translate Problem I.1 with $p = N$ into the following problem^{e)} with the inversion $\tilde{\varphi}_{\mathcal{S}}^{-1}: Q_{N,N} \rightarrow O(N) \setminus E_{N,N}(\mathcal{S})$ of a *localized Cayley transform* $\tilde{\varphi}_{\mathcal{S}}$ for $\mathcal{S} \in O(N)$ (see Def. II.1 for $\tilde{\varphi}_{\mathcal{S}}, \tilde{\varphi}_{\mathcal{S}}^{-1}, Q_{N,N}, E_{N,N}(\mathcal{S})$):

Problem I.2. For an arbitrarily given $\epsilon > 0$, a continuously differentiable function $f: \mathbb{R}^{N \times N} \rightarrow \mathbb{R}$ and $\mathcal{S} \in SO(N) := \{\mathbf{X} \in O(N) | \det(\mathbf{X}) = 1\}$,

$$\text{Find } \mathbf{V}^* \in Q_{N,N} \text{ s.t. } f \circ \tilde{\varphi}_{\mathcal{S}}^{-1}(\mathbf{V}^*) < \min f(SO(N)) + \epsilon. \quad (3)$$

Problem I.2 is a sound relaxation of Problem I.1 because the existence of \mathbf{V}^* in (3) is guaranteed by the denseness of $\tilde{\varphi}_{\mathcal{S}}^{-1}(Q_{N,N}) = O(N) \setminus E_{N,N}(\mathcal{S})$ in $SO(N)$ [21]. Since $Q_{N,N}$ is a vector space, we can enjoy for Problem I.2 various arts, even (2) of NAG, of optimization over a vector space. To overcome the raised issue (see [3], [21]) regarding the performance degradation around the singular point set $E_{N,N}(\mathcal{S})$, we proposed ALCP in [20] by changing a center point \mathcal{S} to a certain $\mathcal{S}' \in O(N)$ after detecting the performance degradation (see Sec. II). We can also exploit ALCP even for general case $p \leq N$ by replacing $\tilde{\varphi}_{\mathcal{S}}^{-1}$ and $SO(N)$ with $\Xi \circ \tilde{\varphi}_{\mathcal{S}}^{-1}$ and $\text{St}(p, N)$ in (3), where $\Xi: O(N) \rightarrow \text{St}(p, N)$ is a canonical projection (see the beginning of Sec. III). However the computational complexity $\mathcal{O}(N^3)$ for $\Xi \circ \tilde{\varphi}_{\mathcal{S}}^{-1}$ is expensive compared with that for retractions of $\text{St}(p, N)$ even if $p \ll N$ as found in typical applications in signal processing and machine learning.

In this paper, we propose (i) a computationally efficient ALCP for $p \leq N$ with a new *Generalized Left Localized Cayley Transform (G-L²CT)* $\Phi_{\mathcal{S}}$; (ii) a NAG combined with the proposed ALCP.

More precisely, for parametrization of $\text{St}(p, N)$, we propose G-L²CT $\Phi_{\mathcal{S}}: \text{St}(p, N) \setminus E_{N,p}(\mathcal{S}) \rightarrow Q_{N,p}$ centered at $\mathcal{S} \in O(N)$ and its inversion $\Phi_{\mathcal{S}}^{-1}: Q_{N,p}(\mathcal{S}) \rightarrow \text{St}(p, N) \setminus E_{N,p}(\mathcal{S})$ (see Def. III.1 and Prop. III.2 for $\Phi_{\mathcal{S}}, \Phi_{\mathcal{S}}^{-1}, Q_{N,p}, E_{N,p}(\mathcal{S})$), where $Q_{N,p}$ is a structured subspace of $Q_{N,N}$. Theorem III.3 ensures the image $\text{St}(p, N) \setminus E_{N,p}(\mathcal{S})$ of $\Phi_{\mathcal{S}}^{-1}$ is dense in $\text{St}(p, N)$, thus Problem I.1 with $p < N$ can be relaxed soundly into the following problem:

Problem I.3. For an arbitrarily given $\epsilon > 0$, a continuously differentiable function $f: \mathbb{R}^{N \times p} \rightarrow \mathbb{R}$ and $\mathcal{S} \in O(N)$,

$$\text{Find } \mathbf{V}^* \in Q_{N,p} \text{ s.t. } f \circ \Phi_{\mathcal{S}}^{-1}(\mathbf{V}^*) < \min f(\text{St}(p, N)) + \epsilon.$$

By imposing a special structure $O_p(N) \subset O(N)$ (see (4)) on center points \mathcal{S} , the computational complexity for $\Phi_{\mathcal{S}}^{-1}$ achieves $\mathcal{O}(Np^2 + p^3)$, which is more efficient than that for $\Xi \circ \tilde{\varphi}_{\mathcal{S}}^{-1}$ and is competitive to that for retractions [22] of $\text{St}(p, N)$.

^{e)}Strictly speaking, Problem I.2 is a relaxation of optimization over $SO(N)$. For a relaxation of optimization over $O(N)$, we need to solve additionally another Problem I.2 replaced $SO(N)$ with $O(N) \setminus SO(N)$.

In Sec. IV, we propose NAG with ALCP for (1) by synchronizing the change of center points of ALCP and the adaptive restart scheme of NAG [16]. We also evaluate a Lipschitz constant of the gradient of f in the CP domain (see Lemma III.6) because the proposed NAG requires the Lipschitz constant. Unlike [16] for Problem I.1, the proposed NAG (i) has a guarantee to generate the gradient sequence of f , in the CP domain, which converges to zero (see Theorem. IV.1); (ii) does not require the inversion of neither the exponential map nor retractions.

Numerical experiments in Sec. V show that the proposed NAG outperforms the standard optimization methods provided in *Manopt* [23], e.g., the Barzilai-Borwein method and quasi-Newton's method (L-BFGS), for (1) in the scenarios of joint diagonalization and orthogonal Procrustes problem.

