# A formula for the Frechet derivative of a generalized matrix function

Noferini, Vanni

2016

Manchester Institute for Mathematical Sciences
School of Mathematics

The University of Manchester

# A FORMULA FOR THE FRÉCHET DERIVATIVE OF A GENERALIZED MATRIX FUNCTION

VANNI NOFERINI[*]

**Abstract.** We state and prove an analogue of the Daleckiĭ-Kreĭn theorem, thus obtaining an explicit formula for the Fréchet derivative of generalized matrix functions. Moreover, we prove the differentiability of generalized matrix functions of real matrices under very mild assumptions. For complex matrices, we argue that, under the same assumptions, generalized matrix functions are real differentiable but generally not complex differentiable. Finally, we discuss the application of our results to the study of the condition number of generalized matrix functions. Along our way, we also derive generalized matrix functional analogues of a few classical theorems on polynomial interpolation of classical matrix functions and their derivatives.

**1. Introduction.** Matrix functions are a central subject in matrix theory and in numerical linear algebra [12, Chapter 9], [18], [20, Chapter 6]. There are several equivalent definitions of matrix functions, based on, for example, the Jordan canonical form, polynomial interpolation, or Cauchy integrals [18].

However, all these equivalent definitions can only be applied to square matrices. Linear algebraists have therefore considered possible extensions of the classical concept of a matrix function that allow for rectangular matrices as their argument. Hawkins and Ben-Israel introduced a definition based on the singular value decomposition, and developed some basic theory [14] . They forged the name "generalized matrix functions"[1] for their singular value-based definition, and showed that generalized matrix functions satisfy four of the so-called Fantappié properties [9, 10]. In other areas of mathematics, the study of generalized matrix functions has been called "singular value functional calculus" [1]. Recently, Arrigo, Benzi and Fenu explored the computational aspects of bilinear forms involving generalized matrix functions in the context of numerical linear algebra [4]. They introduced the notation $f^\diamond$, that we adopt in this paper, for generalized matrix functions. The reader may find in [1, Section 1] and in [4, Section 4] a survey of applications of generalized matrix functions, including complex network analysis, computer vision, finance, control system, computation of classical functions of large skew-symmetric matrices, solution of Hamiltonian differential systems, and filter factors.

An important part of the theory of classical matrix functions is devoted to the study of their Gâteaux and Fréchet derivatives. This has intrinsic theoretical interest, and has also relevant implications in numerical analysis, namely, it is important for the analysis of the condition number of a matrix function [18, Chapter 3]. In particular, a basic result in matrix theory is the Daleckiĭ-Kreĭn theorem [8], that we review in Section 2.4 and that gives a formula for the derivative of the classical matrix function of any diagonalizable matrix.

The main goal of this paper is to study the differentiability of generalized matrix functions, both developing a theoretical framework and analyzing the implications on

---
[*]Department of Mathematical Sciences, University of Essex, Wivenhoe Park, Colchester, UK, CO4 3SQ. (`vnofer@essex.ac.uk`)
[1]In spite of this name, generalized matrix functions do not generally reduce to classical matrix functions when the argument matrix is square [4, 18]. We adhere to the original terminology of [14], following also the more recent paper [4].

numerical conditioning. Our main result is a "generalized Daleckiĭ-Kreĭn theorem": an explicit formula for the derivative of a generalized matrix function $f^\diamond(A)$. Unlike the classical case, where a closed-form expression for the Fréchet derivative is not known for a generic function and a matrix with nontrivial Jordan form, our theorem holds in full generality. Among the applications of a formula for the derivative of generalized matrix functions is the study of their conditioning: we will discuss this matter in the present paper. More generally, our "generalized Daleckiĭ-Kreĭn theorem" may, at least potentially, be useful whenever there is an interest in studying how $f^\diamond(A)$ changes when $A$ is perturbed. This happens, for example, in complex network analysis [3].

The paper is structured as follows. Section 2 summarizes the mathematical background that we need: singular value decompositions, generalized matrix functions, Fréchet and Gâteaux derivatives of functions between Banach spaces, and the Daleckiĭ-Kreĭn theorem. Section 3 investigates the existence of the real Fréchet and Gâteaux derivatives of generalized matrix functions, shows that they are always equal to each other, and states and proves our main result: an explicit formula for them. We will also explain why, for complex matrices, generalized matrix functions are generally not complex-differentiable. Finally, Section 4 discusses the application of our results to the study of the condition number of generalized matrix functions.

## 2. Background.

**2.1. Singular value decompositions.** Let $A \in \mathbb{C}^{m \times n}$ have rank $r$, and throughout the paper we denote

$$\nu := \min\{m, n\}.$$

A singular value decomposition (SVD) [12] of $A$ is a factorization $A = USV^*$ such that $S \in \mathbb{R}^{m \times n}$ is diagonal, i.e., $S_{ij} = 0$ if $i \neq j$, and $U \in \mathbb{C}^{m \times m}$ and $V \in \mathbb{C}^{n \times n}$ are unitary, i.e., $UU^* = I_m$, $VV^* = I_n$. Moreover, the diagonal entries $S_{ii} = \sigma_i$ are called the singular values of $A$ and appear in nonincreasing order: $\sigma_1 \geq \cdots \geq \sigma_r > \sigma_{r+1} = \cdots = \sigma_\nu = 0$. The columns of $U$ and $V$ are called the left and right singular vectors of $A$, respectively. The matrix $S$ is uniquely determined by $A$, but there exist degrees of freedom in the choice of $U$ and $V$, which is why one speaks of "an SVD", rather than "the SVD". However, if $A \in \mathbb{R}^{m \times n}$, then $U$ and $V$ can always be chosen to be real and orthogonal, i.e., $UU^T = I_m$, $VV^T = I_n$, and we will always implicitly make this assumption whenever we refer to an SVD of a real matrix.

Following [4], given an SVD of $A$ we define the partial isometries $U_r \in \mathbb{C}^{m \times r}$ and $V_r \in \mathbb{C}^{n \times r}$ as the matrices whose columns are equal to the $r$ leftmost columns of $U$ and $V$, respectively, and $S_r \in \mathbb{R}^{r \times r}$ as the $r \times r$ top-left block of $S$. The resulting *compact SVD (CSVD)* of the matrix $A$ is the factorization

$$A = U_r S_r V_r^*,$$

whose existence can be immediately deduced from the SVD. For the definition of the CSVD to make sense when $r = 0$ and $U_r, S_r, V_r$ are empty matrices, we tacitly understand here (and throughout the paper) that, if $X \in \mathbb{C}^{m \times 0}$ and $Y \in \mathbb{C}^{0 \times n}$, then $XY = 0 \in \mathbb{C}^{m \times n}$.

