

1 **ITERATIVE COMPUTATION OF THE FRÉCHET DERIVATIVE OF**
2 **THE POLAR DECOMPOSITION***

3 EVAN S. GAWLIK[†] AND MELVIN LEOK[†]

4 **Abstract.** We derive iterative methods for computing the Fréchet derivative of the map which
5 sends a full-rank matrix A to the factor U in its polar decomposition $A = UH$, where U has
6 orthonormal columns and H is Hermitian positive definite. The methods apply to square matrices
7 as well as rectangular matrices having more rows than columns. Our derivation relies on a novel
8 identity that relates the Fréchet derivative of the polar decomposition to the matrix sign function
9 $\text{sign}(X) = X(X^2)^{-1/2}$ applied to a certain block matrix X .

10 **Key words.** Polar decomposition, Fréchet derivative, matrix function, matrix iteration, Newton
11 iteration, Newton-Schulz iteration, matrix sign function

12 **AMS subject classifications.** 65F30, 15A23, 15A24

13 **1. Introduction.** The polar decomposition theorem asserts that every matrix
14 $A \in \mathbb{C}^{m \times n}$ ($m \geq n$) can be written as the product $A = UH$ of a matrix $U \in \mathbb{C}^{m \times n}$
15 having orthonormal columns times a Hermitian positive semidefinite matrix $H \in$
16 $\mathbb{C}^{n \times n}$ [18, Theorem 8.1]. If A is full-rank, then this decomposition is unique and H
17 is positive definite, allowing one to define a map \mathcal{P} which sends a full-rank matrix
18 $A \in \mathbb{C}^{m \times n}$ to the factor $\mathcal{P}(A) = U \in \mathbb{C}^{m \times n}$ in its polar decomposition $A = UH$. We
19 refer to U as the unitary factor in the polar decomposition of A , bearing in mind that
20 this is a slight abuse of terminology when A (and hence U) is rectangular. The aim of
21 this paper is to derive iterative algorithms for computing the Fréchet derivative of \mathcal{P} .

22 Our interest in differentiating the polar decomposition stems from several sources.
23 First, differentiating the polar decomposition gives precise information about the sen-
24 sitivity of the polar decomposition to perturbations. This is a topic of longstanding
25 interest in numerical analysis [25, 26, 6, 21, 28, 7, 4, 24], where much of the literature
26 has focused on bounding the deviations in the perturbed factors in the polar decom-
27 position of A after a small-normed perturbation of A . These analyses often rely on a
28 formula for the Fréchet derivative of \mathcal{P} that involves the singular value decomposition
29 of A [21, Equation 2.18]. While theoretically useful, such a formula loses some of its
30 appeal in the numerical setting, where computing the singular value decomposition
31 tends to be costly. As a second source of motivation, differentiating the polar de-
32 composition has proven necessary in the design of certain schemes for interpolating
33 functions which take values in the special orthogonal group [12], the group of real
34 square matrices with orthonormal columns and positive determinant. These interpo-
35 lation schemes have applications in computer animation, mechanics, and other areas
36 in which continuously varying rotation matrices play a role.

37 A number of authors have addressed the computation of the Fréchet derivatives
38 of other functions of matrices, such as the matrix exponential [1, 27, 30], the matrix
39 logarithm [3, 23], the matrix square root [1, Section 2], the matrix p^{th} root [19, 9,
40 8], and the matrix sign function $\text{sign}(X) = X(X^2)^{-1/2}$ [21]. The aforementioned

*Submitted to the editors August 15, 2016.

Funding: EG has been supported in part by NSF under grants DMS-1411792, DMS-1345013. ML has been supported in part by NSF under grants DMS-1010687, CMMI-1029445, DMS-1065972, CMMI-1334759, DMS-1411792, DMS-1345013.

[†]Department of Mathematics, University of California, San Diego (egawlik@ucsd.edu, mleok@math.ucsd.edu).

41 functions, unlike the map \mathcal{P} , are examples of *primary matrix functions*. Roughly
 42 speaking, a primary matrix function is a scalar function that has been extended to
 43 square matrices in a canonical way; for a precise definition, see [18, Section 1.2]
 44 and [20]. The polar decomposition is not a primary matrix function, which is perhaps
 45 the main reason that the computation of its Fréchet derivative (a quantity whose
 46 existence is justified in Section 3.1) has largely evaded scrutiny until now.

47 Formally, iterative schemes for computing the Fréchet derivatives of matrix func-
 48 tions (be they primary or nonprimary) can be derived as follows. Let $f : \mathbb{C}^{m \times n} \rightarrow$
 49 $\mathbb{C}^{m \times n}$ be a function with Fréchet derivative L_f . That is, given $X \in \mathbb{C}^{m \times n}$, the map
 50 $L_f(X, \cdot) : \mathbb{C}^{m \times n} \rightarrow \mathbb{C}^{m \times n}$ is a linear map satisfying

$$51 \quad (1) \quad f(X + E) - f(X) - L_f(X, E) = o(\|E\|)$$

52 for every $E \in \mathbb{C}^{m \times n}$, where $\|\cdot\|$ denotes any matrix norm. Let $A \in \mathbb{C}^{m \times n}$, and
 53 suppose that

$$54 \quad (2) \quad X_{k+1} = g(X_k), \quad X_0 = A$$

55 is an iterative scheme for computing $f(A)$; that is, $X_k \rightarrow f(A)$ as $k \rightarrow \infty$. Then,
 56 as noted in [1, Section 2], differentiation of (2) with respect to A in the direction
 57 $E \in \mathbb{C}^{m \times n}$ yields the coupled iteration

$$58 \quad (3) \quad X_{k+1} = g(X_k), \quad X_0 = A,$$

$$59 \quad (4) \quad E_{k+1} = L_g(X_k, E_k), \quad E_0 = E,$$

61 for computing $f(A)$ and $L_f(A, E)$. The validity of this formal derivation, of course,
 62 depends on the commutativity of $\lim_{k \rightarrow \infty}$ with differentiation, which is generally
 63 nontrivial to establish.

64 For a primary matrix function f , proving the validity of this formal derivation is
 65 greatly simplified by the following identity. For any primary matrix function f and
 66 any square matrices A and E ,

$$67 \quad (5) \quad f \begin{pmatrix} A & E \\ 0 & A \end{pmatrix} = \begin{pmatrix} f(A) & L_f(A, E) \\ 0 & f(A) \end{pmatrix},$$

68 provided that f is $2p - 1$ times continuously differentiable on an open subset of \mathbb{C}
 69 containing the spectrum of A , where p is the size of the largest Jordan block of A [29].
 70 From this it follows that if (2) is an iterative scheme for computing $f(A)$, and if g
 71 maps block upper triangular matrices to block upper triangular matrices, then

$$72 \quad (6) \quad \begin{pmatrix} X_{k+1} & E_{k+1} \\ 0 & X_{k+1} \end{pmatrix} = g \begin{pmatrix} X_k & E_k \\ 0 & X_k \end{pmatrix}, \quad \begin{pmatrix} X_0 & E_0 \\ 0 & X_0 \end{pmatrix} = \begin{pmatrix} A & E \\ 0 & A \end{pmatrix}$$

73 defines an iterative scheme for computing $\begin{pmatrix} f(A) & L_f(A, E) \\ 0 & f(A) \end{pmatrix}$, provided that it con-
 74 verges and provided that f has the requisite regularity to apply (5). Using (5) again to
 75 isolate each block of the iteration (6), one obtains the coupled iteration (3-4). Details
 76 behind this argument, as well as an example of its application, can be found in [1,
 77 Section 2].

78 Our main result in this paper, Theorem 1, establishes the validity of schemes
 79 like (3-4) when the function f under consideration is the function \mathcal{P} which sends A

80 to the unitary factor U in its polar decomposition, *even though \mathcal{P} is not a primary*
 81 *matrix function.* In particular,

$$82 \quad \mathcal{P} \begin{pmatrix} A & E \\ 0 & A \end{pmatrix} \neq \begin{pmatrix} \mathcal{P}(A) & L_{\mathcal{P}}(A, E) \\ 0 & \mathcal{P}(A) \end{pmatrix},$$

83 so the argument in the preceding paragraph does not apply. For example, if $A = 2$
 84 and $E = 3$, then it is not hard to check that $\mathcal{P} \begin{pmatrix} A & E \\ 0 & A \end{pmatrix} = \begin{pmatrix} 4/5 & 3/5 \\ -3/5 & 4/5 \end{pmatrix}$, but

$$85 \quad \begin{pmatrix} \mathcal{P}(A) & L_{\mathcal{P}}(A, E) \\ 0 & \mathcal{P}(A) \end{pmatrix} = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}. \text{ Instead of using (5), our derivation relies on a novel}$$

86 identity that relates the Fréchet derivative of \mathcal{P} to the matrix sign function $\text{sign}(X) =$
 87 $X(X^2)^{-1/2}$ applied to a certain block matrix X ; see Theorem 2.

88 One notable corollary of Theorem 1 is that the popular Newton iteration [16]

$$89 \quad X_{k+1} = \frac{1}{2}(X_k + X_k^{-*}), \quad X_0 = A$$

90 for computing the unitary factor $\mathcal{P}(A) = U$ in the polar decomposition $A = UH$ of a
 91 square matrix A extends to a coupled iteration for computing $\mathcal{P}(A)$ and its Fréchet
 92 derivative. In particular, Corollary 3 shows that for any nonsingular $A \in \mathbb{C}^{n \times n}$ and
 93 any $E \in \mathbb{C}^{n \times n}$, the scheme

$$94 \quad X_{k+1} = \frac{1}{2}(X_k + X_k^{-*}), \quad X_0 = A,$$

$$95 \quad E_{k+1} = \frac{1}{2}(E_k - X_k^{-*} E_k^* X_k^{-*}), \quad E_0 = E,$$

$$96$$

97 produces iterates X_k and E_k that converge to $\mathcal{P}(A) = U$ and $L_{\mathcal{P}}(A, E)$, respectively,
 98 as $k \rightarrow \infty$.

99 The fact that the matrix sign function will play a role in our study of Fréchet
 100 derivatives of the polar decomposition should come as no surprise, given the sign
 101 function's intimate connection with the polar decomposition. The sign function and
 102 polar decomposition are linked via the identity

$$103 \quad (7) \quad \text{sign} \begin{pmatrix} 0 & A \\ A^* & 0 \end{pmatrix} = \begin{pmatrix} 0 & \mathcal{P}(A) \\ \mathcal{P}(A)^* & 0 \end{pmatrix},$$

104 which holds for any square nonsingular matrix A [18]. This identity has been used,
 105 among other things, to derive iterative schemes for computing the polar decomposi-
 106 tion. The essence of this approach is to write down an iterative scheme for computing
 107 $\text{sign} \begin{pmatrix} 0 & A \\ A^* & 0 \end{pmatrix}$, check that its iterates retain the relevant block structure, and read
 108 off the (1,2)-block of the resulting algorithm. In principle, one can adopt a similar
 109 strategy to derive iterative schemes for computing the Fréchet derivatives of the polar
 110 decomposition. Indeed, any iterative scheme that computes

$$111 \quad \text{sign} \begin{pmatrix} 0 & A & 0 & E \\ A^* & 0 & E^* & 0 \\ 0 & 0 & 0 & A \\ 0 & 0 & A^* & 0 \end{pmatrix}$$

112 while retaining its block structure will suffice, owing to the following observation. By
 113 appealing to the definition (1) of the Fréchet derivative, the identity (7) can be used

114 to verify that

$$115 \quad (8) \quad L_{\text{sign}} \left(\begin{pmatrix} 0 & A \\ A^* & 0 \end{pmatrix}, \begin{pmatrix} 0 & E \\ E^* & 0 \end{pmatrix} \right) = \begin{pmatrix} 0 & L_{\mathcal{P}}(A, E) \\ L_{\mathcal{P}}(A, E)^* & 0 \end{pmatrix}.$$

116 Now since the sign function is a primary matrix function, (5), (7), and (8) imply that

$$117 \quad \text{sign} \begin{pmatrix} 0 & A & 0 & E \\ A^* & 0 & E^* & 0 \\ 0 & 0 & 0 & A \\ 0 & 0 & A^* & 0 \end{pmatrix} = \begin{pmatrix} \text{sign} \begin{pmatrix} 0 & A \\ A^* & 0 \end{pmatrix} & L_{\text{sign}} \left(\begin{pmatrix} 0 & A \\ A^* & 0 \end{pmatrix}, \begin{pmatrix} 0 & E \\ E^* & 0 \end{pmatrix} \right) \\ 0 & \text{sign} \begin{pmatrix} 0 & A \\ A^* & 0 \end{pmatrix} \end{pmatrix}$$

$$118 \quad = \begin{pmatrix} 0 & \mathcal{P}(A) & 0 & L_{\mathcal{P}}(A, E) \\ \mathcal{P}(A)^* & 0 & L_{\mathcal{P}}(A, E)^* & 0 \\ 0 & 0 & 0 & \mathcal{P}(A) \\ 0 & 0 & \mathcal{P}(A)^* & 0 \end{pmatrix}.$$

119

120 A drawback of this approach is that it is valid only for square matrices A . The
 121 strategy we adopt in the present paper will be quite different, and will be valid not
 122 just for square matrices A but also for rectangular matrices A having more rows than
 123 columns.