Notation The matrices $\mathbf{X}_{\text{up}} \in \mathbb{R}^{p \times p}$ and $\mathbf{X}_{\text{lo}} \in \mathbb{R}^{(N-p) \times p}$ denote the upper and the lower block matrices of $\mathbf{X} \in \mathbb{R}^{N \times p}$ respectively. The matrices $\mathbf{S}_{\text{le}} \in \mathbb{R}^{N \times p}$ and $\mathbf{S}_{\text{ri}} \in \mathbb{R}^{N \times (N-p)}$ denote the left and right block matrices of $\mathbf{S} \in \mathbb{R}^{N \times N}$ respectively. The matrix $\mathbf{I}_{N \times p} \in \text{St}(p, N)$ denotes the first p columns of the identity matrix $\mathbf{I} \in \mathbb{R}^{N \times N}$. The norms $\|\cdot\|_2$ and $\|\cdot\|_F$ denote the spectral norm and the Frobenius norm respectively. The function $\sigma_{\min}(\cdot)$ denotes the smallest singular value of a given matrix. The function $o(\epsilon)$ denotes a matrix-valued function that satisfies $\lim_{\epsilon \rightarrow 0} \|o(\epsilon)\|_F / \epsilon = 0$.

II. PRELIMINARY

Definition II.1 (Right localized Cayley transform for $O(N)$ [20]). Let $\mathcal{S} \in O(N)$, $Q_{N,N}(\mathcal{S}) := \{\mathbf{V} \in \mathbb{R}^{N \times N} | \mathbf{V}^T = -\mathbf{V}\}$, $E_{N,N}(\mathcal{S}) := \{\mathbf{X} \in O(N) | \det(\mathbf{I} + \mathbf{S}^T \mathbf{X}) = 0\}$. The right localized Cayley transform $\tilde{\varphi}_{\mathcal{S}}: O(N) \setminus E_{N,N}(\mathcal{S}) \rightarrow Q_{N,N}(\mathcal{S})$ centered at \mathcal{S} is defined by

$$(\mathbf{X} \in O(N) \setminus E_{N,N}(\mathcal{S})) \quad \tilde{\varphi}_{\mathcal{S}}(\mathbf{X}) = (\mathcal{S} - \mathbf{X})(\mathcal{S} + \mathbf{X})^{-1},$$

its inversion $\tilde{\varphi}_{\mathcal{S}}^{-1}: Q_{N,N}(\mathcal{S}) \rightarrow O(N) \setminus E_{N,N}(\mathcal{S})$ is given as ($\mathbf{V} \in Q_{N,N} := Q_{N,N}(\mathcal{S})$) $\tilde{\varphi}_{\mathcal{S}}^{-1}(\mathbf{V}) = (\mathbf{I} - \mathbf{V})(\mathbf{I} + \mathbf{V})^{-1} \mathcal{S}$.

Remark II.2 (On Def. II.1). For $\mathbf{V} \in Q_{N,N}$, $\mathbf{I} + \mathbf{V}$ is invertible because every eigenvalue of \mathbf{V} is pure imaginary [21]. A specialization $\tilde{\varphi}_{\mathcal{I}}$ of $\tilde{\varphi}_{\mathcal{S}}$ is the classical Cayley transform [20]. $Q_{N,N}$ is a vector space because it is closed under matrix addition and scalar multiplication. For every $\mathcal{S} \in O(N)$, $Q_{N,N}(\mathcal{S})$ is the same set as $Q_{N,N}$, however we distinguish them later from the viewpoint of the parametrization with \mathcal{S} .

One of the first applications of the Cayley transform to Problem I.1 with $p = N$ is found in [21] where Problem I.1 is relaxed to Problem I.2. Since $Q_{N,N}(\mathcal{S})$ is a vector space, we can enjoy for Problem I.2 various optimization methods designed specially on a vector space. However, the convergence speed of optimization algorithms combined with CP tends to slow down severely in the case where the minimizer $\mathbf{X}^* \in O(N)$ is close to $E_{N,N}(\mathcal{S})$ [3], [21]. Since it is a priori unclear whether \mathbf{X}^* is close to $E_{N,N}(\mathcal{S})$ or not, the center point \mathcal{S} was used as a hyper parameter in [3].

Through a quantitative analysis on the mobility of $\tilde{\varphi}_{\mathcal{S}}^{-1}$ [20, Eq. (7)], we recently revealed that $\tilde{\varphi}_{\mathcal{S}}^{-1}$ can be insensitive to the change on $Q_{N,N}(\mathcal{S})$ at points distant from zero, which is

certainly a major reason of the performance degradation of optimization algorithms combined with CP. We also proposed, for $p = N$, the Adaptive Localized Cayley Parametrization Technique (ALCP), which changes a center point \mathbf{S} to $\mathbf{S}' := \tilde{\varphi}_{\mathbf{S}'}^{-1}(\mathbf{V}_n)$ adaptively and then restarts iterative suppression of $f \circ \tilde{\varphi}_{\mathbf{S}'}^{-1}$ over $Q_{N,N}(\mathbf{S}')$ with the initial point $\mathbf{V}'_n := \tilde{\varphi}_{\mathbf{S}'} \circ \tilde{\varphi}_{\mathbf{S}'}^{-1}(\mathbf{V}_n) = \mathbf{0}$. More precisely, we change \mathbf{S} in the case where the latest estimate $\mathbf{V}_n \in Q_{N,N}(\mathbf{S})$ is distant from zero. Thanks to this adaptive strategy, a center point \mathbf{S} is no longer a hyperparameter because we can use the center point strategically not to lose a sufficient mobility in the update of the estimation.

