A MATRIX-ALGEBRAIC ALGORITHM FOR THE RIEMANNIAN LOGARITHM ON THE STIEFEL MANIFOLD UNDER THE CANONICAL METRIC

RALF ZIMMERMANN*

Abstract. We derive a numerical algorithm for evaluating the Riemannian logarithm on the Stiefel manifold with respect to the canonical metric. In contrast to the existing optimization-based approach, we work from a purely matrix-algebraic perspective. Moreover, we prove that the algorithm converges locally and exhibits a linear rate of convergence.

Key words. Stiefel manifold, Riemannian logarithm, Riemannian exponential, Dynkin series, Goldberg series, Baker-Campbell-Hausdorff series

AMS subject classifications. 15A16, 15B10, 15B57, 33B30, 33F05, 53-04, 65F60

1. Introduction. Consider an arbitrary Riemannian manifold \mathcal{M} . Geodesics on \mathcal{M} are locally shortest curves that are parametrized by the arc length. Because they satisfy an initial value problem, they are uniquely determined by specifying a starting point $p_0 \in \mathcal{M}$ and a starting velocity $\Delta \in T_{p_0}\mathcal{M}$ from the tangent space at p_0 . Geodesics give rise to the *Riemannian exponential function* that maps a tangent vector $\Delta \in T_{p_0}\mathcal{M}$ to the endpoint $\mathcal{C}(1)$ of a geodesic path $\mathcal{C}: [0,1] \to \mathcal{M}$ starting at $\mathcal{C}(0) = p_0 \in \mathcal{M}$ with velocity $\Delta = \dot{\mathcal{C}}(0) \in T_{p_0}\mathcal{M}$. It thus depends on the base point p_0 and is denoted by

$$Exp_{p_0}: T_{p_0}\mathcal{M} \to \mathcal{M}, Exp_{p_0}(\Delta) := \mathcal{C}(1).$$
(1.1)

The Riemannian exponential is a local diffeomorphism, [13, §5]. This means that it is locally invertible and that its inverse, called the *Riemannian logarithm* is also differentiable. Moreover, the exponential is radially isometric, i.e., the Riemannian distance between the starting point p_0 and the endpoint $p_1 := Exp_{p_0}(\Delta)$ on \mathcal{M} is the same as the length of the velocity vector Δ of the geodesic $t \mapsto Exp_{p_0}(t\Delta)$ when measured on the tangent space $T_{p_0}\mathcal{M}$, [13, Lem. 5.10 & Cor. 6.11]. In this way, the exponential mapping gives a local parametrization from the (flat, Euclidean) tangent space to the (possibly curved) manifold. This is also referred to as to representing the manifold in *normal coordinates* [12, §III.8].

The Riemannian exponential and logarithm are important both from the theoretical perspective as well as in practical applications. The latter fact holds true in particular, when \mathcal{M} is a *matrix manifold* [2]. Examples range from data analysis and signal processing [7, 17, 3, 18] over computer vision [4, 14] to adaptive model reduction and subspace interpolation [5] and, more generally speaking, optimization techniques on manifolds [6, 1, 2]. This list is far from being exhaustive.

Original contribution. In the work at hand, we present a matrix-algebraic derivation of an algorithm for computing the Riemannian logarithm on the *Stiefel manifold*. The matrix-algebraic perspective allows us to prove local linear convergence. The approach is based on an iterative inversion of the closed formula for the associated Riemannian exponential that has been derived in [6, §2.4.2]. Our main tools are Dynkin's explicit Baker-Campbell-Hausdorff formula [19] and Goldberg's exponential

^{*}Department of Mathematics and Computer Science, University of Southern Denmark (SDU) Odense, (zimmermann@imada.sdu.dk).

series [8], both of which represent a solution Z to the matrix equation

$$\exp_m(Z(X,Y)) = \exp_m(X) \exp_m(Y) \iff Z(\log_m(V), \log_M(W)) = \log_m(VW))$$

where $V = \exp_m(X), W = \exp_m(Y)$ and \exp_m, \log_m are the standard matrix exponential and matrix logarithm [11, §10, §11]. As an aside, we improve Thompson's norm bound from [20] on ||Z(X,Y)|| for the Goldberg series by a factor of 2, where $|| \cdot ||$ is any submultiplicative matrix norm.

The Stiefel log algorithm can be implemented in $\mathcal{O}(10)$ lines of (commented) MATLAB [15] code, which we include in Appendix E.1.

Comparison with previous work. To the best of our knowledge, up to now, the only algorithm for evaluating the Stiefel logarithm appeared in Q. Rentmeesters' thesis [18, Alg. 4, p. 91]. This algorithm is based on a Riemannian optimization problem. It turns out that this approach and the ansatz that is pursued here, though very different in their course of action, lead to essentially the same numerical scheme. Rentmeesters observes numerically a linear rate of convergence [18, p.83, p.100]. Proving linear convergence for [18, Alg. 4, p. 91] would require estimates on the Hessian, see [18, §5.2.1], [2, Thm. 4.5.6]. In contrast, the derivation presented here uses only elementary matrix algebra and the convergence proof given here formally avoids the requirements of computing/estimating step sizes, gradients and Hessians that are inherent to analyzing the convergence of optimization approaches. In fact, the convergence proof applies to [18, Alg. 4, p. 91] and yields the linear convergence of this optimization approach when using a fixed *unit step size*, but only on a sufficiently small domain. The thesis [18] was published under a two-years access embargo and the fundamentals of the work at hand were developed independently before [18] was accessible.

Transition to the complex case. The basic geometric concepts of the Stiefel manifold, the algorithm for the Riemannian log mapping developed here and its convergence proof carry over to complex matrices, where orthogonal matrices have to be replaced with unitary matrices and skew-symmetric matrices with skew-Hermitian matrices and so forth, see also [6, §2.1]. The thus adjusted log mapping algorithm was also confirmed numerically to work in the complex case.

Organization. Background information on the Stiefel manifold are reviewed in Section 2. The new derivation for the Stiefel log algorithm is in Section 3, convergence analysis is performed in Section 4, experimental results are in Section 5, and the conclusions follow in Section 6.

Notational specifics. The $(p \times p)$ -identity matrix is denoted by $I_p \in \mathbb{R}^{p \times p}$. If the dimension is clear, we will simply write I. The $(p \times p)$ -orthogonal group, i.e., the set of all square orthogonal matrices is denoted by

$$O_{p \times p} = \{ \Phi \in \mathbb{R}^{p \times p} | \Phi^T \Phi = \Phi \Phi^T = I_p \}.$$

The standard matrix exponential and matrix logarithm are denoted by

$$\exp_m(X) := \sum_{j=0}^{\infty} \frac{X^j}{j!}, \quad \log_m(I+X) := \sum_{j=1}^{\infty} (-1)^{j+1} \frac{X^j}{j}.$$

We use the symbols Exp^{St}, Log^{St} for the Riemannian counterparts on the Stiefel manifold.

When we employ the qr-decomposition of a rectangular matrix $A \in \mathbb{R}^{n \times p}$, we implicitly assume that $n \geq p$ and refer to the 'economy size' qr-decomposition A = QR, with $Q \in \mathbb{R}^{n \times p}$, $R \in \mathbb{R}^{p \times p}$.

2. The Stiefel manifold in numerical representation. This section reviews the essential aspects of the numerical treatment of Stiefel manifolds, where we rely heavily on the excellent references [2, 6]. The *Stiefel manifold* is the compact homogeneous matrix manifold of all column-orthogonal rectangular matrices

$$St(n,p) := \{ U \in \mathbb{R}^{n \times p} | \quad U^T U = I_p \}.$$

The tangent space $T_U St(n, p)$ at a point $U \in St(n, p)$ can be thought of as the space of velocity vectors of differentiable curves on St(n, p) passing through U:

$$T_U St(n,p) = \{ \mathcal{C}(t_o) | \mathcal{C} : (t_0 - \epsilon, t_0 + \epsilon) \to St(n,p), \mathcal{C}(t_0) = U \}.$$

For any matrix representative $U \in St(n,p)$, the tangent space of St(n,p) at U is represented by

$$T_U St(n,p) = \left\{ \Delta \in \mathbb{R}^{n \times p} | \quad U^T \Delta = -\Delta^T U \right\} \subset \mathbb{R}^{n \times p}$$

Every tangent vector $\Delta \in T_U St(n, p)$ may be written as

$$\Delta = UA + (I - UU^T)T, \quad A \in \mathbb{R}^{p \times p} \text{ skew}, \quad T \in \mathbb{R}^{n \times p} \text{ arbitrary}.$$
(2.1)

The dimension of both $T_U St(n,p)$ and St(n,p) is $np - \frac{1}{2}p(p+1)$.

Each tangent space carries an inner product $\langle \Delta, \tilde{\Delta} \rangle_U = tr \left(\Delta^T (I - \frac{1}{2} U U^T) \tilde{\Delta} \right)$

with corresponding norm $\|\Delta\|_U = \sqrt{\langle \Delta, \Delta \rangle_U}$. This is called the *canonical metric* on $T_U St(n, p)$. It is derived from the quotient space representation $St(n, p) = O_{n \times n}/O_{(n-p) \times (n-p)}$ that identifies two square orthogonal matrices in $O_{n \times n}$ as the same point on St(n, p), if their first p columns coincide [6, §2.4]. Endowing each tangent space with this metric (that varies differentiably in U) turns St(n, p) into a Riemannian manifold.

We now turn to the Riemannian exponential (1.1) but for $\mathcal{M} = St(n, p)$. An efficient algorithm for evaluating the Stiefel exponential was derived in [6, §2.4.2]. The algorithm starts with decomposing an input tangent vector $\Delta \in T_U St(n, p)$ into its horizontal and vertical components with respect to the base point U,

$$\Delta = UU^T \Delta + (I - UU^T) \Delta \stackrel{(\text{qr of } (I - UU^T)\Delta)}{=} UA + Q_E R_E.$$

Because Δ is tangent, $A \in \mathbb{R}^{p \times p}$ is skew. Then the matrix exponential is invoked to compute

$$\binom{M}{N_E} := \exp_m \left(\begin{pmatrix} A & -R_E^T \\ R_E & 0 \end{pmatrix} \right) \begin{pmatrix} I_p \\ 0 \end{pmatrix}.$$
(2.2)

The final output is¹

$$\tilde{U} := Exp_U^{St}(\Delta) = UM + Q_E N_E \in St(n, p).$$
(2.3)

(A MATLAB function for the Stiefel exponential is in the supplement in Appendix H.) The matrix exponential in (2.2) is related with the solution of the initial value problem that defines a geodesic on St(n, p), see [6, §2.4.2] for details. It turns out that the main obstacle in computing the inverse of the Stiefel exponential and thus the Stiefel logarithm is inverting (2.2), i.e. finding A, R_E given M, N_E , compare to [18, eq. (5.21)].

¹The index in Q_E, R_E, N_E is used to emphasize that these matrices stem from the Stiefel exponential as opposed to the closely related matrices Q, R, N that will appear in the procedure for the Stiefel logarithm.

3. Derivation of the Stiefel log algorithm. Let $U, \tilde{U} \in St(n, p)$ and assume that \tilde{U} is contained in a neighborhood \mathcal{D} of U such that Exp_U^{St} is a diffeomorphism from a neighborhood of $0 \in T_USt(n, p)$ onto \mathcal{D} . The central objective is to find $\Delta \in T_USt(n, p)$ such that $Exp_U^{St}(\Delta) = \tilde{U}$.