124 *Organization.* This paper is organized as follows. We begin in Section 2 by giving
 125 statements of our main results, deferring their proof to Section 3. In Section 4, we
 126 discuss several practical aspects of the iterative schemes, including stability, scaling,
 127 and termination criteria. We compare the iterative schemes to other methods for
 128 computing the Fréchet derivative of the polar decomposition in Section 5. We finish
 129 with some numerical experiments in Section 6.

130 **2. Statement of Results.** In this section, we give a presentation of this paper's
 131 main result, which is a theorem that details a class of iterative schemes for computing
 132 the Fréchet derivative $L_{\mathcal{P}}$ of the map \mathcal{P} which sends a matrix A to the unitary factor
 133 U in its polar decomposition $A = UH$. A proof of the theorem is given in Section 3.

134 The class of iterative schemes to be considered comprises schemes of the form (3-
 135 4), with a mild constraint on the form of the function g . To understand this constraint,
 136 it is helpful to develop some intuition concerning iterative schemes for computing the
 137 polar decomposition and their relationship to iterative schemes for computing the
 138 matrix sign function. Fundamental to that intuition are the identities

$$139 \quad (9) \quad \text{sign}(A) = A(A^2)^{-1/2}, \quad \mathcal{P}(A) = A(A^*A)^{-1/2},$$

140 and the integral representation formulas [17, Equations 6.2 and 6.3]

$$141 \quad \text{sign}(A) = \frac{2}{\pi} A \int_0^\infty (t^2 I + A^2)^{-1} dt, \quad \mathcal{P}(A) = \frac{2}{\pi} A \int_0^\infty (t^2 I + A^* A)^{-1} dt,$$

142 which hint at two rules of thumb. First, iterative schemes for computing the matrix
 143 sign function tend to have the form $X_{k+1} = X_k h(X_k^2)$, where h is a primary matrix
 144 function. Second, to each iterative scheme $X_{k+1} = X_k h(X_k^2)$ for computing the
 145 matrix sign function, there corresponds an iterative scheme $X_{k+1} = X_k h(X_k^* X_k)$
 146 for computing the polar decomposition. The first of these rules of thumb appears
 147 to hold empirically to our knowledge. The second is made precise in [18, Theorem
 148 8.13]. The theorem below extends [18, Theorem 8.13] by showing, in essence, that
 149 to each iterative scheme $X_{k+1} = X_k h(X_k^2)$ for computing the matrix sign function,

150 there corresponds an iterative scheme for computing the polar decomposition *and its*
 151 *Fréchet derivative*. This iterative scheme is given by (3-4) with $g(X) = Xh(X^*X)$.

152 In what follows, we denote by $\text{skew}(B) = \frac{1}{2}(B - B^*)$ and $\text{sym}(B) = \frac{1}{2}(B + B^*)$ the
 153 skew-Hermitian and Hermitian parts, respectively, of a square matrix B . We denote
 154 the spectrum of B by $\Lambda(B)$.

155 **THEOREM 1.** *Let $A \in \mathbb{C}^{m \times n}$ ($m \geq n$) be a full-rank matrix having polar de-*
 156 *composition $A = UH$, where $U \in \mathbb{C}^{m \times n}$ has orthonormal columns and $H \in \mathbb{C}^{n \times n}$*
 157 *is Hermitian positive definite. Let $E \in \mathbb{C}^{m \times n}$, and define $\Omega = \text{skew}(U^*E)$ and*
 158 *$S = \text{sym}(U^*E)$. Let h be a primary matrix function satisfying $h(Z^*) = h(Z)^*$ for*
 159 *every Z , and suppose that the iteration $Z_{k+1} = Z_k h(Z_k^2)$ produces iterates Z_k that*
 160 *converge to $\text{sign}(Z_0)$ as $k \rightarrow \infty$ when the initial condition is*

$$161 \quad (10) \quad Z_0 = \begin{pmatrix} H & \Omega \\ 0 & -H \end{pmatrix},$$

162 *as well as when the initial condition is*

$$163 \quad (11) \quad Z_0 = \begin{pmatrix} H & S \\ 0 & H \end{pmatrix}.$$

164 *Assume that in both cases, h is smooth on an open subset of \mathbb{C} containing $\cup_{k=0}^{\infty} \Lambda(Z_k)$.*
 165 *Let $g(X) = Xh(X^*X)$. Then the iteration*

$$166 \quad (12) \quad X_{k+1} = g(X_k), \quad X_0 = A,$$

$$167 \quad (13) \quad E_{k+1} = L_g(X_k, E_k), \quad E_0 = E,$$

169 *produces iterates X_k and E_k that converge to $\mathcal{P}(A) = U$ and $L_{\mathcal{P}}(A, E)$, respectively,*
 170 *as $k \rightarrow \infty$.*

171 *Remark.* Taking $E = 0$ in the preceding theorem, one recovers [18, Theorem 8.13],
 172 up to the following modification: Instead of requesting that h is a primary matrix
 173 function satisfying $h(Z^*) = h(Z)^*$, [18, Theorem 8.13] makes the weaker assumption
 174 that the function $\tilde{g}(Z) = Zh(Z^2)$ satisfies $\tilde{g}(Z^*) = \tilde{g}(Z)^*$ for every Z . It is easily
 175 checked using elementary properties of primary matrix functions [18, Theorem 1.13]
 176 that the latter is implied by the former.

177 Note that it is sometimes the case that the convergence of the matrix sign function
 178 iteration $Z_{k+1} = Z_k h(Z_k^2)$ referenced in Theorem 1 is dictated by the spectrum of Z_0 .
 179 If this is the case, then the hypothesis that the iteration converges when Z_0 is given
 180 by (10) or (11) is equivalent to the simpler hypothesis that the iteration converges
 181 when $Z_0 = \pm H$. This follows from the fact that the spectra of (10) and (11) are given
 182 by $\Lambda(H) \cup \Lambda(-H)$ and $\Lambda(H)$, respectively.

183 Central to the proof of Theorem 1 is an identity that relates the Fréchet derivative
 184 of the polar decomposition to the sign of the block matrix Z_0 appearing in (10). We
 185 state the identity below to emphasize its importance. A proof is given in Section 3.1.

186

187 **THEOREM 2.** *Let A, U, H, E , and Ω be as in Theorem 1. Then*

$$188 \quad (14) \quad \text{sign} \begin{pmatrix} H & \Omega \\ 0 & -H \end{pmatrix} = \begin{pmatrix} I & U^* L_{\mathcal{P}}(A, E) \\ 0 & -I \end{pmatrix}.$$

189

190 *In particular, if U^*E is skew-Hermitian, then*

$$191 \quad (15) \quad \text{sign} \left(\begin{pmatrix} U^* & 0 \\ 0 & -U^* \end{pmatrix} \begin{pmatrix} A & E \\ 0 & A \end{pmatrix} \right) = \begin{pmatrix} U^* & 0 \\ 0 & -U^* \end{pmatrix} \begin{pmatrix} \mathcal{P}(A) & L_{\mathcal{P}}(A, E) \\ 0 & \mathcal{P}(A) \end{pmatrix}.$$

192 In addition to being useful in the proof of Theorem 1, the identity (15) bears an
193 interesting resemblance to (5).

194 Theorem 1 has several corollaries, each corresponding to a different choice of iter-
195 ative scheme $Z_{k+1} = Z_k h(Z_k^2)$ for computing the matrix sign function. The simplest
196 is the well-known Newton iteration

$$197 \quad (16) \quad Z_{k+1} = \frac{1}{2}(Z_k + Z_k^{-1}),$$

198 which corresponds to the choice $h(Z) = \frac{1}{2}(I + Z^{-1})$. It is known that this iteration
199 converges quadratically to $\text{sign}(Z_0)$ for any Z_0 having no pure imaginary eigenval-
200 ues [18, Theorem 5.6]. Since (10) and (11) have eigenvalues equal to plus or minus
201 the eigenvalues of H , all of which are nonzero real numbers, we obtain the following
202 corollary. In it, we restrict the discussion to square matrices, since this leads to a
203 particularly simple iterative scheme.

204 **COROLLARY 3.** *Let $A \in \mathbb{C}^{n \times n}$ be a nonsingular matrix having polar decomposition*
205 *$A = UH$, where $U \in \mathbb{C}^{n \times n}$ is unitary and $H \in \mathbb{C}^{n \times n}$ is Hermitian positive definite.*
206 *Let $E \in \mathbb{C}^{n \times n}$. Then the iteration*

$$207 \quad (17) \quad X_{k+1} = \frac{1}{2}(X_k + X_k^{-*}), \quad X_0 = A,$$

$$208 \quad (18) \quad E_{k+1} = \frac{1}{2}(E_k - X_k^{-*} E_k^* X_k^{-*}), \quad E_0 = E,$$

209
210 *produces iterates X_k and E_k that converge to $\mathcal{P}(A) = U$ and $L_{\mathcal{P}}(A, E)$, respectively,*
211 *as $k \rightarrow \infty$.*

212 *Remark.* If A is rectangular, then the iteration obtained from (16) reads

$$213 \quad (19) \quad X_{k+1} = \frac{1}{2}X_k(I + (X_k^* X_k)^{-1}), \quad X_0 = A,$$

$$214 \quad (20) \quad E_{k+1} = \frac{1}{2} \left[E_k(I + (X_k^* X_k)^{-1}) \right. \\ 215 \quad \left. - X_k(X_k^* X_k)^{-1}(E_k^* X_k + X_k^* E_k)(X_k^* X_k)^{-1} \right], \quad E_0 = E.$$

216
217 This scheme simplifies to (17-18) when A is square.

218 A second corollary of Theorem 1 is obtained by considering the Newton-Schulz
219 iteration

$$220 \quad (21) \quad Z_{k+1} = \frac{1}{2}Z_k(3I - Z_k^2),$$

221 which corresponds to the choice $h(Z) = \frac{1}{2}(3I - Z)$. It is known that this iteration
222 converges to $\text{sign}(Z_0)$ provided that (i) Z_0 has no pure imaginary eigenvalues and
223 (ii) the eigenvalues of $I - Z_0^2$ all have magnitude strictly less than one [22, Theorem
224 5.2]. Note that [22, Theorem 5.2] replaces the latter condition with $\|I - Z_0^2\| < 1$,
225 but it is evident from their proof that this condition can be relaxed to what we have
226 written here. Since the eigenvalues of

$$227 \quad \begin{pmatrix} I & 0 \\ 0 & I \end{pmatrix} - \begin{pmatrix} H & \Omega \\ 0 & -H \end{pmatrix}^2 = \begin{pmatrix} I - H^2 & \Omega H - H\Omega \\ 0 & I - H^2 \end{pmatrix}$$

228 and

$$229 \quad \begin{pmatrix} I & 0 \\ 0 & I \end{pmatrix} - \begin{pmatrix} H & S \\ 0 & H \end{pmatrix}^2 = \begin{pmatrix} I - H^2 & -HS - SH \\ 0 & I - H^2 \end{pmatrix}$$

230 coincide with those of $I - H^2 = I - A^*A$, we obtain the following corollary.

231 COROLLARY 4. Let $A \in \mathbb{C}^{m \times n}$ ($m \geq n$) be a full-rank matrix having polar de-
 232 composition $A = UH$, where $U \in \mathbb{C}^{m \times n}$ has orthonormal columns and $H \in \mathbb{C}^{n \times n}$ is
 233 Hermitian positive definite. Let $E \in \mathbb{C}^{m \times n}$. If all of the singular values of A lie in
 234 the interval $(0, \sqrt{2})$, then the iteration

$$235 \quad (22) \quad X_{k+1} = \frac{1}{2}X_k(3I - X_k^*X_k), \quad X_0 = A,$$

$$236 \quad (23) \quad E_{k+1} = \frac{1}{2}E_k(3I - X_k^*X_k) - \frac{1}{2}X_k(E_k^*X_k + X_k^*E_k), \quad E_0 = E,$$

238 produces iterates X_k and E_k that converge to $\mathcal{P}(A) = U$ and $L_{\mathcal{P}}(A, E)$, respectively,
 239 as $k \rightarrow \infty$.

240 *Remark.* A more direct analysis of (22), without appealing to its relationship
 241 to a matrix sign function iteration, shows that $X_k \rightarrow U$ under the less stringent
 242 requirement that all of the singular values of A lie in the interval $(0, \sqrt{3})$ [18, Problem
 243 8.20]. Our numerical experiments suggest that the coupled iteration (22-23) enjoys
 244 convergence under the same condition, but Theorem 1 alone appears inadequate to
 245 conclude such a claim.

246 Other corollaries to Theorem 1 can be derived in a similar fashion. For instance,
 247 iterative schemes based on Padé approximations of $\text{sign}(Z) = Z(I - (I - Z^2))^{-1/2}$ (of
 248 which (21) is a special case) can be used; see [18, Chapter 5.4] for further details.

249 **3. Proofs.** In this section, we present proofs of Theorems 1 and 2. Our presenta-
 250 tion is divided into two parts. First, in Section 3.1, we derive a few identities involving
 251 the Fréchet derivative of the polar decomposition, proving Theorem 2 in the process.
 252 Then, in Section 3.2, we use the aforementioned identities to prove convergence of the
 253 iteration (12-13), thereby proving Theorem 1.