**2.2. Generalized matrix functions.** In [4, 14] the following definition, based on the CSVD, is given, albeit in a slightly different form. Here and below, $[0, \infty)$ denotes the set of nonnegative real numbers, while $\mathcal{S}$ denotes an arbitrary, but fixed, subset of $[0, \infty)$.

DEFINITION 2.1. *Let $A \in \mathbb{C}^{m \times n}$ be a rank-$r$ matrix and let $A = U_r S_r V_r^*$ be a CSVD. Let $\mathcal{S} \subseteq [0, \infty)$ be such that $\sigma_i \in \mathcal{S}$ for all $i = 1, \ldots, r$ and $f : \mathcal{S} \to \mathbb{R}$ be a*

*scalar function.* The generalized matrix function $f^\diamond : \mathbb{C}^{m \times n} \to \mathbb{C}^{m \times n}$ *induced by* $f$ *is defined as*

$$f^\diamond(A) := U_r f(S_r) V_r^*,$$

*where* $f(S_r)$ *is the* $r \times r$ *diagonal matrix such that*

$$(f(S_r))_{ii} = f((S_r)_{ii}) = f(\sigma_i) \quad for \quad i = 1, \ldots, r.$$

REMARK 2.2. *Definition 2.1 is well posed, in the sense that it does not depend on the particular choice of an SVD (and hence of the resulting CSVD).*

*Indeed, suppose that* $A$ *has* $k$ *distinct singular values, denoted by*

$$\sigma_1 = \sigma_{i_1} > \sigma_{i_2} > \cdots > \sigma_{i_k} = \sigma_\nu,$$

*and suppose that*

$$A = \sum_{j=1}^{k} \sigma_{i_j} U_{(j)} V_{(j)}^*$$

*Here,* $U_{(j)}, V_{(j)}$ *are matrices with orthonormal columns spanning the* $j$*th left and right singular spaces. Note that this implies that the block matrices* $U_r = \begin{bmatrix} U_{(1)} & \ldots & U_{(\ell)} \end{bmatrix}$ *and* $V_r = \begin{bmatrix} V_{(1)} & \ldots & V_{(\ell)} \end{bmatrix}$ *appear in one arbitray, but fixed, CSVD* $A = U_r S_r V_r$. *Noting that* $\sigma_{i_j}$ *are all positive, with the one possible exception of* $\sigma_{i_k}$ *which is allowed to be zero, we conclude that the index* $\ell$ *is equal to either* $k$ *or* $k - 1$, *according to whether* $A$ *has full rank or not.*

*There are some degrees of freedom in the SVD (and hence in the CSVD). In particular, we may map* $U_{(j)} \mapsto U_{(j)} Q_j$ *and* $V_{(j)} \mapsto V_{(j)} Z_j$, *for* $j = 1, \ldots, \ell$ *and where* $Q_j, Z_j$ *are arbitrary unitary matrices of the appropriate size. Moreover,* $Q_j$ *and* $Z_j$ *cannot be chosen independently for nonzero singular values, namely,* $\sigma_{i_j} \neq 0 \Rightarrow Q_j = Z_j$. *The latter observation implies that* $f^\diamond(A)$ *does not depend on which particular CSVD one starts from: indeed, by Definition 2.1*

$$f^\diamond(A) = \sum_{j=1}^{\ell} f(\sigma_{i_j}) U_{(j)} Q_j Q_j^* V_{(j)} = \sum_{j=1}^{\ell} f(\sigma_{i_j}) U_{(j)} V_{(j)},$$

*where again* $\ell = k$ *if* $\sigma_\nu > 0$ *or* $\ell = k - 1$ *if* $\sigma_\nu = 0$.

*The argument above holds even when* $A$ *is rank deficient and* $f(0) \neq 0$. *The crucial observation is that the zero singular values are always mapped to 0 by* $f^\diamond(A)$, *regardless of the value of* $f(0)$.

REMARK 2.3. *The observations in Remark 2.2 imply that it holds* $\mathrm{rank}\, f^\diamond(A) \leq \mathrm{rank}\, A$, *with equality if and only if* $f(\sigma_i) \neq 0$ *for any singular value* $\sigma_i$ *of* $A$.

Observe also that, if we restrict the domain of $f^\diamond$ to real matrices, then clearly $f^\diamond$ maps $\mathbb{R}^{m \times n}$ to itself.

Several other elementary properties of generalized matrix functions are discussed in [4, 14, 18]; below we collect a few results that are useful to us, and we add some new observations of our own.

PROPOSITION 2.4. ([4, Proposition 3.2]) *For any* $A \in \mathbb{C}^{m \times n}$, *and any generalized matrix function* $f^\diamond$ *such that* $f^\diamond(A)$ *is defined:*

3

(i) $[f^\diamond(A)]^* = f^\diamond(A^*)$;

(ii) if $U_1 \in \mathbb{C}^{m \times m}$ and $U_2 \in \mathbb{C}^{n \times n}$ are unitary, then $f^\diamond(U_1 A U_2) = U_1 f^\diamond(A) U_2$.

THEOREM 2.5. ([4, Theorem 3.4]). *For any $A \in \mathbb{C}^{m \times n}$, and any generalized matrix function $f^\diamond$, induced by the scalar function $f$ and such that $f^\diamond(A)$ is defined, it holds*

$$f^\diamond(A) = f(\sqrt{AA^*})(\sqrt{AA^*})^\dagger A = A(\sqrt{A^*A})^\dagger f(\sqrt{A^*A}),$$

*where $X^\dagger$ denotes the Moore-Penrose pseudoinverse of the matrix $X$ [24] and $f(Y)$ denotes the classical matrix function [18], induced by the same scalar function $f$, of the matrix $Y$.*

We focus now on generalized polynomial functions $p^\diamond(A)$. Note that, even when $A$ is square, $p^\diamond(A)$ is not a polynomial in $A$ in the classical sense, but a generalized polynomial, whose explicit form is clarified in the next Remark.

REMARK 2.6. *[14] Let first $p = x^{2k+1}$. Then the induced generalized odd powers of $A$ are equal to $p^\diamond(A) = (AA^*)^k A = A(A^*A)^k$. If $p = x^{2k}$, given a CSVD $A = U_r S_r V_r^*$, then the generalized even powers of $A$ are $p^\diamond(A) = (AA^*)^k Q = Q(A^*A)^k$, where $Q = U_r V_r^*$. Formulae for a generic generalized polynomial $p^\diamond(A)$ can be obtained by linearity.*

Definition 2.1, being based on the CSVD, is advantageous for computational purposes. In this paper, we will sometimes find more convenient to use the next, equivalent, definition, based on the SVD.