Because of Alg. 2.3, we know that \tilde{U} allows for a representation $\tilde{U} = UM + Q_E N_E$. Hence, we have to determine the unknown matrices $M, N_E \in \mathbb{R}^{p \times p}, Q_E \in \mathbb{R}^{n \times p}$, which feature the following properties: $Q_E^T U = 0$ and $M^T M + N_E^T N_E = I_p$. (Note that by (2.2), M and N_E are the left upper and lower $p \times p$ blocks of a $2p \times 2p$ orthogonal matrix.) We directly obtain

$$M = U^T \tilde{U}, \quad Q_E N_E = (I - U U^T) \tilde{U}.$$

We compute candidates for Q_E, N_E via a qr-decomposition

$$QN \stackrel{q_T}{=} (I - UU^T)U, \quad Q \in St(n, p).$$

The set of all orthogonal matrices with M, N as an upper diagonal and lower off-diagonal block is parametrized via

$$\left\{ \begin{pmatrix} M & X \\ N & Y \end{pmatrix} \mid \quad \begin{pmatrix} X \\ Y \end{pmatrix} = \begin{pmatrix} X_0 \\ Y_0 \end{pmatrix} \Phi, \quad \Phi \in O_{p \times p} \right\},$$

where $(X_0^T, Y_0^T)^T$ is a specific orthogonal completion, computed, say, via the Gram-Schmidt process.

Thus, the objective is reduced to solving the following nonlinear matrix equation

$$0 = \begin{pmatrix} 0 & I_p \end{pmatrix} \log_m \left(\begin{pmatrix} M & X_0 \\ N & Y_0 \end{pmatrix} \begin{pmatrix} I_p & 0 \\ 0 & \Phi \end{pmatrix} \right) \begin{pmatrix} 0 \\ I_p \end{pmatrix}, \quad \Phi \in O_{p \times p}.$$
(3.1)

Writing $\log_m \left(\begin{pmatrix} M & X_0 \\ N & Y_0 \end{pmatrix} \begin{pmatrix} I_p & 0 \\ 0 & \Phi \end{pmatrix} \right) = \begin{pmatrix} A & -B^T \\ B & C \end{pmatrix}$, this means finding a rotation Φ such that C = 0.

The first result is that solving (3.1) indeed leads to the Riemannian logarithm on the Stiefel manifold.

THEOREM 3.1. Let $U, \tilde{U} \in St(n, p)$ and assume that \tilde{U} is contained in a neighborhood \mathcal{D} of U such that Exp_U^{St} is a diffeomorphism from a neighborhood of $0 \in T_USt(n, p)$ onto \mathcal{D} .

Let M, Q_E , N_E , Q, N, X_0 , Y_0 as introduced in the above setting. Suppose that $\Phi \in O_{p \times p}$ solves (3.1), i.e.,

$$\log_m \left(\begin{pmatrix} M & X_0 \Phi \\ N & Y_0 \Phi \end{pmatrix} \right) = \begin{pmatrix} A & -B^T \\ B & 0 \end{pmatrix}.$$

Define $\Delta := UA + QB \in T_USt(n, p)$. Then $Exp_U^{St}(\Delta) = \tilde{U}$, i.e., $\Delta = Log_U^{St}(\tilde{U})$. Proof. By construction, it holds $QN = (I - UU^T)\tilde{U}$ and hence

$$U^{T}Q = 0, \quad (I - UU^{T})Q = Q, \quad QQ^{T}\tilde{U} = QQ^{T}(I - UU^{T})\tilde{U} = (I - UU^{T})\tilde{U}.$$
 (3.2)

Now, we apply the Stiefel exponential Alg. 2.3 to $\Delta = UA + QB$. This gives $(I - UU^T)\Delta = QB$ and

$$Q_E R_E \stackrel{qr}{=} QB \Leftrightarrow R_E = \Psi B$$
, where $\Psi := (Q_E^T Q) \in O_{p \times p}$.²

²The matrices Q_E and Q differ by an orthogonal rotation but span the same subspace.

With $U^T \Delta = A$, we obtain

$$\begin{pmatrix} M\\N_E \end{pmatrix} := \exp_m \left(\begin{pmatrix} A & -R_E^T\\R_E & 0 \end{pmatrix} \right) \begin{pmatrix} I_p\\0 \end{pmatrix}$$

$$= \begin{pmatrix} I & 0\\0 & \Psi \end{pmatrix} \exp_m \left(\begin{pmatrix} A & -B^T\\B & 0 \end{pmatrix} \right) \begin{pmatrix} I & 0\\0 & \Psi^T \end{pmatrix} \begin{pmatrix} I_p\\0 \end{pmatrix}$$

$$= \begin{pmatrix} I & 0\\0 & \Psi \end{pmatrix} \begin{pmatrix} M & X_0 \Phi\\N & Y_0 \Phi \end{pmatrix} \begin{pmatrix} I_p\\0 \end{pmatrix} = \begin{pmatrix} M\\\Psi N \end{pmatrix}.$$

Keeping in mind that $Q_E \Psi = Q_E Q_E^T Q = Q$, this leads to an output of

$$Exp_U^{St}(\Delta) = UM + Q_E N_E = UM + Q_E \Psi N = UM + QN = U.$$

Thus, Δ is a valid tangent vector in $T_U St(n,p)$ such that $Exp_U^{St}(\Delta) = \tilde{U} \in St(n,p)$. From abstract differential geometry, we know that $Log_U^{St}(\tilde{U}) \in T_U St(n,p)$ is the unique tangent with $Exp_U^{St}(Log_U^{St}(\tilde{U})) = \tilde{U}$. We arrive at the claim

$$\Delta = Log_U^{St}(\tilde{U}).$$

 \square Having established Theorem 3.1, we now focus on solving (3.1). Let

$$V_0 := \begin{pmatrix} M & X_0 \\ N & Y_0 \end{pmatrix}, \quad \log_m(V_0) := \begin{pmatrix} A_0 & -B_0^T \\ B_0 & C_0 \end{pmatrix},$$
$$W_0 := \begin{pmatrix} I_p & 0 \\ 0 & \Phi_0 \end{pmatrix}, \quad \log_m(W_0) = \begin{pmatrix} 0 & 0 \\ 0 & \log_m(\Phi_0) \end{pmatrix}.$$

Up to terms of first order, it holds $\log_m(V_0W_0) = \log_m(V_0) + \log_m(W_0)$. Hence, the choice

$$\Phi_0 = \exp_m(-C_0)$$

gives an approximate solution to (3.1). We define

$$V_1 := \begin{pmatrix} M & X_0 \\ N & Y_0 \end{pmatrix} \begin{pmatrix} I_p & 0 \\ 0 & \Phi_0 \end{pmatrix}, \quad \log_m(V_1) := \begin{pmatrix} A_1 & -B_1^T \\ B_1 & C_1 \end{pmatrix}$$
(3.5)

and iterate. This is the essential idea of Alg. 3.1 for the Riemannian logarithm.³

In Section 4 we make use of the Baker-Campbell-Hausdorff formula [19, §1.3, p. 22] that corrects for the misfit in the approximative matrix relation $\log_m(VW) \approx \log_m(V) + \log_m(W)$ for two non-commuting matrices V, W in order to show that the above procedure leads to

$$\|C_{k+1}\|_2 < \alpha \|C_k\|_2$$

for all $k \in \mathbb{N}_0$ and a constant $\alpha < 1$ and is thus convergent.

Since the Riemannian exponential is a local diffeomorphism, we have to postulate a suitable bound on the distance between the input matrices U and \tilde{U} . Suppose that $\|U - \tilde{U}\|_2 < \varepsilon$. Recalling the definitions $M = U^T \tilde{U}$ and $(I - UU^T)\tilde{U} = QN$, this

 $^{^{3}}$ This is the same algorithm as [18, Alg. 4, p. 91] that Rentmeesters obtains from his geometrical perspective when a fixed unit step length is employed and when [18, §5.3] is taken into account.

gives the following bounds for the horizontal and the vertical component of $U - \tilde{U}$ with respect to the subspace spanned by U:

$$||UU^{T}(U-\tilde{U})||_{2} = ||I_{p}-M||_{2} < \varepsilon, \quad ||(I-UU^{T})(U-\tilde{U})||_{2} = ||QN||_{2} = ||N||_{2} < \varepsilon.$$

However, it turns out that for the convergence proof, estimates on the norms of X_0 , Y_0 and $Y_0 - I_p$ are also required. By the CS-decomposition of orthonormal matrices [9, Thm 2.6.3, p. 78], the diagonal blocks M and Y_0 share the same singular values and so do the off-diagonal blocks N, X_0 . Hence, $||N||_2 = ||X_0||_2 < \varepsilon$. Let $D_1 \Sigma R_1^T$ be the SVD of M and $D_2 \Sigma R_2^T$ be the SVD of Y_0 . An estimate for the singular values of M can be obtained as follows:

$$\varepsilon^{2} > \|N\|_{2}^{2} = \|N^{T}N\|_{2} = \|I - M^{T}M\|_{2} = \|I - \Sigma^{2}\|_{2} = \max_{\sigma_{k}}(1 - \sigma_{k}^{2}), \qquad (3.6)$$

where we have used that $\sigma_1 = ||M||_2 \leq 1$. Now, we replace the Y_0 that has been obtained via, say, Gram-Schmidt by $Y_0R_2D_2^T = D_2\Sigma D_2^T$ (and, correspondingly, X_0 by $X_0R_2D_2^T$). Essentially, this is the orthogonal Procrustes method, [9, §12.4.1, p.601], applied to $\min_{\Psi \in O_{p \times p}} ||I - Y_0\Psi||_2$. This operation preserves the orthogonality of $V_0 = \begin{pmatrix} M & X_0 \\ N & Y_0 \end{pmatrix}$, but the new Y_0 is symmetric with eigenvalue decomposition $Y_0 = D_2\Sigma D_2^T$. This gives

$$||Y_0 - I_p||_2 = ||\Sigma - I_p||_2 = \max_{\sigma_k} |1 - \sigma_k| < \max_{\sigma_k} (1 - \sigma_k^2) < \varepsilon^2.$$

In summary, if $||U - \tilde{U}||_2 < \varepsilon$ and if we start the iterations indicated by (3.5) with the Procrustes orthogonal completion X_0, Y_0 rather than the standard Gram-Schmidt process, we obtain Alg. 3.1 with the starting conditions

$$||I_p - M||_2 < \varepsilon, \quad ||N||_2 = ||X_0||_2 < \varepsilon, \quad ||Y_0 - I_p||_2 < \varepsilon^2.$$
(3.7)

Computational costs. W.l.o.g. suppose that $n \geq p$. In fact the most important case in practical applications is $n \gg p$. Because of the matrix product in step 1 and the qr-decomposition in step 2 of Alg. 3.1, the preparatory steps 1–3 require $\mathcal{O}(np^2)$ FLOPS. The dominating costs in the iterative loop, steps 5–10, are the evaluation of the matrix logarithm for a 2*p*-by-2*p* orthogonal matrix and the matrix exponential for a *p*-by-*p* skew-symmetric matrix in every iteration, both of which can be achieved efficiently via the Schur decomposition. The costs are $\mathcal{O}(p^3)$, see [9, Alg. 7.5.2].

A MATLAB function for Alg. 3.1 is in Appendix E.1.

4. Convergence proof. In this section, we establish the convergence of Alg. 3.1 under suitable conditions. We state the main result as Theorem 4.1; the proof is subdivided into the auxiliary results Lemma 4.2, and Lemma 4.3 as well as Lemma A.1 that appears in Appendix A. An essential requirement is that the point $\tilde{U} \in St(n, p)$ that is to be mapped to the tangent space $T_U St(n, p)$ is sufficiently close to the base point $U \in St(n, p)$ in the sense that $||U - \tilde{U}||_2 < \varepsilon$. Throughout, we will make extensive use of Dynkin's explicit BCH formula [19, §1.3, p. 22].

THEOREM 4.1. Let $U, \tilde{U} \in St(n, p)$. Assume that $||U - \tilde{U}||_2 < \varepsilon$. Let $(V_k)_k$ be the sequence of orthogonal matrices generated by Alg. 3.1.