254 **3.1. Identities Involving the Fréchet Derivative of the Polar Decom-**
 255 **position.** This section studies the Fréchet derivative of the polar decomposition and
 256 its relationship to the matrix sign function, culminating in a proof of Theorem 2. A
 257 couple of main observations will be made. First, as will be seen in Lemma 7, the task
 258 of evaluating $L_{\mathcal{P}}(A, E)$ can essentially be reduced to the case in which A is Hermitian
 259 positive definite and E is skew-Hermitian. This is relatively simple to show when A
 260 is square, but the rectangular case turns out to be more subtle, requiring that some
 261 attention be paid to the relationship between the column space of A and that of E .
 262 This observation will be followed with a proof of Theorem 2, which reveals that the
 263 value of $U^*L_{\mathcal{P}}(A, E)$ can be read off of the (1, 2)-block of the matrix sign function
 264 applied to a certain block matrix.

265 Before studying the derivatives of \mathcal{P} in detail, it is worth pointing out that \mathcal{P}
 266 is a smooth map from the set of full-rank $m \times n$ ($m \geq n$) matrices to the set of
 267 $m \times n$ matrices with orthonormal columns. This follows from two facts: (1) the
 268 latter set of matrices constitutes a smooth, compact manifold, the Stiefel manifold
 269 $V_n(\mathbb{C}^m) = \{U \in \mathbb{C}^{m \times n} \mid U^*U = I\}$, and (2) the map \mathcal{P} coincides with the closest
 270 point projection onto $V_n(\mathbb{C}^n)$. That is, in the Frobenius norm $\|\cdot\|_F$,

$$271 \quad \mathcal{P}(A) = \arg \min_{U \in V_n(\mathbb{C}^m)} \|A - U\|_F$$

272 for any full-rank $A \in \mathbb{C}^{m \times n}$ [18, Theorem 8.4]. It is a classical result from differential
 273 geometry that the closest point projection onto a smooth, compact manifold embedded
 274 in Euclidean space is a smooth map [11]. In particular, \mathcal{P} is Fréchet differentiable at

275 any full-rank $A \in \mathbb{C}^{m \times n}$. (For a different justification of this fact, see [10, Section
276 2.3(c)].)

277 We now turn our attention to the differentiation of \mathcal{P} . We begin by recording a
278 useful formula for the Fréchet derivative of a function of the form $g(X) = Xh(X^*X)$.
279 Along the way, we make some observations concerning the column space $\mathcal{R}(A)$ of
280 a matrix $A \in \mathbb{C}^{m \times n}$ and the column space $\mathcal{R}(L_g(A, E))$ of the Fréchet derivative
281 $L_g(A, E)$ of g at A in a direction $E \in \mathbb{C}^{m \times n}$. We denote by $\mathcal{N}(A^*)$ the null space of
282 A^* ; equivalently, $\mathcal{N}(A^*)$ is the orthogonal complement to $\mathcal{R}(A)$ in \mathbb{C}^m .

283 LEMMA 5. *Let $A \in \mathbb{C}^{m \times n}$ ($m \geq n$), let $h : \mathbb{C}^{n \times n} \rightarrow \mathbb{C}^{n \times n}$ be Fréchet differentiable
284 at A^*A , and define $g(X) = Xh(X^*X)$. Then for any $E \in \mathbb{C}^{m \times n}$,*

$$285 \quad (24) \quad L_g(A, E) = Eh(A^*A) + AL_h(A^*A, A^*E + E^*A).$$

286 *In particular, if $\mathcal{R}(E) \subseteq \mathcal{R}(A)$, then $\mathcal{R}(L_g(A, E)) \subseteq \mathcal{R}(A)$. On the other hand, if
287 $\mathcal{R}(E) \subseteq \mathcal{N}(A^*)$, then*

$$288 \quad (25) \quad L_g(A, E) = Eh(A^*A),$$

289 *and hence $\mathcal{R}(L_g(A, E)) \subseteq \mathcal{N}(A^*)$.*

290 *Proof.* The formula (24) is a consequence of the product rule and the chain
291 rule [18, Theorems 3.3 & 3.4]. The implication $\mathcal{R}(E) \subseteq \mathcal{R}(A) \implies \mathcal{R}(L_g(A, E)) \subseteq$
292 $\mathcal{R}(A)$ is immediate since the columns of $L_g(A, E)$ are linear combinations of the
293 columns of A and E . Equation (25) follows from the fact that $A^*E + E^*A = 0$
294 whenever $\mathcal{R}(E) \subseteq \mathcal{N}(A^*)$. \square

295 The preceding lemma has several important consequences. The first of these is
296 an application of Lemma 5 to the function $g(X) = \mathcal{P}(X)$, which has the requisite
297 functional form in view of (9).

298 LEMMA 6. *Let $A \in \mathbb{C}^{m \times n}$ ($m \geq n$) be a full-rank matrix having polar decomposi-
299 tion $A = UH$, where $U \in \mathbb{C}^{m \times n}$ has orthonormal columns and $H \in \mathbb{C}^{n \times n}$ is
300 Hermitian positive definite. Let $E \in \mathbb{C}^{m \times n}$, and write*

$$301 \quad (26) \quad E = E^{\parallel} + E^{\perp}, \quad E^{\parallel} = UU^*E, \quad E^{\perp} = (I - UU^*)E.$$

302 *Then*

$$303 \quad (27) \quad UU^*L_{\mathcal{P}}(A, E^{\parallel}) = L_{\mathcal{P}}(A, E^{\parallel})$$

304 *and*

$$305 \quad (28) \quad L_{\mathcal{P}}(A, E^{\perp}) = E^{\perp}H^{-1}.$$

306 *Proof.* Apply Lemma 5 with the choice $h(X) = X^{-1/2}$, so that
307 $g(X) = X(X^*X)^{-1/2} = \mathcal{P}(X)$. Equation (27) is a restatement of the fact that
308 $\mathcal{R}(L_{\mathcal{P}}(A, E^{\parallel})) \subseteq \mathcal{R}(A) = \mathcal{R}(U)$, while (28) follows from (25) together with the iden-
309 tity $H = (A^*A)^{1/2}$. \square

310 We will now show, with the help of Lemma 6, that the task of evaluating $L_{\mathcal{P}}(A, E)$
311 can essentially be reduced to the case in which A is Hermitian positive definite and
312 E is skew-Hermitian.

313 LEMMA 7. *Let $A \in \mathbb{C}^{m \times n}$ ($m \geq n$) be a full-rank matrix having polar decom-
314 position $A = UH$, where $U \in \mathbb{C}^{m \times n}$ has orthonormal columns and $H \in \mathbb{C}^{n \times n}$ is
315 Hermitian positive definite. Then for any $E \in \mathbb{C}^{m \times n}$,*

$$316 \quad (29) \quad \text{skew}(U^*L_{\mathcal{P}}(A, E)) = L_{\mathcal{P}}(H, \Omega),$$

$$317 \quad (30) \quad \text{sym}(U^*L_{\mathcal{P}}(A, E)) = L_{\mathcal{P}}(H, S) = 0,$$

319 where $\Omega = \text{skew}(U^*E)$ and $S = \text{sym}(U^*E)$. Hence,

$$320 \quad U^*L_{\mathcal{P}}(A, E) = L_{\mathcal{P}}(H, \Omega).$$

321 *Proof.* Let E^{\parallel} and E^{\perp} be as in (26). The linearity of the Fréchet derivative
322 implies that

$$323 \quad L_{\mathcal{P}}(A, E) = L_{\mathcal{P}}(A, E^{\parallel}) + L_{\mathcal{P}}(A, E^{\perp}).$$

325 The formula (28) and the identities $A = UH$ and $UU^*E^{\parallel} = E^{\parallel}$ then give

$$326 \quad L_{\mathcal{P}}(A, E) = L_{\mathcal{P}}(UH, UU^*E^{\parallel}) + E^{\perp}H^{-1}.$$

327 Now note that the map \mathcal{P} clearly satisfies $\mathcal{P}(VB) = V\mathcal{P}(B)$ for any nonsingular
328 $B \in \mathbb{C}^{n \times n}$ and any $V \in \mathbb{C}^{m \times n}$ ($m \geq n$) with orthonormal columns. From this it
329 follows that for any such V and B , and any $F \in \mathbb{C}^{n \times n}$,

$$330 \quad (31) \quad L_{\mathcal{P}}(VB, VF) = VL_{\mathcal{P}}(B, F).$$

331 Applying this identity to the case in which $B = H$, $V = U$, and $F = U^*E^{\parallel}$, we obtain

$$332 \quad L_{\mathcal{P}}(UH, UU^*E^{\parallel}) = UL_{\mathcal{P}}(H, U^*E^{\parallel}) \\ 333 \quad \quad \quad = UL_{\mathcal{P}}(H, U^*E),$$

335 where the second line follows from the fact that $U^*E^{\perp} = 0$. Thus,

$$336 \quad L_{\mathcal{P}}(A, E) = UL_{\mathcal{P}}(H, U^*E) + E^{\perp}H^{-1}.$$

337 Multiplying from the left by U^* gives

$$338 \quad U^*L_{\mathcal{P}}(A, E) = L_{\mathcal{P}}(H, U^*E).$$

339 since $U^*U = I$ and $U^*E^{\perp} = 0$. Equivalently, in terms of $\Omega = \text{skew}(U^*E)$ and
340 $S = \text{sym}(U^*E)$,

$$341 \quad U^*L_{\mathcal{P}}(A, E) = L_{\mathcal{P}}(H, \Omega) + L_{\mathcal{P}}(H, S).$$

342 The proof will be complete if we can show that $L_{\mathcal{P}}(H, \Omega)$ is skew-Hermitian and

$$343 \quad (32) \quad L_{\mathcal{P}}(H, S) = 0.$$

344 In fact, (32) holds for any Hermitian matrix S since, for all sufficiently small ε , $H + \varepsilon S$
345 is Hermitian positive definite, showing that $\mathcal{P}(H + \varepsilon S) = I$. The skew-Hermiticity of
346 $L_{\mathcal{P}}(H, \Omega)$ follows from differentiating the identity

$$347 \quad \mathcal{P}(H + \varepsilon\Omega)^*\mathcal{P}(H + \varepsilon\Omega) = I$$

348 with respect to ε and using the fact that $\mathcal{P}(H) = I$. □

349 Another consequence of Lemma 5 is the following identity that relates the Fréchet
350 derivative of the polar decomposition of a Hermitian positive definite matrix to the
351 matrix sign function applied to a certain block matrix.

352 **LEMMA 8.** *Let $H \in \mathbb{R}^{n \times n}$ be Hermitian positive definite, and let $\Omega \in \mathbb{R}^{n \times n}$ be*
353 *skew-Hermitian. Then*

$$354 \quad (33) \quad \text{sign} \begin{pmatrix} H & \Omega \\ 0 & -H \end{pmatrix} = \begin{pmatrix} I & L_{\mathcal{P}}(H, \Omega) \\ 0 & -I \end{pmatrix},$$

$$355 \quad (34) \quad \text{sign} \begin{pmatrix} H & S \\ 0 & H \end{pmatrix} = \begin{pmatrix} I & L_{\mathcal{P}}(H, S) \\ 0 & I \end{pmatrix} = \begin{pmatrix} I & 0 \\ 0 & I \end{pmatrix}.$$

356

357 *Proof.* By definition,

$$\begin{aligned}
358 \quad \text{sign} \begin{pmatrix} H & \Omega \\ 0 & -H \end{pmatrix} &= \begin{pmatrix} H & \Omega \\ 0 & -H \end{pmatrix} \begin{pmatrix} H^2 & H\Omega - \Omega H \\ 0 & H^2 \end{pmatrix}^{-1/2} \\
359 \quad &= \begin{pmatrix} H & \Omega \\ 0 & -H \end{pmatrix} \begin{pmatrix} H^2 & H\Omega + \Omega^* H \\ 0 & H^2 \end{pmatrix}^{-1/2}. \\
360
\end{aligned}$$

361 Now apply (5) to the primary matrix function $f(X) = X^{-1/2}$ to obtain

$$\begin{aligned}
362 \quad \begin{pmatrix} H^2 & H\Omega + \Omega^* H \\ 0 & H^2 \end{pmatrix}^{-1/2} &= \begin{pmatrix} H^{-1} & L_{x^{-1/2}}(H^2, H\Omega + \Omega^* H) \\ 0 & H^{-1} \end{pmatrix}, \\
363
\end{aligned}$$

364 where the identity $(H^2)^{-1/2} = H^{-1}$ follows from the positive-definiteness of H . Thus,

$$\begin{aligned}
365 \quad \text{sign} \begin{pmatrix} H & \Omega \\ 0 & -H \end{pmatrix} &= \begin{pmatrix} H & \Omega \\ 0 & -H \end{pmatrix} \begin{pmatrix} H^{-1} & L_{x^{-1/2}}(H^2, H\Omega + \Omega^* H) \\ 0 & H^{-1} \end{pmatrix} \\
366 \quad &= \begin{pmatrix} I & HL_{x^{-1/2}}(H^2, H\Omega + \Omega^* H) + \Omega H^{-1} \\ 0 & -I \end{pmatrix}. \\
367
\end{aligned}$$

368 The identity (33) follows upon observing that, by (24),

$$\begin{aligned}
369 \quad L_{\mathcal{P}}(H, \Omega) &= \Omega H^{-1} + HL_{x^{-1/2}}(H^2, H\Omega + \Omega^* H). \\
370
\end{aligned}$$

371 The proof of (34) is simpler, since, by (5) and the identity $L_{\text{sign}}(H, S) = 0$,

$$\begin{aligned}
372 \quad \text{sign} \begin{pmatrix} H & S \\ 0 & H \end{pmatrix} &= \begin{pmatrix} I & L_{\text{sign}}(H, S) \\ 0 & I \end{pmatrix} = \begin{pmatrix} I & 0 \\ 0 & I \end{pmatrix}. \quad \square
\end{aligned}$$

373 We remark that an alternative proof of (33) exists. It is based on the observation
374 that $L_{\mathcal{P}}(H, \Omega)$ is the solution of a Lyapunov equation which can be solved by reading
375 off the (1, 2)-block of $\text{sign} \begin{pmatrix} H & \Omega \\ 0 & -H \end{pmatrix}$. For details, see Section 5.