III. ADAPTIVE LOCALIZED CAYLEY PARAMETRIZATION FOR STIEFEL MANIFOLD

We can apply ALCP for $O(N)$ to Problem I.1 with $p < N$ by replacing $\tilde{\varphi}_{\mathbf{S}}^{-1}$ and $SO(N)$ in (3) with $\Xi \circ \tilde{\varphi}_{\mathbf{S}}^{-1}$ and $\text{St}(p, N)$, where $\Xi: O(N) \rightarrow \text{St}(p, N): \mathbf{X} \rightarrow \mathbf{X}\mathbf{I}_{N \times p}$. However, the computational complexity for $\Xi \circ \tilde{\varphi}_{\mathbf{S}}^{-1}$ is $\mathcal{O}(N^3)$ that is expensive even in comparison with $\mathcal{O}(Np^2 + p^3)$ required for retractions of $\text{St}(p, N)$.

For a computationally efficient extension of ALCP to $\text{St}(p, N)$, we newly propose G-L²CT.

Definition III.1 (G-L²CT). Let $\mathbf{S} \in O(N)$, $E_{N,p}(\mathbf{S}) := \{\mathbf{X} \in \text{St}(p, N) \mid \det(\mathbf{I} + \mathbf{S}_{\text{le}}^T \mathbf{X}) = 0\}$, and

$$Q_{N,p}(\mathbf{S}) := Q_{N,p} := \left\{ \begin{bmatrix} \mathbf{A} & -\mathbf{B}^T \\ \mathbf{B} & \mathbf{0} \end{bmatrix} \middle| \begin{matrix} \mathbf{A}^T = -\mathbf{A} \in \mathbb{R}^{p \times p} \\ \mathbf{B} \in \mathbb{R}^{(N-p) \times p} \end{matrix} \right\} \subset Q_{N,N}.$$

The G-L²CT $\Phi_{\mathbf{S}}: \text{St}(p, N) \setminus E_{N,p}(\mathbf{S}) \rightarrow Q_{N,p}(\mathbf{S})$ centered at \mathbf{S} is defined by

$$(\mathbf{X} \in \text{St}(p, N) \setminus E_{N,p}(\mathbf{S})) \quad \Phi_{\mathbf{S}}(\mathbf{X}) = \begin{bmatrix} \mathbf{A} & -\mathbf{B}^T \\ \mathbf{B} & \mathbf{0} \end{bmatrix}$$

with $\mathbf{A} := (\mathbf{I} + \mathbf{S}_{\text{le}}^T \mathbf{X})^{-T} (\mathbf{X}^T \mathbf{S}_{\text{le}} - \mathbf{S}_{\text{le}}^T \mathbf{X}) (\mathbf{I} + \mathbf{S}_{\text{le}}^T \mathbf{X})^{-1} \in \mathbb{R}^{p \times p}$ and $\mathbf{B} := -\mathbf{S}_{\text{ri}}^T \mathbf{X} (\mathbf{I} + \mathbf{S}_{\text{le}}^T \mathbf{X})^{-1} \in \mathbb{R}^{(N-p) \times p}$.

Proposition III.2 (Inversion of G-L²CT). *The G-L²CT $\Phi_{\mathbf{S}}$ centered at $\mathbf{S} \in O(N)$ is a diffeomorphism between $\text{St}(p, N) \setminus E_{N,p}(\mathbf{S})$ and $Q_{N,p}(\mathbf{S})$, and its inversion $\Phi_{\mathbf{S}}^{-1}: Q_{N,p}(\mathbf{S}) \rightarrow \text{St}(p, N) \setminus E_{N,p}(\mathbf{S})$ is given by*

$$(\mathbf{V} \in Q_{N,p}(\mathbf{S})) \quad \Phi_{\mathbf{S}}^{-1}(\mathbf{V}) = \mathbf{S}(\mathbf{I} - \mathbf{V})(\mathbf{I} + \mathbf{V})^{-1} \mathbf{I}_{N \times p}.$$

We call $E_{N,p}(\mathbf{S})$ the singular point set of $\Phi_{\mathbf{S}}$. In the following, $\mathbf{A}(\mathbf{V}) \in \mathbb{R}^{p \times p}$ and $\mathbf{B}(\mathbf{V}) \in \mathbb{R}^{(N-p) \times p}$ denote the left upper and the left lower block matrices of $\mathbf{V} \in Q_{N,p}(\mathbf{S})$ respectively. When there is no possibility of confusion, we will use \mathbf{A} and \mathbf{B} instead of them respectively.

Theorem III.3 (Denseness of $\text{St}(p, N) \setminus E_{N,p}(\mathbf{S})$). *For $p < N$ and $\mathbf{S} \in O(N)$, it holds $\text{St}(p, N) \setminus E_{N,p}(\mathbf{S}) = \Xi(O(N) \setminus E_{N,N}(\mathbf{S}))$ and $\text{St}(p, N) \setminus E_{N,p}(\mathbf{S})$ is dense in $\text{St}(p, N)$.*

Remark III.4 (On G-L²CT).

- (a) (Specialization of $\Phi_{\mathbf{S}}$). A specialization of $\Phi_{\mathbf{S}}$ with $p = N$ and $\mathbf{S} = \mathbf{I}$ is the classical Cayley transform.
- (b) (Sound relaxation of Problem I.1). By the denseness of $\text{St}(p, N) \setminus E_{N,p}(\mathbf{S})$ in $\text{St}(p, N)$, Problem I.3 is a sound relaxation of Problem I.1 with $p < N$.

- (c) (Mobility of $\Phi_{\mathbf{S}}^{-1}$). As [20], for $\mathbf{V}, \mathbf{F} \in Q_{N,p}(\mathbf{S})$ and $\epsilon > 0$, we evaluate a mobility of $\Phi_{\mathbf{S}}^{-1}$ with $p < N$ at \mathbf{V} : $\|\Phi_{\mathbf{S}}^{-1}(\mathbf{V} + \epsilon \mathbf{F}) - \Phi_{\mathbf{S}}^{-1}(\mathbf{V})\|_F \leq 2\epsilon r(\mathbf{B}) \|\mathbf{F}\|_F + o(\epsilon)\|_F$, where $r(\mathbf{B}) := (1 + \|\mathbf{B}\|_2^2)^{1/2} / (1 + \sigma_{\min}^2(\mathbf{B}))$. Since $r(\mathbf{B})$ is bounded as $(1 + \|\mathbf{B}\|_2^2)^{-1/2} \leq r(\mathbf{B})$ with holding equality when $\sigma_{\min}(\mathbf{B}) = \|\mathbf{B}\|_2$, the mobility shows $\Phi_{\mathbf{S}}^{-1}$ can be insensitive to change on $Q_{N,p}(\mathbf{S})$ in the case where $\|\mathbf{B}\|_2$ is large. This leads to slow convergence of naive optimization methods applied to Problem I.3.