If $\varepsilon < 0.0912$, then Alg. 3.1 converges to a limit matrix $V_{\infty} := \lim_{k \to \infty} V_k$ such that

$$\log_m(V_\infty) := \begin{pmatrix} A_\infty & -B_\infty^T \\ B_\infty & C_\infty \end{pmatrix} = \begin{pmatrix} A_\infty & -B_\infty^T \\ B_\infty & 0 \end{pmatrix}.$$

Algorithm 3.1 Stiefel logarithm, iterative procedure

Input: base point $U \in St(n,p)$ and $\tilde{U} \in St(n,p)$ 'close' to base point, $\tau > 0$ convergence threshold 1: $M := U^T \tilde{U} \in \mathbb{R}^{p \times p}$ 2: $QN := \tilde{U} - UM \in \mathbb{R}^{n \times p}$ # (thin) gr-decomp. of normal component of \tilde{U} 3: $V_0 := \begin{pmatrix} M & X_0 \\ N & Y_0 \end{pmatrix} \in O_{2p \times 2p}$ 4: **for** $k = 0, 1, 2, \dots$ **do** # orthogonal completion and Procrustes $\begin{pmatrix} A_k & -B_k^T \\ B_k & C_k \end{pmatrix} := \log_m(V_k) \quad \# \text{ matrix } \log, A_k, C_k \text{ skew}$ 5: if $||C_k||_2 \leq \tau$ then 6: 7: break end if 8: $\Phi_{k} = \exp_{m} (-C_{k}) \qquad \text{# matrix exp, } \Phi_{k} \text{ orthogonal}$ $V_{k+1} := V_{k}W_{k}, \text{ where } W_{k} := \begin{pmatrix} I_{p} & 0\\ 0 & \Phi_{k} \end{pmatrix}$ 9: 10: # update 11: end for **Output:** $\Delta := Log_U^{St}(\tilde{U}) = UA_k + QB_k \in T_USt(n,p)$

Given a numerical convergence threshold $\tau > 0$, see Alg. 3.1, line 7, the algorithm requires at most $k = \left\lceil \frac{\log(||C_0||_2) - \log(\tau)}{\log(2)} \right\rceil - 1$ iteration steps to meet the convergence criterion under the above conditions.

REMARK 1. Alg. 3.1 generates a sequence of orthonormal matrices

$$V_{k+1} = V_k W_k = V_0(W_0 W_1 \dots W_k) = V_0 \left(\begin{pmatrix} I_p & 0\\ 0 & \Phi_0 \end{pmatrix} \dots \begin{pmatrix} I_p & 0\\ 0 & \Phi_k \end{pmatrix} \right) \in O_{2p \times 2p}.$$
(4.1)

The proof of Theorem 4.1 will show that $\lim_{k\to\infty} W_k = I_{2p}$, see (4.12). Therefore, the sequence of orthogonal products $\Phi_0 \dots \Phi_k$ converges to a limit Φ_∞ for $k \to \infty$. The limit Φ_∞ solves (3.1). However, it is not required to actually form Φ_∞ . In pursuit of the proof of Theorem 4.1, we first show that if the norm of the matrix logarithm of the orthogonal matrix V_k produced by Alg. 3.1 at iteration k is sufficiently small, then the norm of the lower p-by-p diagonal block of the matrix logarithm of the next iterate V_{k+1} is strictly decreasing by a constant factor.

LEMMA 4.2. Let $U, U \in St(n, p)$. Let $(V_k)_k \subset O_{2p \times 2p}$ be the sequence of orthogonal matrices generated by Alg. 3.1. Suppose that at stage k, it holds

$$\|\log_m(V_k)\|_2 := \| \begin{pmatrix} A_k & -B_k^T \\ B_k & C_k \end{pmatrix} \|_2 < \frac{1}{2}(\sqrt{5} - 1).$$
(4.2)

 $Then \log_m(V_{k+1}) := \begin{pmatrix} A_{k+1} & -B_{k+1}^T \\ B_{k+1} & C_{k+1} \end{pmatrix} features \ a \ lower \ (p \times p) - diagonal \ block \ of \ norm$

$$||C_{k+1}||_2 < \alpha ||C_k||_2, \quad 0 < \alpha < \frac{1}{2}$$

Proof. Given $V_k = \begin{pmatrix} M & X_k \\ N & Y_k \end{pmatrix} = \exp_m \left(\begin{pmatrix} A_k & -B_k^T \\ B_k & C_k \end{pmatrix} \right)$, Alg. 3.1 computes the next iterate V_{k+1} via

$$V_{k+1} := V_k W_k$$

where $W_k := \begin{pmatrix} I_p & 0 \\ 0 & \exp_m(-C_k) \end{pmatrix}$. For brevity, we introduce the notation $L_V := \log_m(V)$ for the matrix logarithm. Recall that [V, W] := VW - WV denotes the commutator or Lie-bracket of the matrices V, W. From Dynkin's formula for the Baker-Campbell-Hausdorff series, see [19, §1.3, p. 22], we obtain

$$\begin{split} L_{V_{k+1}} &= \log_m(V_k W_k) \\ &= L_{V_k} + L_{W_k} + \frac{1}{2} [L_{V_k}, L_{W_k}] \\ &+ \frac{1}{12} \left(\left[L_{V_k}, [L_{V_k}, L_{W_k}] \right] + \left[L_{W_k}, [L_{W_k}, L_{V_k}] \right] \right) \\ &- \frac{1}{24} \left[L_{W_k}, \left[L_{V_k}, [L_{V_k}, L_{W_k}] \right] \right] + \sum_{l=5}^{\infty} z_l (L_{V_k}, L_{W_k}), \end{split}$$

where $\sum_{l=5}^{\infty} z_l(L_{V_k}, L_{W_k}) =: h.o.t.(5)$ are the terms of fifth order and higher in the series. In the case at hand, it holds

$$L_{V_k} + L_{W_k} = \begin{pmatrix} A_k & -B_k^T \\ B_k & C_k \end{pmatrix} + \begin{pmatrix} 0 & 0 \\ 0 & -C_k \end{pmatrix} = \begin{pmatrix} A_k & -B_k^T \\ B_k & 0 \end{pmatrix},$$
$$[L_{V_k}, L_{W_k}] = \begin{pmatrix} 0 & B_k^T C_k \\ C_k B_k & 0 \end{pmatrix}.$$

(Note that the basic idea in designing Alg. 3.1 was exactly to choose W_k such that the lower diagonal block in the BCH-series cancels in the first order terms.)

The third and fourth order terms are

$$\frac{1}{12} \begin{pmatrix} -2B_k^T C_k B_k & A_k B_k^T C_k - 2B_k^T C_k^2 \\ 2C_k^2 B_k - C_k B_k A_k & B_k B_k^T C_k + C_k B_k B_k^T \end{pmatrix}, \text{ and} \\
\frac{1}{24} \begin{pmatrix} 0 & -B_k^T C_k^3 + A_k B_k^T C_k^2 \\ -C_k^3 B_k + C_k^2 B_k A_k & B_k B_k^T C_k^2 - C_k^2 B_k B_k^T \end{pmatrix}, \text{ respectively.}$$

Therefore, the series expansion for the lower diagonal block in $\log_m(V_{k+1})$ starts with the terms of third order:

$$\|C_{k+1}\|_{2} = \|\frac{1}{12}(B_{k}B_{k}^{T}C_{k} + C_{k}B_{k}B_{k}^{T}) - \frac{1}{24}(B_{k}B_{k}^{T}C_{k}^{2} - C_{k}^{2}B_{k}B_{k}^{T}) + h.o.t.(5)\|_{2}$$

$$\leq \left(\frac{1}{6}\|B_{k}\|_{2}^{2} + \frac{1}{12}\|B_{k}\|_{2}^{2}\|C_{k}\|_{2} + \frac{\|h.o.t.(5)\|_{2}}{\|C_{k}\|_{2}}\right)\|C_{k}\|_{2}.$$
(4.6a)

We tackle the higher order terms via Lemma A.1 from the appendix. The lemma applies because $||C_k||_2 = ||L_{W_k}||_2 \le ||L_{V_k}||_2 < \frac{1}{2}(\sqrt{5}-1) < 1$. In this setting, it gives

$$\|h.o.t.(5)\|_{2} \leq \sum_{l=5}^{\infty} \|z_{l}(L_{V_{k}}, L_{W_{k}})\|_{2} < \sum_{l=5}^{\infty} \|L_{V_{k}}\|^{l-1} \|L_{W_{k}}\|_{2},$$

since each of the "letters" L_{V_k}, L_{W_k} appears at least once in every "word" that contributes to $z_k(L_{V_k}, L_{W_k})$, see Appendix A and [20, 16, 21].

Writing $s := ||L_{V_k}||_2$ and substituting in (4.6a) leads to

$$\|C_{k+1}\|_2 < \left(\frac{1}{6}s^2 + \frac{1}{12}s^3 + \sum_{l=4}^{\infty}s^l\right)\|C_k\|_2 =: \alpha \|C_k\|_2.$$
(4.7)

The proof is complete, if we can show that $\alpha < 1$. Note that $\sum_{l=4}^{\infty} s^l = \frac{1}{1-s} - 1 - s - s^2 - s^3$. As a consequence

$$\alpha < 1 \Leftrightarrow \frac{s^2}{1-s} - \frac{5}{6}s^2 - \frac{11}{12}s^3 < 1.$$

An obvious bound on the size of s is obtained via observing that $\frac{s^2}{1-s} < 1$, if $s < \frac{1}{2}(\sqrt{5}-1) \approx 0.618$. The corresponding α is $0.4653 < \frac{1}{2}$. A sharper bound can be obtained via solving the associated quartic equation. This shows that the inequality even holds for all s < 0.7111. \Box In order to make use of Lemma 4.2, we establish conditions such that $\|\log_m(V_k)\|_2 < \frac{1}{2}(\sqrt{5}-1)$ holds throughout the iterations of Alg. 3.1.

This is the goal of the the next lemma. It relies on the auxiliary results Proposition B.1, Proposition B.2 and Lemma B.3 from Appendix B. Proposition B.1 shows that $\|\exp_m(C)-I\|_2 < \|C\|_2$ for C skew-symmetric; Proposition B.2 establishes a bound in the opposite direction: if V is orthogonal such that $\|V-I\|_2 < r$, then $\|\log_m(V)\|_2 < r\sqrt{1-\frac{r^2}{4}}/(1-\frac{r^2}{2})$. Finally, Lemma B.3 shows that $\|V_0-I\|_2 < 2\varepsilon$ for the first iterate V_0 of Alg. 3.1, provided that $\|U-\tilde{U}\|_2 < \varepsilon$.