376 Combining Lemma 8 with Lemma 7 proves Theorem 2.

377 **3.2. Convergence of the Iteration.** We now focus our efforts on proving convergence
378 of the iteration (12-13), thereby proving Theorem 1. The cornerstone of the
379 proof is Lemma 11, where a relationship is established between certain blocks of the
380 matrices Z_k defined by the matrix sign function $Z_{k+1} = Z_k h(Z_k)^2$ and the matrices
381 X_k and E_k defined by the iteration (12-13). Once this has been shown, convergence
382 of the iteration (12-13) will follow from the convergence of Z_k to $\text{sign}(Z_0)$, together
383 with the knowledge (from Theorem 2) that the Fréchet derivative of the polar decom-
384 position is related to the (1, 2)-block of $\text{sign}(Z_0)$ for certain values of Z_0 .

385 We begin by examining the block structure of the iterates Z_k .

386 LEMMA 9. *The iterates Z_k produced by the iteration $Z_{k+1} = Z_k h(Z_k^2)$ with initial*
387 *condition (10) have the form*

$$\begin{aligned}
388 \quad Z_k &= \begin{pmatrix} H_k & \Omega_k \\ 0 & -H_k \end{pmatrix},
\end{aligned}$$

389 where H_k is Hermitian and Ω_k is skew-Hermitian.

390 *Proof.* Assume the statement is true at iteration k . Then by (5),

$$\begin{aligned}
391 \quad Z_{k+1} &= \begin{pmatrix} H_k & \Omega_k \\ 0 & -H_k \end{pmatrix} h \begin{pmatrix} H_k^2 & H_k \Omega_k - \Omega_k H_k \\ 0 & H_k^2 \end{pmatrix} \\
392 &= \begin{pmatrix} H_k & \Omega_k \\ 0 & -H_k \end{pmatrix} \begin{pmatrix} h(H_k^2) & L_h(H_k^2, H_k \Omega_k - \Omega_k H_k) \\ 0 & h(H_k^2) \end{pmatrix} \\
393 \quad (35) &= \begin{pmatrix} H_k h(H_k^2) & H_k L_h(H_k^2, H_k \Omega_k - \Omega_k H_k) + \Omega_k h(H_k^2) \\ 0 & -H_k h(H_k^2) \end{pmatrix}. \\
394
\end{aligned}$$

395 By the remark following Theorem 1, $H_k h(H_k^2) = [H_k h(H_k^2)]^*$, showing that $H_{k+1} =$
396 $H_k h(H_k^2)$ is Hermitian. On the other hand, the fact that h is a primary matrix
397 function implies that Z_k commutes with $h(Z_k^2)$, so, by a calculation similar to that
398 above, we also have

$$\begin{aligned}
399 \quad (36) \quad Z_{k+1} &= \begin{pmatrix} h(H_k^2)H_k & h(H_k^2)\Omega_k - L_h(H_k^2, H_k \Omega_k - \Omega_k H_k)H_k \\ 0 & -h(H_k^2)H_k \end{pmatrix}. \\
400
\end{aligned}$$

401 Denote $C_k = H_k \Omega_k - \Omega_k H_k$. Since H_k is Hermitian and Ω_k is skew-Hermitian, C_k is
402 Hermitian. Hence, since $h(Z^*) = h(Z)^*$ for every Z ,

$$\begin{aligned}
403 \quad L_h(H_k^2, C_k)^* &= L_h((H_k^2)^*, C_k^*) \\
404 &= L_h(H_k^2, C_k).
\end{aligned}$$

406 Comparing the (1, 2) blocks of (35) and (36) then shows that

$$\begin{aligned}
407 \quad 0 &= H_k L_h(H_k^2, C_k) + \Omega_k h(H_k^2) - h(H_k^2)\Omega_k + L_h(H_k^2, C_k)H_k \\
408 &= H_k L_h(H_k^2, C_k) + \Omega_k h(H_k^2) + h(H_k^2)^* \Omega_k^* + L_h(H_k^2, C_k)^* H_k^* \\
409 &= \Omega_{k+1} + \Omega_{k+1}^*.
\end{aligned}$$

411 It follows that $\Omega_k = -\Omega_k^*$ for every k . □

412 The proof above also reveals a recursion satisfied by H_k and Ω_k , namely,

$$\begin{aligned}
413 \quad (37) \quad H_{k+1} &= H_k h(H_k^2) \\
414 \quad (38) \quad \Omega_{k+1} &= \Omega_k h(H_k^2) + H_k L_h(H_k^2, H_k \Omega_k - \Omega_k H_k).
\end{aligned}$$

416 Next, we examine the block structure of the iterates Z_k with initial condition (11).

417 LEMMA 10. *The iterates Z_k produced by the iteration $Z_{k+1} = Z_k h(Z_k^2)$ with initial*
418 *condition (11) have the form*

$$\begin{aligned}
419 \quad Z_k &= \begin{pmatrix} H_k & S_k \\ 0 & H_k \end{pmatrix},
\end{aligned}$$

420 where H_k is the same Hermitian matrix as in Lemma 9 and S_k is Hermitian.

421 *Proof.* We omit the proof, which is very similar to the proof of Lemma 9. □

422 In analogy with (38), the iterates S_k satisfy the recursion

$$\begin{aligned}
423 \quad (39) \quad S_{k+1} &= S_k h(H_k^2) + H_k L_h(H_k^2, S_k H_k + H_k S_k).
\end{aligned}$$

424 We now relate the matrices H_k , Ω_k , and S_k defined in the preceding pair of
425 lemmas to the matrices X_k and E_k defined by the coupled iteration (12-13).

426 LEMMA 11. *The iterates H_k , Ω_k , and S_k are related to X_k and E_k via*

$$427 \quad (40) \quad UH_k = X_k,$$

$$428 \quad (41) \quad \Omega_k = \text{skew}(U^*E_k),$$

$$429 \quad (42) \quad S_k = \text{sym}(U^*E_k).$$

431 *Proof.* The first of these equalities follows easily by induction, for if it holds at
432 iteration k , then

$$\begin{aligned} 433 \quad X_{k+1} &= g(X_k) \\ 434 \quad &= X_k h(X_k^* X_k) \\ 435 \quad &= UH_k h(H_k^* U^* UH_k) \\ 436 \quad &= UH_k h(H_k^2) \\ 437 \quad &= UH_{k+1}. \end{aligned}$$

439 Furthermore, $X_0 = A = UH = UH_0$, which proves (40). To prove (41) and (42),
440 we will show that if $\Omega_k = \text{skew}(U^*E_k)$ and $S_k = \text{sym}(U^*E_k)$ for a given k , and if
441 E_{k+1} , Ω_{k+1} , and S_{k+1} are given by (13), (38), and (39), respectively, then $\Omega_{k+1} =$
442 $\text{skew}(U^*E_{k+1})$ and $S_{k+1} = \text{sym}(U^*E_{k+1})$. Recalling (24), we have

$$\begin{aligned} 443 \quad U^*E_{k+1} &= U^*L_g(X_k, E_k) \\ 444 \quad &= U^*E_k h(X_k^* X_k) + U^*X_k L_h(X_k^* X_k, E_k^* X_k + X_k^* E_k) \\ 445 \quad &= U^*E_k h(H_k^2) + H_k L_h(H_k^2, E_k^* UH_k + H_k U^* E_k) \\ 446 \quad &= \Omega_k h(H_k^2) + H_k L_h(H_k^2, \Omega_k^* H_k + H_k \Omega_k) + S_k h(H_k^2) + H_k L_h(H_k^2, S_k H_k + H_k S_k) \\ 447 \quad &= \Omega_{k+1} + S_{k+1}, \end{aligned}$$

449 where we have used (38), (39), and the decomposition $U^*E_k = \Omega_k + S_k$. By Lemmas 9
450 and 10, Ω_{k+1} is skew-Hermitian and S_{k+1} is Hermitian, proving (41) and (42). \square

451 The proof of Theorem 1 is now almost complete, since by Lemma 9 and Theorem 2,

$$452 \quad \begin{pmatrix} H_k & \Omega_k \\ 0 & -H_k \end{pmatrix} \rightarrow \text{sign} \begin{pmatrix} H & \Omega \\ 0 & -H \end{pmatrix} = \begin{pmatrix} I & U^*L_{\mathcal{P}}(A, E) \\ 0 & -I \end{pmatrix}$$

453 as $k \rightarrow \infty$. Likewise, by (34) and Lemma 10,

$$454 \quad \begin{pmatrix} H_k & S_k \\ 0 & H_k \end{pmatrix} \rightarrow \text{sign} \begin{pmatrix} H & S \\ 0 & H \end{pmatrix} = \begin{pmatrix} I & 0 \\ 0 & I \end{pmatrix},$$

455 as $k \rightarrow \infty$. These observations, together with (40-42), show that

$$\begin{aligned} 456 \quad X_k &\rightarrow U, \\ 457 \quad \text{skew}(U^*E_k) &\rightarrow U^*L_{\mathcal{P}}(A, E), \\ 458 \quad \text{sym}(U^*E_k) &\rightarrow 0 \end{aligned}$$

460 as $k \rightarrow \infty$. In other words,

$$461 \quad (43) \quad X_k \rightarrow \mathcal{P}(A),$$

$$462 \quad (44) \quad U^*E_k \rightarrow U^*L_{\mathcal{P}}(A, E)$$

464 as $k \rightarrow \infty$. The latter limit implies that $E_k \rightarrow L_{\mathcal{P}}(A, E)$ when U is square, but
 465 not when U is rectangular. To handle the rectangular case, consider the decomposi-
 466 tions (26) and

$$467 \quad E_k = E_k^{\parallel} + E_k^{\perp}, \quad E_k^{\parallel} = UU^*E_k, \quad E_k^{\perp} = (I - UU^*)E_k.$$

469 By Lemma 6 and the linearity of the Fréchet derivative, the statement (44) is equiv-
 470 alent to the statement that

$$471 \quad U^*E_k^{\parallel} \rightarrow U^*L_{\mathcal{P}}(A, E^{\parallel}) + U^*L_{\mathcal{P}}(A, E^{\perp}) \\ 472 \quad = U^*L_{\mathcal{P}}(A, E^{\parallel}).$$

474 Multiplying from the left by U and recalling that $UU^*E_k^{\parallel} = E_k^{\parallel}$ and $UU^*L_{\mathcal{P}}(A, E^{\parallel}) =$
 475 $L_{\mathcal{P}}(A, E^{\parallel})$ (by (27)), we conclude that

$$476 \quad E_k^{\parallel} \rightarrow L_{\mathcal{P}}(A, E^{\parallel}).$$

477 The proof of Theorem 1 will be complete if we can show that

$$478 \quad E_k^{\perp} \rightarrow L_{\mathcal{P}}(A, E^{\perp}).$$

479 This is carried out in the following lemma.

480 LEMMA 12. As $k \rightarrow \infty$, $E_k^{\perp} \rightarrow L_{\mathcal{P}}(A, E^{\perp})$.

481 *Proof.* By (28), it suffices to show that

$$482 \quad E_k^{\perp} \rightarrow E^{\perp}H^{-1}.$$

483 Using Lemma 5, it is straightforward to see that E_k^{\parallel} and E_k^{\perp} satisfy independent
 484 recursions of the form

$$485 \quad E_{k+1}^{\parallel} = L_g(X_k, E_k^{\parallel}), \\ 486 \quad E_{k+1}^{\perp} = L_g(X_k, E_k^{\perp}).$$

488 Now since $\mathcal{R}(E_k^{\perp})$ is orthogonal to $\mathcal{R}(U) \supseteq \mathcal{R}(UH_k) = \mathcal{R}(X_k)$, it follows from (25)
 489 that

$$490 \quad L_g(X_k, E_k^{\perp}) = E_k^{\perp}h(X_k^*X_k),$$

491 so

$$492 \quad E_{k+1}^{\perp} = E_k^{\perp}h(X_k^*X_k).$$

494 If we introduce the matrix $B_k \in \mathbb{C}^{n \times n}$ defined by the recursion

$$495 \quad B_{k+1} = B_k h(X_k^*X_k), \quad B_0 = I,$$

496 then an inductive argument shows that

$$497 \quad E_k^{\perp} = E^{\perp}B_k.$$

498 We claim that $B_k \rightarrow H^{-1}$ as $k \rightarrow \infty$. To see this, observe that (12) implies that

$$499 \quad X_k = X_0B_k = AB_k.$$

500 Since $X_k \rightarrow U$ as $k \rightarrow \infty$, we conclude that

$$501 \quad I = U^*U = U^* \lim_{k \rightarrow \infty} X_k = U^*A \lim_{k \rightarrow \infty} B_k = H \lim_{k \rightarrow \infty} B_k.$$

502 It follows that $E_k^{\perp} = E^{\perp}B_k \rightarrow E^{\perp}H^{-1}$ as $k \rightarrow \infty$. □

503 **4. Practical Considerations.** This section discusses several practical consid-
504 erations concerning the iterative schemes detailed in Theorem 1.