DEFINITION 2.7. *Let $A \in \mathbb{C}^{m \times n}$ be a rank-r matrix and let $A = USV^*$ be an SVD. Let $\mathcal{S} \subseteq [0, \infty)$ be such that $\sigma_i \in \mathcal{S}$ for all $i = 1, \ldots, r$, and let $f : \mathcal{S} \to \mathbb{R}$. Then, we define the scalar function $f^\diamond(\sigma) : (\mathcal{S} \cup \{0\}) \to \mathbb{R}$ as*

$$f^\diamond(\sigma) = \begin{cases} f(\sigma) & \text{if } \sigma > 0; \\ 0 & \text{if } \sigma = 0. \end{cases} \tag{2.1}$$

*The generalized matrix function $f^\diamond : \mathbb{C}^{m \times n} \to \mathbb{C}^{m \times n}$ induced by $f$ is defined as*

$$f^\diamond(A) := U f^\diamond(S) V^*,$$

*where $f^\diamond(S)$ is defined as the $m \times n$ diagonal matrix such that*

$$(f^\diamond(S))_{ii} = f^\diamond(S_{ii}) = f^\diamond(\sigma_i) \quad \text{for} \quad i = 1, \ldots, \nu.$$

It is immediate that Definitions 2.1 and 2.7 are equivalent. Indeed, if $A = USV^*$ is an SVD and labelling by $U_i$ (resp. $V_i$) the $i$th column of $U$ (resp. $V$),

$$U_r f(S_r) V_r^* = \sum_{i=1}^{r} f(\sigma_i) U_i V_i^* = \sum_{i=1}^{\nu} f^\diamond(\sigma_i) U_i V_i^* = U f^\diamond(S) V^*.$$

A third characterization is also possible, as briefly mentioned in [18, Solution to Problem 1.53]. If $m \geq n$ and $A = QH$ is a polar decomposition [18, Chapter 8], Definition 2.7 and Proposition 2.4 yield $f^\diamond(A) = Q f^\diamond(H)$, where $f^\diamond(H)$ is the classical matrix function of $H$ induced by the scalar function (2.1). If either $A$ has full rank or $f(0) = 0$, then we have the stronger property $f^\diamond(A) = Q f(H)$. Note, however, that the latter statement is not true if $f(0) \neq 0$ and $A$ is rank deficient: indeed, in this

scenario $Qf(H)$ is not even uniquely defined – unlike $Qf^\diamond(H) = f^\diamond(A)$ – because of the nonuniqueness of $Q$.

Definition 2.7 makes it manifest that the scalar functions of the form (2.1) cannot be continuous at 0 unless they are induced by a continuous function $f$ satisfying $f(0) = 0$. Generalized matrix functions are built upon the modified scalar functions (2.1), and hence the same observation holds for rank deficient matrices.

REMARK 2.8. *Suppose that $0 \in S$ but $f(0) \neq 0$. If $A$ does not have full rank, i.e., if $\operatorname{rank} A < \nu$, then $f^\diamond(X)$ is not continuous (let alone differentiable) at $X = A$.*

EXAMPLE 2.9. *Suppose that $S$ contains a right neighbourhood of $0$. For $t > 0$ let $A(t) = \begin{bmatrix} 0 & t \end{bmatrix}$. It is easy to show that $f^\diamond(A(t)) = \begin{bmatrix} 0 & f(t) \end{bmatrix}$. Therefore,*

$$\lim_{t \to 0} f^\diamond(A(t)) = f^\diamond(A(0)) = \begin{bmatrix} 0 & 0 \end{bmatrix} \iff f(0) = 0.$$

**2.3. Fréchet derivatives, Gâteaux derivatives, and their relation.** In this subsection we review some basic notions in functional analysis. A more detailed treatment can be found, e.g., in [23], or in [25] for the finite dimensional case.

Suppose that $X, Y$ are Banach spaces and let $f : X \to Y$. Then $f$ is said to be Fréchet differentiable at $x \in X$ if there exists a bounded linear map $L_f(x, \cdot)$ such that

$$\lim_{\|h\|_X \to 0} \frac{\|f(x+h) - f(x) - L_f(x, h)\|_Y}{\|h\|_X} = 0 \qquad \forall\, h \in X.$$

The map $L_f$ is required to be real-linear for real Banach spaces and complex-linear for complex Banach spaces. Since any complex Banach space is also a real Banach space, it is possible for a function $f$ defined on a complex Banach space to be Fréchet real-differentiable but not Fréchet complex-differentiable.

Under the assumptions above, $L_f(x, \cdot)$ is called the Fréchet derivative of $f$ at $x$. It is, by definition, linear in $h$. When it exists, the Fréchet derivative is equal to the Gâteaux derivative, defined as[2]

$$G_f(x, h) = \lim_{t \to 0} \frac{f(x+th) - f(x)}{t}$$

where $t \in \mathbb{R}$ for real Banach spaces, and $t \in \mathbb{C}$ for complex Banach spaces.

The existence of the Gâteaux derivative alone does not imply Fréchet differentiability. However, if the Gâteaux derivative exists, additional sufficient conditions are known that imply that $f$ is Fréchet differentiable and the two derivatives coincide, for instance: (i) the Gâteaux derivative is linear in $h$, and is continuous in $x$ [18, Chapter 3], or (ii) $f$ is jointly (in $x$ and $h$) continuously Gâteaux differentiable [13, Section 3], or (iii) $X$ is finite dimensional, $f$ is Lipschitz continuous, and the Gâteaux derivative is linear in $h$ [2, Proposition A.4].

EXAMPLE 2.10. *Let $f : \mathbb{R}^2 \to \mathbb{R}, (x, y) \mapsto \frac{x^3}{x^2 + y^2}$ if $(x, y) \neq (0, 0)$ and $f(0, 0) = 0$. Then $f$ is Gâteaux differentiable for any $(x, y) \in \mathbb{R}^2$, with*

$$G_f((x, y), (h_x, h_y)) = \begin{cases} (x^2 + y^2)^{-2}(-2xy h_y + (x^2 + 3y^2)h_x) \ for \ (x, y) \neq (0, 0) \\ f(h_x, h_y) \ for \ (x, y) = (0, 0). \end{cases}$$

---

[2]Some authors require that the Gâteaux derivative is linear in $h$, using the term "Gâteaux differential" if this condition is dropped. We do not insist on linearity, but, as a warning against potential confusion, we note that both customs are common in the literature.

