

Vector Transport Free Riemannian LBFGS for Optimization on Symmetric Positive Definite Matrix Manifolds

Reza Godaz*

REZA.GODAZ@MAIL.UM.AC.IR

Department of Computer Engineering, Ferdowsi University of Mashhad, Mashhad, Iran

Benyamin Ghojogh*

BGHOJOGH@UWATERLOO.CA

Department of Electrical and Computer Engineering, University of Waterloo, ON, Canada

Reshad Hosseini

RESHAD.HOSSEINI@UT.AC.IR

Department of Electrical and Computer Engineering, University of Tehran, Tehran, Iran

Reza Monsefi

MONSEFI@UM.AC.IR

Department of Computer Engineering, Ferdowsi University of Mashhad, Mashhad, Iran

Fakhri Karray

KARRAY@UWATERLOO.CA

Mark Crowley

MCROWLEY@UWATERLOO.CA

Department of Electrical and Computer Engineering, University of Waterloo, ON, Canada

Editors: Vineeth N Balasubramanian and Ivor Tsang

Abstract

This work concentrates on optimization on Riemannian manifolds. The Limited-memory Broyden-Fletcher-Goldfarb-Shanno (LBFGS) algorithm is a commonly used quasi-Newton method for numerical optimization in Euclidean spaces. Riemannian LBFGS (RLBFGS) is an extension of this method to Riemannian manifolds. RLBFGS involves computationally expensive vector transports as well as unfolding recursions using adjoint vector transports. In this article, we propose two mappings in the tangent space using the inverse second root and Cholesky decomposition. These mappings make both vector transport and adjoint vector transport identity and therefore isometric. Identity vector transport makes RLBFGS less computationally expensive and its isometry is also very useful in convergence analysis of RLBFGS. Moreover, under the proposed mappings, the Riemannian metric reduces to Euclidean inner product, which is much less computationally expensive. We focus on the Symmetric Positive Definite (SPD) manifolds which are beneficial in various fields such as data science and statistics. This work opens a research opportunity for extension of the proposed mappings to other well-known manifolds.

Keywords: LBFGS, Riemannian optimization, positive definite manifolds, isometric vector transport, quasi-Newton's method

1. Introduction

Various numerical optimization methods have appeared, for the Euclidean spaces, which can be categorized into first-order and second-order methods (Nocedal and Wright, 2006). Examples for the former category are steepest descent and gradient descent and for the latter group are Newton's method. Computation of the Hessian matrix is usually expensive in

* The first two authors contributed equally to this work.

Newton’s method encouraging many practical problems to either use quasi-Newton’s methods for approximating the Hessian matrix or use non-linear conjugate gradient (Hestenes and Stiefel, 1952). The most well-known algorithm for quasi-Newton optimization is Broyden-Fletcher-Goldfarb-Shanno (BFGS) (Fletcher, 2013). Limited-memory BFGS (LBFGS) is a simplified version of BFGS which utilizes less memory (Nocedal, 1980; Liu and Nocedal, 1989). It has recursive unfoldings which approximate the descent directions in optimization.

Unlike Euclidean spaces in which the optimization direction lie in a linear coordinate system, Riemannian spaces have curvature in coordinates. Recently, extension of optimization methods from Euclidean spaces to Riemannian manifolds has been extensively noticed in the literature (Absil et al., 2009; Boumal, 2020; Hu et al., 2020). For example, Euclidean BFGS has been extended to Riemannian manifolds, named Riemannian BFGS (RBFGS), (Qi et al., 2010), its convergence has been proven (Ring and Wirth, 2012; Huang et al., 2015), and its properties have been analyzed in the literature (Seibert et al., 2013). As vector transport is computationally expensive in RBFGS, cautious RBFGS was proposed (Huang et al., 2016) which ignores the curvature condition in the Wolfe conditions (Wolfe, 1969) and only checks the Armijo condition (Armijo, 1966). Since the curvature condition guarantees that the approximation of Hessian remains positive definite, it compensates by checking a cautious condition (Li and Fukushima, 2001) before updating the approximation of Hessian. This cautious RBFGS has been used in the Manopt optimization toolbox (Boumal et al., 2014). Another approach is an extension of the Euclidean LBFGS to Riemannian manifolds, named Riemannian LBFGS (RLBFGS), using both Wolfe conditions (Wolfe, 1969) in linesearch can be found in (Sra and Hosseini, 2015, 2016; Hosseini and Sra, 2020). Some other direct extensions of Euclidean BFGS to Riemannian spaces exist (e.g., see (Ji, 2007, Chapter 7)).

In this paper, we address the computationally expensive parts of RLBFGS algorithm, which are computation of vector transports, their adjoints, and Riemannian metrics. To achieve this, we propose two mappings, in the tangent space, which make RLBFGS free of vector transport. We name the obtained algorithm Vector Transport Free (VTF)-RLBFGS. One mapping uses inverse second root and the other uses Cholesky decomposition which is very efficient (Golub and Van Loan, 2013). The proposed mappings make both vector transport and adjoint vector transport, which are used in RLBFGS, identity. This reduction of transports to identity makes optimization much less expensive computationally. Moreover, as the vector transports become identity, they are isometric which is a suitable property mostly used in the convergence proofs of RBFGS and RLBFGS algorithms (Ring and Wirth, 2012; Huang et al., 2015). Furthermore, under the proposed mappings, the Riemannian metric reduces to Euclidean inner product which is much less computationally expensive. In this paper, we concentrate on the Symmetric Positive Definite (SPD) manifolds (Sra and Hosseini, 2015, 2016; Bhatia, 2009) which are very useful in data science and machine learning, such as in mixture models (Hosseini and Sra, 2020; Hosseini and Mash’al, 2015). This paper opens a new research path for extension of the proposed mappings to other well-known manifolds such as Grassmann and Stiefel (Edelman et al., 1998).

The remainder of this paper is organized as follows. Section 2 reviews the notations and technical background on Euclidean LBFGS, Wolfe conditions, Riemannian LBFGS, and SPD manifold. The proposed mappings using inverse second root and Cholesky decompo-

sition are introduced in Sections 3 and 4, respectively. Simulation results are reported in Section 6. Finally, Section 7 concludes the paper and proposes the possible future directions.

2. Background and Notations

2.1. Euclidean BFGS and LBFGS

Consider minimization of the cost function $f(\boldsymbol{\Sigma})$ where the point $\boldsymbol{\Sigma}$ belongs to some domain. In Newton's method, the descent direction \mathbf{p}_k at the iteration k is calculated as $\mathbf{B}_k \mathbf{p}_k = -\nabla f(\boldsymbol{\Sigma}_k) \implies \mathbf{p}_k = -\mathbf{B}_k^{-1} \nabla f(\boldsymbol{\Sigma}_k)$ (Nocedal and Wright, 2006), where \mathbf{B}_k is the Hessian or approximation of Hessian and $\nabla f(\boldsymbol{\Sigma}_k)$ is the gradient of function at iteration k . The Euclidean BFGS method is a quasi-Newton's method which approximates the Hessian matrix as $\mathbf{B}_{k+1} := \mathbf{B}_k + (\mathbf{y}_k \mathbf{y}_k^\top) / (\mathbf{y}_k^\top \mathbf{s}_k) - (\mathbf{B}_k \mathbf{s}_k \mathbf{s}_k^\top \mathbf{B}_k^\top) / (\mathbf{s}_k^\top \mathbf{B}_k \mathbf{s}_k)$ (Fletcher, 2013; Nocedal and Wright, 2006), where $\mathbf{s}_k := \boldsymbol{\Sigma}_{k+1} - \boldsymbol{\Sigma}_k$ and $\mathbf{y}_k := \nabla f(\boldsymbol{\Sigma}_{k+1}) - \nabla f(\boldsymbol{\Sigma}_k)$. The descent direction is \mathbf{p}_k whose expression was provided above.