LEMMA 4.3. Let $U, \tilde{U} \in St(n, p)$ with $||U - \tilde{U}||_2 < \varepsilon$. Let $(V_k)_k \subset O_{2p \times 2p}$ be the sequence of orthogonal matrices generated by Alg. 3.1, where $V_k = \begin{pmatrix} M & X_k \\ N & Y_k \end{pmatrix}$. Let $\tilde{\varepsilon} = 2\varepsilon \frac{\sqrt{1-\varepsilon^2}}{1-2\varepsilon^2}$ and $\hat{\varepsilon} := (e^{2\tilde{\varepsilon}} - 1) + \varepsilon + \varepsilon^2$. If $0 < \varepsilon$ is small enough such that $\hat{\varepsilon} \frac{\sqrt{1-\varepsilon^2}}{1-\varepsilon^2} < \frac{1}{2}(\sqrt{5}-1)$, then

$$\|\log_m(V_k)\|_2 = \|\begin{pmatrix} A_k & -B_k^T \\ B_k & C_k \end{pmatrix}\|_2 < \frac{1}{2}(\sqrt{5}-1) \text{ for all } k$$

Proof. Let $\delta_0 := \frac{1}{2}(\sqrt{5}-1)$. By Lemma B.3 from Appendix B, it holds

$$\|\log_m(V_0)\|_2 < 2\varepsilon \frac{\sqrt{1-\varepsilon^2}}{1-2\varepsilon^2} = \tilde{\varepsilon} \quad (<\delta_0)$$

In particular, $\tilde{\varepsilon} > \| \begin{pmatrix} A_0 & -B_0^T \\ B_0 & C_0 \end{pmatrix} \|_2 \ge \|C_0\|_2$. By Alg. 3.1, $\Phi_0 = \exp_m(-C_0)$, where Φ_0 is orthogonal. By Proposition B.1 from Appendix B

$$\begin{aligned} \|\Phi_0 - I\|_2 &\leq \|C_0\|_2 < \tilde{\varepsilon}. \\ \end{aligned}$$

Writing $V_1 = I + (V_1 - I) =: I + E_1$, this leads to the estimate

$$\begin{split} \|E_1\|_2 &= \| \begin{pmatrix} M-I & X_0 \Phi_0 \\ N & Y_0 \Phi_0 - I \end{pmatrix} \|_2 \\ &= \| \begin{pmatrix} M-I & 0 \\ 0 & Y_0(\Phi_0 - I) \end{pmatrix} + \begin{pmatrix} 0 & 0 \\ 0 & Y_0 - I \end{pmatrix} + \begin{pmatrix} 0 & X_0 \Phi_0 \\ N & 0 \end{pmatrix} \|_2 \\ &\leq \max\{\|M-I\|_2, \|Y_0(\Phi_0 - I)\|_2\} + \|Y_0 - I\|_2 + \max\{\|N\|_2, \|X_0 \Phi_0\|_2\} \\ &\leq \max\{\varepsilon, \|Y_0\|_2 \| (\Phi_0 - I) \|_2\} + \varepsilon^2 + \varepsilon \leq \tilde{\varepsilon} + \varepsilon^2 + \varepsilon \\ &\leq (e^{2\tilde{\varepsilon}} - 1) + \varepsilon + \varepsilon^2 = \hat{\varepsilon}, \end{split}$$

where we have used (3.7) and the fact that $||Y_0||_2 \leq 1$, see (B.2a), (B.2b). By Lemma B.3, $||\log_m(V_1)||_2 < \hat{\varepsilon}\sqrt{1-\frac{\hat{\varepsilon}^2}{4}}/(1-\frac{\hat{\varepsilon}^2}{2}) < \delta_0$. Thus, the claim holds for k = 0, 1.

Lemma 4.2 applies to $\|\log_m(V_0)\|_2$ and leads to $\|C_1\|_2 < \frac{1}{2}\|C_0\|_2 < \frac{\tilde{\varepsilon}}{2}$ for the lower diagonal block C_1 of the next iterate $\log_m(V_1)$. Therefore, by using Proposition B.1 once more, we see that

$$\|\Phi_1 - I\|_2 \le \|C_1\|_2 < \frac{\tilde{\varepsilon}}{2}.$$

By induction, we obtain $V_k = I + (V_k - I) =: I + E_k$ with

$$\begin{split} \|E_k\|_2 &= \| \begin{pmatrix} M & X_0 \hat{\Phi}_{k-1} \\ N & Y_0 \hat{\Phi}_{k-1} \end{pmatrix} - I \|_2 \\ &= \| \begin{pmatrix} M - I & 0 \\ 0 & Y_0 (\hat{\Phi}_{k-1} - I) \end{pmatrix} + \begin{pmatrix} 0 & 0 \\ 0 & Y_0 - I \end{pmatrix} + \begin{pmatrix} 0 & X_0 \hat{\Phi}_{k-1} \\ N & 0 \end{pmatrix} \|_2 \\ &\leq \max\{ \|M - I\|_2, \|Y_0 (\hat{\Phi}_{k-1} - I)\|_2 \} + \|Y_0 - I\|_2 + \max\{\|N\|_2, \|X_0 \hat{\Phi}_{k-1}\|_2 \} \\ &\leq \max\{\varepsilon, \|Y_0\|_2 \| \hat{\Phi}_{k-1} - I\|_2 \} + \varepsilon^2 + \varepsilon. \end{split}$$
(4.9a)

where $\hat{\Phi}_{k-1} = \Phi_0 \cdots \Phi_{k-1}$.

We can estimate $\|\hat{\Phi}_{k-1} - I\|_2$ as follows: By the induction hypothesis, we assume that we have checked that $\|\log_m(V_j)\|_2 < \delta_0$ for $j = 0, \ldots, k-1$. Hence, Lemma 4.2 ensures that $\|C_j\|_2 < \frac{1}{2}\|C_{j-1}\|_2 < \ldots < \frac{1}{2^j}\|C_0\|_2 < \frac{\tilde{\varepsilon}}{2^j}$ for the lower diagonal block of $\log_m(V_j), j = 0, \ldots, k-1$. As above, this gives $\|\Phi_j - I\|_2 \leq \|C_j\|_2 < \frac{\tilde{\varepsilon}}{2^j}$. We thus may write $\Phi_j = I + (\Phi_j - I) =: I + \Gamma_j$ with $\|\Gamma_j\|_2 =: g_j < \frac{\tilde{\varepsilon}}{2^j}$. This gives

$$\|\hat{\Phi}_{k-1} - I\|_2 = \|(I + \Gamma_1) \cdots (I + \Gamma_{k-1}) - I\|_2 \le (1 + g_1) \cdots (1 + g_{k-1}) - 1.$$
(4.10)

It holds

$$\ln\left(\prod_{j=0}^{k-1}\left(1+g_j\right)\right) = \sum_{j=0}^{k-1}\ln(1+g_j) \le \sum_{j=0}^{k-1}g_j \le \sum_{j=0}^{\infty}\frac{\tilde{\varepsilon}}{2^j} = 2\tilde{\varepsilon}.$$

Using this estimate in (4.10) gives

$$\|\hat{\Phi}_{k-1} - I\|_2 < e^{2\tilde{\varepsilon}} - 1$$

and we finally arrive at

$$||E_k||_2 \le (e^{2\tilde{\varepsilon}} - 1) + \varepsilon^2 + \varepsilon = \hat{\varepsilon}.$$
10

Recalling (4.9a), we have $V_k = I + E_k$ with $||E_k||_2 < \hat{\varepsilon}$ at every iteration k. By Lemma B.3, $||\log_m(V_k)||_2 < \hat{\varepsilon} \frac{\sqrt{1-\frac{\hat{\varepsilon}^2}{4}}}{1-\frac{\hat{\varepsilon}^2}{2}}$ and we see that the postulate on the size of ε is such that $||\log_m(V_k)||_2 < \delta_0$. Thus Lemma 4.2 indeed applies at iteration k, which closes the induction. \square **Remark:** The inequality $\hat{\varepsilon} \frac{\sqrt{1-\frac{\hat{\varepsilon}^2}{4}}}{1-\frac{\hat{\varepsilon}^2}{2}} < \delta_0$ holds precisely for $\hat{\varepsilon} < \sqrt{2} \left(1 - \frac{1}{\sqrt{1+\delta_0^2}}\right)^{\frac{1}{2}} =: \hat{\varepsilon}_0$. A further calculations shows that if $\varepsilon < 0.0912$, then

 $\hat{\varepsilon} = (e^{2\tilde{\varepsilon}} - 1) + \varepsilon^2 + \varepsilon < \hat{\varepsilon}_0$, i.e., the conditions of Lemma 4.3 hold, for all $\varepsilon < 0.0912$. With the tools established above at hand, we are now in a position to prove Theorem 4.1.

Proof. [Theorem 4.1] Let $(V_k)_{k \in \mathbb{N}_0}$ be the sequence of orthogonal matrices generated by Alg. 3.1. By Lemma 4.2 and Lemma 4.3, it holds

$$\|\log_m V_k\|_2 := \| \begin{pmatrix} A_k & -B_k^T \\ B_k & C_k \end{pmatrix} \|_2 < \frac{1}{2}(\sqrt{5}-1), \quad \|C_{k+1}\|_2 < \alpha^{k+1} \|C_0\|_2$$
(4.11)

for all $k \ge 0$, where $0 < \alpha < \frac{1}{2}$. From this equation and the continuity of the matrix exponential, we obtain

$$\lim_{k \to \infty} W_k = \lim_{k \to \infty} \begin{pmatrix} I_p & 0\\ 0 & \exp_m(-C_k) \end{pmatrix} = \begin{pmatrix} I_p & 0\\ 0 & I_p \end{pmatrix}.$$
 (4.12)

The convergence result is now an immediate consequence of Alg. 3.1, step 10. The upper bound on the iteration count required for numerical convergence is also obvious from (4.11). \Box

5. Examples and experimental results. In this section, we discuss a special case that can be treated analytically. Following, we present numerical results on the performance of Alg. 3.1.

5.1. A special case. Here, we consider the special situation, where the two points $U, \tilde{U} \in St(n, p)$ are such that their columns span the same subspace.⁴ Hence, there exists an orthogonal matrix $M \in O_{p \times p}$ such that $\tilde{U} = UM = UM + (I - UU^T)0$. In this case, Alg. 3.1 produces the initial matrices $V_0 = \begin{pmatrix} M & 0 \\ 0 & Y_0 \end{pmatrix}$ and $\Phi_0 = \exp_m(-\log_m(Y_0)) = Y_0^{-1}$. Note that the corresponding $W_0 = \begin{pmatrix} I_p & 0 \\ 0 & Y_0^{-1} \end{pmatrix}$ commutes with V_0 . Thus, we have the reduced BCH formula $\log_m(V_0W_0) = \log_m(V_0) + \log_m(W_0) = \begin{pmatrix} \log_m(M) & 0 \\ 0 & 0 \end{pmatrix}$, i.e., Alg. 3.1 converges after a single iteration and gives

$$Log_U^{St}(UM) = U\log_m(M).$$
(5.1)

(Of course, it is also straight forward to show this directly without invoking Alg. 3.1.) Let $\sigma(M) = \{e^{i\varphi_1}, \ldots, e^{i\varphi_p}\}$ be the spectrum of $M \in O_{p \times p}$ and suppose that M is such that none of its eigenvalues is on the negative real axis, i.e., $\varphi_j \in (-\pi, \pi)$. Then,

⁴We may alternatively express this by saying that $[U] := \operatorname{colspan}(U)$ and $[\tilde{U}] := \operatorname{colspan}(\tilde{U})$ are the same points on the Grassmann manifold $[U] = [\tilde{U}] \in Gr(n, p)$.

the maximal Riemannian distance between two points U and UM is bounded by

$$\begin{aligned} \operatorname{dist}(U, UM) &= \sqrt{\langle U \log_m(M), U \log_m(M) \rangle_U} \\ &= \left(\frac{1}{2} \operatorname{tr}(\log_m(M)^T \log_m(M))\right)^{\frac{1}{2}} = \left(\frac{1}{2} \sum_{j=1}^p \varphi_j^2\right)^{\frac{1}{2}}. \end{aligned}$$

As a consequence

$$dist(U, UM) < \begin{cases} \sqrt{\frac{p}{2}}\pi, & p \text{ even}, \\ \sqrt{\frac{p-1}{2}}\pi, & p \text{ odd}. \end{cases}$$

The latter fact holds, because the eigenvalues of M come in complex conjugate pairs. Hence, if p is odd, there is at least one real eigenvalue $\lambda_j = e^{i\varphi_j}$ and because $\varphi_j \in (-\pi, \pi)$, there is at least one zero argument $\varphi_j = 0$. Related is [6, eq. (2.15)].

5.2. Numerical performance. First, we try to mimic the experiments featured in [18, §5.4]. Fig. 5.5 (lower left) of the aforementioned reference shows the average iteration count when applying the optimization-based Stiefel logarithm to matrices within a Riemannian annulus of inner radius 0.4π and outer radius 0.44π around $(I_p, 0)^T \in St(n, p)$ for dimensions of n = 10, p = 2. Convergence is detected, if $\|C_k\|_F < \tau = 10^{-7}$, where C_k is the same as in Alg. 3.1. ([18, Alg. 4, p. 91] uses $\tau^2 < 10^{-14}$). Since [18, §5.4] does not list the precise input data, we create comparable data randomly. To this end, we fix an arbitrary point $U \in St(10, 2)$ and create artificially but randomly another point $\tilde{U} \in St(10, 2)$ such that the Riemannian distance from U to \tilde{U} is exactly 0.44π . For full comparability, we replace the 2-norm in Alg. 3.1, line 7 with the Frobenius norm. We average over 1000 random experiments and arrive at an average iteration count of $\bar{k} = 7.83$. A MATLAB script that performs the required computations is available in Appendix F. When the distance of U and \tilde{U} is lowered to 0.4π , the average iteration count drops to a value of $\bar{k} = 6.92$.