- (d) (Comparisons to other Cayley transforms on $\text{St}(p, N)$). In the literature, we have found different ideas (e.g., [24], [25]) for parametrization of $\text{St}(p, N)$ with the Cayley transform. [24] proposed to parameterize $\text{St}(p, N)$ with a diffeomorphism $\tilde{\Phi}_{\mathbf{S}}^{-1}: \tilde{Q}_{N,p}(\mathbf{S}) \rightarrow \text{St}(p, N) \setminus \tilde{E}_{N,p}(\mathbf{S})$: $\mathbf{V} \mapsto \tilde{\varphi}_{\mathbf{S}}^{-1}(\mathbf{V}) \mathbf{I}_{N \times p}$, where $\mathbf{S} \in O(N)$, $\tilde{E}_{N,p}(\mathbf{S}) := \{\mathbf{X} \in \text{St}(p, N) \mid \det([\mathbf{S}_{\text{le}}]_{\text{up}} + \mathbf{X}_{\text{up}}) = 0\}$ and $\tilde{Q}_{N,p}(\mathbf{S}) \subset Q_{N,p}(\mathbf{S})$ is the domain of $\tilde{\Phi}_{\mathbf{S}}^{-1}$. Although this map seems to be a natural extension of Def. II.1, the condition $\tilde{Q}_{N,p}(\mathbf{S}) \ni \mathbf{0}$ for some \mathbf{S} is not guaranteed^d, which implies $\tilde{Q}_{N,p}(\mathbf{S})$ is not a vector subspace of $Q_{N,p}(\mathbf{S})$. To see this fact, choose $\mathbf{S} \in O(N)$ such that $\det([\mathbf{S}_{\text{le}}]_{\text{up}}) = 0$, implying thus $\det([\mathbf{S}_{\text{le}}]_{\text{up}} + [\mathbf{S}_{\text{le}}]_{\text{up}}) = 0$ and $\mathbf{S}_{\text{le}} \in \tilde{E}_{N,p}(\mathbf{S})$. Assume also that $\mathbf{0} \in \tilde{Q}_{N,p}(\mathbf{S})$. Then, we have $\mathbf{S}_{\text{le}} = \tilde{\Phi}_{\mathbf{S}}^{-1}(\mathbf{0}) \in \text{St}(p, N) \setminus \tilde{E}_{N,p}(\mathbf{S})$, which is absurd. On the other hand, $\Phi_{\mathbf{S}}^{-1}$ for every $\mathbf{S} \in O(N)$ is a diffeomorphism between the vector space $Q_{N,p}(\mathbf{S})$ and $\text{St}(p, N) \setminus E_{N,p}(\mathbf{S})$, thus we can enjoy various optimization methods designed over a vector space for minimization of $f_{\mathbf{S}} := f \circ \Phi_{\mathbf{S}}^{-1}$ over $Q_{N,p}(\mathbf{S})$.

We consider the computational complexities for $\Phi_{\mathbf{S}}$ and $\Phi_{\mathbf{S}}^{-1}$. From the fact $(\mathbf{I} - \mathbf{V})(\mathbf{I} + \mathbf{V})^{-1} = 2(\mathbf{I} + \mathbf{V})^{-1} - \mathbf{I}$ for $\mathbf{V} \in Q_{N,p}(\mathbf{S})$ and the Schur complement formula [26, Sec. 0.8.5], the inversion $\Phi_{\mathbf{S}}^{-1}(\mathbf{V})$ can be expressed as $2(\mathbf{S}_{\text{le}} - \mathbf{S}_{\text{ri}} \mathbf{B}) \mathbf{M}^{-1} - \mathbf{S}_{\text{le}}$, where $\mathbf{M} := \mathbf{I} + \mathbf{A} + \mathbf{B}^T \mathbf{B} \in \mathbb{R}^{p \times p}$. The complexities for $\Phi_{\mathbf{S}}$ and $\Phi_{\mathbf{S}}^{-1}$ are $\mathcal{O}(N^2 p + p^3)$, which are dominated by computations of $\mathbf{S}_{\text{ri}}^T \mathbf{X}$ and $\mathbf{S}_{\text{ri}} \mathbf{B}$ respectively.

To reduce these complexities, we impose a structure

$$O_p(N) := \left\{ \Gamma_N(\mathbf{T}) := \begin{bmatrix} \mathbf{T} & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_{N-p} \end{bmatrix} \middle| \mathbf{T} \in O(p) \right\} \subset O(N) \quad (4)$$

to center points. For $\mathbf{S} \in O_p(N)$, we obtain $\mathbf{S}_{\text{ri}}^T \mathbf{X} = \mathbf{X}_{\text{lo}}$ and $\mathbf{S}_{\text{ri}} \mathbf{B} = [\mathbf{0}_p \quad \mathbf{B}^T]^T$, thus the complexities for $\Phi_{\mathbf{S}}$ and $\Phi_{\mathbf{S}}^{-1}$ centered at $\mathbf{S} \in O_p(N)$ can be reduced to $\mathcal{O}(Np^2 + p^3)$, which is competitive to that for retractions of $\text{St}(p, N)$ [22].