505 **4.1. Scaling.** Scaling the iterates X_k in the Newton iteration (17) often reduces
506 the number of iterations required to achieve convergence [18, Chapter 8.6]. If this
507 strategy is generalized to the coupled iteration (17-18), then the resulting iteration
508 reads

$$509 \quad (45) \quad X_{k+1} = \frac{1}{2}(\mu_k X_k + \mu_k^{-1} X_k^{-*}), \quad X_0 = A,$$

$$510 \quad (46) \quad E_{k+1} = \frac{1}{2}(\mu_k E_k - \mu_k^{-1} X_k^{-*} E_k^* X_k^{-*}), \quad E_0 = E,$$

512 where $\mu_k > 0$ is a scaling factor chosen heuristically. Practical choices for μ_k in-
513 clude [18]

$$514 \quad (47) \quad \mu_k = \left(\frac{\|X_k^{-1}\|_1 \|X_k^{-1}\|_\infty}{\|X_k\|_1 \|X_k\|_\infty} \right)^{1/4}$$

515 and

$$516 \quad (48) \quad \mu_k = \left(\frac{\|X_k^{-1}\|_F}{\|X_k\|_F} \right)^{1/2},$$

517 where $\|\cdot\|_1$, $\|\cdot\|_\infty$, and $\|\cdot\|_F$ denote the matrix 1-, ∞ - and Frobenius norms,
518 respectively.

519 More generally, scaling can be applied to other iterative schemes of the form (12-
520 13), leading to iterative schemes of the form

$$521 \quad X_{k+1} = g(\mu_k X_k), \quad X_0 = A,$$

$$522 \quad E_{k+1} = L_g(\mu_k X_k, \mu_k E_k), \quad E_0 = E.$$

524 Note that if A is rectangular, then (47) and (48) are inapplicable. We have found

$$525 \quad (49) \quad \mu_k = \left(\frac{\|(X_k^* X_k)^{-1}\|_1 \|(X_k^* X_k)^{-1}\|_\infty}{\|X_k^* X_k\|_1 \|X_k^* X_k\|_\infty} \right)^{1/8}$$

526 and

$$527 \quad (50) \quad \mu_k = \left(\frac{\|(X_k^* X_k)^{-1}\|_F}{\|X_k^* X_k\|_F} \right)^{1/4}$$

528 to be effective alternatives to (47) and (48) in our numerical experiments with rect-
529 angular A .

530 **4.2. Termination Criteria.** Determining when to terminate the iteration (12-
531 13) is a delicate task. Termination criteria for (12) by itself are, of course, well-studied,
532 but the accuracy of E_k should be taken into account when choosing termination
533 criteria for the coupled iteration (12-13).

534 One possibility is to appeal to the relationship between X_k and E_k and the sign
535 function iterates Z_k referenced in the statement of Theorem 1. Convergence of the
536 sign function iterates to $\text{sign}(Z_0) = \lim_{k \rightarrow \infty} Z_k$ can be readily verified with the aid of
537 the inequality

$$538 \quad \frac{\|Z_k^2 - I\|}{\|\text{sign}(Z_0)\|(\|Z_k\| + \|\text{sign}(Z_0)\|)} \leq \frac{\|Z_k - \text{sign}(Z_0)\|}{\|\text{sign}(Z_0)\|} \leq \|Z_k^2 - I\|,$$

539 which holds in any submultiplicative matrix norm, so long as $\|\text{sign}(Z_0)(Z_k$
 540 $-\text{sign}(Z_0))\| < 1$ and Z_0 has no pure imaginary eigenvalues [18, Lemma 5.12]. In
 541 other words, $\|Z_k^2 - I\|$ provides an estimate for the accuracy of Z_k .

542 For the iterates Z_k with initial condition (10), we have, in the notation of
 543 Lemma 9,

$$544 \quad Z_k^2 - I = \begin{pmatrix} H_k^2 - I & H_k \Omega_k - \Omega_k H_k \\ 0 & H_k^2 - I \end{pmatrix}.$$

545 Likewise, for the iterates Z_k with initial condition (11), we have, in the notation of
 546 Lemma 10,

$$547 \quad Z_k^2 - I = \begin{pmatrix} H_k^2 - I & H_k S_k + S_k H_k \\ 0 & H_k^2 - I \end{pmatrix}.$$

548 Thus, accuracy is assured when the quantities $\|H_k^2 - I\|$, $\|H_k \Omega_k - \Omega_k H_k\|$, and $\|H_k S_k +$
 549 $S_k H_k\|$ are small. Of course, H_k , Ω_k , and S_k are never computed explicitly in the
 550 iteration (12-13), so we must relate these quantities to X_k and E_k using Lemma 11.
 551 By (40), we have

$$552 \quad H_k^2 - I = X_k^* X_k - I.$$

553 The quantities $H_k \Omega_k - \Omega_k H_k$ and $H_k S_k + S_k H_k$ are more difficult to relate to X_k and
 554 E_k in a computable way (i.e., a way that does not involve knowing U in advance).
 555 However, second-order accurate approximations to $H_k \Omega_k - \Omega_k H_k$ and $H_k S_k + S_k H_k$
 556 are available. As shown in Appendix A, we have

$$557 \quad (51) \quad H_k \Omega_k - \Omega_k H_k = \frac{1}{2} (X_k^* X_k X_k^* E_k - X_k^* E_k X_k^* X_k) + F_k,$$

$$558 \quad (52) \quad H_k S_k - S_k H_k = X_k^* E_k + E_k^* X_k - \frac{1}{2} (X_k^* X_k X_k^* E_k - X_k^* E_k X_k^* X_k) - F_k,$$

560 where

$$561 \quad \|F_k\| = O(\|H_k^2 - I\|^2 + \|H_k^2 - I\| \|H_k S_k + S_k H_k\|).$$

562 Roughly speaking, (51) arises from the approximations $H_k \approx \frac{1}{2}(I + X_k^* X_k)$ and $\Omega_k \approx$
 563 $X_k^* E_k$. It turns out that only the first of these approximations is second-order accurate
 564 (see Lemma 15), but delicate cancellations detailed in Appendix A lead to the validity
 565 of (51). One then deduces (52) by noting that $X_k^* E_k + E_k^* X_k = (H_k \Omega_k - \Omega_k H_k) +$
 566 $(H_k S_k - S_k H_k)$ (see Lemma 13).

567 In summary, the quantities

$$568 \quad (53) \quad \mathcal{A}_k = X_k^* X_k - I,$$

$$569 \quad (54) \quad \mathcal{B}_k = \frac{1}{2} (X_k^* X_k X_k^* E_k - X_k^* E_k X_k^* X_k),$$

$$570 \quad (55) \quad \mathcal{C}_k = X_k^* E_k + E_k^* X_k - \mathcal{B}_k$$

572 are computable approximations to $H_k^2 - I$, $H_k \Omega_k - \Omega_k H_k$, and $H_k S_k + S_k H_k$, re-
 573 spectively. These are small in norm if and only if $\|Z_k - \text{sign}(Z_0)\|$ is small (for each
 574 of the initial conditions (10) and (11)), which is true if and only if $\|X_k - U\|$ and
 575 $\|E_k - L_{\mathcal{P}}(A, E)\|$ are small. As a practical note, these arguments appear to break
 576 down if A is very ill-conditioned, as illustrated in Section 6.

577 Based on these considerations, we propose that the iterations be terminated when

$$578 \quad (56) \quad \|\mathcal{A}_k\| \leq \delta \|X_k\| \quad \text{and} \quad \|\mathcal{B}_k\| + \|\mathcal{C}_k\| \leq \varepsilon \|E_k\|,$$

580 where δ and ε are relative error tolerances for $\|X_k - U\|$ and $\|E_k - L_{\mathcal{P}}(A, E)\|$,
 581 respectively.

582 As an alternative approach to terminating the iterations, one could consider basing
 583 the decision to terminate on the smallness of the step lengths $\|X_{k+1} - X_k\|$ and
 584 $\|E_{k+1} - E_k\|$. Details of this approach, for the case in which E_k is absent, can be
 585 found in [18, Chapter 8.7].

586 **4.3. Stability.** Stability of the iterative schemes detailed in Theorem 1 is rela-
 587 tively easy to establish. Indeed, the map

$$588 \quad (57) \quad \mathcal{F} \begin{pmatrix} A \\ E \end{pmatrix} = \begin{pmatrix} \mathcal{P}(A) \\ L_{\mathcal{P}}(A, E) \end{pmatrix}$$

589 is idempotent, since $\mathcal{P}(\mathcal{P}(A)) = \mathcal{P}(A)$ and $L_{\mathcal{P}}(\mathcal{P}(A), L_{\mathcal{P}}(A, E)) = L_{\mathcal{P}}(A, E)$ by
 590 the chain rule. It follows that any superlinearly convergent iteration for computing
 591 $\begin{pmatrix} \mathcal{P}(A) \\ L_{\mathcal{P}}(A, E) \end{pmatrix}$ is automatically stable [18, Theorem 4.19]. More precisely, if

$$592 \quad (58) \quad \begin{pmatrix} X_{k+1} \\ E_{k+1} \end{pmatrix} = \begin{pmatrix} g(X_k) \\ L_g(X_k, E_k) \end{pmatrix}$$

593 converges superlinearly to $\begin{pmatrix} \mathcal{P}(X_0) \\ L_{\mathcal{P}}(X_0, E_0) \end{pmatrix}$ for all X_0 and E_0 sufficiently close to A
 594 and E , respectively, then the iteration is stable in the sense of [18, Definition 4.17].
 595 Moreover, the Fréchet derivative of the map (57) coincides with the Fréchet derivative
 596 of the map (58) at the fixed point $\begin{pmatrix} \mathcal{P}(A) \\ L_{\mathcal{P}}(A, E) \end{pmatrix}$ [18, Theorem 4.19].

597 As an example, the Newton iteration (17-18) is superlinearly convergent by virtue
 598 of the superlinear (indeed, quadratic) convergence of the corresponding matrix sign
 599 function iteration (16). The Newton-Schulz iteration (22-23) is likewise superlinearly
 600 (indeed, quadratically) convergent, provided that the singular values of A lie in the
 601 interval $(0, \sqrt{2})$. Thus, both iterations are stable. Using, for instance, (22-23), we
 602 find that the Fréchet derivative of the map (58) (and hence of the map (57)) at
 603 $\begin{pmatrix} U \\ K \end{pmatrix} = \begin{pmatrix} \mathcal{P}(A) \\ L_{\mathcal{P}}(A, E) \end{pmatrix}$ is given by

$$604 \quad L_{\mathcal{F}} \left(\begin{pmatrix} U \\ K \end{pmatrix}, \begin{pmatrix} F \\ G \end{pmatrix} \right) = \begin{pmatrix} F - \frac{1}{2}U(U^*F + F^*U) \\ G - \frac{1}{2}[U(U^*G + G^*U + K^*F + F^*K) + K(U^*F + F^*U)] \end{pmatrix}.$$

605 Note that when U is square, the identities $UU^* = I$ and $U^*K = -K^*U$ (by (30))
 606 imply that this formula reduces to

$$607 \quad L_{\mathcal{F}} \left(\begin{pmatrix} U \\ K \end{pmatrix}, \begin{pmatrix} F \\ G \end{pmatrix} \right) = \begin{pmatrix} \frac{1}{2}(F - UF^*U) \\ \frac{1}{2}(G - UG^*U - UF^*K - KF^*U) \end{pmatrix},$$

608 in agreement with [18, Theorem 8.19].