As a second experiment, we now return to the 2-norm and lower the convergence threshold to $\|C_k\|_2 < \tau = 10^{-13}$ in the convergence criterion of Alg. 3.1. We create randomly points $U, \tilde{U} \in St(n, p)$ that are also a Riemannian distance of 0.44π away from each other, where we consider various different dimensions (n, p), see Table 5.1. We apply Alg. 3.1 to compute $\Delta = Log_U^{St}(\tilde{U})$.

(n,p)	$\operatorname{dist}\left(U,\tilde{U}\right)$	$\ U - \tilde{U}\ _2$	iters.	$\ \Delta - Log_U^{St}(\tilde{U})\ _2$	time
(10,2)	0.44π	1.0179	16	8.7903e-15	0.01s
(10,2)	0.89π	1.7117	95	4.1934e-13	0.06s
(1,000, 200)	0.44π	0.1616	5	1.5119e-14	0.7s
(1,000, 200)	0.89π	0.3256	7	1.7272e-14	0.8s
(1,000, 900)	0.44π	0.1234	4	9.6999e-14	16.1s
(1,000, 900)	0.89π	0.2491	5	7.9052e-14	21.0s
(100,000, 500)	0.44π	0.0875	4	5.9857 e-14	13.1s
(100,000, 500)	0.89π	0.1768	5	6.1041e-14	14.0s

 $\begin{array}{c} \text{TABLE 5.1}\\ \text{Convergence of Alg. 3.1 for random data to an accuracy of } \|C_k\|_2 \leq 10^{-13}. \end{array}$



FIG. 5.1. Convergence of Alg. 3.1 for random data $U, \tilde{U} \in St(n, p)$ for various n and p. Convergence accuracy is set to $\|C_k\|_2 \leq 10^{-13}$. Left: convergence graphs for $dist(U, \tilde{U}) = 0.44\pi$; right: for $dist(U, \tilde{U}) = 0.89\pi$.

Fig. 5.1 shows the associated convergence histories. The associated computation times⁵ are listed in Table 5.1. As can be seen from the figure and the table, Alg. 3.1 converges slowest (in terms of the iteration count) in the case of St(10, 2). Note that in this case, the constant $||U - \tilde{U}||_2$ that played a major role in the convergence analysis of Alg. 3.1 is largest. Moreover, we observe that the algorithm converges in all test cases even though in only one of the experiments the theoretical convergence guarantee $||U_0 - \tilde{U}||_2 < 0.09$ is satisfied, so that the theoretical bound derived here can probably be improved. Table 5.1 suggests that the impact of the size of $||U - \tilde{U}||_2$ on the iteration count is more direct than that of the actual Riemannian distance.

We repeat the exercise with random data $U, \tilde{U} \in St(n, p)$ that are a distance of 0.89π apart, which is the lower bound for the injectivity radius on the Stiefel manifold given in [18, eq. (5.14)]. In the case of St(10, 2), we hit a random matrix pair U, \tilde{U} , where the associated value $||U - \tilde{U}||_2$ is so large that the conditions of Theorem 4.1 and Lemma 4.2, Lemma 4.3 do not hold. In fact, we have $||\log_m(V_0)||_2 = 3.141$ for the starting point of Alg. 3.1 in this case, which is close to π . Yet, the algorithm converges, but very slowly so, see Table 5.1, second row and Fig. 5.1, right side. In all of the other cases, Alg. 3.1 converges in well under ten iterations, even for the larger test cases.

A MATLAB script that performs the required computations is available in Appendix F.

5.3. Dependence of the convergence on the Riemannian and the Euclidean distance. In this section, we examine the convergence of Alg. 3.1 depending on the Riemannian distance $dist(U, \tilde{U})$ and the distance $||U - \tilde{U}||_2$ in the Euclidean

 $^{^5 \}rm as$ measured on a Dell desktop computer endowed with six processors of type Intel(R) Core(TM) i7-3770 CPU@3.40GHz

operator-2-norm. To this end, we create a random point $U \in St(n, p)$ with MATLAB by computing the thin qr-decomposition of an $(n \times p)$ matrix with entries sampled uniformly from (0, 1). Likewise, we create a random tangent vector $\Delta \in T_U St(n, p)$ by chosing randomly a skew-symmetric matrix $A = \tilde{A} - \tilde{A}^T \in \mathbb{R}^{p \times p}$ and a matrix $T \in \mathbb{R}^{n \times p}$, where the entries of \tilde{A} and T are again uniformly sampled from (0, 1), and setting $\tilde{\Delta} = UA + (I - UU^T)T$. We normalize $\tilde{\Delta}$ according to the canonical metric $\Delta = \frac{\tilde{\Delta}}{\sqrt{\langle \tilde{\Delta}, \tilde{\Delta} \rangle_U}}$, see Section 2. In this way, we obtain for every $t \in [0, \pi)$ a point $\tilde{U} = U(t)$ the composition of L is the same for U = U(t).

U = U(t) that is a Riemannian distance of $dist(U, U(t))) = ||t\Delta||_U = t$ away from U. We discretize the interval $[0.1, 0.9\pi)$ by 100 equidistant points $\{x_k | k = 1, ..., 100\}$ and compute

- the number of iterations until convergence when computing $\log_U^{St}(U(t_k))$ with Alg. 3.1 for k = 1, ..., 100.
- the distance in spectral norm $||U U(t_k)||_2$, k = 1, ..., 100.
- the norm of the matrix logarithm of the first iterate $\|\log_m(V_0)\|_2$ from Alg. 3.1, step 3.

The results are displayed in Figures 5.2 – 5.4 for dimensions of St(10,000,400), St(100,10) and St(4,2), respectively. In all cases, the convergence threshold was set to $\|C_l\|_2 < \tau = 10^{-13}$. The algorithm converged in all cases, where $\|\log_m(V_0)\|_2 < \pi$ and produced a tangent vector $\Delta(t_k) := \log_U^{St}(U(t_k))$ of accuracy $\|\Delta(t_k) - t_k\Delta\|_2 < 10^{-13}$. A MATLAB script that performs the required computations is available in Appendix G. In the case of St(4,2), the algorithm starts to fail for $t_k \approx \frac{\pi}{2}$, where $\|\log_m(V_0)\|_2$



FIG. 5.2. Convergence of Alg. 3.1 for $U, \tilde{U} = U(t_k) = Exp_U^{St}(t_k\Delta) \in St(n,p)$, where Δ is a random tangent vector of canonical norm 1 and n = 10,000, p = 400. Convergence accuracy is set to $\|C_k\|_2 \leq 10^{-13}$. Left: number of iterations until convergence vs. $dist(U, \tilde{U})$; middle: $\|U - \tilde{U}\|_2$ vs. $dist(U, \tilde{U})$; right: $\|\log_m(V_0)\|_2$ vs. $dist(U, \tilde{U})$.

jumps to a value of π . This indicates that V_0 features (up to numerical errors) an eigenvalue $\lambda = -1$ so that the standard principal matrix logarithm is no longer well-defined. In all the experiments that were conducted, this behavior was observed only for small values of p < 8, while there was never produced a random data set where

Alg. 3.1 failed for $t < 0.9\pi$ and p > 10. The figures suggest that for small column-



FIG. 5.3. Same as Fig. 5.2, but for n = 100, p = 10.



FIG. 5.4. Same as Fig. 5.2, but for n = 4, p = 2.

numbers p, the ratio between the Riemannian distance $dist(U, \tilde{U})$ and the spectral distance $||U - \tilde{U}||_2$ is smaller than in higher dimensions. Moreover, for smaller p, it seems to be more likely to hit a random tangent direction along which Alg. 3.1 fails early than for higher p. This may partly be explained by the star-shaped nature of the domain of injectivity of the Riemannian exponential, [13, Lemma 5.7], and the

richer variety of directions in higher dimensions.

From these observations, it is tempting to conjecture that Alg. 3.1 will converge, whenever $\|\log_m(V_0)\|_2 < \pi$. However, these results are based on a limited notion of randomness and a more thorough examination of the numerical behavior of Alg. 3.1 is required to obtained conclusive results, which is beyond the scope of this work. Note that the domain of convergence of Alg. 3.1 is related to the injectivity radius of St(n,p) but it does not have to be the same. In Appendix C from the supplement, we state an explicit example in St(4,2), where Alg. 3.1 produces a first iterate V_0 with $\lambda = -1$ for an input pair $U, \tilde{U} \in St(4,2)$ with $dist(U,\tilde{U}) = \frac{\pi}{2}$, while the injectivity radius is estimated to be $\approx 0.71\pi$ in [18, §5]. An analytical investigation in St(4,2)might be possible and may shed more light on the precise value of the Stiefel manifold's injectivity radius.

6. Conclusions and outlook. We have presented a matrix-algebraic derivation of an algorithm for evaluating the Riemannian logarithm $Log_U^{St}(\tilde{U})$ on the Stiefel manifold. In contrast to [18, Alg. 4, p. 91], the construction here is not based on an optimization procedure but on an iterative solution to a non-linear matrix equation. Yet, it turns out that both approaches lead to essentially the same numerical scheme. More precisely, our Alg. 3.1 coincides with [18, Alg. 4, p. 91], when a unit step size is employed in the optimization scheme associated with the latter method. Apart from its comparatively simplicity, a key benefit is that our matrix-algebraic approach allows for a convergence analysis that does not require estimates on gradients nor Hessians and we are able to prove that the convergence rate of Alg. 3.1 is at least linear. This, in turn, proves the local linear convergence of [18, Alg. 4, p. 91] when using a unit step size. The algorithm shows a very promising performance in numerical experiments, even when the dimensions n, p become large.

So far, we have carried out a theoretical *local* convergence analysis. Open questions to be tackled in the future include estimates on how large the convergence domain of Alg. 3.1 is in terms of the Riemannian distance of the input points $dist(U, \tilde{U})$. This is related with the question of determining the injectivity radius of the Stiefel manifold. Estimates on the injectivity radius are featured in [18, §5.2.1].

Appendix A. A sharper majorizing series for Goldberg's Exponential series. As an alternative to Dynkin's BCH formula of nested commutators, Goldberg has shown in [8] that the solution to the exponential equation

$$\exp_m(X)\exp_m(Y) = \exp_m(Z)$$

can be written as a formal series

$$Z = X + Y + \sum_{k=2}^{\infty} z_k(X, Y), \quad z_k(X, Y) = \sum_{w, |w|=k} g_w w.$$
(A.1)

Each term $z_k(X, Y)$ in (A.1) is the sum over all *words* of length k in the alphabet $\{X, Y\}$. For example, $YXYX^2$ and X^2YXY^2 are such words of length 5 and 6 and thus contributing to $z_5(X, Y)$ and $z_6(X, Y)$, respectively. The coefficients are rational numbers $g_w \in \mathbb{Q}$, called Goldberg coefficients.

Thompson [20] has shown that the series converges provided that $||X||, ||Y|| \le \mu < 1$ for any submultiplicative norm $||\cdot||$. More precisely, his result is that $||z_k(X,Y)|| = ||\sum_{w,|w|=k} g_w w|| \le 2\mu^k$ for $k \ge 2$, see also [16, eq. 2]. In the next lemma, we improve this bound by cutting the factor 2.