*Hence $f$ is Fréchet differentiable for $(x, y) \neq (0, 0)$, with $L_f(x, h) = G_f(x, h)$, but it is not Fréchet differentiable at $(0, 0)$.*

In the following, we will take $X = Y = \mathbb{R}^{m \times n}$ or $\mathbb{C}^{m \times n}$ (the latter seen as a *real* Banach space, for technical reasons to be discussed later on). In particular, $X$ will always be a finite dimensional Banach space, so that any linear map defined on $X$ is necessarily bounded, and the definition of the Fréchet derivative can be slightly simplified accordingly.

**2.4. The Daleckiĭ-Kreĭn theorem.** Suppose that a square matrix $A \in \mathbb{C}^{n \times n}$ is diagonalizable by similarity, i.e., that there exists an invertible matrix $Z \in \mathbb{C}^{n \times n}$ such that $A = ZDZ^{-1}$ and $D$ is diagonal. Then, given a scalar function $f$ defined on the eigenvalues of $A$, the classical matrix function $f(D)$ is the diagonal matrix satisfying $(f(D))_{ii} = f(D_{ii})$, and $f(A)$ is defined as $f(A) = Zf(D)Z^{-1}$ (one can check that the definition is well posed in the sense that it does not depend on the choice of $Z$). This definition can be extended to any square matrix, including those whose Jordan canonical form is not diagonal. The details can be found in classical references such as [12, 18, 20].

Since $\mathbb{C}^{n \times n}$ is a Banach space, it makes sense to study the Fréchet derivative of the classical matrix function $f(A)$. Besides the intrinsic theoretical interest, the main application is the study of the condition number of matrix functions, see [18, Chapter 3]. An explicit formula is known for diagonalizable matrices, and was first formulated by Daleckiĭ and Kreĭn. We recall [19] that the Schur (or Hadamard) product of two matrices $A, B \in \mathbb{C}^{m \times n}$ is denoted by $(A \circ B) \in \mathbb{C}^{m \times n}$ and it is defined entrywise as $(A \circ B)_{ij} = A_{ij}B_{ij}$.

THEOREM 2.11. [**Daleckiĭ–Kreĭn Theorem**][8]. *Let $A = ZDZ^{-1} \in \mathbb{C}^{n \times n}$ be a diagonalizable matrix, with $D$ diagonal, and let $f$ be continuously differentiable on the spectrum of $A$. Then the Fréchet derivative of the classical matrix function $f(X)$ at $X = A$, applied to the perturbation $E$, is equal to*

$$L_f(A, E) = Z(F \circ (Z^{-1}EZ))Z^{-1}$$

*where the symbol $\circ$ denotes the Schur product and the matrix $F \in \mathbb{C}^{n \times n}$ is defined as*

$$F_{ij} = \frac{f(D_{ii}) - f(D_{jj})}{D_{ii} - D_{jj}} \ \ if \ \ D_{ii} \neq D_{jj}, \qquad F_{ij} = f'(D_{ii}) \ \ otherwise.$$

**3. Main result.** The Daleckiĭ-Kreĭn theorem only applies to diagonalizable (by similarity) square matrices. In this section, we derive an analogous result, valid for *any matrix*, either square or not, and generalized matrix functions. Namely, we first prove the Gâteaux and Fréchet real-differentiability, under suitable assumptions, of generalized matrix functions. Then, we give explicit formulae for the derivatives.

**3.1. Existence of the real Gâteaux and Fréchet derivatives.** Generalized matrix functions of a complex matrix $A$ are generally[3] not complex-differentiable (neither in the Gâteaux nor in the Fréchet sense), not even in the scalar case.

EXAMPLE 3.1. *Let us compute $f^\diamond(\rho + z) - f^\diamond(\rho)$ for $0 \neq \rho \in \mathbb{C}$ and $z = \epsilon e^{\iota \zeta} \in \mathbb{C}$. By item (ii) in Proposition 2.4, without loss of generality we may take $\rho$ real and positive. Defining*

$$\mu(\epsilon, \zeta) = \sqrt{\rho^2 + \epsilon^2 + 2\rho\epsilon \cos \zeta}, \qquad \theta(\epsilon, \zeta) = \arctan \frac{\epsilon \sin \zeta}{\rho + \epsilon \cos \zeta}$$

---

[3]That is, except a few trivial exceptions, e.g., linear functions.

*we get (assuming $\epsilon < \rho$)*

$$f^\diamond(\rho + z) - f^\diamond(\rho) = e^{\iota\theta(\epsilon,\zeta)} f(\mu(\epsilon,\zeta)) - f(\rho),$$

*and expanding in a power series in $\epsilon$,*

$$f^\diamond(\rho + z) - f^\diamond(\rho) = \epsilon\left(f'(\rho)\cos\zeta + \iota\frac{f(\rho)}{\rho}\sin\zeta\right) + O(\epsilon^2).$$

*This shows that, for a generic $f$, the complex Gâteaux derivative of the generalized matrix function $f^\diamond$ does not exist. Indeed, unless $\rho f'(\rho) = f(\rho)$, letting $\epsilon \to 0^+$ with $\zeta = $const. in $z = \epsilon e^{\iota\zeta}$ yields different results of the limit*

$$\lim_{z \to 0} \frac{f^\diamond(\rho + z) - f^\diamond(\rho)}{z}$$

*depending on $\zeta$.*

Therefore, from now on we will always consider *real* derivatives of $f^\diamond(A)$, even when $A$ is complex.

Let us start by considering Gâteaux differentiability.

THEOREM 3.2. [**Gâteaux real-differentiability of generalized matrix functions**] *Let $A \in \mathbb{C}^{m\times n}$ and let $f : \mathcal{S} \to \mathbb{R}$ be differentiable on an open subset of $\mathcal{S}$ containing the positive singular values of $A$. Moreover, if $A$ is rank deficient, i.e., if rank $A < \nu$, suppose further that $f(0) = 0$ and that $f$ is right differentiable at $0$. Then $f^\diamond(X)$, defined as in Definitions 2.1 or 2.7, is Gâteaux differentiable at $X = A$.*