The Euclidean LBFGS calculates the descent direction recursively where it uses the approximation of the inverse Hessian as $\mathbf{H}_{k+1} := \mathbf{V}_k^\top \mathbf{H}_k \mathbf{V}_k + \rho_k \mathbf{s}_k \mathbf{s}_k^\top$ (Nocedal, 1980; Liu and Nocedal, 1989), where \mathbf{H}_k denotes the approximation of the inverse of the Hessian at iteration k , $\rho_k := 1 / (\mathbf{y}_k^\top \mathbf{s}_k)$, and $\mathbf{V}_k := \mathbf{I} - \rho_k \mathbf{y}_k \mathbf{s}_k^\top$ in which \mathbf{I} denotes the identity matrix. The LBFGS algorithm updates the approximation of the inverse of the Hessian matrix recursively and for that it always stores a memory window of pairs $\{\mathbf{y}_k, \mathbf{s}_k\}$ (Liu and Nocedal, 1989).

2.2. Linesearch and Wolfe Conditions

After finding the descent direction \mathbf{p}_k at each iteration k of optimization, one needs to know what step size α_k should be taken in that direction. Linesearch should be performed to find the largest step, for faster progress, satisfying Wolfe conditions which are $f(\boldsymbol{\Sigma}_k + \alpha_k \mathbf{p}_k) \leq f(\boldsymbol{\Sigma}_k) + c_1 \alpha_k \mathbf{p}_k^\top \nabla f(\boldsymbol{\Sigma}_k)$ and $-\mathbf{p}_k^\top \nabla f(\boldsymbol{\Sigma}_k + \alpha_k \mathbf{p}_k) \leq -c_2 \mathbf{p}_k^\top \nabla f(\boldsymbol{\Sigma}_k)$ (Wolfe, 1969), where the parameters $0 < c_1 < c_2 < 1$ are recommended to be $c_1 = 10^{-1}$ and $c_2 = 0.9$ (Nocedal and Wright, 2006). The former condition is the Armijo condition to check if cost decreases sufficiently (Armijo, 1966) while the latter is the curvature condition making sure that the slope is reduced sufficiently in a way that the approximation of Hessian remains positive definite. Note that there also exists a strong curvature condition, i.e., $|\mathbf{p}_k^\top \nabla f(\boldsymbol{\Sigma}_k + \alpha_k \mathbf{p}_k)| \leq c_2 |\mathbf{p}_k^\top \nabla f(\boldsymbol{\Sigma}_k)|$.

2.3. Riemannian Notations

Consider a Riemannian manifold denoted by \mathcal{M} . At every point $\boldsymbol{\Sigma} \in \mathcal{M}$, there is a tangent space to the manifold, denoted by $T_{\boldsymbol{\Sigma}} \mathcal{M}$. A tangent space includes tangent vectors. We denote a tangent vector by $\boldsymbol{\xi}$. For $\boldsymbol{\xi}, \boldsymbol{\eta} \in T_{\boldsymbol{\Sigma}} \mathcal{M}$, a metric on this manifold is the inner product defined on manifold and is denoted by $g_{\boldsymbol{\Sigma}}(\boldsymbol{\xi}, \boldsymbol{\eta})$. Note that the gradient of a cost function, whose domain is a manifold, is a tangent vector in the tangent space and is denoted by $\nabla f(\boldsymbol{\Sigma})$ for point $\boldsymbol{\Sigma} \in \mathcal{M}$. Vector transport is an operator which maps a tangent vector $\boldsymbol{\xi} \in T_{\boldsymbol{\Sigma}_1} \mathcal{M}$ from its tangent space at point $\boldsymbol{\Sigma}_1$ to another tangent space at another point $\boldsymbol{\Sigma}_2$. We denote this vector transport by $\mathcal{T}_{\boldsymbol{\Sigma}_1, \boldsymbol{\Sigma}_2}(\boldsymbol{\xi}) : T_{\boldsymbol{\Sigma}_1} \mathcal{M} \mapsto T_{\boldsymbol{\Sigma}_2} \mathcal{M}$. Now, consider a point $\boldsymbol{\Sigma}$ on a manifold \mathcal{M} and a descent direction $\boldsymbol{\xi}$ in the tangent space $T_{\boldsymbol{\Sigma}} \mathcal{M}$. The

retraction $\text{Ret}_{\Sigma}(\xi) : T_{\Sigma}\mathcal{M} \mapsto \mathcal{M}$ retracts or maps the direction ξ in the tangent space onto the manifold \mathcal{M} . The operator exponential map, denoted by $\text{Exp}_{\Sigma}(\xi)$, is also capable of this mapping by moving along the geodesic. One can use the second-order Taylor expansion of exponential map, which is positive-preserving (Jeuris et al., 2012) for the case of SPD manifold, for approximating the exponential map. In this paper, $\text{tr}(\cdot)$ denotes the trace of matrix and $\|\cdot\|_F$ denotes the Frobenius norm.

2.4. Riemannian BFGS and LBFGS

The Riemannian extension of Euclidean BFGS (RBFGS) (Qi et al., 2010; Ring and Wirth, 2012; Huang et al., 2015) performs updates of Hessian approximation by the \mathbf{B}_{k+1} (see Section 2.1) using Riemannian operators. RBFGS methods check both Wolfe conditions (see Section 2.2) which are computationally expensive. Cautious RBFGS (Huang et al., 2016) ignores the curvature condition and only checks the Armijo condition for linesearch. However, for ensuring that the Hessian approximation remains positive definite, it checks a cautious condition (Li and Fukushima, 2001) before updating the Hessian approximation.

Riemannian LBFGS (Sra and Hosseini, 2015, 2016; Hosseini and Sra, 2020) performs the recursions of LBFGS in the Riemannian space for finding the descent direction $\xi_k \in T_{\Sigma_k}\mathcal{M}$ and checks both Wolfe conditions in linesearch. In every iteration k of optimization, recursion starts with the direction $\mathbf{p} = -\nabla f(\Sigma_k)$. Let the recursive function $\text{GetDirection}(\mathbf{p}, k)$ returns the descent direction. Inside every step of this recursion, we have (Hosseini and Sra, 2020, Algorithm 3):

$$\tilde{\mathbf{p}} := \mathbf{p} - \rho_k \mathbf{g}_{\Sigma_k}(\mathbf{s}_k, \mathbf{p}) \mathbf{y}_k, \quad (1)$$

$$\hat{\mathbf{p}} := \mathcal{T}_{\Sigma_{k-1}, \Sigma_k}(\text{GetDirection}(\mathcal{T}_{\Sigma_{k-1}, \Sigma_k}^*(\tilde{\mathbf{p}}), k-1)), \quad (2)$$

$$\text{return } \xi_k := \hat{\mathbf{p}} - \rho_k \mathbf{g}_{\Sigma_k}(\mathbf{y}_k, \hat{\mathbf{p}}) \mathbf{s}_k + \rho_k \mathbf{g}_{\Sigma_k}(\mathbf{s}_k, \mathbf{s}_k) \mathbf{p}, \quad (3)$$

where, $\rho_k := 1/\mathbf{g}_{\Sigma_k}(\mathbf{y}_k, \mathbf{s}_k)$, and inspired by the introduced \mathbf{s}_k and \mathbf{y}_k for Euclidean spaces, we have:

$$\Sigma_{k+1} := \text{Exp}_{\Sigma_k}(\alpha_k \xi_k) \text{ or } \text{Ret}_{\Sigma_k}(\alpha_k \xi_k), \quad (4)$$

$$\mathbf{s}_{k+1} := \mathcal{T}_{\Sigma_k, \Sigma_{k+1}}(\alpha_k \xi_k), \quad (5)$$

$$\mathbf{y}_{k+1} := \nabla f(\Sigma_{k+1}) - \mathcal{T}_{\Sigma_k, \Sigma_{k+1}}(\nabla f(\Sigma_k)). \quad (6)$$

Note that the new point in every iteration is found by retraction, or an exponential map, of the searched point along the descent direction onto manifold. According to Eq. (2) in recursion and Eq. (5), RLFBFGS involves both adjoint vector transport and vector transport which are computationally expensive. Moreover, Eqs. (1) and (3) show that Riemannian metric is utilized many times inside recursions. Our proposed mappings simplify all vector transport, adjoint vector transport, and metric which are used in the RLFBFGS algorithm.

2.5. Symmetric Positive Definite (SPD) Manifold

Consider the SPD manifold (Sra and Hosseini, 2015, 2016) whose every point is a SPD matrix, i.e., $\Sigma \in \mathcal{M}$ and $\mathbb{S}_{++}^n \ni \Sigma \succ \mathbf{0}$. It can be shown that the tangent space to the SPD manifold is the space of symmetric matrices, i.e., $T_{\mathcal{M}}(\Sigma) \subset \mathbb{S}_{++}^n$ (Bhatia, 2009).