We present in Lemma III.5 a computational choice of center points $\mathbf{S} \in O_p(N)$ from $\mathbf{X} \in \text{St}(p, N)$ such that the performance degradation hardly appears around $\Phi_{\mathbf{S}}(\mathbf{X})$ as in Remark III.4 (c). Since the complexity for the Singular Value Decomposition (SVD) of $\mathbf{X}_{\text{up}} \in \mathbb{R}^{p \times p}$ is $\mathcal{O}(p^3)$, thus the choice of \mathbf{S} in Lemma III.5 is computationally efficient.

Lemma III.5 (Choice of center points). *Define a set $\tau(\mathbf{X}) := \{\mathbf{Q}_1 \mathbf{Q}_2^T \in O(p) \mid (\mathbf{Q}_1, \Sigma, \mathbf{Q}_2) \in \text{SVD}(\mathbf{X}_{\text{up}})\}$, where $\text{SVD}(\mathbf{X}_{\text{up}})$ is the set of all tuples $(\mathbf{Q}_1, \Sigma, \mathbf{Q}_2)$ such that*

^dAnother Cayley transform on $\text{St}(p, N)$ is found in [25], however its image Γ^{π} in [25, Theorem A] is not guaranteed as well to be a vector space.

Algorithm 1 Nesterov accelerated gradient based on ALCP

Input: $\mathbf{X}_0 \in \text{St}(p, N)$, $c \in (0, 2^{-1}]$, $T > 0$, L_{CP}
 $\mathbf{S}_0 \in \Gamma_N(\tau(\mathbf{X}_0))$, $\mathbf{Y}_0 = \mathbf{Z}_0 = \Phi_{\mathbf{S}_0}(\mathbf{X}_0)$, $k \leftarrow 0$
for $n = 0, 1, \dots, \hat{n} - 1$ **do**
 $\mathbf{Y}_{n+1} = \mathbf{Z}_n - \gamma_n \nabla f_{\mathbf{S}_n}(\mathbf{Z}_n)$ (γ_n is the output of Alg. 2)
 $\mathbf{X}_{n+1} = \Phi_{\mathbf{S}_n}^{-1}(\mathbf{Y}_{n+1})$
 if (6) or (7) are satisfied **then**
 $\mathbf{X}_{n+1} = \mathbf{X}_n$, $\mathbf{S}_{n+1} \in \Gamma_N(\tau(\mathbf{X}_{n+1}))$
 $\mathbf{Y}_{n+1} = \mathbf{Z}_{n+1} = \Phi_{\mathbf{S}_{n+1}}(\mathbf{X}_{n+1})$, $k \leftarrow 0$ (reset β_n)
 else
 $\mathbf{Z}_{n+1} = \mathbf{Y}_{n+1} + \beta_n(\mathbf{Y}_{n+1} - \mathbf{Y}_n)$, $\mathbf{S}_{n+1} = \mathbf{S}_n$,
 $k \leftarrow k + 1$ (see Remark IV.2 for $\beta_n \in \mathbb{R}$)
 end if
end for
Output: $\mathbf{X}_{\hat{n}}$

$\mathbf{X}_{up} = \mathbf{Q}_1 \Sigma \mathbf{Q}_2^T$ is an SVD of \mathbf{X}_{up} with $\mathbf{Q}_1, \mathbf{Q}_2 \in \text{O}(p)$ and a non-negative diagonal matrix $\Sigma \in \mathbb{R}^{p \times p}$. For $\mathbf{X} \in \text{St}(p, N)$, $\mathbf{T} \in \tau(\mathbf{X})$ and $\mathbf{S} := \Gamma_N(\mathbf{T})$, we have $\mathbf{X} \in \text{St}(p, N) \setminus E_{N,p}(\mathbf{S})$, $\mathbf{A}(\Phi_{\mathbf{S}}(\mathbf{X})) = \mathbf{0}$, $\|\mathbf{B}(\Phi_{\mathbf{S}}(\mathbf{X}))\|_2 \leq 1$, and $\|\Phi_{\mathbf{S}}(\mathbf{X})\|_F \leq \sqrt{2p}$.

From the discussion above, we propose a computationally efficient ALCP for Problem I.1 with $\Phi_{\mathbf{S}}^{-1}$. By letting $\mathcal{A}_n : Q_{N,p}(\mathbf{S}_n) \rightarrow Q_{N,p}(\mathbf{S}_n)$ be an update scheme at n th iteration and $T > 0$, ALCP generates estimates $(\mathbf{X}_n)_{n=0}^{\infty} \subset \text{St}(p, N)$ of solutions to Problem I.1, center points $(\mathbf{S}_n)_{n=0}^{\infty} \subset \text{O}_p(N)$ and $(\mathbf{V}_n)_{n=0}^{\infty} \subset Q_{N,p}$ iteratively as

$$\left. \begin{aligned} \mathbf{X}_{n+1} &:= \Phi_{\mathbf{S}_n}^{-1}(\mathcal{A}_n(\mathbf{V}_n)) \\ \mathbf{S}_{n+1} &:= \begin{cases} \mathbf{S}_n & (\|\mathcal{A}_n(\mathbf{V}_n)\|_F \leq T) \\ \mathbf{S} \in \Gamma_N(\tau(\mathbf{X}_{n+1})) & (\text{otherwise}) \end{cases} \\ \mathbf{V}_{n+1} &:= \Phi_{\mathbf{S}_{n+1}}(\mathbf{X}_{n+1}). \end{aligned} \right\} (5)$$

We note that $\mathbf{V}_{n+1} = \mathcal{A}_n(\mathbf{V}_n)$ when $\mathbf{S}_{n+1} = \mathbf{S}_n$. A simple specialization of ALCP is a steepest descent method by $\mathcal{A}_n(\mathbf{V}) := \mathbf{V} - \gamma_n \nabla f_{\mathbf{S}_n}(\mathbf{V})$ with a suitable stepsize $\gamma_n > 0$. In the same way as [3, Theorem 3.2], we obtain the gradient $\nabla f_{\mathbf{S}}(\mathbf{V}) = 2(\mathbf{W}_{\mathbf{S}}^f(\mathbf{V}) - \mathbf{W}_{\mathbf{S}}^f(\mathbf{V})^T)$, where

$$[\mathbf{W}_{\mathbf{S}}^f(\mathbf{V})]_{ij} := \begin{cases} 0 & (p+1 \leq i, j \leq N) \\ [\overline{\mathbf{W}}_{\mathbf{S}}^f(\mathbf{V})]_{ij} & (\text{otherwise}), \end{cases}$$