*Proof.* Recall that [6, Theorem 1] any real-valued real-analytic matrix function admits an analytic SVD. Note that, although [6, Theorem 1] only states this result for a real-valued real-analytic matrix function $A(t)$, its proof relies on specializing the analysis of [21, Section II.6.2] to the matrix $\begin{bmatrix} 0 & A(t) \\ A(t)^T & 0 \end{bmatrix}$. However, the theory of [21] applies to any matrix-valued function, possibly complex, which is Hermitian for any value of the *real* parameter $t$. Therefore, slightly modifying the proof of [6, Theorem 1] by considering instead the matrix $\begin{bmatrix} 0 & A(t) \\ A(t)^* & 0 \end{bmatrix}$, we see that every real-analytic matrix-valued function (possibly complex) admits an analytic SVD. In particular, letting $A(t) = A + tE$, it holds

$$A + tE = U(t)S(t)V(t)^* \tag{3.1}$$

where $U(t), V(t)$ are analytic and unitary and $S(t)$ is real, analytic and diagonal for all real $t$, and in particular in some neighbourhood of $0$. These facts immediately yield that the *real* Gâteaux derivative

$$G_{f^\diamond}(A, E) = \lim_{t \in \mathbb{R}, t \to 0} \frac{f^\diamond(A + tE) - f^\diamond(A)}{t} \tag{3.2}$$

exists provided that $f^\diamond(S(t))$ is differentiable at $t = 0$. The latter condition is satisfied if and only if the scalar function $f$ is differentiable on an open set containing the singular values of $A$, with the additional conditions that $f(0) = 0$ and that $f$ is right differentiable at $0$ if $A$ is not full rank. Indeed, expanding $U(t) = U_0 + tU_1 + O(t^2)$,

7

$V(t) = V_0 + tV_1 + O(t^2)$, $S(t) = S_0 + tS_1 + O(t^2)$, from (3.1) we get that $A = U_0 S_0 V_0^*$ is an SVD , that $U_0 U_1^* + U_1 U_0^* = 0 = V_0 V_1^* + V_1 V_0^*$, that $S_1$ is diagonal, and that

$$G_{f^\diamond}(A,E) = U_1 f^\diamond(S_0)V_0^* + U_0 f^\diamond(S_0)V_1^* + U_0 \left.\frac{df^\diamond(S(t))}{dt}\right|_{t=0} V_0^*, \qquad (3.3)$$

where by the chain rule

$$\left.\frac{df^\diamond(S(t))}{dt}\right|_{t=0} = (f')^\diamond(S_0) \circ S_1.$$

□

The computation of the Gâteaux derivative from (3.3) is, in principle, not impossible employing the sophisticated techniques of [6]; however, this may be very challenging in practice. We will give a much more explicit formula in Theorem 3.8.

If $A$ is not full rank and $f(0) \neq 0$, then by Remark 2.8 $f^\diamond(X)$ cannot be differentiable at $X = A$. If we assume $f(0) = 0$, then generalized matrix functions induced by a Lipschitz continuous function $f$ are Lipschitz continuous [1, Theorem 1.1]. Hence, by Rademacher's Theorem [15, Theorem 3.1], they must be Fréchet differentiable almost everywhere; yet, in principle, there might exist a measure zero subset of $\mathbb{C}^{m \times n}$ on which they are not.

To fill this gap, we follow a different approach based on polynomial interpolation. The following theorem is a generalized matrix functional analogue of [20, Theorem 6.6.14]. Its first part is new, while the second part was already mentioned (without giving details) in [14].

THEOREM 3.3. *Let $A \in \mathbb{C}^{m \times n}$ and let $f : \mathcal{S} \to \mathbb{R}$ be differentiable on an open subset of $\mathcal{S}$ containing the positive singular values of $A$. Moreover, if $A$ is rank deficient suppose further that $f(0) = 0$ and that $f$ is right differentiable at $0$. Let $A$ have $k$ distinct singular values, denoted by*

$$\sigma_1 = \sigma_{i_1} > \sigma_{i_2} > \cdots > \sigma_{i_k} = \sigma_\nu,$$

*and let $q$ be the unique polynomial of degree $2k - 1$ satisfying*

$$q(\sigma_{i_j}) = f(\sigma_{i_j}) \qquad and \qquad q'(\sigma_{i_j}) = f'(\sigma_{i_j}), \qquad j = 1, \dots, k.$$

*Then the real Gâteaux derivatives of $f^\diamond(X)$ and $q^\diamond(X)$ coincide at $X = A$:*

$$G_{f^\diamond}(A,E) = G_{q^\diamond}(A,E) \qquad \forall E \in \mathbb{C}^{m \times n}.$$

*Moreover, let $p$ be the unique polynomial of degree $h - 1$ satisfying*

$$p(\sigma_{i_j}) = f(\sigma_{i_j}), \qquad j = 1, \dots, h$$

*where $h = k$ if $A$ has full rank and $h = k - 1$ if $A$ is rank deficient. (If $h = 0$, set $p = 0$). Then $f^\diamond(A) = p^\diamond(A)$.*

*Proof.* The first part of the statement is a consequence of (3.3). Indeed, $U_0, U_1, V_0, V_1, S_0, S_1$ depend on $A$ and $E$, but not on $f$. On the other hand, the definition of $q$ guarantees that $f^\diamond(S_0) = q^\diamond(S_0)$ and that $(f')^\diamond(S_0) = (q')^\diamond(S_0)$.

The second part of the statement is straightforward from Definition 2.1. □

If $f^\diamond$ is the generalized matrix power induced by $f(x) = 1 = x^0$, then $f^\diamond(X)$ is differentiable at $X = A$ if and only if $A$ is full rank. In contrast, positive generalized matrix powers are always differentiable, as we next show.

LEMMA 3.4. *Let $f(x) = x^h$, $h = 1, 2, 3, \ldots$, and let $A \in \mathbb{C}^{m \times n}$. Then the generalized power matrix $f^\diamond(X)$, defined as in Definitions 2.1 or 2.7, is Fréchet real-differentiable at $X = A$.*

*Proof.* By Remark 2.6, if $h = 2k + 1$ is odd then $f^\diamond(A)$ is manifestly Fréchet differentiable: indeed, each entry of $(AA^*)^k A$ is a polynomial function of the entries of $A$. It remains to argue that the same is true for even and positive $h = 2k \geq 2$. By Remark 2.6, $f^\diamond(A) = (A^*A)^k U_r V_r^*$, where $U_r$ and $V_r$ are the partial isometries defining a CSVD of $A = U_r S_r V_r^*$. Suppose first that $r = m = n$, i.e., $A$ is square and nonsingular. Observe that