Table 1: Operators on SPD manifold under the proposed mappings.

Operator	No mapping
Metric, $g_{\Sigma}(\xi, \eta)$	$\text{tr}(\Sigma^{-1}\xi\Sigma^{-1}\eta)$
Gradient, $\nabla f(\Sigma)$	$\frac{1}{2}\Sigma(\nabla_E f(\Sigma) + (\nabla_E f(\Sigma))^{\top})\Sigma$
Exponential map, $\text{Exp}_{\Sigma}(\xi)$	$\Sigma \exp(\Sigma^{-1}\xi) = \Sigma^{\frac{1}{2}} \exp(\Sigma^{-\frac{1}{2}}\xi\Sigma^{-\frac{1}{2}}) \Sigma^{\frac{1}{2}}$
Vector transport, $\mathcal{T}_{\Sigma_1, \Sigma_2}(\xi)$	$\Sigma_2^{\frac{1}{2}} \Sigma_1^{-\frac{1}{2}} \xi \Sigma_1^{-\frac{1}{2}} \Sigma_2^{\frac{1}{2}}$ or $L_2 L_1^{-1} \xi L_1^{-\top} L_2^{\top}$
Approx. Euclidean retraction, $\text{Ret}_{\Sigma}(\xi)$	$\Sigma + \xi + \frac{1}{2}\xi\Sigma^{-1}\xi$
Operator	Mapping by inverse second root
Mapping	$\xi' := \Sigma^{-\frac{1}{2}} \xi \Sigma^{-\frac{1}{2}}$
Metric, $g'_{\Sigma}(\xi', \eta')$	$\text{tr}(\xi'\eta')$
Gradient, $\nabla' f(\Sigma)$	$\frac{1}{2}\Sigma^{\frac{1}{2}}(\nabla_E f(\Sigma) + (\nabla_E f(\Sigma))^{\top})\Sigma^{\frac{1}{2}}$
Exponential map, $\text{Exp}_{\Sigma}(\xi')$	$\Sigma^{\frac{1}{2}} \exp(\xi') \Sigma^{\frac{1}{2}}$
Vector transport, $\mathcal{T}'_{\Sigma_1, \Sigma_2}(\xi')$	ξ'
Approx. Euclidean retraction, $\text{Ret}_{\Sigma}(\xi')$	$\Sigma + \Sigma^{\frac{1}{2}}\xi'\Sigma^{\frac{1}{2}} + \frac{1}{2}\Sigma^{\frac{1}{2}}\xi'^2\Sigma^{\frac{1}{2}}$
Operator	Mapping by Cholesky decomposition
Mapping	$\xi' := L^{-1} \xi L^{-\top}$
Metric, $g'_{\Sigma}(\xi', \eta')$	$\text{tr}(\xi'\eta')$
Gradient, $\nabla' f(\Sigma)$	$\frac{1}{2}L^{\top}(\nabla_E f(\Sigma) + (\nabla_E f(\Sigma))^{\top})L$
Exponential map, $\text{Exp}_{\Sigma}(\xi')$	$\Sigma \exp(L^{-\top} \xi' L^{\top})$
Vector transport, $\mathcal{T}'_{\Sigma_1, \Sigma_2}(\xi')$	ξ'
Approx. Euclidean retraction, $\text{Ret}_{\Sigma}(\xi')$	$\Sigma + L \xi' L^{\top} + \frac{1}{2}L \xi'^2 L^{\top}$

In this paper, we focus on SPD manifolds which are widely used in data science. The operators for metric, gradient, exponential map, and vector transport on SPD manifolds are listed in Table 1 (Sra and Hosseini, 2016; Hosseini and Sra, 2020). In this table, $\nabla_E f(\Sigma)$ denotes the Euclidean gradient and L_1 and L_2 are the lower-triangular matrices in Cholesky decomposition of points Σ_1 and Σ_2 , respectively.

3. Vector Transport Free Riemannian LBFGS Using Mapping by Inverse Second Root

We propose two mappings on tangent vectors in the tangent space where the first mapping is by inverse second root. Our first proposed mapping in the tangent space of every point $\Sigma \in \mathcal{M}$ is:

$$\xi' := \Sigma^{-\frac{1}{2}} \xi \Sigma^{-\frac{1}{2}} \implies \xi = \Sigma^{\frac{1}{2}} \xi' \Sigma^{\frac{1}{2}}, \quad (7)$$

where the mapped tangent vector still remains in the tangent space, i.e. $\xi, \xi' \in T_{\Sigma}\mathcal{M} \subset \mathbb{S}_{++}^n$. It is important the proposed mapping is bijective and keeps the tangent vector in the tangent space while it simplifies vector transport, adjoint vector transport, and metric.

Under the proposed mapping (7), the Riemannian operators on a SPD manifold are modified and mostly simplified. These operators are listed in Table 1. In the following, we provide proofs for these modifications.

Proposition 1 *After mapping (7), we have*

- *Metric: the metric on SPD manifold is reduced to the Euclidean inner product, i.e., $g_{\Sigma}(\xi', \eta') = \text{tr}(\xi' \eta')$.*
- *Gradient: the gradient on a SPD manifold is changed to $\nabla' f(\Sigma) = \frac{1}{2} \Sigma^{\frac{1}{2}} (\nabla_E f(\Sigma) + (\nabla_E f(\Sigma))^{\top}) \Sigma^{\frac{1}{2}}$, where $\nabla_E f(\Sigma)$ denotes the Euclidean gradient.*
- *Vector transport: the vector transport is changed to identity, i.e., $\mathcal{T}'_{\Sigma_1, \Sigma_2}(\xi') = \xi'$, hence, optimization becomes vector transport free.*
- *Exponential map: the exponential map on a SPD manifold becomes $\text{Exp}_{\Sigma}(\xi') = \Sigma^{\frac{1}{2}} \exp(\xi') \Sigma^{\frac{1}{2}}$.*
- *Adjoint vector transport: the adjoint of vector transport on a SPD manifold remains the same. In other words, if before mapping we have the definition of adjoint vector transport as $g_{\Sigma_1}(\xi, \mathcal{T}_{\Sigma_1, \Sigma_2}^* \eta) = g_{\Sigma_2}(\mathcal{T}_{\Sigma_1, \Sigma_2} \xi, \eta)$, $\forall \xi \in T_{\Sigma_1} \mathcal{M}, \forall \eta \in T_{\Sigma_2} \mathcal{M}$ (Ring and Wirth, 2012), we will have $g_{\Sigma_1}(\xi', \mathcal{T}_{\Sigma_1, \Sigma_2}^* \eta') = g_{\Sigma_2}(\mathcal{T}'_{\Sigma_1, \Sigma_2} \xi', \eta')$, $\forall \xi' \in T_{\Sigma_1} \mathcal{M}, \forall \eta' \in T_{\Sigma_2} \mathcal{M}$.*
- *Retraction: the approximation of Euclidean retraction, using second-order Taylor expansion, on a SPD manifold becomes $\text{Ret}_{\Sigma}(\xi') = \Sigma + \Sigma^{\frac{1}{2}} \xi' \Sigma^{\frac{1}{2}} + \frac{1}{2} \Sigma^{\frac{1}{2}} \xi'^2 \Sigma^{\frac{1}{2}}$.*

Proof

• **Metric:** $g_{\Sigma}(\xi', \eta') \stackrel{(a)}{=} \text{tr}(\Sigma^{-1} \xi \Sigma^{-1} \eta) = \text{tr}(\Sigma^{-\frac{1}{2}} \Sigma^{-\frac{1}{2}} \xi \Sigma^{-\frac{1}{2}} \Sigma^{-\frac{1}{2}} \eta) \stackrel{(b)}{=} \text{tr}(\underbrace{\Sigma^{-\frac{1}{2}} \xi \Sigma^{-\frac{1}{2}}}_{=\xi'}) \underbrace{\Sigma^{-\frac{1}{2}} \eta \Sigma^{-\frac{1}{2}}}_{=\eta'} \stackrel{(7)}{=} \text{tr}(\xi' \eta')$ where (a) is because of definition of metric on SPD manifolds (see Table 1) and (b) is thanks to the cyclic property of trace.