$\overline{\mathbf{W}}_{\mathbf{S}}^f(\mathbf{V}) := (\mathbf{I} + \mathbf{V})^{-1} \mathbf{I}_{N \times p} \nabla f(\mathbf{X})^T \mathbf{S} (\mathbf{I} + \mathbf{V})^{-1}$ and $\mathbf{X} := \Phi_{\mathbf{S}}^{-1}(\mathbf{V})$. For $\mathbf{T} \in \text{O}(p)$, $\mathbf{W}_{\Gamma_N(\mathbf{T})}^f$ can be simplified from the Schur complement formula as

$$\mathbf{W}_{\Gamma_N(\mathbf{T})}^f(\mathbf{V}) := \begin{bmatrix} \mathbf{P} & \mathbf{P}\mathbf{B}^T + \mathbf{M}^{-1}[\nabla f(\mathbf{X})]_{\text{lo}}^T \\ -\mathbf{B}\mathbf{P} & \mathbf{0} \end{bmatrix} \in \mathbb{R}^{N \times N},$$

where $\mathbf{P} := \mathbf{M}^{-1} \nabla f(\mathbf{X})^T [\mathbf{S}_p^T \quad -\mathbf{B}^T]^T \mathbf{M}^{-1} \in \mathbb{R}^{p \times p}$ and $\mathbf{M} := \mathbf{I} + \mathbf{A} + \mathbf{B}^T \mathbf{B} \in \mathbb{R}^{p \times p}$. Therefore, the complexity for $\nabla f_{\mathbf{S}}(\mathbf{V})$ with $\mathbf{S} \in \text{O}_p(N)$ is $\mathcal{O}(Np^2 + p^3)$.

Moreover, we evaluate a Lipschitz constant $L_{\text{CP}} > 0$ of $\nabla f_{\mathbf{S}}$ on $Q_{N,p}(\mathbf{S})$ under a mild assumption because the proposed NAG combined with ALCP in Section IV requires L_{CP} . Many applications of Problem I.1 satisfy our assumption, e.g., Sec. V, thus Lemma III.6 is useful for plugging optimization methods, exploiting the Lipschitz constant, into ALCP.

Lemma III.6 (Lipschitz constant L_{CP} of $\nabla f_{\mathbf{S}}$). *We assume*

Algorithm 2 Backtracking algorithm

Input: $c \in (0, 2^{-1}]$, $\rho \in (0, 1)$, $\gamma > 0$, L_{CP}
while γ satisfies the condition (6) and $\gamma \geq L_{\text{CP}}^{-1}$ **do**
 $\gamma \leftarrow \rho\gamma$
end while
Output: γ

that ∇f is Lipschitz continuous with $L > 0$ on $\text{St}(p, N)$, i.e., $(\mathbf{X}_1, \mathbf{X}_2 \in \text{St}(p, N)) \|\nabla f(\mathbf{X}_1) - \nabla f(\mathbf{X}_2)\|_F \leq L \|\mathbf{X}_1 - \mathbf{X}_2\|_F$. Let $\mu := \max \|\nabla f(\text{St}(p, N))\|_2^{\circ}$. Then, it holds with $L_{\text{CP}} := 8(L + \mu)$ for $\mathbf{S} \in \text{O}(N)$ and $\mathbf{V}_1, \mathbf{V}_2 \in Q_{N,p}(\mathbf{S})$ that $\|\nabla f_{\mathbf{S}}(\mathbf{V}_1) - \nabla f_{\mathbf{S}}(\mathbf{V}_2)\|_F \leq L_{\text{CP}} \|\mathbf{V}_1 - \mathbf{V}_2\|_F$.

IV. NESTEROV ACCELERATION BASED ON ALCP

To achieve a faster convergence speed for Problem I.1, by assuming a Lipschitz constant $L_{\text{CP}} > 0$ of $\nabla f_{\mathbf{S}}$ is available, we present a NAG with a restart scheme based on ALCP. In the proposed NAG illustrated in Alg. 1, we synchronize the change of center points $\mathbf{S} \in \text{O}_p(N)$ of ALCP in (5) and the adaptive restart scheme of NAG [16]. It is known that ripples in the sequence of the values of f in NAG lead to slow convergence of NAG, thus to avoid the ripples, adaptive restart schemes for NAG have been studied and their remarkable efficacies have been examined numerically [16], [27].

We change a center point for Alg. 1 when one of the following conditions hold true:

$$f_{\mathbf{S}_n}(\mathbf{Y}_{n+1}) > f_{\mathbf{S}_n}(\mathbf{Y}_n) - c\gamma_n \|\nabla f_{\mathbf{S}_n}(\mathbf{Z}_n)\|_F^2, \quad (6)$$

$$\|\mathbf{Y}_{n+1}\|_F > T \quad \text{and} \quad k > 0, \quad (7)$$

where $c \in (0, 2^{-1}]$, $T > 0$, $\gamma_n > 0$ is a stepsize and k is the number of iterations that the current center point has not been changed. The first criterion (6), which is proposed by [16], ensures for NAG to decrease monotonically the value of $f_{\mathbf{S}_n}$. The former of the second criterion (7) is derived from the mobility analysis of $\Phi_{\mathbf{S}}^{-1}$ in Remark III.4 (c). We need the latter condition $k > 0$ in (7) to work Alg. 1 well even if T is small. Indeed, if we do not use the condition $k > 0$, then the estimate \mathbf{X}_n has not been changed in a case where the former condition in (7) is satisfied at every iteration due to small T .