LEMMA A.1. Let $||X||, ||Y|| \le \mu < 1$. The Goldberg series is majorized by

$$||Z|| < ||X|| + ||Y|| + \sum_{k=2}^{\infty} \mu^k$$

Proof. One ingredient of Thompson's proof is the following basic estimate on binomial terms:

$$m\begin{pmatrix} m-1\\ \lfloor \frac{m}{2} \rfloor \end{pmatrix} \ge 2^{m-1}.$$
 (A.2)

Here, $\lfloor x \rfloor$ denotes the largest integer smaller or equal to x. Thompson's argument is that $2^{m-1} = (1+1)^{m-1} = \sum_{l=0}^{m-1} \binom{m-1}{l}$ and that $\binom{m-1}{\lfloor \frac{m}{2} \rfloor}$ is the largest out of the m terms in the binomial sum. (It appears twice, if m-1 is odd.) In the following, we prefer to write this term with using the ceil-operator as $\binom{m-1}{\lfloor \frac{m}{2} \rfloor} = \binom{m-1}{\lceil \frac{m-1}{2} \rceil}$, because in this way, the same index m-1 appears in the upper and lower entry of the binomial coefficient.

For larger m, the inequality (A.2) can in fact be improved by a factor of 2:

Claim:
$$m\begin{pmatrix} m-1\\ \lceil \frac{m-1}{2} \rceil \end{pmatrix} > 2^m$$
 for all $m \ge 7$. (A.3)

For m = 7, we have $7 \begin{pmatrix} 7-1 \\ \lceil \frac{7-1}{2} \rceil \end{pmatrix} = 7 \cdot 20 = 140 > 128 = 2^7$; for m = 8, the inequality evaluates to $280 > 256 = 2^8$. To prove the claim, we proceed by induction.

Case 1: "m even". In this case, $\lceil \frac{m}{2} \rceil = \frac{m}{2} = \lceil \frac{m-1}{2} \rceil$ and

$$(m+1) \binom{m}{\lceil \frac{m}{2} \rceil} = (m+1) \left(\binom{m-1}{\frac{m}{2}-1} + \binom{m-1}{\frac{m}{2}} \right)$$
$$= 2(m+1) \binom{m-1}{\lceil \frac{m-1}{2} \rceil} > 2(m+1) \frac{2^m}{m} > 2^{m+1},$$
(A.4a)

where we have used the symmetry in the Pascal triangle (m - 1 is odd) and the induction hypothesis to arrive at (A.4a).

Case 2: "m odd". In this case, $\lceil \frac{m}{2} \rceil = \frac{m+1}{2}$ and

$$(m+1) \binom{m}{\left\lceil \frac{m}{2} \right\rceil} = (m+1) \left(\binom{m-1}{\frac{m+1}{2} - 1} + \binom{m-1}{\frac{m+1}{2}} \right)$$
$$= (m+1) \left(\binom{m-1}{\left\lceil \frac{m-1}{2} \right\rceil} + \binom{m-1}{\left\lceil \frac{m}{2} \right\rceil} \right).$$
(A.5a)

Note that $\binom{m-1}{\lceil \frac{m}{2} \rceil}$ is the second-to-largest term in the binomial expansion of $(1+1)^{m-1}$. Moreover, since m-1 is even, the relation to the largest term is

$$\binom{m-1}{\lceil \frac{m}{2} \rceil} = \frac{m-1}{m+1} \binom{m-1}{\lceil \frac{m-1}{2} \rceil}$$
17

Substituting in (A.5a) and applying the induction hypothesis gives

$$(m+1)\binom{m}{\lceil \frac{m}{2} \rceil} > (m+1)\left(\frac{2^m}{m} + \frac{m-1}{m+1}\frac{2^m}{m}\right) = \left(\frac{m+1}{m} + \frac{m-1}{m}\right)2^m = 2^{m+1}.$$

Using (A.3) rather than (A.2) in Thompson's original proof leads to the improved bound of $||z_k(X;Y)|| \le \mu^k$ for $k \ge 7$.

We tackle the terms involving words of lengths k = 2, 3, ..., 6 manually. The reference [21] lists explicit expressions of the summands in the Goldberg BCH series up to z_8 . The first three of them read

$$z_{2}(X,Y) = \frac{1}{2}(XY - YX) \Rightarrow ||z_{2}(X,Y)|| \le \frac{2}{2}\mu^{2}.\checkmark$$

$$z_{3}(X,Y) = \frac{1}{12}(X^{2}Y - 2XYX + XY^{2} + YX^{2} - 2YXY + Y^{2}X)$$

$$\Rightarrow ||z_{3}(X,Y)|| \le \frac{8}{12}\mu^{3}.\checkmark$$

$$z_{4}(X,Y) = \frac{1}{24}(X^{2}Y^{2} - 2XYXY + 2YXYX - Y^{2}X^{2}) \Rightarrow ||z_{4}(X,Y)|| \le \frac{6}{24}\mu^{4}.\checkmark$$

The expressions for $z_5(X, Y)$ and $z_6(X, Y)$ are too cumbersome to be restated here. However, for our purposes, a very rough counting argument is sufficient: The expression for $z_5(X, Y)$ features 30 length-5 words with non-zero Goldberg coefficient and the largest Goldberg coefficient is $\frac{1}{30}$. Hence, $||z_5(X,Y)|| = ||\sum_{w,|w|=5} g_w w|| < \frac{30}{30} \mu^5 \cdot \sqrt{(A more careful consideration reveals <math>||z_5(X,Y)|| \le \frac{176}{720} \mu^5$.) The expression for $z_6(X,Y)$ features 28 length-6 words with non-zero Goldberg Coefficient and the complete the second seco

The expression for $z_6(X, Y)$ features 28 length-6 words with non-zero Goldberg coefficient and the largest Goldberg coefficient is $\frac{1}{60}$. Hence, $||z_6(X, Y)|| = ||\sum_{w,|w|=6} g_w w|| \leq \frac{28}{60} \mu^6 \sqrt{\Box}$

Appendix B. Norm bound for the matrix logarithm.

PROPOSITION B.1. Let $C \in \mathbb{R}^{p \times p}$ be skew-symmetric with $\|C\|_2 < \pi$. Then

$$\|\exp_m(C) - I\|_2 < \|C\|_2.$$

Proof. Since C is skew-symmetric, it features an EVD $C = Q\Lambda Q^H$ with $\Lambda = diag(\lambda_1, \ldots, \lambda_p) = diag(i\varphi_1, \ldots, i\varphi_p)$, where $\varphi \in (-\pi, \pi)$ and $\max_j |i\varphi_j| = ||C||_2$. Therefore, $\exp_m(C) = Q \exp_m(\Lambda)Q^H$ with $\exp_m(\Lambda) = diag(e^{i\varphi_1}, \ldots, e^{i\varphi_p})$ and

$$\|\exp_m(C) - I\|_2 = \max_j |e^{i\varphi_j} - 1| < \max_j |\varphi_j| = \|C\|_2.$$

(The latter estimate may also be deduced from Fig. B.1.) \Box

PROPOSITION B.2. Let $V \in O_{n \times n}$ be such that $||V - I||_2 < r < 1$. Then

$$\|\log_m(V)\|_2 < r \frac{\sqrt{1 - \frac{r^2}{4}}}{1 - \frac{r^2}{2}}.$$

Proof. Let E = V - I. The matrices V and E share the same (orthonormal) basis of eigenvectors Q and the spectrum of V is precisely the spectrum of E shifted by +1. By assumption, $r > ||E||_2 = \max_{\mu \in \sigma(E)} |\mu|$. Hence, the eigenvalues $\lambda \in \sigma(V)$ are complex numbers of modulus one of the form $\lambda = e^{i\alpha} = 1 + \mu$, with $|\mu| < r$. Thus, λ lies on the unit circle but within a ball of radius r around $1 \in \mathbb{C}$, see Fig. B.1. The maximal angle α for such a λ is bounded by the slope of the line that starts in $0 \in \mathbb{C}$ and crosses the points of intersection of the two circles $\{|z| < 1\}$ and $\{|z - 1| < r\}$.

The intersection points are $(x_s, \pm y_s) = \left(1 - \frac{r^2}{2}, \pm r\sqrt{1 - \frac{r^2}{4}}\right)$. Therefore

$$\alpha| < \arctan\left(\frac{y_s}{x_s}\right) = \arctan\left(\frac{r\sqrt{1-\frac{r^2}{4}}}{1-\frac{r^2}{2}}\right) < r\frac{\sqrt{1-\frac{r^2}{4}}}{1-\frac{r^2}{2}}$$

As a consequence,

$$\|\log_m(V)\|_2 = \|Q\log_m(\Lambda)Q^H\|_2 = \max_{\lambda \in \sigma(V)} |\ln(\lambda)| = \max_{\lambda = e^{i\alpha} \in \sigma(V)} |i\alpha| < r \frac{\sqrt{1 - \frac{r^2}{4}}}{1 - \frac{r^2}{2}}$$



FIG. B.1. Geometrical illustration of Proposition B.2 in the complex plane.

LEMMA B.3. Let $U, \tilde{U} \in St(n, p)$ with $||U - \tilde{U}||_2 < \epsilon$. Let M, N, X_0, Y_0 and $V_0 := \begin{pmatrix} M & X_0 \\ N & Y_0 \end{pmatrix} \in O_{2p \times 2p}$ be as constructed in the first steps of Alg. 3.1. Then

$$\|\log_m(V_0)\|_2 < 2\varepsilon \frac{\sqrt{1-\varepsilon^2}}{1-2\varepsilon^2}.$$
(B.1)

Proof. Because V_0 is orthogonal,

$$1 = \|V_0\|_2 \ge \nu(V_0) = \max_{\|w\|_2 = 1} |w^H \begin{pmatrix} M & X_0 \\ N & Y_0 \end{pmatrix} w|$$
(B.2a)

$$\geq \max_{\|v\|_{2}=1} |(0, v^{H}) \begin{pmatrix} M & X_{0} \\ N & Y_{0} \end{pmatrix} \begin{pmatrix} 0 \\ v \end{pmatrix} | = \|Y_{0}\|_{2},$$
(B.2b)

where $\nu(V_0)$ denotes the numerical radius of V_0 , see [10, eq. 1.21, p. 21]. Likewise, $||M||_2 \leq 1$ so that the singular values of M and Y_0 range between 0 and 1. Moreover, by the Procrustes preprocessing outlined at the end of Section 3,

$$||N||_2 = ||X_0||_2 < \varepsilon, \quad ||M - I||_2 < \epsilon, \quad ||Y_0 - I_p||_2 < \varepsilon^2,$$

see (3.7). Combining these facts, we obtain V = I + (V - I) = I + E, where

$$||E||_{2} = || \begin{pmatrix} M - I & X_{0} \\ N & Y_{0} - I \end{pmatrix} ||_{2} \le || \begin{pmatrix} M - I & 0 \\ 0 & Y_{0} - I \end{pmatrix} + \begin{pmatrix} 0 & X_{0} \\ N & 0 \end{pmatrix} ||_{2} \le \max\{||M - I||_{2}, ||Y_{0} - I||_{2}\} + \max\{||N||_{2}, ||X_{0}||_{2}\} < \max\{\varepsilon, \varepsilon^{2}\} + \varepsilon = 2\varepsilon.$$

Applying Proposition B.2 to V = I + E proves the claim.

Appendix C. A critical special case. We present an example that shows that Alg. 3.1 may fail at computing $Log_U^{St}(\tilde{U})$ even for $U, \tilde{U} \in St(n, p)$ that are only a Riemannian distance of $dist(U, \tilde{U}) = \frac{\pi}{2}$ apart.