$$(AA^*)^k U_r V_r^* = (AA^*)^{k-1} A V_r S_r V_r^* =: (AA^*)^{k-1} AH,$$

where $H = V_r S_r V_r^*$ is the Hermitian factor of the polar decomposition of $A$ [18, Chapter 8], which is uniquely defined for an invertible square matrix [18, Theorem 8.1]. That $H$ is Gâteaux real-differentiable can be argued as in the proof of Theorem 3.2 via the existence of the analytic SVD (3.1). Indeed, observe that if $A(t) = A + tE = U(t)S(t)V(t)^*$ is an SVD then the Hermitian factor of any polar decomposition of $A(t)$ is $H(A(t)) = \sqrt{V(t)S(t)^T S(t) V(t)^*}$. Note also that $H$ is a Lipschitz continuous function of $A$ [18, Theorem 8.9]. Therefore, if we can show that the Gâteaux derivative is real-linear in $E$, it will follow that $H$ is a Fréchet real-differentiable function of $A$ [2, Proposition A.4]. Let $G_H(A, E)$ be the real Gâteaux derivative of $H$ as a function of $A$, applied to the direction $E$. Following [17, Proof of Theorem 2.5], one can show by differentiating the equation $H(t)^2 = A(t)^* A(t)$ and evaluating at $t = 0$ that it holds

$$HG_H(A, E) + G_H(A, E)H = A^* E + E^* A.$$

The latter is a Sylvester equation in the unknown $G_H(A, E)$, whose right hand side depends linearly on the real and imaginary parts $E$. Hence, it displays the real-linearity of $G_H(A, E)$ in $E$ (It is known [20, Section 4.4] that the Sylvester equation $AX - XB = C$ has a unique solution for any right hand side $C$ if and only if $A$ and $B$ have no eigenvalues in common. Note that, in our case, $A = -B = H$. On the other hand, by assumption, $H$ is Hermitian positive definite, thus ensuring existence and uniqueness of the solution of the Sylvester equation). This concludes the proof for a square and invertible $A$.

For a general $A$, we will give a proof assuming for simplicity of exposition that $m \geq n$: the case $n > m$ is similar, or it can be argued that it follows applying the argument to $A^*$ and invoking item (i) in Proposition 2.4. Let $A = QH$, $Q = U_r V_r^*$, $H = V_r S_r V_r^*$, be a polar decomposition. When $A$ does not have full rank, it is not any more true that the Hermitian factor $H$ is Fréchet real-differentiable at $A$. However, we will argue that $B = AH$ is, implying that

$$(AA^*)^k U_r V_r^* = (AA^*)^{k-1} B$$

is differentiable as well. Observe that $B = AH = f^\diamond(A)$ is the generalized matrix function of $A$ induced by the locally Lipschitz continuous (on any bounded interval containing the singular values of $A$) scalar function $f(x) = x^2$, that satisfies $f(0) = 0$. Hence, $B$ is locally Lipschitz continuous in some neighbourhood of $A$ by [1, Theorem 1.1], and in view of Theorem 3.2 it suffices to show that its real Gâteaux derivative

is real-linear in the perturbation $E$. Note that by item (ii) in Proposition 2.4 the generalized matrix function $B$ is differentiable at $A$ if and only if it is differentiable at $S$, where $A = USV^*$ is an SVD. Therefore, with no loss of generality, we may take $A$ to be of the fom

$$A = \begin{bmatrix} S_r & 0_{r \times (n-r)} \\ 0_{(m-r) \times r} & 0_{(m-r) \times (n-r)} \end{bmatrix} \Rightarrow B = \begin{bmatrix} S_r^2 & 0_{r \times (n-r)} \\ 0_{(m-r) \times r} & 0_{(m-r) \times (n-r)} \end{bmatrix}.$$

We will suppose that $E$ is partitioned coherently with $A$:

$$E = \begin{bmatrix} E_{11} & E_{12} \\ E_{21} & E_{22} \end{bmatrix}.$$

Letting $A(t) = A + tE$, suppose that $A(t) = U(t)S(t)V(t)^*$ is an analytic SVD (3.1). Partition

$$S(t) = \begin{bmatrix} S_r(t) & 0 \\ 0 & O(t) \end{bmatrix},$$

$$U(t) = \begin{bmatrix} U_{11}^{(0)} + tU_{11}^{(1)} & tU_{12} \\ tU_{21} & U_{22}^{(0)} + tU_{22}^{(1)} \end{bmatrix} + O(t)^2, \qquad V(t) = \begin{bmatrix} V_{11}^{(0)} + tV_{11}^{(1)} & tV_{12} \\ tV_{21} & V_{22}^{(0)} + tV_{22}^{(1)} \end{bmatrix} + O(t^2),$$

where the top-left blocks are all $r \times r$ and the fact that the off-diagonal blocks of $U(t)$ and $V(t)$, as well as the bottom-right block of $S(t)$, are 0 at $t = 0$ is a consequence of the zero pattern of $A(0) = A$. Observe that, if $f^\diamond$ is the generalized matrix function induced by $f(x) = x^2$, one has

$$f^\diamond(S(t)) = \begin{bmatrix} S_r^2(t) & 0 \\ 0 & O(t^2) \end{bmatrix}.$$

Hence, for $B(t) = U(t)f^\diamond(S(t))V(t)^*$, and with $B = B(0)$, we obtain

$$B(t) = B + t \begin{bmatrix} X & U_{11}^{(0)} S_r^2 V_{21}^* \\ U_{21} S_r^2 (V_{11}^{(0)})^* & 0 \end{bmatrix} + O(t^2), \tag{3.4}$$

where $X := U_{11}^{(0)}(S_r^2)'(0)(V_{11}^{(0)})^* + U_{11}^{(0)} S_r^2 (V_{11}^{(1)})^* + U_{11}^{(1)} S_r^2 (V_{11}^{(0)})^*$. We deduce that the real Gâteaux derivative of $B$ at $A$, applied to the perturbation $E$, has the form

$$G_B(A, E) = \begin{bmatrix} X & Y \\ Z & 0 \end{bmatrix} \in \mathbb{C}^{m \times n},$$

and by (3.4) it is immediate to check that $X \in \mathbb{C}^{r \times r}$ is precisely $G_B(S_r, E_{11})$, i.e., the real Gâteaux derivative of the generalized matrix function $f^\diamond$ induced by $f(x) = x^2$ (but seen as a function defined on $\mathbb{C}^{r \times r}$ rather than on $\mathbb{C}^{m \times n}$ as elsewhere in this proof) at $S_r$, applied to the perturbation $E_{11}$. Since $S_r$ is square and invertible, by the first part of the proof $X$ is also a real Fréchet derivative, and therefore it is real-linear in $E_{11}$, and hence, in $E$.