• **Gradient:** $\nabla' f(\Sigma) \stackrel{(7)}{=} \Sigma^{-\frac{1}{2}} \nabla f(\Sigma) \Sigma^{-\frac{1}{2}} \stackrel{(a)}{=} \frac{1}{2} \Sigma^{-\frac{1}{2}} \Sigma (\nabla_E f(\Sigma) + (\nabla_E f(\Sigma))^{\top}) \Sigma \Sigma^{-\frac{1}{2}} = \frac{1}{2} \Sigma^{\frac{1}{2}} (\nabla_E f(\Sigma) + (\nabla_E f(\Sigma))^{\top}) \Sigma^{\frac{1}{2}}$, where (a) is due to definition of gradient on a SPD manifold (see Table 1).

• **Vector transport:** $\mathcal{T}'_{\Sigma_1, \Sigma_2}(\xi') \stackrel{(7)}{=} \Sigma_2^{-\frac{1}{2}} \mathcal{T}_{\Sigma_1, \Sigma_2}(\xi) \Sigma_2^{-\frac{1}{2}} \stackrel{(a)}{=} \underbrace{\Sigma_2^{-\frac{1}{2}} (\Sigma_2^{\frac{1}{2}})}_{=I} \Sigma_1^{-\frac{1}{2}} \xi \Sigma_1^{-\frac{1}{2}} \underbrace{(\Sigma_2^{\frac{1}{2}} \Sigma_2^{-\frac{1}{2}})}_{=I} = \Sigma_1^{-\frac{1}{2}} \xi \Sigma_1^{-\frac{1}{2}} \stackrel{(7)}{=} \xi'$, where (a) is due to definition of vector transport on a SPD manifold (see Table 1).

• **Exponential map:** $\text{Exp}_{\Sigma}(\xi') \stackrel{(a)}{=} \Sigma \exp(\Sigma^{-1} \xi) \stackrel{(b)}{=} \Sigma^{\frac{1}{2}} \exp(\Sigma^{-\frac{1}{2}} \xi \Sigma^{-\frac{1}{2}}) \Sigma^{\frac{1}{2}} \stackrel{(7)}{=} \Sigma^{\frac{1}{2}} \exp(\xi') \Sigma^{\frac{1}{2}}$ where (a) is shown in (Sra and Hosseini, 2015, Eq. 3.3) and (b) is shown in (Sra and Hosseini, 2015, Eq. 3.2). Also see Table 1.

- **Retraction:** $\text{Ret}_{\Sigma}(\xi') \stackrel{(a)}{=} \Sigma + \xi + \frac{1}{2}\xi\xi^{-1}\xi \stackrel{(7)}{=} \Sigma + \Sigma^{\frac{1}{2}}\xi'\Sigma^{\frac{1}{2}} + \frac{1}{2}\Sigma^{\frac{1}{2}}\xi'\underbrace{\Sigma^{\frac{1}{2}}\Sigma^{-1}\Sigma^{\frac{1}{2}}}_{=I}\xi'\Sigma^{\frac{1}{2}} = \Sigma + \Sigma^{\frac{1}{2}}\xi'\Sigma^{\frac{1}{2}} + \frac{1}{2}\Sigma^{\frac{1}{2}}\xi'^2\Sigma^{\frac{1}{2}}$, where (a) is because of approximation of Euclidean retraction, using the second-order Taylor expansion, on SPD manifolds (see Table 1). ■

4. Vector Transport Free Riemannian LBFGS Using Mapping by Cholesky Decomposition

Our second proposed mapping is by Cholesky decomposition which is very efficient computationally. Consider the Cholesky decomposition of point $\Sigma \in \mathcal{M}$ (Golub and Van Loan, 2013):

$$\mathbf{0} \prec \Sigma = \mathbf{L}\mathbf{L}^{\top} \implies \Sigma^{-1} = \mathbf{L}^{-\top}\mathbf{L}^{-1}, \quad (8)$$

where $\mathbf{L} \in \mathbb{R}^{n \times n}$ is the lower-triangular matrix in Cholesky decomposition. It is noteworthy that many of the MATLAB matrix multiplication operators, which the Manopt toolbox (Boumal et al., 2014) also uses, apply Cholesky decomposition internally due to its efficiency.

In the tangent space of every point $\Sigma \in \mathcal{M}$, the proposed mapping is:

$$\xi' := \mathbf{L}^{-1}\xi\mathbf{L}^{-\top} \implies \xi = \mathbf{L}\xi'\mathbf{L}^{\top}, \quad (9)$$

where $\xi, \xi' \in T_{\Sigma}\mathcal{M} \subset \mathbb{S}_{++}^n$. Note that, similar to the previous mapping, the tangent matrix is still symmetric under this mapping; hence, it remains in the tangent space of the SPD manifold (Bhatia, 2009). Similar to the previous mapping, under the second proposed mapping (7), the Riemannian operators on SPD manifold are simplified. These operators can be found in Table 1. In the following, we provide proofs for these operators.

Proposition 2 *After mapping (9), we have:*

- *Metric: the metric on a SPD manifold is reduced to the Euclidean inner product, i.e., $g_{\Sigma}(\xi', \eta') = \text{tr}(\xi'\eta')$.*
- *Gradient: the gradient on SPD manifold is changed to $\nabla'f(\Sigma) = \frac{1}{2}\mathbf{L}^{\top}(\nabla_E f(\Sigma) + (\nabla_E f(\Sigma))^{\top})\mathbf{L}$, where $\nabla_E f(\Sigma)$ denotes the Euclidean gradient.*
- *Vector transport: the vector transport is changed to identity, i.e., $\mathcal{T}'_{\Sigma_1, \Sigma_2}(\xi') = \xi'$, hence, optimization becomes vector transport free.*
- *Exponential map: the exponential map becomes $\text{Exp}_{\Sigma}(\xi') = \Sigma \exp(\mathbf{L}^{-\top}\xi'\mathbf{L}^{\top})$.*
- *Adjoint vector transport: the adjoint vector transport becomes $g_{\Sigma_1}(\xi', \mathcal{T}'_{\Sigma_1, \Sigma_2}^* \eta') = g_{\Sigma_2}(\mathcal{T}'_{\Sigma_1, \Sigma_2} \xi', \eta')$, $\forall \xi' \in T_{\Sigma_1}\mathcal{M}, \forall \eta' \in T_{\Sigma_2}\mathcal{M}$ as we had in Proposition 1.*
- *Retraction: the approximation of Euclidean retraction by second-order Taylor expansion on a SPD manifold becomes $\text{Ret}_{\Sigma}(\xi') = \Sigma + \mathbf{L}\xi'\mathbf{L}^{\top} + \frac{1}{2}\mathbf{L}\xi'^2\mathbf{L}^{\top}$.*