Consider n = 4, p = 2 and set

$$U = \frac{1}{2} \begin{pmatrix} 1 & 1 & 1 & 1 \\ 1 & 1 & -1 & -1 \end{pmatrix}^T \in St(4,2), \quad \Delta = \frac{1}{2} \begin{pmatrix} -1 & 1 & -1 & 1 \\ 0 & 0 & 0 & 0 \end{pmatrix}^T \in T_U St(4,2).$$

Note that $\Delta^T U = A = 0$ and that $\Delta = QR$ with $R = \begin{pmatrix} 1 & 0 \\ 0 & 0 \end{pmatrix}$ is the qr-decomposition of the tangent vector Δ . Hence, the Stiefel exponential (3.1) applied to this data set yields

$$\tilde{U}(t) = Exp_U^{St}(t\Delta) = (U,Q)\exp_m\left(\begin{pmatrix} 0 & -tR\\ tR & 0 \end{pmatrix}\right)\begin{pmatrix} I_2\\ 0 \end{pmatrix}.$$

Because of the simple structure of R, the matrix exponential can be computed explicitly

$$\exp_m\left(\begin{pmatrix} 0 & -tR\\ tR & 0 \end{pmatrix}\right) = \begin{pmatrix} \cos(t) & 0 & -\sin(t) & 0\\ 0 & 1 & 0 & 0\\ \sin(t) & 0 & \cos(t) & 0\\ 0 & 0 & 0 & 1 \end{pmatrix}$$

Recall from Section 2 that $dist(U, \tilde{U}(t)) = \sqrt{\langle t\Delta, t\Delta \rangle}_U$, which in this setting evaluates to t, since Δ is of unit norm also with respect to the canonical metric. For $t = \frac{\pi}{2}$, we obtain

$$\tilde{U} := \tilde{U} \begin{pmatrix} \frac{\pi}{2} \end{pmatrix} = U \begin{pmatrix} 0 & 0 \\ 0 & 1 \end{pmatrix} + Q \begin{pmatrix} 1 & 0 \\ 0 & 0 \end{pmatrix} = \frac{1}{2} \begin{pmatrix} -1 & 1 & -1 & 1 \\ 1 & 1 & -1 & -1 \end{pmatrix}^T \in St(4, 2).$$

If we now apply Alg. 3.1 to the matrix pair U, \tilde{U} , then we obtain in step 3 of the algorithm a corresponding

$$V_0 = \begin{pmatrix} 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix},$$
20

which features -1 as an eigenvalue and thus leads to a failure in the principal matrix logarithm. The problem here is the ambiguity in the orthogonal completion. If we replace the first row of the above V_0 with its negative, then we have still a valid orthogonal completion, and the method works. This example suggests that in a practical implementation of Alg. 3.1, one should try and explore strategies to compute a suitable starting iterate V_0 with small $\|\log_m(V_0)\|_2$.

Appendix D. Why is the Grassmann case simpler than the Stiefel case?. An important matrix manifold that is related with the Stiefel manifold and that arises frequently in applications is the *Grassmann manifold*. It is defined as the set of all p-dimensional subspaces $\mathcal{U} \subset \mathbb{R}^n$, i.e.,

$$Gr(n,p) := \{ \mathcal{U} \subset \mathbb{R}^n | \quad \mathcal{U} \text{ subspace, } \dim(\mathcal{U}) = p \}.$$

In this supplementary section, I give sketches for derivations for the Riemannian exponential and logarithm on the Grassmannian. Closed-form expressions for both mappings are known from the literature and I try to explain why the Stiefel case is more difficult. For background theory, the reader is referred to [2], [6].

The Grassmann manifold can be realized as a quotient manifold of the Stiefel manifold under actions of the orthogonal group via

$$Gr(n,p) = St(n,p)/O_{p \times p} = \{ [U] | \quad U \in St(n,p) \}.$$
 (D.1)

The quotient view point allows for using points $U \in St(n, p)$ as representatives for points $[U] \in Gr(n, p)$, i.e., subspaces, see [6] for details. For any matrix representative $U \in St(n, p)$ of $\mathcal{U} = [U] \in Gr(n, p)$, the tangent space at \mathcal{U} is represented by

$$T_{\mathcal{U}}Gr(n,p) = \left\{ \Delta \in \mathbb{R}^{n \times p} | \quad U^T \Delta = 0 \right\} \subset \mathbb{R}^{n \times p}.$$

This representation also stems from considering Gr(n, p) as a quotient manifold with St(n, p) as the total space. In fact, the tangent space of the Stiefel manifold can be decomposed into the so-called vertical space and the horizontal space with respect to the quotient mapping, $T_USt(n, p) = \mathcal{V}_U \oplus \mathcal{H}_U$, see [13, Problem 3.8], [2, §3.5.8], [6, §2.3.2]. The explicit representation of vectors in $T_{[U]}Gr(n, p)$ that we have introduced above corresponds to the identification of the actual abstract tangent space $T_{[U]}Gr(n, p)$ with the horizontal space \mathcal{H}_U .

From the quotient perspective, Grassmann tangent vectors are special Stiefel tangent vectors $\Delta = U_0 A + (I - UU^T)T$, namely those associated with the special skew-symmetric matrix $A = 0 \in \mathbb{R}^{p \times p}$, cf. (2.1). Hence, we may use the Stiefel exponential to compute the Grassmann exponential:

- Given $\Delta \in T_{[U]}Gr(n,p)$, compute the qr-decomposition $(I U_0 U_0^T)\Delta = \Delta = Q_E R_E$.
- Compute the matrix exponential

$$\begin{pmatrix} M \\ N_E \end{pmatrix} := \exp_m \left(\begin{pmatrix} 0 & -R_E^T \\ R_E & 0 \end{pmatrix} \right) \begin{pmatrix} I_p \\ 0 \end{pmatrix}.$$
(D.2)

• Return $\tilde{U} = Exp^{Gr}_{[U_0]}(\Delta) = [U_0M + Q_EN_E].$

It is precisely the extra upper-left zero-block in the matrix exponential in (D.2), that makes the Grassmann case easier to tackle than the Stiefel case: By using the SVD $R_E = \Phi \Sigma D^T$ and the series expansion of \exp_m , it is straight-forward to show that

$$\exp_m \left(\begin{pmatrix} 0 & -R_E^T \\ R_E & 0 \end{pmatrix} \right) = \begin{pmatrix} D & 0 \\ 0 & \Phi \end{pmatrix} \begin{pmatrix} \cos(\Sigma) & -\sin(\Sigma) \\ \sin(\Sigma) & \cos(\Sigma) \end{pmatrix} \begin{pmatrix} D^T & 0 \\ 0 & \Phi^T \end{pmatrix}, \quad (D.3)$$
²¹

which gives

$$\begin{pmatrix} M \\ N_E \end{pmatrix} = \begin{pmatrix} D\cos(\Sigma)D^T \\ \Phi\sin(\Sigma)D^T \end{pmatrix}.$$
 (D.4)

(In the above formulae, it is understood that sin and cos are to be applied pointwise to the diagonal elements of the diagonal matrix Σ .) Eventually, we arrive at

$$Exp_{\mathcal{U}_0}^{Gr}(\Delta) = [U_0M + QN_E] = [U_0D\cos(\Sigma)D^T + Q\Phi\sin(\Sigma)D^T]$$

Instead of starting with the qr-decomposition $\Delta = Q_E R_E$, we now see that we could have directly worked with the SVD $\Delta = \hat{Q} \Sigma D^T (= (Q\Phi) \Sigma D^T)$, which yields $Exp_{\mathcal{U}_0}^{Gr}(\Delta) = [U_0 D \cos(\Sigma) D^T + \hat{Q} \sin(\Sigma) D^T].$

This is exactly the expression that Edelman et al. have found in [6, Thm. 2.3] for the Riemannian exponential on Gr(n, p) and the derivation above can be considered as a 'thin SVD'-version of [6, Thm. 2.3, Proof 2, p. 320].

The inverse of this mapping, i.e., the Riemannian logarithm on Gr(n, p) can be deduced as follows: Consider $[U_0], [\tilde{U}] \in Gr(n, p)$. Under the assumption that $[\tilde{U}]$ is sufficiently close to [U], it holds that $\tilde{U} = U_0M + QN$ and the task is to find $M, N \in \mathbb{R}^{p \times p}$ and $Q \in St(n, p)$ such that $Q^T U_0 = 0$. The first matrix factor M is uniquely determined by $U_0^T \tilde{U} = M$. We obtain candidates for Q, N by computing the qr-decomposition $(I - U_0 U_0^T) \tilde{U} = QN$. Yet, in order to reverse (D.4), it is require to work with consistent coordinates. Taking (D.3), (D.2) into account, this is established by setting $N_L = NM^{-1}$ and computing the SVD $N_L = \Phi SD^T$, because by defining $\Sigma = \arctan(S)$, we can decompose

$$N_L = \Phi S D^T = \Phi \tan(\Sigma) D^T = \Phi \sin(\Sigma) D^T D(\cos(\Sigma))^{-1} D^T.$$

This shows that the choice $R_L := \Phi \Sigma D^T$ yields a tangent vector $\Delta := QR_L$ such that

$$Exp_{\mathcal{U}_0}^{Gr}(\Delta) = \left[(U_0, Q) \exp_m \left(\begin{pmatrix} 0 & -R_L^T \\ R_L & 0 \end{pmatrix} \right) \begin{pmatrix} I_p \\ 0 \end{pmatrix} \right]$$
$$= \left[U_0 D \cos(\Sigma) D^T + Q \Phi \sin(\Sigma) D^T \right] = \left[\tilde{U} \right]$$

Note that $QN_L = Q\Phi SD^T \stackrel{\text{SVD}}{=} QNM^{-1} = (I - U_0U_0^T)\tilde{U}M^{-1}$. Hence, we now see that we could have directly started with the SVD of $\hat{Q}SD^T = (I - U_0U_0^T)\tilde{U}M^{-1}$ to arrive at

$$Log_{\mathcal{U}_0}^{Gr}(\tilde{\mathcal{U}}) = \Delta = \hat{Q} \arctan(S) D^T.$$

This is the well-known closed-form of the Grassmann logarithm. Unfortunatley, I was not able to track down the original derivation. The earliest appearance in the literature that I found was [4, Alg. 3]. However, this reference only mentions the above formula but does not cite a source. In summary, the Grassmann case is easier to deal with because of the extra off-diagonal block structure in the associated matrix exponential (D.2), which leads to a CS-decomposition in (D.3) by a *similarity transformation*; compare this to [9, Thm. 2.6.3, p.78].

Appendix E. MATLAB code.

```
E.1. Alg. 3.1.
%
function [Delta, k, conv_hist, norm_logV0] = ...
                          Stiefel_Log_supp(U0, U1, tau)
%_-----
%@author: Ralf Zimmermann, IMADA, SDU Odense
%
% Input arguments
% UO, U1 : points on St(n,p)
% tau : convergence threshold
% Output arguments
% Delta : Log<sup>{</sup>St}_UO(U1),
%
           i.e. tangent vector such that Exp^St_UO(Delta) = U1
%
       k : iteration count upon convergence
% supplementary output
% conv_hist : convergence history
% norm_logVO : norm of matrix log of first iterate VO
%-----
% get dimensions
[n,p] = size(U0);
% store convergence history
conv_hist = [0];
% step 1
M = UO'*U1;
% step 2
[Q,N] = qr(U1 - U0*M,0); % thin qr of normal component of U1
% step 3
[V, ~] = qr([M;N]);
                                  % orthogonal completion
% "Procrustes preprocessing"
[D,S,R] = svd(V(p+1:2*p,p+1:2*p));
V(:,p+1:2*p) = V(:,p+1:2*p)*(R*D');
           = [[M;N], V(:,p+1:2*p)]; %
V
                                      |M XO|
                                   \% now, V = |N YO|
% just for the record
norm_logV0 = norm(logm(V),2);
% step 4: FOR-Loop
for k = 1:10000
   % step 5
   [LV, exitflag] = logm(V);
                              % standard matrix logarithm
                              % |Ak -Bk'|
                              \% now, LV = |Bk Ck |
   C = LV(p+1:2*p, p+1:2*p);
                              % lower (pxp)-diagonal block
   % steps 6 - 8: convergence check
   normC = norm(C, 2);
   conv_hist(k) = normC;
```

```
if normC<tau;</pre>
        disp(['Stiefel log converged after ', num2str(k),...
              ' iterations.']);
        break;
   end
   % step 9
   Phi = expm(-C);
                                 % standard matrix exponential
   % step 10
   V(:,p+1:2*p) = V(:,p+1:2*p)*Phi; % update last p columns
end
                                         |A −B'|
% prepare output
% upon convergence, we have logm(V) = |B \ 0| = LV
     A = LV(1:p, 1:p);
                         B = LV(p+1:2*p, 1:p)
%
% Delta = UO*A+Q*B
Delta = U0*LV(1:p,1:p) + Q*LV(p+1:2*p, 1:p);
return:
end
```

Note: The performance of this method may be enhanced by computing $\exp_m,\,\log_m$ via a Schur decomposition.