609 **4.4. Condition Number Estimation.** A seemingly natural application of
 610 Theorem 1 is to leverage the iterative scheme (12-13) to estimate the (absolute) con-
 611 dition number

$$612 \quad \kappa(\mathcal{P}, A) = \|L_{\mathcal{P}}(A, \cdot)\| = \sup_{\substack{E \in \mathbb{C}^m \times \mathbb{C}^n, \\ E \neq 0}} \frac{\|L_{\mathcal{P}}(A, E)\|}{\|E\|}$$

613 of the map \mathcal{P} at A (or its counterpart $\kappa_{rel}(\mathcal{P}, A) = \kappa(\mathcal{P}, A)\|A\|/\|\mathcal{P}(A)\|$, the relative
 614 condition number of A). As tempting as it may seem, a much simpler (and undoubtedly
 615 more efficient) algorithm is available for estimating $\kappa(\mathcal{P}, A)$. As explained in [18,
 616 Theorem 8.9], the value of $\kappa(\mathcal{P}, A)$ at $A \in \mathbb{C}^{m \times n}$ ($m \geq n$) is σ_n^{-1} , where σ_n denotes
 617 the smallest singular value of A . This quantity can be estimated efficiently by apply-
 618 ing the power method [18, Algorithm 3.19] to $(A^*A)^{-1}$. In most iterative algorithms
 619 for computing the polar decomposition, this matrix (or A^{-1}) is computed in the first
 620 iteration, so the additional cost of computing $\kappa(\mathcal{P}, A)$ is negligible.

621 Before finishing our discussion of condition number estimation, it is worth pointing
 622 out a subtlety that arises when considering the polar decomposition of a real square
 623 matrix. If A is real and square ($m = n$), then it can be shown that the (absolute)
 624 condition number of A with respect to *real* perturbations is $2(\sigma_n + \sigma_{n-1})^{-1}$ [18,
 625 Theorem 8.9]. This fact will play a role in our interpretation of certain numerical
 626 experiments in Section 6.

627 **5. Comparison with Other Methods.** There are several other methods that
 628 can be used to compute the Fréchet derivative of the polar decomposition. Below, we
 629 describe a few and compare them with iterative schemes of the form (12-13).

630 One alternative is to recognize that $L_{\mathcal{P}}(A, E)$ is the solution to a Lyapunov equa-
 631 tion. Indeed, upon noting that $\mathcal{P}(A)^*A = U^*A = H$ is Hermitian, one can differenti-
 632 ate the relation

$$633 \quad \text{skew}(\mathcal{P}(A)^*A) = 0$$

634 with the aid of the product rule to obtain

$$635 \quad \text{skew}(L_{\mathcal{P}}(A, E)^*A + \mathcal{P}(A)^*E) = 0$$

636 for any $E \in \mathbb{C}^{m \times n}$. Substituting $A = UH$ and $\mathcal{P}(A) = U$, and denoting $Y :=$
 637 $U^*L_{\mathcal{P}}(A, E) = -L_{\mathcal{P}}(A, E)^*U$, we obtain

$$638 \quad (59) \quad HY + YH = U^*E - E^*U.$$

639 Given H , U , and E , this is a Lyapunov equation in the unknown Y , which, by the
 640 positive-definiteness of H , has a unique solution. It can be solved using standard
 641 algorithms for the solution of Lyapunov and Sylvester equations [5, 13]. It also has
 642 theoretical utility, offering an alternative proof of part of Theorem 2, owing to a well-
 643 known connection between the solution of Lyapunov and Sylvester equations and the
 644 matrix sign function [31, 18, Chapter 2.4]. Indeed, (59) is equivalent to the equation

$$645 \quad \begin{pmatrix} H & E^*U - U^*E \\ 0 & -H \end{pmatrix} = \begin{pmatrix} I & Y \\ 0 & I \end{pmatrix} \begin{pmatrix} H & 0 \\ 0 & -H \end{pmatrix} \begin{pmatrix} I & Y \\ 0 & I \end{pmatrix}^{-1}.$$

646 Taking the sign of both sides, noting that $\text{sign}(H) = I$, and using the fact that the
 647 matrix sign function commutes with similarity transformations, we conclude that

$$648 \quad \text{sign} \begin{pmatrix} H & E^*U - U^*E \\ 0 & -H \end{pmatrix} = \begin{pmatrix} I & Y \\ 0 & I \end{pmatrix} \left[\text{sign} \begin{pmatrix} H & 0 \\ 0 & -H \end{pmatrix} \right] \begin{pmatrix} I & Y \\ 0 & I \end{pmatrix}^{-1} \\ 649 \quad \quad \quad = \begin{pmatrix} I & -2Y \\ 0 & -I \end{pmatrix}. \\ 650$$

651 This is precisely the identity (14), up to a rescaling of E . Its connection with the
 652 Lyapunov equation (59) reveals that the coupled iteration (12-13) is effectively solv-
 653 ing (59) and computing the polar decomposition simultaneously. In comparison to a

naive approach in which (59) is solved after first computing the polar decomposition, the coupled iteration (12-13) is attractive, as it computes $L_{\mathcal{P}}(A, E)$ at the expense of a few extra matrix-matrix multiplications and additions on top of the computation of $\mathcal{P}(A)$.

When A and E are real, another method for computing Fréchet derivative of a matrix function f is to use the complex step approximation [2]

$$(60) \quad L_f(A, E) \approx \operatorname{Im} \left(\frac{f(A + ihE) - f(A)}{h} \right),$$

where h is a small positive scalar and $\operatorname{Im}(B)$ denotes the imaginary part of a matrix B . By using a pure imaginary step ih , this approximation does not suffer from cancellation errors that plague standard finite differencing, allowing h to be taken arbitrarily small [2]. This approximation can be applied to the polar decomposition, but care must be exercised in order to do so correctly. In particular, a meaningful approximation is obtained only if the conjugate transposes X_k^* appearing in the algorithm are interpreted as transposes X_k^T when evaluating the “polar decomposition” of $A + ihE$. We have put “polar decomposition” in quotes since the result of such a computation is the matrix $(A + ihE) [(A + ihE)^T (A + ihE)]^{-1/2}$, not $\mathcal{P}(A + ihE) = (A + ihE) [(A + ihE)^* (A + ihE)]^{-1/2}$. The cost of this approximation is close to the cost of computing two polar decompositions.

Another approach is to appeal to the relation $\mathcal{P}(A) = A(A^*A)^{-1/2}$. By (24), the Fréchet derivative of \mathcal{P} at A in the direction E is given by

$$\begin{aligned} L_{\mathcal{P}}(A, E) &= E(A^*A)^{-1/2} + AL_{x^{-1/2}}(A^*A, E^*A + A^*E) \\ &= EH^{-1} + AL_{x^{-1/2}}(A^*A, E^*A + A^*E). \end{aligned}$$

Evaluating the second term, the Fréchet derivative of the inverse square root, can be reduced to the task of solving a Lyapunov equation, so this approach is essentially of the same complexity as the one based on (59).

Any of the aforementioned methods, including our own, can be applied in two different ways when A is rectangular ($m \times n$ with $m > n$). One way is to apply the methods verbatim, working at all times with rectangular matrices. The alternative is to first compute a reduced QR decomposition $A = QR$, where $Q \in \mathbb{C}^{m \times n}$ has orthonormal columns and $R \in \mathbb{C}^{n \times n}$ is upper triangular. Then, one can compute $\mathcal{P}(R)$ and $L_{\mathcal{P}}(R, Q^*E)$ (which are square matrices) and invoke the identities

$$U = \mathcal{P}(A) = Q\mathcal{P}(R), \quad H = \mathcal{P}(R)^*R$$

and

$$\begin{aligned} L_{\mathcal{P}}(A, E) &= L_{\mathcal{P}}(A, QQ^*E) + L_{\mathcal{P}}(A, (I - QQ^*)E) \\ &= QL_{\mathcal{P}}(R, Q^*E) + (I - QQ^*)EH^{-1} \end{aligned}$$

to recover $\mathcal{P}(A)$ and $L_{\mathcal{P}}(A, E)$. The validity of the latter identity is a consequence of (31), (28), and the fact that $QQ^* = UU^*$. In summary, computations for rectangular A can be reduced to the square case by performing a reduced QR decomposition of A at the outset.

Finally, when A is square, one more method for computing $L_{\mathcal{P}}(A, E)$ is available, as noted in, for instance, [21]. The idea is to make use of the singular value decomposition $A = P\Sigma Q^*$, where $P, Q \in \mathbb{C}^{n \times n}$ are unitary and $\Sigma \in \mathbb{C}^{n \times n}$ is diagonal.

698 The singular value decomposition is related to the polar decomposition $A = UH$ via
 699 the relations $U = PQ^*$ and $H = Q\Sigma Q^*$. Moreover, the Lyapunov equation (59) is
 700 equivalent to

$$701 \quad \Sigma G + G\Sigma = F - F^*,$$

702 where $F = P^*EQ$ and $G = P^*L_{\mathcal{P}}(A, E)Q$ [21, Equation 2.18]. Given Σ and F , this
 703 equation admits an explicit solution for the components of G . Namely,

$$704 \quad G_{ij} = \frac{1}{\sigma_i + \sigma_j} (F_{ij} - \overline{F_{ji}}),$$

705 where σ_i denotes the i^{th} diagonal entry of Σ , and $\overline{F_{ji}}$ denotes the complex conjugate of
 706 F_{ji} . One then obtains $L_{\mathcal{P}}(A, E)$ from $L_{\mathcal{P}}(A, E) = PGQ^*$. This method is attractive
 707 if the singular value decomposition of A has already been computed, but otherwise it
 708 is an expensive approach in general.

709 **5.1. Floating Point Operations.** Relative to the methods listed above, the
 710 iterative schemes derived in this paper are distinguished by their efficiency, at least
 711 when n is large and the columns of A are close to being orthonormal. To see this, con-
 712 sider the number of floating point operations needed to compute $\mathcal{P}(A)$ and $L_{\mathcal{P}}(A, E)$.
 713 For simplicity, assume that A and E are real and of size $n \times n$. Then, to leading or-
 714 der in n , and excluding the costs associated with termination criteria in the iterative
 715 schemes, the methods have the following computational costs:

- 716 • The iteration (17-18) requires n_{iter} matrix inversions (each requiring $2n^3$
 717 flops [18, Appendix C]) and $2n_{\text{iter}}$ matrix multiplications (each requiring $2n^3$
 718 flops), where n_{iter} denotes the number of iterations used. Its computational
 719 cost is thus $n_{\text{iter}}(2n^3) + 2n_{\text{iter}}(2n^3) = 6n_{\text{iter}}n^3$ flops.
- 720 • Solving the Lyapunov equation (59) with a direct method involves diagonaliz-
 721 ing H ($9n^3$ flops [18, Appendix C]) and performing 4 matrix multiplications,
 722 for a total of $9n^3 + 4(2n^3) = 17n^3$ flops. The additional cost of computing U ,
 723 $H = U^*A$, $L_{\mathcal{P}}(A, E) = UY$, and U^*E (assuming that (17) is used to com-
 724 pute U) is dominated by the cost of performing n_{iter} matrix inversions and
 725 3 matrix multiplications, bringing the total to $17n^3 + n_{\text{iter}}(2n^3) + 3(2n^3) =$
 726 $(23 + 2n_{\text{iter}})n^3$ flops.
- 727 • The complex step approximation (assuming that (17) is used to compute the
 728 polar decomposition of A and $A + ihE$) requires $2n_{\text{iter}}$ matrix inversions, of
 729 which n_{iter} involve complex arithmetic. Since each inversion of a complex
 730 matrix requires n^3 additions of complex scalars (2 real flops) and n^3 mul-
 731 tiplications of complex scalars (6 real flops), the computational cost of the
 732 complex step approximation is $n_{\text{iter}}(2n^3) + n_{\text{iter}}(8n^3) = 10n_{\text{iter}}n^3$ flops.
- 733 • The method based on the singular value decomposition requires 5 matrix
 734 multiplications plus the computation of the SVD. Assuming, for instance,
 735 that the Golub-Reinsch algoirthm ($22n^3$ flops [14]) is used to compute the
 736 SVD, this method's total cost is $5(2n^3) + 22n^3 = 32n^3$ flops.

737 We conclude from this analysis that, for sufficiently large n , the iteration (17-18)
 738 requires fewer floating point operations than its competitors whenever $n_{\text{iter}} \leq 5$. Note
 739 that this is no longer the case if the costs of computing the residual estimates (53-
 740 55) are taken into account. However, if efficiency is the primary objective, then
 741 cheaper termination criteria (based, for instance, on $\|X_k - X_k^{-*}\|$, $\|X_{k+1} - X_k\|$,
 742 and/or $\|E_{k+1} - E_k\|$) may be appropriate.