Proof

- **Metric:** $g_{\Sigma}(\xi', \eta') \stackrel{(a)}{=} \text{tr}(\Sigma^{-1}\xi\Sigma^{-1}\eta) = \text{tr}(\mathbf{L}^{-\top}\mathbf{L}^{-1}\xi\mathbf{L}^{-\top}\mathbf{L}^{-1}\eta) \stackrel{(b)}{=} \underbrace{\text{tr}(\mathbf{L}^{-1}\xi\mathbf{L}^{-\top}}_{=\xi'} \underbrace{\mathbf{L}^{-1}\eta\mathbf{L}^{-\top}}_{=\eta'} \stackrel{(9)}{=} \text{tr}(\xi'\eta')$, where (a) is because of definition of metric on SPD manifolds (see Table 1) and (b) is thanks to the cyclic property of trace.
- **Gradient:** $\nabla' f(\Sigma) \stackrel{(9)}{=} \mathbf{L}^{-1}\nabla f(\Sigma)\mathbf{L}^{-\top} \stackrel{(a)}{=} \frac{1}{2}\mathbf{L}^{-1}\Sigma(\nabla_E f(\Sigma) + (\nabla_E f(\Sigma))^{\top})\Sigma\mathbf{L}^{-\top} \stackrel{(8)}{=} \frac{1}{2}\mathbf{L}^{-1}\mathbf{L}\mathbf{L}^{\top}(\nabla_E f(\Sigma) + (\nabla_E f(\Sigma))^{\top})\mathbf{L}\mathbf{L}^{\top}\mathbf{L}^{-\top} = \frac{1}{2}\mathbf{L}^{\top}(\nabla_E f(\Sigma) + (\nabla_E f(\Sigma))^{\top})\mathbf{L}$, where (a) is due to definition of gradient on a SPD manifold (see Table 1).
- **Vector transport:** $\mathcal{T}'_{\Sigma_1, \Sigma_2}(\xi') \stackrel{(9)}{=} \mathbf{L}_2^{-1}\mathcal{T}_{\Sigma_1, \Sigma_2}(\xi)\mathbf{L}_2^{-\top} \stackrel{(a)}{=} \mathbf{L}_2^{-1}(\mathbf{L}_2\mathbf{L}_1^{-1}\xi\mathbf{L}_1^{-\top}\mathbf{L}_2^{\top})\mathbf{L}_2^{-\top} = \mathbf{L}_1^{-1}\xi\mathbf{L}_1^{-\top} \stackrel{(9)}{=} \xi'$, where (a) is due to definition of vector transport on SPD manifold (see Table 1).
- **Exponential map:** The exponential map is changed to $\text{Exp}_{\Sigma}(\xi') \stackrel{(a)}{=} \Sigma \exp(\Sigma^{-1}\xi) \stackrel{(b)}{=} \Sigma \exp(\mathbf{L}^{-\top}\mathbf{L}^{-1}\mathbf{L}\xi\mathbf{L}^{\top}) = \Sigma \exp(\mathbf{L}^{-\top}\xi'\mathbf{L}^{\top})$, where (a) is shown in (Sra and Hosseini, 2015, Eq. 3.3) and (b) is because of Eqs. (8) and (9).
- **Retraction:** $\text{Ret}_{\Sigma}(\xi') \stackrel{(a)}{=} \Sigma + \xi + \frac{1}{2}\xi\Sigma^{-1}\xi \stackrel{(b)}{=} \Sigma + \mathbf{L}\xi'\mathbf{L}^{\top} + \frac{1}{2}\mathbf{L}\xi'\mathbf{L}^{\top}\mathbf{L}^{-\top}\mathbf{L}^{-1}\mathbf{L}\xi'\mathbf{L}^{\top} = \Sigma + \mathbf{L}\xi'\mathbf{L}^{\top} + \frac{1}{2}\mathbf{L}\xi'^2\mathbf{L}^{\top}$, where (a) is because of approximation of Euclidean retraction, using second-order Taylor expansion, on SPD manifolds (see Table 1), and (b) is because of Eqs. (8) and (9). ■

Noticing Eq. (9), the approximation of retraction under mapping by Cholesky decomposition (see Proposition 2), can be restated as $\text{Ret}_{\Sigma}(\xi') = 0.5\Sigma + 0.5\mathbf{L}(\mathbf{I} + \xi')^2\mathbf{L}^{\top}$. Defining $\Psi := \mathbf{L}(\mathbf{I} + \xi')$ restates this retraction as $\text{Ret}_{\Sigma}(\xi') = 0.5\Sigma + 0.5\Psi\Psi^{\top}$ because ξ' is symmetric. The term $\Psi\Psi^{\top}$ is very efficient and fast to compute because it is symmetric.

5. Analytical Discussion and Complexity Analysis

Corollary 3 *Propositions 1 and 2 show that under mapping (7) or (9), both vector transport and adjoint vector transport are identity. As these transforms become identity, they also become isometric because inner products of vectors do not change under these transforms. As they are identity, these transforms also preserve the length of vectors under the proposed mappings.*

Propositions 1 and 2 and Corollary 3 show the two proposed mappings simplify vector transport and adjoint vector transport to isometric identity and reduce the Riemannian metric to Euclidean inner product. These reductions and simplifications reduce computations significantly during optimization on the manifold. The VTF-RLBFGS algorithm using either of the proposed mappings is shown in Algorithm 1. As the algorithm shows, at every new point Σ_k , the entire parameters are computed in the paradigm of mapping because the manifold operators, calculated as in Table 1, are in that paradigm. This simplifies operators such as metric and removes vector transports from RL BFGS. In case the Riemannian gradient is given directly by the user to RL BFGS, the Riemannian gradient,

Algorithm 1: The VTF-RLBFGS algorithm

Input: Initial point Σ_0

$$\mathbf{H}_0 := \frac{1}{\sqrt{g'_{\Sigma_0}(\nabla' f(\Sigma_0), \nabla' f(\Sigma_0))}} \mathbf{I}$$

for $k = 0, 1, \dots$ **do**

 Compute $\nabla' f(\Sigma_k)$ from Euclidean gradient by one of the mappings in Table 1

$$\xi'_k := \text{GetDirection}(-\nabla' f(\Sigma_k), k)$$

 $\alpha_k :=$ Line search with Wolfe conditions

$$\Sigma_{k+1} := \text{Exp}_{\Sigma_k}(\alpha_k \xi'_k) \text{ or } \text{Ret}_{\Sigma_k}(\alpha_k \xi'_k)$$

$$s'_{k+1} := \alpha_k \xi'_k$$

$$y'_{k+1} := \nabla' f(\Sigma_{k+1}) - \nabla' f(\Sigma_k)$$

$$\mathbf{H}_{k+1} := \frac{g'_{\Sigma_{k+1}}(s'_{k+1}, y'_{k+1})}{g'_{\Sigma_{k+1}}(y'_{k+1}, y'_{k+1})}$$

 Store y'_{k+1} , s'_{k+1} , $g'_{\Sigma_{k+1}}(s'_{k+1}, y'_{k+1})$, $g'_{\Sigma_{k+1}}(s'_{k+1}, s'_{k+1})$, and \mathbf{H}_{k+1}
end
return Σ_{k+1}
Function $\text{GetDirection}(p', k)$
if $k > 0$ **then**

$$\rho_k := \frac{1}{g'_{\Sigma_k}(y'_k, s'_k)}$$

$$\tilde{p}' := p' - \rho_k g'_{\Sigma_k}(s'_k, p') y'_k$$

$$\hat{p}' := \text{GetDirection}(\tilde{p}', k - 1)$$

$$\text{return } \hat{p}' - \rho_k g'_{\Sigma_k}(y'_k, \hat{p}') s'_k + \rho_k g'_{\Sigma_k}(s'_k, s'_k) p'$$

else

$$\text{return } \mathbf{H}_0 p'$$

end

which is in the tangent space, should be mapped explicitly by Eqs. (7) and (9) at every iteration. However, if the Riemannian gradient is calculated from the Euclidean gradient, it should not be mapped explicitly, since it is already in the paradigm of mapping implicitly because of the used operators of Table 1 in that paradigm.

Lemma 4 *Vector transport is valid under both proposed mappings (7) and (9) because they preserve the properties of vector transport.*

Proof A valid vector transport should satisfy three properties (Hosseini and Sra, 2020) (also see (Absil et al., 2009, Definition 8.1.1) and (Boumal, 2020, Definition 10.62)):

(1) The following property $\exists v \in T_{\Sigma_1} \mathcal{M} : \mathcal{T}_{\Sigma_1, \Sigma_2}(\xi) \in T_{\text{Ret}_{\Sigma}(v)} \mathcal{M}, \forall \xi \in T_{\Sigma_1} \mathcal{M}$ holds because the tangent space $T_{\text{Ret}_{\Sigma}(v)} \mathcal{M}$ is isomorphic to the space of symmetric matrices and the identity vector transports return the tangent vector itself which is in the space of symmetric matrices.

(2) The vector transport of a tangent vector from one point to itself should be the same tangent vector. This holds because vector transport is identity under the proposed mappings: $\mathcal{T}_{\Sigma, \Sigma}(\xi) = \xi, \forall \xi \in T_{\Sigma_1} \mathcal{M}$.

(3) The vector transport $\mathcal{T}_{\Sigma_1, \Sigma_2}(\xi)$ should be linear which is because it is equal to ξ under the proposed mappings.