Appendix F. MATLAB code corresponding to Section 5.2.

First experiment discribed in Section 5.2.

```
%_____
% script_Stiefel_Log_supp52.m
% %Qauthor: Ralf Zimmermann, IMADA, SDU Odense
%_____
clear;
% set dimensions
n = 10;
p = 2;
% fix stream of random numbers for reproducability
s = RandStream('mt19937ar', 'Seed',1);
% set number of random experiments
runs = 100;
dist = 0.4*pi;
average_iters = 0;
for j=1:runs
   %create random stiefel data
   [U0, U1, Delta] = create_random_Stiefel_data(s, n, p, dist);
   % 'project' Delta onto St(n,p) via the Stiefel exponential
   U1 = Stiefel_Exp_supp(U0, Delta);
   % compute the Stiefel logarithm
   [Delta_rec, k] = Stiefel_Log_supp(U0, U1, 1.0e-13);
                   % uncomment the following lines to check
                      % if Stiefel logarithm recovers Delta
   %norm(Delta_rec - Delta)
   average_iters = average_iters +k;
end
average_iters = average_iters/runs;
                              24
```

```
disp(['The average iteration count of the Stiefel log is ',...
    num2str(average_iters)]);
% EOF: script_Stiefel_Log_supp52.m
%_____
  Second experiment discribed in Section 5.2.
%_____
% script_Stiefel_Log_supp52b.m
% %Qauthor: Ralf Zimmermann, IMADA, SDU Odense
%-----
clear; close all;
dist = 0.44*pi;
<u>%_____</u>
% set dimensions
n = 10;
p = 2;
% fix stream of random numbers for reproducability
s = RandStream('mt19937ar', 'Seed',1);
%create random stiefel matrix:
[U0, U1, Delta] = create_random_Stiefel_data(s, n, p, dist);
norm_UO_U1 = norm(U0 - U1, 2)
% compute the Stiefel logarithm
tic;
[Delta_rec, k, conv_hist1, norm_logV01] = ...
                   Stiefel_Log_supp(U0, U1, 1.0e-13);
toc;
norm_recon11 = norm(Delta_rec - Delta)
%-----
%_____
\% reset dimensions
n = 1000;
p = 200;
%create random stiefel matrix:
[U0, U1, Delta] = create_random_Stiefel_data(s, n, p, dist);
norm_UO_U1 = norm(UO - U1, 2)
% compute the Stiefel logarithm
tic:
[Delta_rec, k, conv_hist2, norm_logV02] = ...
                   Stiefel_Log_supp(U0, U1, 1.0e-13);
toc:
norm_recon12 = norm(Delta_rec - Delta)
<u>%_____</u>
%_____
% reset dimensions
```

```
n = 1000;
p = 900;
%create random stiefel matrix:
[U0, U1, Delta] = create_random_Stiefel_data(s, n, p, dist);
norm_UO_U1 = norm(UO - U1, 2)
% compute the Stiefel logarithm
tic;
[Delta_rec, k, conv_hist3, norm_logV03] = ...
                      Stiefel_Log_supp(U0, U1, 1.0e-13);
toc;
norm_recon13 = norm(Delta_rec - Delta)
%_____
%_____
% reset dimensions
n = 100000:
p = 500;
%create random stiefel matrix:
[U0, U1, Delta] = create_random_Stiefel_data(s, n, p, dist);
norm_UO_U1 = norm(UO - U1, 2)
% compute the Stiefel logarithm
tic;
[Delta_rec, k, conv_hist4, norm_logV04] = ...
                      Stiefel_Log_supp(U0, U1, 1.0e-13);
toc;
norm_recon14 = norm(Delta_rec - Delta)
<u>%_____</u>
% plot convergence history
figure;
subplot(1,2,1);
semilogy(1:length(conv_hist1), conv_hist1, 'k-s', ...
```

% set dimensions

```
n = 100;
p = 10;
% fix stream of random numbers for reproducability
s = RandStream('mt19937ar', 'Seed',1);
%create random stiefel data
[U0, U1, Delta] = create_random_Stiefel_data(s, n, p, 1.0);
% discretize the interval [0.1, 0.9pi] with resolution res
res = 100;
start = 0.01;
t = linspace(start, 0.9*pi, res)';
%*********************
% initialize observations
% spectral distance U, Uk
norm_U_Uk = zeros(res,1);
% iterations until convergence
iters_convk = zeros(res,1);
% norm log(VO)
norm_logV0k = zeros(res,1);
\% accuracy of the reconstruction
norm_Delta_Delta_rec_k = zeros(res,1);
for k = 1:res
   % 'project' tDelta onto St(n,p) via the Stiefel exponential
   Uk = Stiefel_Exp_supp(U0, t(k)*Delta);
   % compute spectral norm
   norm_U_Uk(k) = norm(U0-Uk,2);
   % execute the Stiefel logarithm
   disp(['Compute log for t=', num2str(t(k))]);
    [Delta_rec, iters_conv, conv_hist, norm_logV0] = ...
       Stiefel_Log_supp(U0, Uk, 1.0e-13);
   % store data
    iters_convk(k) = iters_conv;
   norm_logVOk(k) = norm_logVO;
   norm_Delta_Delta_rec_k(k) = norm(t(k)*Delta-Delta_rec, 2);
end
% visualize results
figure;
subplot(1,3,1);
plot(t, iters_convk, 'k-');
legend('iters until convergence');
hold on
subplot(1,3,2);
plot(t, norm_U_Uk, 'k-');
```

```
Appendix H. Auxiliary MATLAB functions.
```

```
Stiefel exponential.
%_-----
%file: Stiefel_Exp_supp.m
% Cauthor: Ralf Zimmermann, IMADA, SDU Odense
%_____
function [U1] = Stiefel_Exp_supp(U0, Delta)
%-----
% Input arguments
%
 UO
     : base point on St(n,p)
%
  Delta : tangent vector in T_UO St(n,p)
% Output arguments
% U1 : Exp^{St}_U0(Delta),
%_-----
% get dimensions
[n,p] = size(U0);
A = U0'*Delta;
                           % horizontal component
K = Delta-U0*A;
                              % normal component
[Qe, Re] = qr(K, 0);
                          % qr of normal component
% matrix exponential
MNe = expm([[A, -Re']; [Re, zeros(p)]]);
U1 = [U0, Qe] *MNe(:,1:p);
return;
end
%EOF: Stiefel_Exp_supp.m
%-----
```

Construction of random data on the Stiefel manifold.

%file: create_random_Stiefel_data.m
% @author: Ralf Zimmermann, IMADA, SDU Odense
%-----function [U0, U1, Delta] =...
 create_random_Stiefel_data(s, n, p, dist)
%------% create a random data set
% U0, U1 on St(n,p),

<u>%_____</u>

```
\% Delta on T_U St(n,p) with canonical norm 'dist',
% which is also the Riemannian distance dist(U0,U1)
%
% input arguments
%
      s = random stream (for reproducability)
% (n,p) = dimension of the Stiefel matrices
% dist = Riemannian distance between the points U0,U1
%
          that are to be created
%-
%create random stiefel matrix:
X = rand(s, n, p);
[U0,~] = qr(X, 0);
% create random tangent vector in T_UO St(n,p)
A = rand(s, p, p);
A = A - A';
                     % random p-by-p skew symmetric matrix
T = rand(s, n, p);
Delta = U0*A + T - U0*(U0'*T);
%normalize Delta w.r.t. the canonical metric
norm_Delta = sqrt(trace(Delta'*Delta) - 0.5*trace(A'*A));
Delta = (dist/norm_Delta)*Delta;
% 'project' Delta onto St(n,p) via the Stiefel exponential
U1 = Stiefel_Exp_supp(U0, Delta);
return;
end
%EOF: create_random_Stiefel_data.m
%-----
```

REFERENCES

- P.-A. ABSIL, R. MAHONY, AND R. SEPULCHRE, Riemannian geometry of Grassmann manifolds with a view on algorithmic computation, Acta Applicandae Mathematica, 80 (2004), pp. 199–220,
- P.-A. ABSIL, R. MAHONY, AND R. SEPULCHRE, Optimization Algorithms on Matrix Manifolds, Princeton University Press, Princeton, New Jersey, 2008,
- [3] L. BALZANO, R. NOWAK, AND B. RECHT, Online identification and tracking of subspaces from highly incomplete information, in Proceedings of Allerton, September 2010.
- [4] E. BEGELFOR AND M. WERMAN, Affine invariance revisited, 2012 IEEE Conference on Computer Vision and Pattern Recognition, 2 (2006), pp. 2087–2094,
- [5] P. BENNER, S. GUGERCIN, AND K. WILLCOX, A survey of projection-based model reduction methods for parametric dynamical systems, SIAM Review, 57 (2015), pp. 483–531,
- [6] A. EDELMAN, T. A. ARIAS, AND S. T. SMITH, The geometry of algorithms with orthogonality constraints, SIAM Journal on Matrix Analysis and Application, 20 (1999), pp. 303–353,
- [7] K. GALLIVAN, A. SRIVASTAVA, X. LIU, AND P. VAN DOOREN, Efficient algorithms for inferences on Grassmann manifolds, in Statistical Signal Processing, 2003 IEEE Workshop on, 2003, pp. 315–318,
- [8] K. GOLDBERG, The formal power series for $\log e^x e^y$, Duke Math. J., 23 (1956), pp. 13–21,
- [9] G. H. GOLUB AND C. F. VAN LOAN, Matrix Computations, The John Hopkins University Press, Baltimore – London, 3 ed., 1996.
- [10] A. GREENBAUM, Iterative Methods for solving linear systems, vol. 17 of Frontiers in applied mathematics, SIAM Society for Industrial and Applied Mathematics, Philadelphia, 1997.
- [11] N. J. HIGHAM, Functions of Matrices: Theory and Computation, Society for Industrial and Applied Mathematics, Philadelphia, PA, USA, 2008.
- [12] S. KOBAYASHI AND K. NOMIZU, Foundations of Differential Geometry, vol. I of Interscience

Tracts in Pure and Applied Mathematics no. 15, John Wiley & Sons, New York – London – Sidney, 1963.

- [13] J. LEE, Riemannian Manifolds: an Introduction to Curvature, Springer Verlag, New York Berlin – Heidelberg, 1997.
- [14] Y. MAN LUI, Advances in matrix manifolds for computer vision, Image and Vision Computing, 30 (2012), pp. 380–388,
- [15] MATLAB, version 7.10.0 (R2010a), The MathWorks Inc., Natick, Massachusetts, 2010.
- [16] M. NEWMAN, S. WASIN, AND R. C. THOMPSON, Convergence domains for the Campbell-Baker-Hausdorff formula, Linear and Multilinear Algebra, 24 (1989), pp. 301–310.
- [17] I. U. RAHMAN, I. DRORI, V. C. STODDEN, D. L. DONOHO, AND P. SCHRÖDER, Multiscale representations for manifold-valued data, SIAM J. Mult. Model. Simul., 4 (2005), pp. 1201– 1232.
- [18] Q. RENTMEESTERS, Algorithms for data fitting on some common homogeneous spaces, PhD thesis, Université Catholique de Louvain, Louvain, Belgium, July 2013/2015
- [19] W. ROSSMANN, Lie Groups: An Introduction Through Linear Groups, Oxford graduate texts in mathematics, Oxford University Press, 2006,
- [20] R. C. THOMPSON, Convergence proof for Goldberg's exponential series, Linear Algebra and its Applications, 121 (1989), pp. 3–7.
- [21] A. VAN-BRUNT AND M. VISSER, Simplifying the Reinsch algorithm for the Baker-Campbell-Hausdorff series. arXiv:1501.05034, 2015,