To find γ_n that (6) does not hold with, we employ a backtracking algorithm in Alg. 2. Since the existence of such a stepsize is not always guaranteed, we use L_{CP}^{-1} as a lower bound of acceptable stepsizes. When Alg. 1 restarts at n th iteration, (6) at $n + 1$ st iteration does not hold with $\gamma_{n+1} = L_{\text{CP}}^{-1}$ at least due to $\mathbf{Y}_{n+1} = \mathbf{Z}_{n+1}$ and the Lipschitz continuous of $\nabla f_{\mathbf{S}_{n+1}}$, thus Alg. 1 is well-defined.

Then, we also show that the gradient $\nabla f_{\mathbf{S}_n}(\mathbf{Y}_n)$ converges to zero from any initial point in Theorem IV.1.

Theorem IV.1 (Convergence analysis for Alg. 1). *Suppose that $(\mathbf{S}_n)_{n=0}^{\infty}$ and $(\mathbf{Y}_n)_{n=0}^{\infty}$ are generated by Alg. 1. Then we have $\lim_{n \rightarrow \infty} \|\nabla f_{\mathbf{S}_n}(\mathbf{Y}_n)\| = 0$.*

Remark IV.2. Since we can enjoy the linear structure of $Q_{N,p}(\mathbf{S})$, Alg. 1 essentially requires only computation of (2)

^oThe existence of μ is guaranteed by the compactness of $\text{St}(p, N)$ and the continuity of $\|\nabla f(\cdot)\|_2$.

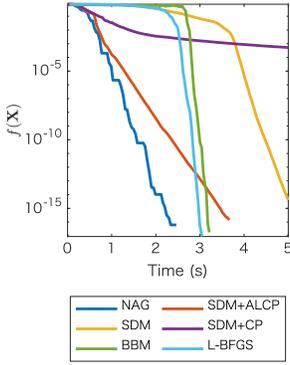


Fig. 1. $f(\mathbf{X})$ versus running time for the methods in the scenario of joint diagonalization.

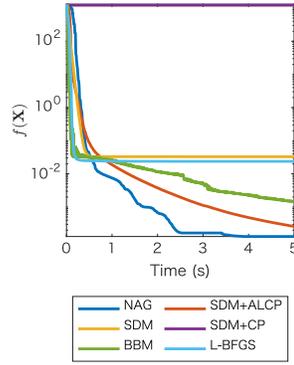


Fig. 2. $f(\mathbf{X})$ versus running time for the methods in the scenario of orthogonal Procrustes problem.

without requiring the inversion of neither the exponential map nor retractions unlike [14]–[16]. Moreover, Theorem IV.1 holds by employing any momentum parameter $\beta_n \in \mathbb{R}$.

V. NUMERICAL EXPERIMENTS

To demonstrate the performance of Alg. 1 (NAG), we compare NAG to the steepest descent method with ALCP (SDM+ALCP) and the standard methods for (1) based on retractions, e.g., the steepest descent method (SDM), the Barzilai-Borwein method (BBM) and L-BFGS method, provided by *Manopt* [23], in two scenarios: A. joint diagonalization ($p = N$) and B. orthogonal Procrustes problem ($p \ll N$). Moreover, to examine the avoidance of the performance degradation of CP, we test the steepest descent method with CP (SDM+CP), i.e., SDM+ALCP with $T = \infty$ in (5). The parameters of NAG are given as $\beta_n = (\alpha_n - 1)/\alpha_{n-1}$ with $\alpha_n = (1 + (1 + 4\alpha_{n-1}^2)^{1/2})/2$ for $n \geq 0$ and $\alpha_{-1} = 1$ as [13], and $T = 5$ for the problem A, $T = \sqrt{2p}$ for the problem B. For the backtracking algorithm, we employ $c = 2^{-13}$, $\rho = 2^{-1}$ and an initial stepsize $\gamma = 1$, which are the default of *Manopt*.

A. Joint diagonalization. We consider minimization of $f(\mathbf{X}) := \sum_{i=1}^I \|\text{off}(\mathbf{X}^T \mathbf{A}_i \mathbf{X})\|_F^2$ subject to $\mathbf{X} \in O(N)$ for symmetric matrices $\mathbf{A}_i \in \mathbb{R}^{N \times N}$, where $\text{off}(\mathbf{X}) := \mathbf{X} - \text{diag}(\mathbf{X})$ and $\text{diag}(\mathbf{X})$ denotes the diagonal matrix whose diagonal entries are those of \mathbf{X} . From $\nabla f(\mathbf{X}) = 4 \sum_{i=1}^I \mathbf{A}_i \mathbf{X} \text{off}(\mathbf{X}^T \mathbf{A}_i \mathbf{X})$ [1], we obtain $L_{CP} = 64(3 + \sqrt{N}) \sum_{i=1}^I \|\mathbf{A}_i\|_2^2$. Each \mathbf{A}_i is generated as $\mathbf{X}^{*T} \mathbf{A}_i \mathbf{X}^*$ with a random orthogonal matrix $\mathbf{X}^* \in O(N)$, where entries of a diagonal $\mathbf{A}_i \in \mathbb{R}^{N \times N}$ are uniformly chosen from $[0, 1]$.

B. Orthogonal Procrustes problem. We consider minimization of $f(\mathbf{X}) := \|\mathbf{C}\mathbf{X} - \mathbf{D}\|_F^2$ subject to $\mathbf{X} \in \text{St}(p, N)$ for $\mathbf{C} \in \mathbb{R}^{q \times N}$ and $\mathbf{D} \in \mathbb{R}^{q \times p}$. A closed-form solution to the problem is still a challenge [2]. From $\nabla f(\mathbf{X}) = 2\mathbf{C}^T(\mathbf{C}\mathbf{X} - \mathbf{D})$, we obtain $L_{CP} = 16(2\|\mathbf{C}^T \mathbf{C}\|_2 + \|\mathbf{C}^T \mathbf{D}\|_2)$. We generate an ill-conditioned matrix \mathbf{C} as [2, Ex. 4 in Sec. 4.3] ($m_1 = N/4, m_2 = N/2, m_3 = N/2 - 10, m_4 = 10$) and $\mathbf{D} := \mathbf{X}^* \mathbf{C}$ with a random $\mathbf{X}^* \in \text{St}(p, N)$.