As after applying the mappings, the vector transport holds the three above properties, it is a valid transport. Q.E.D. \blacksquare

Proposition 5 *The time complexity of the recursion part in RLBFGS is improved from $\Theta(mn^3)$ to $\Theta(mn^2)$ after mapping (7) or (9), where m is the memory limit, i.e., the maximum number of recursions in RLBFGS (proof is available in Supplementary Material). This time improvement shows off better in problems whose computation of gradients in Eq. (6), or line \mathbf{y}'_{k+1} in Algorithm 1, is not dominant in complexity.*

6. Simulations

In this section, we evaluate the effectiveness of the proposed mappings, i.e. (7) and (9), in the tangent space. Here, we show that these mappings often improve the performance and speed of RLBFGS. The code of this article is available in <https://github.com/bghojogh/LBFGS-Vector-Transport-Free>. In our reports, we denote the proposed Vector Transport Free (VTF) RLBFGS with VTF-RLBFGS where ISR and Cholesky (or Chol.) stand for VTF-RLBFGS using mapping by inverse second root and Cholesky decomposition, respectively. For RLBFGS, with and without the proposed mappings, we use both Wolfe linesearch conditions. The programming language, used for experiments, was MATLAB and the hardware was Intel Core-i7 CPU with the base frequency 2.20 GHz and 12 GB RAM. For every experiment, we performed optimization for ten runs and the reported results are the average of performances over the runs. We evaluated our mappings with various application problems, explained below.

6.1. Gaussian Mixture Model

- **Formulation:** An optimization problem, which we selected for evaluation, is the Riemannian optimization for Gaussian Mixture Model (GMM) without the use of expectation maximization. We employ RLBFGS with and without the proposed mappings for fitting the GMM problem whose algorithm can be found in (Hosseini and Sra, 2020; Hosseini and Mash'al, 2015). This is a suitable problem for evaluation of the proposed mappings because the covariance matrices are SPD (Bhatia, 2009). For this, we minimize the negative log-likelihood of GMM where the covariance matrix is constrained to belong to the SPD matrix manifold (Hosseini and Sra, 2020). For n -dimensional GMM, the optimization problem is:

$$\begin{aligned} & \underset{\{\alpha_j, \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j\}_{j=1}^K}{\text{minimize}} && - \sum_{i=1}^N \log \left(\sum_{j=1}^K \alpha_j \mathcal{N}(\mathbf{x}_i; \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j) \right), \\ & \text{subject to} && \boldsymbol{\Sigma}_j \in \mathcal{M} = \mathbb{S}_{++}^n, \quad \forall j \in \{1, \dots, K\}, \end{aligned} \tag{10}$$

where N denotes the sample size, K denotes the number of components in mixture model, and α_j , $\boldsymbol{\mu}_j$, and $\boldsymbol{\Sigma}_j$ are the mixing probability, mean, and covariance of the j -th component, respectively. We use the same reformulation trick of (Hosseini and Sra, 2020) to reformulate the cost function of (10). The mixture parameters were initialized using K-means++

Table 2: Comparison of average results over ten runs where exponential map is used in algorithms and $K \in \{2, 5\}$, $n \in \{2, 10\}$, $N = 10n^2 \in \{40, 1000\}$. The #iters, conv, iter, diff, and std are short for number of iterations, convergence, iteration, difference, and standard deviation, respectively.

K	n	Separation	Algorithm	#iters	conv. time	time diff. std	iter. time	last cost
2	2	Low	VTF (ISR)	53.100±18.248	68.380±52.834	13.405	1.140±0.504	0.364±0.444
			VTF (Chol.)	52.100±17.866	61.096±48.433	17.443	1.030±0.463	0.364±0.444
			RLBFGS	52.500±16.595	65.890±49.208	–	1.125±0.458	0.364±0.444
		Mid	VTF (ISR)	56.400±21.813	76.124±69.189	10.712	1.150±0.574	0.657±0.344
			VTF (Chol.)	54.400±19.156	68.456±64.290	48.016	1.099±0.504	0.638±0.333
			RLBFGS	57.700±19.844	81.957±64.463	–	1.250±0.550	0.657±0.344
		High	VTF (ISR)	25.500±3.064	9.518±2.922	0.333	0.366±0.067	0.341±0.371
			VTF (Chol.)	26.000±4.738	10.254±5.468	3.132	0.377±0.108	0.341±0.371
			RLBFGS	25.500±3.064	10.054±3.141	–	0.386±0.073	0.341±0.371
	10	Low	VTF (ISR)	77.600±54.175	235.615±466.119	12.040	1.936±1.751	4.210±0.889
			VTF (Chol.)	84.100±73.260	313.398±714.908	257.528	2.041±2.150	4.209±0.889
			RLBFGS	77.600±52.754	242.743±458.641	–	2.069±1.733	4.209±0.889
		Mid	VTF (ISR)	44.600±7.260	43.654±16.872	4.139	0.948±0.208	4.262±1.098
			VTF (Chol.)	45.900±8.647	45.661±20.088	5.615	0.955±0.236	4.262±1.098
			RLBFGS	45.200±8.080	48.104±19.981	–	1.025±0.243	4.262±1.098
		High	VTF (ISR)	44.400±9.252	43.684±22.460	11.746	0.936±0.254	3.874±1.395
			VTF (Chol.)	47.100±8.333	48.278±21.112	8.992	0.987±0.240	3.874±1.395
			RLBFGS	43.300±7.150	43.904±17.626	–	0.981±0.225	3.874±1.395
5	2	Low	VTF (ISR)	152.400±62.819	1344.573±999.949	649.595	7.560±3.406	0.260±0.458
			VTF (Chol.)	174.800±114.139	2168.886±3090.820	2347.949	8.680±6.348	0.265±0.466
			RLBFGS	166.300±76.884	1767.676±1406.454	–	8.788±4.427	0.267±0.454
		Mid	VTF (ISR)	120.700±69.620	942.713±1063.075	1120.404	5.807±3.877	0.781±0.207
			VTF (Chol.)	121.600±65.030	897.110±988.826	1315.410	5.720±3.445	0.781±0.207
			RLBFGS	136.300±91.493	1404.654±1986.798	–	7.057±5.378	0.764±0.225
		High	VTF (ISR)	46.400±17.322	94.741±105.045	22.071	1.739±0.902	1.805±0.385
			VTF (Chol.)	46.200±19.803	98.065±122.950	6.436	1.729±1.020	1.805±0.385
			RLBFGS	46.200±18.937	104.934±126.734	–	1.879±1.064	1.805±0.385
	10	Low	VTF (ISR)	295.900±86.577	5524.495±3173.855	2664.627	17.299±5.221	6.175±0.745
			VTF (Chol.)	292.300±83.828	5294.565±3165.735	5106.693	16.737±5.350	6.197±0.718
			RLBFGS	318.800±100.216	7083.082±4801.718	–	20.279±6.859	6.173±0.744
		Mid	VTF (ISR)	134.200±54.956	1139.332±919.917	306.469	7.302±3.227	6.753±0.708
			VTF (Chol.)	133.300±52.415	1100.519±842.840	296.996	7.183±3.045	6.753±0.708
			RLBFGS	135.300±54.965	1268.707±967.678	–	8.070±3.580	6.753±0.708
		High	VTF (ISR)	68.600±12.367	241.487±87.593	18.122	3.398±0.764	6.599±0.836
			VTF (Chol.)	74.000±14.071	279.271±105.828	40.158	3.632±0.838	6.599±0.836
			RLBFGS	68.300±11.982	258.475±93.959	–	3.661±0.790	6.599±0.836

(Arthur and Vassilvitskii, 2007) following (Hosseini and Sra, 2020). Three different levels of separation of Gaussian models, namely low, mid, and high, were used. The reader can refer to (Hosseini and Sra, 2020) for mathematical details of these separation levels.

• **Results:** We compared RLBFGS performances with and without our proposed mappings. The average number of iterations, convergence time, time per iteration, and the cost value of last iteration are reported in Table 2 where exponential map is used for retraction. Results for Taylor approximation of exponential map used as retraction is available in Table 1 of supplementary material. In these experiments, we report performances for $K \in \{2, 5\}$, $n \in \{2, 10\}$, $N = 10n^2 \in \{40, 1000\}$. For the sake of fairness, all the three compared

algorithms start with the same initial points in every run. The reader can see more extensive experiments for more sample size and dimensionality, i.e. $K \in \{2, 5\}$, $n \in \{2, 10, 100\}$, $N = 100n^2 \in \{400, 10000, 1000000\}$, in the Supplementary Material. In these tables, we have also provided the standard deviation (std) of time differences between VTF-RLBFGS and RL BFGS, over the ten runs. This value shows how much the average convergence time improvements over RL BFGS are reliable.