Now, differentiating the equations $B(t)B(t)^* = (A(t)A(t)^*)^2$ and $B(t)^*B(t) = (A(t)^*A(t))^2$ and evaluating them at $t = 0$ we obtain, respectively,

$$BG_B(A, E)^* + G_B(A, E)B^* = EA^*AA^* + AE^*AA^* + AA^*EA^* + AA^*AE^*$$

10

and

$$G_B(A, E)^* B + B^* G_B(A, E) = E^* AA^* A + A^* EA^* A + A^* AE^* A + A^* AA^* E.$$

Computing the $(2, 1)$ block of the first equation yields $ZS_r^2 = E_{21} S_r^3$, while from the $(1, 2)$ block of the second equation we get $S_r^2 Y = S_r^3 E_{12}$. Hence, $Y$ and $Z$ are also both linear (and hence real-linear) in $E$, and this concludes the proof. $\square$

The next theorem is the main result of this subsection.

THEOREM 3.5. [**Fréchet real-differentiability of generalized matrix functions**] *Let $A \in \mathbb{C}^{m \times n}$ and let $f : \mathcal{S} \to \mathbb{R}$ be differentiable on an open subset of $\mathcal{S}$ containing the positive singular values of $A$. Moreover, if $A$ is rank deficient, i.e., if $\operatorname{rank} A < \nu$, suppose further that $f(0) = 0$ and that $f$ is right differentiable at $0$. Then $f^\diamond(X)$, defined as in Definitions 2.1 or 2.7, is Fréchet real-differentiable at $X = A$.*

*Proof.* If $A$ has full rank then the statement follows by Theorem 2.5 and by standard results on the differentiability of standard matrix functions [18, Chapter 3] and pseudoinverses of real full rank matrices [11]. To deal with the case of a rank deficient $A$, we may suppose that $f(0) = 0$. By Theorem 3.3, it suffices to prove the statement for a polynomial of the form (note that the trailing coefficient is 0 by assumption)

$$f(x) = \sum_{i=1}^{\kappa} f_i x^i.$$

The statement then follows by linearity from Lemma 3.4. $\square$

COROLLARY 3.6. *Under the assumptions of Theorem 3.5, the real Gâteaux and Fréchet derivatives coincide, and are real-linear in $E$.*

REMARK 3.7. *Corollary 3.6 has very useful pratical consequences, as Gâteaux derivatives are easy to compute. More in detail, if we have a basis (over the real field) $E_i$, $i = 1, \ldots, 2mn$, of $\mathbb{C}^{m \times n}$, and if we can compute the $2mn$ Gâteaux derivatives $G_{f^\diamond}(A, E_i)$, then for $E = \sum_i \alpha_i E_i$ we can obtain $L_{f^\diamond}(A, E) = \sum_i \alpha_i G_{f^\diamond}(A, E_i)$. This property is crucial for the proof of Theorem 3.8.*

*Note also that, if $f$ is continuously differentiable on the singular values of $A$, it is easy to show the (real) Gâteaux and Fréchet derivatives are also continuous in $A$. Hence, if $f$ is a continuously differentiable function and we have a converging sequence $A_n \to A$, then we can compute $L_{f^\diamond}(A_n, E)$, then we can obtain $L_{f^\diamond}(A, E) = \lim_{n \to \infty} L_{f^\diamond}(A_n, E)$.*

In summary, under mild assumptions on the underlying scalar function $f$, generalized matrix functions on $\mathbb{C}^{m \times n}$ are real-differentiable, but not complex-differentiable. The only way around this obstacle is to see $\mathbb{C}^{m \times n}$ as a real Banach space of dimension $2mn$. We now turn to an explicit formula for the real Fréchet derivative of complex generalized matrix functions in this context. As a special case, we will also recover the Fréchet derivative of real generalized matrix functions on $\mathbb{R}^{m \times n}$.

**3.2. Explicit formulae for the derivative.** The following theorem is our main result, and it gives an explicit formula for the real Fréchet derivative of a generalized matrix function $f^\diamond(X) : \mathbb{C}^{m \times n} \to \mathbb{C}^{m \times n}$.

THEOREM 3.8. [**Daleckiĭ-Kreĭn Theorem for generalized matrix functions**]
*Let $A = USV^*$ be an SVD of $A \in \mathbb{C}^{m \times n}$, where $U \in \mathbb{C}^{m \times m}$, $V \in \mathbb{C}^{n \times n}$, $S \in \mathbb{R}^{m \times n}$, and $S_{ii} =: \sigma_i$, $i = 1, \ldots, \nu$, are the singular values of $A$. Let $f : \mathcal{S} \to \mathbb{R}$ be differentiable on an open subset of $\mathcal{S}$ containing the positive singular values of $A$. Moreover,*

*if $A$ is rank deficient, i.e., if $\operatorname{rank} A < \nu$, suppose further that $f(0) = 0$ and that $f$ is right differentiable at $0$. Then, if we see $\mathbb{C}^{m \times n}$ as a $2mn$-dimensional real vector space, the real Fréchet derivative at $X = A$ of the generalized matrix function $f^\diamond(X)$, applied to the complex perturbation $E$, is*

$$L_{f^\diamond}(A, E) = U\left(F \circ \Re(\widehat{E}) + \iota H \circ \Im(\widehat{E}) + G \circ \Upsilon(\widehat{E})\right) V^*, \qquad (3.5)$$

*where*

- *$\iota$ is the imaginary unit;*
- *the symbol $\circ$ denotes the Schur product;*
- *$\widehat{E} = U^* E V$, $\Re(\widehat{E})$ is its real part, and $\Im(\widehat{E})$ is its imaginary part;*
- *the real-linear operator $\Upsilon$ is the following generalization of the conjugate transposition operator: for any $X \in \mathbb{C}^{m \times n}$, $\Upsilon(X) \in \mathbb{C}^{m \times n}$ and*

$$\text{if } m = n, \qquad \Upsilon(X) = X^*;$$

$$\text{if } m > n \text{ and } X = \begin{bmatrix} X_1 \\ X_2 \end{bmatrix}, X_1 \in \mathbb{C}^{n \times n}, \qquad \Upsilon(X) = \begin{bmatrix} X_1^* \\ X_2 \end{bmatrix};$$

$$\text{if } m < n \text{ and } X = \begin{bmatrix} X_1 & X_2 \end{bmatrix}, X_1 \in \mathbb{C}^{m \times m}, \qquad \Upsilon(X) = \begin{bmatrix} X_1^* & X_2 \end{bmatrix};$$