Fig. 1 and Fig. 2 show the evolution of $f(\mathbf{X})$ of the methods with respect to running time (seconds) for solving the problem A with $N = 100$ and $I = 100$, and the problem B with $N = 1000$ and $p = q = 10$ respectively. In both experiments, we can see that SDM+ALCP outperforms SDM+CP by changing

center points adaptively and the convergence speed of NAG is faster than one of SDM+ALCP. Moreover, we can observe that NAG overwhelms the standard methods, e.g., BBM and L-BFGS, using a retraction in both experiments.

REFERENCES

- [1] M. Nikpour, J. H. Manton, and G. Hori, “Algorithms on the Stiefel manifold for joint diagonalisation,” in *ICASSP*, 2002.
- [2] J. B. Francisco and F. S. V. Bazán, “Nonmonotone algorithm for minimization on closed sets with applications to minimization on Stiefel manifolds,” *J. Comput. Appl. Math.*, vol. 236, no. 10, 2012.
- [3] K. Helfrich, D. Willmott, and Q. Ye, “Orthogonal recurrent neural networks with scaled Cayley transform,” in *ICML*, vol. 80, 2018.
- [4] P.-A. Absil, R. Mahony, and R. Sepulchre, *Optimization Algorithms on Matrix Manifolds*. Princeton University Press, 2008.
- [5] T. E. Abrudan, J. E. Eriksson, and V. Koivunen, “Steepest descent algorithms for optimization under unitary matrix constraint,” *IEEE Trans. Signal Proces.*, vol. 56, no. 3, 2008.
- [6] S. Fiori, T. Kaneko, and T. Tanaka, “Learning on the compact Stiefel manifold by a Cayley-transform-based pseudo-retraction map,” in *IJCNN*, 2012.
- [7] Z. Wen and W. Yin, “A feasible method for optimization with orthogonality constraints,” *Math. Program.*, vol. 142, no. 1, 2013.
- [8] A. Edelman, T. A. Arias, and S. T. Smith, “The geometry of algorithms with orthogonality constraints,” *SIAM J. Matrix Anal. Appl.*, vol. 20, no. 2, 1998.
- [9] W. Ring and B. Wirth, “Optimization methods on Riemannian manifolds and their application to shape space,” *SIAM J. Optim.*, vol. 22, no. 2, 2012.
- [10] H. Sato and T. Iwai, “A new, globally convergent Riemannian conjugate gradient method,” *Optimization*, vol. 64, no. 4, 2015.
- [11] W. Huang, K. A. Gallivan, and P.-A. Absil, “A Broyden class of quasi-Newton methods for Riemannian optimization,” *SIAM J. Optim.*, vol. 25, no. 3, 2015.
- [12] H. Kasai and B. Mishra, “Inexact trust-region algorithms on Riemannian manifolds,” in *NIPS*, 2018.
- [13] Y. Nesterov, “A method for solving the convex programming problem with convergence rate $o(1/k^2)$,” in *Dokl. Akad. nauk SSSR*, vol. 269, 1983.
- [14] Y. Liu, F. Shang, J. Cheng, H. Cheng, and L. Jiao, “Accelerated first-order methods for geodesically convex optimization on Riemannian manifolds,” in *NIPS*, 2017.
- [15] H. Zhang and S. Sra, “An estimate sequence for geodesically convex optimization,” in *ICML*, vol. 75, 2018.
- [16] J. W. Siegel, “Accelerated optimization with orthogonality constraints,” *arXiv preprint arXiv:1903.05204v3 (to appear in J. Comput. Math.)*, Oct 2019.
- [17] Y. Nesterov, *Introductory lectures on convex optimization: A basic course*. Springer Science & Business Media, 2013, vol. 87.
- [18] S. Ghadimi and G. Lan, “Accelerated gradient methods for nonconvex nonlinear and stochastic programming,” *Math. Program.*, vol. 156, 2016.
- [19] R. Zimmermann, “A matrix-algebraic algorithm for the Riemannian logarithm on the Stiefel manifold under the canonical metric,” *SIAM J. Matrix Anal. Appl.*, vol. 38, no. 2, 2017.
- [20] K. Kume and I. Yamada, “Adaptive localized Cayley parametrization technique for smooth optimization over the Stiefel manifold,” in *EU-SIPCO*, 2019.
- [21] I. Yamada and T. Ezaki, “An orthogonal matrix optimization by dual Cayley parametrization technique,” in *ICA*, 2003.
- [22] X. Zhu, “A Riemannian conjugate gradient method for optimization on the Stiefel manifold,” *Comput. Optim. Appl.*, vol. 67, no. 1, 2017.
- [23] N. Boumal, B. Mishra, P.-A. Absil, and R. Sepulchre, “Manopt, a Matlab toolbox for optimization on manifolds,” *J. Mach. Learn. Res.*, vol. 15, 2014.
- [24] C. Fraikin, K. Hüper, and P. V. Dooren, “Optimization over the Stiefel manifold,” in *PAMM*, vol. 7, no. 1, 2007.
- [25] E. Macías-Virgós, M. J. Pereira-Sáez, and D. Tanré, “Cayley transform on Stiefel manifolds,” *J. Geom. Phys.*, vol. 123, 2018.
- [26] R. A. Horn and C. R. Johnson, *Matrix analysis*. Cambridge university press, 2012.
- [27] B. O’Donoghue and E. Candès, “Adaptive restart for accelerated gradient schemes,” *Found. Comput. Math.*, vol. 15, no. 3, 2015.