- **Discussion by Varying Separability:** As Table 2 and the Table 1 of Supplementary Material show, the proposed ISR mapping converges faster than no mapping most often in all separability levels. Both its time per iteration and number of iterations are often less than no mapping. These tables also show that the proposed Cholesky mapping converges faster than no mapping in most of the cases, although not all cases. Overall, we see that the two proposed mappings make RL BFGS faster and more efficient most often. This pacing improvement can be noticed more for low and mid separability levels because they are harder cases to converge. In terms of quality of local minimum, the proposed mappings often find the same local minimum as no mapping, but in a faster way. The equality of the found local optima in the three methods is because of using the same RL BFGS algorithm as their base in addition to having the same initial points. In some cases, the proposed mappings have even found better local minima. Moreover, as expected, the time difference std reduces in higher separability which is a simpler task.

- **Discussion by Varying Dimensionality:** We can discuss the results of Table 2 and the Table 1 of Supplementary Material by varying dimensionality, $n \in \{2, 10\}$. Tables 2 and 3 of Supplementary Material also report performance for dimensions $n \in \{2, 10, 100\}$ and larger sample size. Obviously, by increasing dimensionality, the time of convergence goes up and the faster pacing of the proposed mappings is noticed more, compared to no mapping.

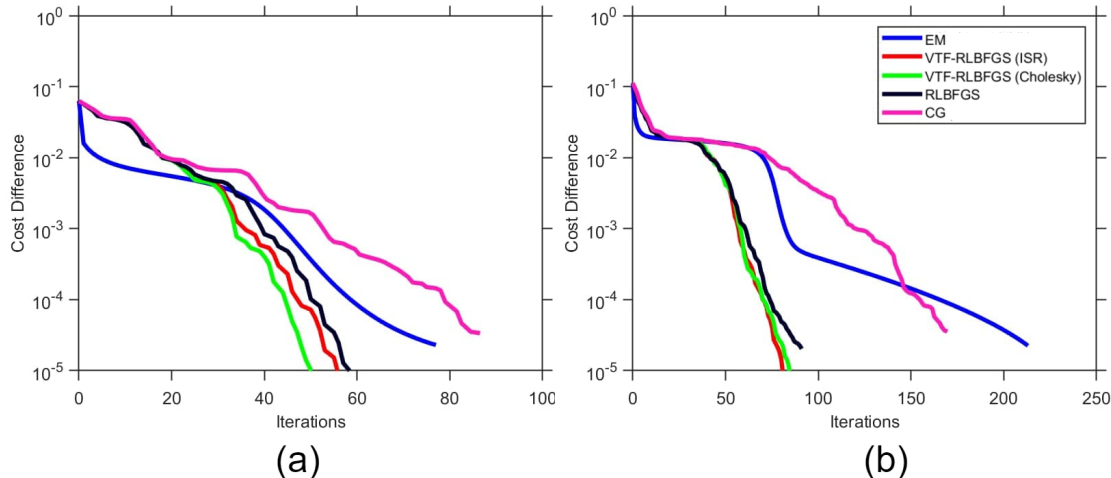
- **Discussion by Varying the Number of Components:** The results of Table 2 and the Table 1 of Supplementary Material can also be interpreted based on varying the number of mixture components, $K \in \{2, 5\}$. The more number of components makes the optimization problem harder. Hence, the difference of speeds of the proposed mappings and no mapping can be noticed more for $K = 5$; although, for both K values, the proposed mappings are often faster than no mapping.

- **Discussion by Varying Retraction Type:** We can also compare the performance of algorithms in terms of type of retraction. Table 2 and the Table 1 of Supplementary Material report performances where exponential map and Taylor approximation of exponential map are used for retraction, respectively. Comparing these tables shows that using Taylor approximation usually converges with less number of iterations compared to using exponential map. It makes sense because exponential map requires passing on geodesics. This difference of pacing is more obvious for larger number of components, i.e., $K = 5$. Our two proposed mappings outperform no mapping for both types retraction. This shows that our mappings are effective regardless of the details of operators.

- **Discussion by Varying the Sample Size:** Table 2 and the Table 1 of Supplementary Material report for $n \in \{2, 10\}$, $N = 10n^2 \in \{40, 1000\}$. More experiments for larger sample size and dimensionality, i.e. $n \in \{2, 10, 100\}$, $N = 100n^2 \in \{400, 10000, 1000000\}$, can be found in Tables 2 and 3 in the Supplementary Material. Comparing those tables with Table 2 and the Table 1 of Supplementary Material shows that larger sample size and/or

Table 3: Comparison of average results over ten runs with various algorithms where exponential map is used in algorithms and $K \in \{2, 5\}$, $n = 2$, $N = 1000n^2 = 4000$.

Algorithm	Low separation		Mid separation		High separation	
	#iters ($K = 2$)	#iters ($K = 5$)	#iters ($K = 2$)	#iters ($K = 5$)	#iters ($K = 2$)	#iters ($K = 5$)
VTF (ISR)	51.500±8.885	123.400±26.069	54.100±20.464	106.500±39.328	24.500±4.353	35.500±6.996
VTF (Chol.)	54.200±11.263	123.000±35.065	58.800±24.943	107.500±44.490	25.100±4.149	35.900±7.622
RLBFGS	52.900±8.825	147.400±35.926	57.900±29.622	110.200±38.064	24.400±4.006	35.600±7.516
CG	83.100±22.684	142.200±44.236	69.500±40.175	155.100±48.732	28.500±9.192	51.600±16.392
EM	94.400±40.114	277.900±142.549	118.100±89.794	224.000±101.576	3.200±0.422	3.600±0.966

Figure 1: The cost difference progress for several runs: (a) $K = 2$, $n = 2$, $N = 4000$, low separation, and (b) $K = 5$, $n = 2$, $N = 4000$, mid separation.

dimensionality takes more time to converge as expected. Still, our proposed mappings often converge faster than no mapping. The difference of pacing is mostly less in larger sample size compared to smaller sample size. This is because very large sample size consumes time on computation of cost function and the difference is not given much chance to show off in that case.

• **Comparison with Other Algorithms:** We compared RLBFGS, with and without the proposed mappings, with some other algorithms, i.e., nonlinear Conjugate Gradient (CG) and Expectation Maximization (EM) for fitting GMM. The log-scale cost difference progress of several insightful runs are illustrated in Fig. 1. The average number of iterations in the algorithms are compared in Table 3. A complete set of plots for cost difference progress can be seen in Figs. 1, 2, 3, and 4 in the Supplementary Material. Figs. 1-a and 1-b show that in some cases, ISR mapping is faster than the Cholesky mapping and in some other cases, we have the other way around. As Fig. 1 and Table 3 show, our proposed mappings are outperforming CG and EM in the number of iterations.

6.2. Geometric Metric Learning

• **Formulation:** As another application, we experiment with geometric metric learning. We implemented an iterative version of (Zadeh et al., 2016) whose regularized problem is:

$$\begin{aligned} \min_{\mathbf{W}} \quad & f := \sum_{(\mathbf{x}_i, \mathbf{x}_j) \in \mathcal{S}} (\mathbf{x}_i - \mathbf{x}_j)^\top \mathbf{W} (\mathbf{x}_i - \mathbf{x}_j) + \sum_{(\mathbf{x}_i, \mathbf{x}_j) \in \mathcal{D}} (\mathbf{x}_i - \mathbf{x}_j)^\top \mathbf{W}^{-1} (\mathbf{x}_i - \mathbf{x}_j) + \frac{1}{2} \|\mathbf{W}\|_F^2, \\ \text{s.t.} \quad & \mathbf{W} \in \mathcal{M} = \mathbb{S}_{++}^n, \end{aligned}$$

where \mathcal{S} and \mathcal{D} denote the sets of similar and dissimilar points, respectively. We used class labels to randomly sample the points for these sets. This metric learning behaves like triplet loss in which the intra- and inter-class variances are decreased and increased, respectively, for better discrimination of classes (Ghojogh et al., 2020). The Euclidean gradient of this problem is $\mathbb{R}^{n \times n} \ni \nabla_E f(\mathbf{W}) = \sum_{(\mathbf{x}_i, \mathbf{x}_j) \in \mathcal{S}} (\mathbf{x}_i - \mathbf{x}_j)(\mathbf{x}_i - \mathbf{x}_j)^\top - \sum_{(\mathbf{x}_i, \mathbf{x}_j) \in \mathcal{D}} (\mathbf{W}^{-1}(\mathbf{x}_i - \mathbf{x}_j))(\mathbf{W}^{-1}(\mathbf{x}_i - \mathbf{x}_j))^\top + \mathbf{W}$.