743 It should be noted that these floating point operation counts are, of course, crude
 744 measures of efficiency that do not account for other factors – parallelizability and

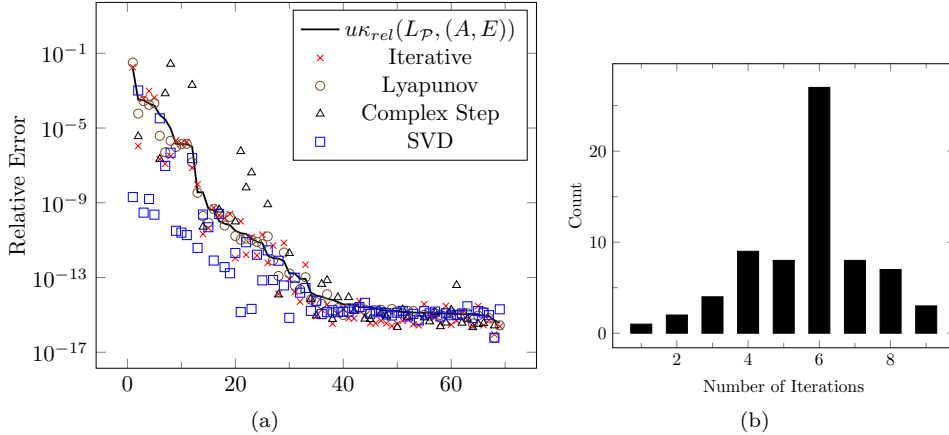


FIG. 1. (a) Relative errors in the computed values of $L_{\mathcal{P}}(A, E)$ for various matrices A and E using four different methods. (b) Histogram showing the number of iterations used by the iterative method in these tests.

k	$\frac{\ X_k - U\ _F}{\ U\ _F}$	$\frac{\ E_k - K\ _F}{\ K\ _F}$	$\ \mathcal{A}_k\ _F$	$\ \mathcal{B}_k\ _F$	$\ \mathcal{C}_k\ _F$	$\ \tilde{\mathcal{B}}_k\ _F$	$\ \tilde{\mathcal{C}}_k\ _F$	μ_k
1	1.3e-5	2.3e-3	1.1e-4	1.6e-4	4.7e-2	1.6e-4	4.7e-2	1.0e0
2	3.0e-10	4.8e-8	2.4e-9	3.5e-9	9.9e-7	3.5e-9	9.9e-7	1.0e0
3	3.4e-16	5.0e-16	1.4e-15	2.3e-15	4.3e-15	5.0e-15	6.2e-15	1.0e0

TABLE 1

Nearly orthogonal matrix, $m = n = 16$, $\sigma_n(A) = 9.9e-1$, $\sigma_{n-1}(A) = 1.0e0$, $\kappa(A) = 1.0e0$.

745 numerical stability, for instance – which may very well render the iterative meth-
 746 ods in this paper even more attractive (for large matrices with nearly orthonormal
 747 columns). Parallelizability is particularly noteworthy, since these iterative methods
 748 require only multiplication, inversion, and addition of matrices. In contrast, meth-
 749 ods based on the SVD or the solution of the Lyapunov equation (59) involve matrix
 750 decompositions (unless the Lyapunov equation (59) is solved iteratively, a strategy
 751 which we have already argued to be inferior to the simultaneous computation of $\mathcal{P}(A)$
 752 and $L_{\mathcal{P}}(A, E)$ via (17-18)). These considerations suggest that for large matrices with
 753 nearly orthonormal columns, the iterative methods in this paper are likely better
 754 suited for parallel computing environments than their competitors.

755 We emphasize that these comparisons are relevant only if the goal is to calcu-
 756 late $L_{\mathcal{P}}(A, E)$ for specific A and E . If condition number estimation is the ultimate
 757 goal, then calculation of the Fréchet derivative of \mathcal{P} is unnecessary, as explained in
 758 Section 4.4. We refer the reader to [12] for an example of an application in which is
 759 desirable to calculate $L_{\mathcal{P}}(A, E)$ and not merely $\kappa(\mathcal{P}, A)$.

760 **6. Numerical Experiments.** To illustrate the performance of the iterative
 761 schemes derived in this paper, we have computed the Fréchet derivative of the polar de-
 762 composition for a collection of 69 matrices: 44 real matrices of size 10×10 from the Ma-
 763 trix Computation Toolbox [15] (we used all of the test matrices in the toolbox except
 764 those that are singular to working precision), as well as 25 complex matrices of size $m \times$
 765 10 generated by the MATLAB command `crandn(1)*gallery('randsvd', [m 10],`
 766 `kappa, mode)`, where $m \in \{2, 4, 6, 8, 10\}$, $\text{kappa} = 10^{(8 * \text{rand}(1))}$, and $\text{mode} \in$
 767 $\{1, 2, 3, 4, 5\}$.

768 We computed $\mathcal{P}(A)$ and $L_{\mathcal{P}}(A, E)$ for each matrix A described above, with E a

k	$\frac{\ X_k - U\ _F}{\ U\ _F}$	$\frac{\ E_k - K\ _F}{\ K\ _F}$	$\ A_k\ _F$	$\ B_k\ _F$	$\ C_k\ _F$	$\ \tilde{B}_k\ _F$	$\ \tilde{C}_k\ _F$	μ_k
1	1.7e1	7.4e1	2.4e3	8.5e4	8.5e4	1.2e2	2.7e3	1.8e-1
2	1.4e0	6.1e0	2.2e1	5.2e0	2.2e1	8.7e-1	2.1e1	5.9e-1
3	1.3e-1	1.2e0	1.1e0	2.4e-2	1.7e0	1.9e-2	1.7e0	9.2e-1
4	2.6e-3	8.7e-2	2.1e-2	5.3e-4	1.0e-1	5.3e-4	1.0e-1	1.0e0
5	1.4e-6	1.8e-4	1.1e-5	4.0e-7	2.1e-4	4.0e-7	2.1e-4	1.0e0
6	3.9e-13	2.1e-10	3.1e-12	6.9e-14	2.4e-10	6.9e-14	2.4e-10	1.0e0
7	2.1e-15	2.4e-15	1.4e-15	9.4e-17	1.3e-16	1.6e-15	1.6e-15	1.0e0

TABLE 2

Binomial matrix, $m = n = 16$, $\sigma_n(A) = 2.6e0$, $\sigma_{n-1}(A) = 2.6e0$, $\kappa(A) = 4.7e3$.

k	$\frac{\ X_k - U\ _F}{\ U\ _F}$	$\frac{\ E_k - K\ _F}{\ K\ _F}$	$\ A_k\ _F$	$\ B_k\ _F$	$\ C_k\ _F$	$\ \tilde{B}_k\ _F$	$\ \tilde{C}_k\ _F$	μ_k
1	2.9e6	8.9e17	8.4e13	2.5e27	2.5e27	1.4e14	6.8e25	1.1e-6
2	1.9e0	4.6e11	4.5e1	1.2e3	1.9e13	4.8e1	1.9e13	4.2e-1
3	2.7e-1	6.7e10	2.5e0	3.2e0	7.5e11	1.8e0	7.5e11	8.3e-1
4	9.5e-3	6.0e8	7.7e-2	9.6e-2	5.7e9	9.4e-2	5.7e9	9.9e-1
5	3.9e-5	3.6e6	3.1e-4	5.3e-4	3.4e7	5.3e-4	3.4e7	1.0e0
6	1.6e-9	4.1e1	1.3e-8	2.6e-8	3.8e2	2.6e-8	3.8e2	1.0e0
7	7.2e-16	4.1e-5	1.1e-15	8.4e-16	4.0e-7	8.4e-15	4.0e-7	1.0e0
8	7.3e-16	4.1e-5	9.3e-16	7.1e-16	1.4e-15	8.5e-15	9.2e-15	1.0e0

TABLE 3

Frank matrix, $m = n = 16$, $\sigma_n(A) = 3.5e-13$, $\sigma_{n-1}(A) = 8.7e-1$, $\kappa(A) = 2.3e14$.

769 matrix (of the same dimensions as A) consisting of random entries sampled from a
770 normal distribution with mean 0 and variance 1. We used the Newton iteration (45-
771 46) with scaling parameter (47) for the square matrices and its generalization (19-20)
772 with scaling parameter (49) for the rectangular matrices. To terminate the iterations,
773 we used (56) with $\delta = \varepsilon = 10^{-14}$ and $\|\cdot\|$ equal to the Frobenius norm. Note that for
774 simplicity, we used scaling throughout the entire iteration, even though the scaling
775 parameter μ_k approaches 1 near convergence. A more efficient approach is to switch
776 to an unscaled iteration after a certain point. A heuristic for deciding when to do so
777 is detailed in [18, Chapter 8.9].

778 We compared the iterative methods with three alternatives described in Section 5:
779 solving the Lyapunov equation (59), using the complex step approximation (60) (when
780 applicable), and using the singular value decomposition. We also computed the “ex-
781 act” values of $\mathcal{P}(A)$ and $L_{\mathcal{P}}(A, E)$ using the singular value decomposition with 100-
782 digit precision using the Symbolic Math Toolbox of MATLAB. Figure 1(a) shows the
783 relative errors (in the Frobenius norm) in the computed values of $L_{\mathcal{P}}(A, E)$ for each
784 of the 69 tests, arranged in order of decreasing condition number $\kappa_{rel}(L_{\mathcal{P}}(A, E))$
785 – the relative condition number of the map $(A, E) \mapsto L_{\mathcal{P}}(A, E)$. We estimated the
786 latter quantity using finite differencing to approximate the derivative of this map in
787 a randomly chosen direction, as advocated in [1]. The solid line in Figure 1(a) shows
788 the estimated value of $u\kappa_{rel}(L_{\mathcal{P}}(A, E))$, where $u = 2^{-53}$ is the unit roundoff. The
789 plot indicates that all of the methods under comparison behave in a forward stable
790 way. Figure 1(b) shows a histogram of the number of iterations used by the iterative
791 methods in these tests.

792 To study the convergence of the iterative methods in more detail, we focus now
793 on a few representative matrices obtained from the MATLAB matrix gallery. Note

k	$\frac{\ X_k - U\ _F}{\ U\ _F}$	$\frac{\ E_k - K\ _F}{\ K\ _F}$	$\ \mathcal{A}_k\ _F$	$\ \mathcal{B}_k\ _F$	$\ \mathcal{C}_k\ _F$	$\ \tilde{\mathcal{B}}_k\ _F$	$\ \tilde{\mathcal{C}}_k\ _F$	μ_k
1	3.4e6	6.9e6	1.0e14	7.9e36	7.9e36	1.9e23	1.5e26	6.5e-7
2	1.2e0	1.4e0	2.2e1	9.8e10	1.6e13	2.0e10	1.6e13	5.1e-1
3	1.4e-1	4.4e-2	1.2e0	5.3e8	4.2e11	5.2e8	4.1e11	9.4e-1
4	7.2e-3	2.7e-3	5.8e-2	2.7e7	2.2e10	2.7e7	2.3e10	1.0e0
5	9.3e-5	3.0e-4	4.9e-4	1.3e4	1.0e7	1.3e4	1.9e9	1.0e0
6	7.0e-5	3.0e-4	5.5e-8	6.2e-2	4.9e1	6.2e-2	1.9e9	1.0e0
7	7.0e-5	3.0e-4	1.2e-15	1.0e-3	1.5e-3	3.4e-3	1.9e9	1.0e0

TABLE 4

Modified Frank matrix, $m = n = 16$, $\sigma_n(A) = 3.5e-13$, $\sigma_{n-1}(A) = 3.5e-13$, $\kappa(A) = 2.3e14$.

794 that the first three matrices are identical to those considered in [18, Chapter 8.9].

- 795 1. A nearly orthogonal matrix, `orth(gallery('moler',16))+ones(16)*1e-3`.
796 2. A binomial matrix, `gallery('binomial',16)`.
797 3. The Frank matrix, `gallery('frank',16)`.
798 4. A modification of the Frank matrix obtained by setting its second smallest
799 singular value equal to its smallest singular value. That is, $A = P\tilde{\Sigma}Q^*$ where
800 $P\Sigma Q^*$ is the singular value decomposition of the Frank matrix, $\tilde{\Sigma}_{ii} = \Sigma_{ii}$ for
801 $i \neq 15$, and $\tilde{\Sigma}_{15,15} = \Sigma_{16,16}$.

802 We computed $\mathcal{P}(A)$ and $L_{\mathcal{P}}(A, E)$ for each A listed above, again with E a matrix
803 consisting of random entries sampled from a normal distribution with mean 0 and
804 variance 1. Tables 1-4 show the values of several quantities monitored during the
805 iterations. The first two columns show the relative errors $\frac{\|X_k - U\|_F}{\|U\|_F}$ and $\frac{\|E_k - K\|_F}{\|K\|_F}$,
806 where $U = \mathcal{P}(A)$ and $K = L_{\mathcal{P}}(A, E)$. The next three columns show the norms
807 of (53-55), which are the quantities we used to determine when to terminate the
808 iterations. Recall that (54) and (55) are computable approximations to $H_k\Omega_k - \Omega_kH_k$
809 and $H_kS_k + S_kH_k$, respectively. We have denoted $\tilde{\mathcal{B}}_k = H_k\Omega_k - \Omega_kH_k$ and $\tilde{\mathcal{C}}_k =$
810 $H_kS_k + S_kH_k$ in the tables and recorded their norms in the seventh and eighth columns.
811 Finally, the last column of the tables shows the value of the scaling parameter μ_k . In
812 the caption of each table, we have made note of the dimensions of the matrix A , the
813 smallest and second smallest singular values $\sigma_n(A)$ and $\sigma_{n-1}(A)$ of A , respectively,
814 and the condition number $\kappa(A)$ of A .

815 Tables 1 and 2 illustrate the effectiveness of the iteration on relatively well-
816 conditioned matrices. In both cases, small relative errors in both X_k and E_k are
817 achieved simultaneously, and convergence is detected appropriately by the termina-
818 tion criteria (56). Comparison of the columns labeled $\|\mathcal{B}_k\|$ and $\|\mathcal{C}_k\|$ with the columns
819 labeled $\|\tilde{\mathcal{B}}_k\|$ and $\|\tilde{\mathcal{C}}_k\|$, respectively, lends credence to the asymptotic accuracy of the
820 approximations $\mathcal{B}_k \approx \tilde{\mathcal{B}}_k$ and $\mathcal{C}_k \approx \tilde{\mathcal{C}}_k$, at least until roundoff errors begin to intervene.