• **Results:** We evaluated our mappings with geometric metric learning optimization on three public datasets, i.e., Fisher Iris, USPS digits, and MNIST digits. The average results over ten runs are reported in Table 4. In Iris data, both mappings have converged faster than RLBFSGS, having a better converged cost function. In USPS and MNIST data, the Cholesky and ISR mappings have outperformed RLBFSGS without mapping, respectively.

Table 4: Comparison of geometric metric learning by RLBFSGS, with and without the proposed mappings, where exponential map is used in algorithms.

Data	Algorithm	#iters	conv. time	iter. time	last cost
Iris	VTF (ISR)	23.500±4.528	2.461±1.174	0.110±0.070	1512.602±347.004
	VTF (Chol.)	25.000±5.676	2.474±1.009	0.096±0.028	1594.975±124.087
	RLBFSGS	25.000±4.714	2.914±1.239	0.113±0.037	1620.207±125.989
USPS	VTF (ISR)	14.700±3.234	1.885±1.024	0.133±0.094	14223.215±0.001
	VTF (Chol.)	13.100±1.524	1.246±0.217	0.095±0.012	14223.216±0.001
	RLBFSGS	13.100±2.234	1.307±0.454	0.098±0.025	14223.215±0.001
MNIST	VTF (ISR)	11.700±3.335	1.288±0.619	0.108±0.041	6254.283±0.001
	VTF (Chol.)	13.100±3.479	1.435±0.793	0.104±0.029	6254.284±0.001
	RLBFSGS	12.100±3.381	1.266±0.754	0.099±0.024	6254.284±0.001

7. Conclusion and Future Direction

In this paper, we proposed two mappings in the tangent space of SPD manifolds by inverse second root and Cholesky decomposition. The proposed mappings simplify the vector transports and adjoint vector transports to identity. These transports are widely used in RLBFSGS quasi-Newton optimization, to identity. They also reduce the Riemannian metric to the Euclidean inner product which is more efficient computationally. Simulation results verified the effectiveness of the proposed mappings for two optimization tasks on SPD matrix manifolds. In this work, we focused on mappings for SPD manifolds which are widely used in machine learning and data science. A possible future direction is to extend the proposed mappings for other well-known Riemannian manifolds, such as Grassmann and

Stiefel (Edelman et al., 1998), as well as other Riemannian optimization methods. This paper opens a new research path for such mappings in the tangent space and we conjecture that such mappings can make numerical Riemannian optimization more efficient.

References

- P-A Absil, Robert Mahony, and Rodolphe Sepulchre. *Optimization algorithms on matrix manifolds*. Princeton University Press, 2009.
- Larry Armijo. Minimization of functions having Lipschitz continuous first partial derivatives. *Pacific Journal of mathematics*, 16(1):1–3, 1966.
- David Arthur and Sergei Vassilvitskii. k-means++: The advantages of careful seeding. In *Proceedings of the eighteenth annual ACM-SIAM symposium on Discrete algorithms*, pages 1027–1035, 2007.
- Rajendra Bhatia. *Positive definite matrices*. Princeton University Press, 2009.
- Nicolas Boumal. *An introduction to optimization on smooth manifolds*. Available online, 2020.
- Nicolas Boumal, Bamdev Mishra, P-A Absil, and Rodolphe Sepulchre. Manopt, a Matlab toolbox for optimization on manifolds. *The Journal of Machine Learning Research*, 15(1):1455–1459, 2014.
- Alan Edelman, Tomás A Arias, and Steven T Smith. The geometry of algorithms with orthogonality constraints. *SIAM journal on Matrix Analysis and Applications*, 20(2): 303–353, 1998.
- Roger Fletcher. *Practical methods of optimization*. John Wiley & Sons, 2013.
- Benyamin Ghojogh, Milad Sikaroudi, Sobhan Shafiei, Hamid R Tizhoosh, Fakhri Karray, and Mark Crowley. Fisher discriminant triplet and contrastive losses for training Siamese networks. In *International joint conference on neural networks*, pages 1–7. IEEE, 2020.
- Gene H Golub and Charles F Van Loan. *Matrix computations*, volume 3. JHU press, 2013.
- Magnus Rudolph Hestenes and Eduard Stiefel. *Methods of conjugate gradients for solving linear systems*, volume 49. NBS Washington, DC, 1952.
- Reshad Hosseini and Mohamadreza Mash’al. Mixest: An estimation toolbox for mixture models. *arXiv preprint arXiv:1507.06065*, 2015.
- Reshad Hosseini and Suvrit Sra. An alternative to EM for Gaussian mixture models: batch and stochastic Riemannian optimization. *Mathematical Programming*, 181(1):187–223, 2020.
- Jiang Hu, Xin Liu, Zai-Wen Wen, and Ya-Xiang Yuan. A brief introduction to manifold optimization. *Journal of the Operations Research Society of China*, 8(2):199–248, 2020.

- Wen Huang, Kyle A Gallivan, and P-A Absil. A Broyden class of quasi-Newton methods for Riemannian optimization. *SIAM Journal on Optimization*, 25(3):1660–1685, 2015.
- Wen Huang, P-A Absil, and Kyle A Gallivan. A Riemannian BFGS method for non-convex optimization problems. In *Numerical Mathematics and Advanced Applications ENUMATH 2015*, pages 627–634. Springer, 2016.
- Ben Jeuris, Raf Vandebril, and Bart Vandereycken. A survey and comparison of contemporary algorithms for computing the matrix geometric mean. *Electronic Transactions on Numerical Analysis*, 39(ARTICLE):379–402, 2012.
- Huibo Ji. *Optimization approaches on smooth manifolds*. PhD thesis, Australian National University, 2007.
- Dong-Hui Li and Masao Fukushima. On the global convergence of the BFGS method for nonconvex unconstrained optimization problems. *SIAM Journal on Optimization*, 11(4):1054–1064, 2001.
- Dong C Liu and Jorge Nocedal. On the limited memory BFGS method for large scale optimization. *Mathematical programming*, 45(1):503–528, 1989.
- Jorge Nocedal. Updating quasi-Newton matrices with limited storage. *Mathematics of computation*, 35(151):773–782, 1980.
- Jorge Nocedal and Stephen Wright. *Numerical optimization*. Springer Science & Business Media, 2006.
- Chunhong Qi, Kyle A Gallivan, and P-A Absil. Riemannian BFGS algorithm with applications. In *Recent advances in optimization and its applications in engineering*, pages 183–192. Springer, 2010.
- Wolfgang Ring and Benedikt Wirth. Optimization methods on Riemannian manifolds and their application to shape space. *SIAM Journal on Optimization*, 22(2):596–627, 2012.
- Matthias Seibert, Martin Kleinsteuber, and Knut Hüper. Properties of the BFGS method on Riemannian manifolds. *Mathematical System Theory C Festschrift in Honor of Uwe Helmke on the Occasion of his Sixtieth Birthday*, pages 395–412, 2013.
- Suvrit Sra and Reshad Hosseini. Conic geometric optimization on the manifold of positive definite matrices. *SIAM Journal on Optimization*, 25(1):713–739, 2015.
- Suvrit Sra and Reshad Hosseini. Geometric optimization in machine learning. In *Algorithmic Advances in Riemannian Geometry and Applications*, pages 73–91. Springer, 2016.
- Philip Wolfe. Convergence conditions for ascent methods. *SIAM Review*, 11(2):226–235, 1969.
- Pourya Zadeh, Reshad Hosseini, and Suvrit Sra. Geometric mean metric learning. In *International conference on machine learning*, pages 2464–2471. PMLR, 2016.