821 Tables 3 and 4 illustrate what can go wrong when A is ill-conditioned. In the
822 case of Table 4, the matrix A (the modified Frank matrix) has condition number
823 $\kappa(A) = 2.3e14$, and its two smallest singular values are both close to zero: $\sigma_n(A) =$
824 $\sigma_{n-1}(A) = 3.5e-13$. As a consequence, the (absolute) condition number of \mathcal{P} with
825 respect to real perturbations (as explained in Section 4.4) is $2(\sigma_n + \sigma_{n-1})^{-1} = 2.9e12$,
826 and we cannot expect much more than 3 or 4 digits of relative accuracy in double
827 precision arithmetic when approximating $\mathcal{P}(A)$, much less $L_{\mathcal{P}}(A, E)$. This expectation
828 is born out in Table 4. A more subtle phenomenon occurs in Table 3. There, the
829 matrix A (the Frank matrix) has condition number $\kappa(A) = 2.3e14$ as well, but only one
830 of its singular values is close to zero. Namely, $\sigma_n(A) = 3.5e-13$, but $\sigma_{n-1}(A) = 8.7e-$

831 1. As a consequence, \mathcal{P} is very well-conditioned with respect to real perturbations,
 832 having (absolute) condition number $2(\sigma_n + \sigma_{n-1})^{-1} = 1.2e0$. Curiously, the result is
 833 that $\mathcal{P}(A)$ is approximated very accurately, but $L_{\mathcal{P}}(A, E)$ is not. The fact that the
 834 performance of the Newton iteration (45) is largely unaffected by poorly conditioned
 835 A (unless A has two singular values close to zero) has been noted in [18, Chapter
 836 8.9]. The observation that, in contrast, it takes only one near-zero singular value to
 837 corrupt the computation of $L_{\mathcal{P}}(A, E)$ via the iteration (45-46) deserves further study.

838 **7. Conclusion.** In this paper, we have derived iterative schemes for computing
 839 the Fréchet derivative of the polar decomposition. The structure of these iterative
 840 schemes lends credence to the mantra that differentiating an iteration for computing
 841 $f(A)$ leads to an iteration for computing $L_f(A, E)$. It would be interesting to deter-
 842 mine what conditions on a matrix function f ensure that this mantra bears out in
 843 practice. Certainly being a primary matrix function suffices, but the results of the
 844 present paper suggest that such a construction might work in a more general setting.

845 On a more specific level, several aspects of this paper warrant further consider-
 846 ation. While the termination criteria devised in Section 4.2 appear to work well in
 847 practice, a more careful analysis of their effectiveness is lacking. In addition, it would
 848 be of interest to better understand the behavior of the iterative scheme (45-46) on
 849 ill-conditioned matrices.

850 **Appendix A. Approximate Residuals.** In this section, we prove the validity
 851 of (51-52). Suppressing the subscript k for the remainder of this section, our goal is
 852 to show that if

$$853 \quad (61) \quad \mathcal{B} = \frac{1}{2} (X^* X X^* E - X^* E X^* X),$$

$$854 \quad (62) \quad \mathcal{C} = (X^* E + E^* X) - \mathcal{B},$$

856 then

$$857 \quad \mathcal{B} = H\Omega - \Omega H + O(\|H^2 - I\|^2 + \|H^2 - I\| \|HS + SH\|),$$

$$858 \quad \mathcal{C} = HS + SH + O(\|H^2 - I\|^2 + \|H^2 - I\| \|HS + SH\|).$$

860 Now since

$$861 \quad H^2 - I = 2(H - I) + (H - I)^2,$$

862 the norms of $H^2 - I$ and $H - I$ are asymptotically equal, up to a factor of 2. Thus,
 863 it is enough to show that

$$864 \quad (63) \quad \mathcal{B} = H\Omega - \Omega H + O(\|H - I\|^2 + \|H - I\| \|HS + SH\|),$$

$$865 \quad (64) \quad \mathcal{C} = HS + SH + O(\|H - I\|^2 + \|H - I\| \|HS + SH\|).$$

867 The following lemma reduces this task to the verification of (63).

868 **LEMMA 13.** *We have*

$$869 \quad (H\Omega - \Omega H) + (HS + SH) = X^* E + E^* X.$$

870 *Proof.* By (40) and the equalities $H = H^*$, $U^* E = \Omega + S$, $\Omega^* = -\Omega$, and $S^* = S$,
 871 we have

$$\begin{aligned} 872 \quad X^* E + E^* X &= HU^* E + E^* UH \\ 873 &= H(\Omega + S) + (\Omega + S)^* H \\ 874 &= H(\Omega + S) + (-\Omega + S)H \\ 875 &= (H\Omega - \Omega H) + (HS + SH). \quad \square \end{aligned}$$

877 It follows from the preceding lemma that if \mathcal{B} satisfies (63), then $\mathcal{C} = (X^*E + E^*X) - \mathcal{B}$
 878 automatically satisfies (64).

879 To prove (63), we begin by noting a few useful relations.

880 LEMMA 14. For any $B \in \mathbb{C}^{n \times n}$,

$$881 \quad H(HB - BH) = HB - BH + O(\|H - I\|^2),$$

$$882 \quad (HB - BH)H = HB - BH + O(\|H - I\|^2).$$

884 *Proof.* These relations follow from the identities

$$885 \quad H(HB - BH) = HB - BH - (H - I)B(H - I) + (H - I)^2B,$$

$$886 \quad (HB - BH)H = HB - BH + (H - I)B(H - I) - B(H - I)^2. \quad \square$$

888 LEMMA 15. We have

$$889 \quad X^*X = 2H - I + O(\|H - I\|^2).$$

890 *Proof.* Use the identity

$$891 \quad H^2 = 2H - I + (H - I)^2$$

892 together with the fact that $X^*X = HU^*UH = H^2$. \square

893 Now consider (61). Substituting $X^*X = 2H - I + O(\|H - I\|^2)$ and $X^*E =$
 894 $HU^*E = H(\Omega + S)$ gives, after simplification,

$$895 \quad \mathcal{B} = H(H(\Omega + S) - (\Omega + S)H) + O(\|H - I\|^2).$$

897 Applying Lemma 14 with $B = \Omega + S$ gives

$$898 \quad \mathcal{B} = H(\Omega + S) - (\Omega + S)H + O(\|H - I\|^2)$$

$$899 \quad = (H\Omega - \Omega H) + (HS - SH) + O(\|H - I\|^2).$$

901 We will finish the proof of (63) by showing that

$$902 \quad HS - SH = O(\|H - I\|^2 + \|H - I\|\|HS + SH\|).$$

903 Averaging the two equalities in Lemma 14 with $B = S$ gives

$$904 \quad HS - SH = \frac{1}{2} [H(HS - SH) + (HS - SH)H] + O(\|H - I\|^2)$$

906 Finally, an algebraic manipulation shows that the term in brackets above is equal to

$$907 \quad H(HS - SH) + (HS - SH)H = (H - I)(HS + SH) - (HS + SH)(H - I),$$

908 and so it is of order $\|H - I\|\|HS + SH\|$.

909 REFERENCES

- 910 [1] A. H. AL-MOHY AND N. J. HIGHAM, *Computing the Fréchet derivative of the matrix exponential, with an application to condition number estimation*, SIAM Journal on Matrix Analysis and Applications, 30 (2009), pp. 1639–1657.
- 911
 912
 913 [2] A. H. AL-MOHY AND N. J. HIGHAM, *The complex step approximation to the Fréchet derivative of a matrix function*, Numerical Algorithms, 53 (2010), pp. 133–148.
- 914

- 915 [3] A. H. AL-MOHY, N. J. HIGHAM, AND S. D. RELTON, *Computing the Fréchet derivative of*
916 *the matrix logarithm and estimating the condition number*, SIAM Journal on Scientific
917 Computing, 35 (2013), pp. C394–C410.
- 918 [4] A. BARRLUND, *Perturbation bounds on the polar decomposition*, BIT Numerical Mathematics,
919 30 (1990), pp. 101–113.
- 920 [5] R. H. BARTELS AND G. STEWART, *Solution of the matrix equation $AX + XB = C$* , Communi-
921 cations of the ACM, 15 (1972), pp. 820–826.
- 922 [6] R. BHATIA, *Matrix factorizations and their perturbations*, Linear Algebra and its applications,
923 197 (1994), pp. 245–276.
- 924 [7] R. BHATIA, *Matrix analysis*, vol. 169, Springer Science & Business Media, 2013.
- 925 [8] J. R. CARDOSO, *Evaluating the Fréchet derivative of the matrix p^{th} root*, Electronic Transac-
926 tions on Numerical Analysis, 38 (2011), pp. 202–217.
- 927 [9] J. R. CARDOSO, *Computation of the matrix p^{th} root and its Fréchet derivative by integrals*,
928 Electronic Transactions on Numerical Analysis, 39 (2012), pp. 414–436.
- 929 [10] L. DIECI AND T. EIROLA, *On smooth decompositions of matrices*, SIAM Journal on Matrix
930 Analysis and Applications, 20 (1999), pp. 800–819.
- 931 [11] R. L. FOOTE, *Regularity of the distance function*, Proceedings of the American Mathematical
932 Society, 92 (1984), pp. 153–155.
- 933 [12] E. S. GAWLIK AND M. LEOK, *Embedding-based interpolation on the special orthogonal group*,
934 (Preprint), (2016).
- 935 [13] G. GOLUB, S. NASH, AND C. VAN LOAN, *A Hessenberg-Schur method for the problem $AX +$*
936 *$XB = C$* , IEEE Transactions on Automatic Control, 24 (1979), pp. 909–913.
- 937 [14] G. H. GOLUB AND C. F. VAN LOAN, *Matrix computations*, vol. 3, JHU Press, 2012.
- 938 [15] N. J. HIGHAM, *The Matrix Computation Toolbox*. [http://www.ma.man.ac.uk/~higham/](http://www.ma.man.ac.uk/~higham/mctoolbox)
939 [mctoolbox](http://www.ma.man.ac.uk/~higham/mctoolbox).
- 940 [16] N. J. HIGHAM, *Computing the polar decomposition – with applications*, SIAM Journal on Sci-
941 entific and Statistical Computing, 7 (1986), pp. 1160–1174.
- 942 [17] N. J. HIGHAM, *The matrix sign decomposition and its relation to the polar decomposition*,
943 Linear Algebra and its Applications, 212 (1994), pp. 3–20.
- 944 [18] N. J. HIGHAM, *Functions of matrices: theory and computation*, SIAM, 2008.
- 945 [19] N. J. HIGHAM AND L. LIN, *An improved Schur–Padé algorithm for fractional powers of a*
946 *matrix and their Fréchet derivatives*, SIAM Journal on Matrix Analysis and Applications,
947 34 (2013), pp. 1341–1360.
- 948 [20] R. A. HORN AND C. R. JOHNSON, *Matrix analysis*, Cambridge University Press, 2012.
- 949 [21] C. KENNEY AND A. J. LAUB, *Polar decomposition and matrix sign function condition estimates*,
950 SIAM Journal on Scientific and Statistical Computing, 12 (1991), pp. 488–504.
- 951 [22] C. KENNEY AND A. J. LAUB, *Rational iterative methods for the matrix sign function*, SIAM
952 Journal on Matrix Analysis and Applications, 12 (1991), pp. 273–291.
- 953 [23] C. S. KENNEY AND A. J. LAUB, *A Schur–Fréchet algorithm for computing the logarithm and*
954 *exponential of a matrix*, SIAM Journal on Matrix Analysis and Applications, 19 (1998),
955 pp. 640–663.
- 956 [24] P. KUNKEL AND V. MEHRMANN, *Smooth factorizations of matrix valued functions and their*
957 *derivatives*, Numerische Mathematik, 60 (1991), pp. 115–131.
- 958 [25] R.-C. LI, *Relative perturbation bounds for the unitary polar factor*, BIT Numerical Mathemat-
959 ics, 37 (1997), pp. 67–75.
- 960 [26] W. LI AND W. SUN, *New perturbation bounds for unitary polar factors*, SIAM Journal on
961 Matrix Analysis and Applications, 25 (2003), pp. 362–372.
- 962 [27] R. MATHIAS, *Evaluating the Fréchet derivative of the matrix exponential*, Numerische Mathe-
963 matik, 63 (1992), pp. 213–226.
- 964 [28] R. MATHIAS, *Perturbation bounds for the polar decomposition*, SIAM Journal on Matrix Anal-
965 ysis and Applications, 14 (1993), pp. 588–597.
- 966 [29] R. MATHIAS, *A chain rule for matrix functions and applications*, SIAM Journal on Matrix
967 Analysis and Applications, 17 (1996), pp. 610–620.
- 968 [30] I. NAJFELD AND T. F. HAVEL, *Derivatives of the matrix exponential and their computation*,
969 Advances in Applied Mathematics, 16 (1995), pp. 321–375.
- 970 [31] J. D. ROBERTS, *Linear model reduction and solution of the algebraic Riccati equation by use*
971 *of the sign function*, International Journal of Control, 32 (1980), pp. 677–687.