- *$F, G \in \mathbb{R}^{m \times n}$ are defined as follows:*

$$F_{ij} = \begin{cases} \frac{\sigma_i f(\sigma_i) - \sigma_j f(\sigma_j)}{\sigma_i^2 - \sigma_j^2} & \text{if } i \neq j,\ \max(i,j) \leq \nu,\ \text{and } \sigma_i \neq \sigma_j; \\ \frac{\sigma_i f'(\sigma_i) + f(\sigma_i)}{2\sigma_i} & \text{if } i \neq j\ ,\ \max(i,j) \leq \nu,\ \text{and } \sigma_i = \sigma_j \neq 0; \\ \frac{f(\sigma_j)}{\sigma_j} & \text{if } i > n,\ \text{and } \sigma_j \neq 0; \\ \frac{f(\sigma_i)}{\sigma_i} & \text{if } j > m,\ \text{and } \sigma_i \neq 0; \\ f'(\sigma_i) & \text{otherwise;} \end{cases} \qquad (3.6)$$

$$G_{ij} = \begin{cases} \frac{\sigma_j f(\sigma_i) - \sigma_i f(\sigma_j)}{\sigma_i^2 - \sigma_j^2} & \text{if } i \neq j,\ i, j \leq \nu,\ \text{and } \sigma_i \neq \sigma_j; \\ \frac{\sigma_i f'(\sigma_i) - f(\sigma_i)}{2\sigma_i} & \text{if } i \neq j,\ i, j \leq \nu,\ \text{and } \sigma_i = \sigma_j \neq 0; \\ 0 & \text{otherwise.} \end{cases} \qquad (3.7)$$

- *and $H \in \mathbb{R}^{m \times n}$ is such that with $H_{ii} = f(\sigma_i)/\sigma_i$ if $\sigma_i \neq 0$, $H_{ii} = f'(0)$ (the right derivative of $f$ at $x = 0$) if $\sigma_i = 0$, whereas $H_{ij} = F_{ij}$ for $i \neq j$.*

*Proof.* Item (ii) in Proposition 2.4 yields

$$f^\diamond(A + tE) = U f(S + tU^* E V) V^* \simeq f^\diamond(A) + t U L_f(S, U^* E V) V^*,$$

where the last approximate equality is exact up to additive terms of higher order in $tE$. Therefore, without loss of generality we can assume that $A$ has zero off-diagonal elements and real positive diagonal elements.

The strategy of the proof is to first prove the result when $E$ is zero except for one element, equal to either $1$ or $\iota$. Using Corollary 3.6, we will compute $L_f(A, E)$ for such an $E$ as the limit at the right hand side of (3.2). The result for a general $E$ will then follow by linearity (over the real field). We now examine a few separate cases according to the value and the exact position of the unique nonzero element of $E$.

We assume first that the nonzero element of $E$ is equal to $1$. There are three cases:

- *Case 1.* If the unique nonzero element of $E$ is its $i$th diagonal element, then $f^\diamond(A + tE) - f^\diamond(A) = \mathrm{diag}(0, \ldots, 0, f(\sigma_i + t) - f(\sigma_i), 0, \ldots, 0)$, where the nonzero element appears in the $i$th position. Dividing by $t$ and going to the limit $t \to 0$ we obtain $\mathrm{diag}(0, \ldots, 0, f'(\sigma_i), 0, \ldots, 0)$, thus proving the theorem in this case.

- *Case 2a.* Suppose now that the unique nonzero element of $E$ is in the position $(i, j)$ with $i < j \leq \nu$. In this case, $A + tE$ is not diagonal, and hence, we need to compute its singular value decomposition to estimate $f^\diamond(A + tE)$. Let us take, without loss of generality (modulo applying a permutation equivalence), $i = 1$, $j = 2$. Then, $A + tE = (U' \oplus I_{m-2})S'(V' \oplus I_{n-2})^*$ is a singular value decomposition if

$$(U')^* \begin{bmatrix} \sigma_i & t \\ 0 & \sigma_j \end{bmatrix} V' \tag{3.8}$$

is real and diagonal and $U', V'$ are unitary.

We now need to distinguish three subcases.

*Subcase 2a-(i).* Assume that $\sigma_i > \sigma_j$. Then, to compute $U'$ and $V'$, let us expand them as $U' = I_2 + t \begin{bmatrix} 0 & -u^* \\ u & 0 \end{bmatrix} + O(t^2)$ and $V' = I_2 + t \begin{bmatrix} 0 & -v^* \\ v & 0 \end{bmatrix} + O(t^2)$, observing that, at the identity matrix, the tangent space to the smooth manifold of unitary matrices is the subspace of skew-Hermitian matrices, as can be easily seen by differentiating the equation $XX^* = I_n$ and evaluating at $X = I_n$ (see also [26, Section 5.4] and [5]). Imposing that (3.8) is diagonal and retaining only the $O(t)$ terms leads to the linear system

$$\begin{bmatrix} -\sigma_j & \sigma_i \\ -\sigma_i & \sigma_j \end{bmatrix} \begin{bmatrix} u \\ v \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \end{bmatrix},$$

which for $\sigma_i \neq \sigma_j$ yields

$$u = \frac{\sigma_j}{\sigma_i^2 - \sigma_j^2}, \qquad v = \frac{\sigma_i}{\sigma_i^2 - \sigma_j^2}.$$

Moreover, with this choice of $u$ and $v$ we have $S' = S + O(t^2)$. At this point, observe that $f^\diamond(A + tE) = (U' \oplus I_{n-2})f(S')(V' \oplus I_{n-2})^*$, and hence, by a direct computation,

$$f^\diamond(A + tE) - f^\diamond(A) = t\left( \begin{bmatrix} 0 & \alpha \\ \beta & 0 \end{bmatrix} \oplus 0_{(m-2)\times(n-2)} \right) + O(t^2),$$

with

$$\alpha = \frac{\sigma_i f(\sigma_i) - \sigma_j f(\sigma_j)}{\sigma_i^2 - \sigma_j^2}$$

and

$$\beta = \frac{\sigma_j f(\sigma_i) - \sigma_i f(\sigma_j)}{\sigma_i^2 - \sigma_j^2}.$$

*Subcase 2a-(ii).* The argument given in Subcase 2a-(i) clearly fails when $\sigma_i = \sigma_j$. However, in this case there are more degrees of freedom in the expansion

13

of $U'$ and $V'$ in (3.8). Indeed, we may have $U' = (I_2 + t \begin{bmatrix} 0 & -u^* \\ u & 0 \end{bmatrix})Q$,

$V' = (I_2 + t \begin{bmatrix} 0 & -v^* \\ v & 0 \end{bmatrix})Q$ for *any* $2 \times 2$ unitary matrix $Q$. (Here, we are using the fact that a matrix in the tangent space at $X = Q$ of the smooth manifold of unitary matrices can always be written as a skew-Hermitian matrix times $Q$, as can be easily seen by differentiating the equation $XX^* = I_2$ and evaluating at $X = Q$; see also [5].) Let us now suppose $\sigma_i = \sigma_j =: \sigma > 0$. We will show that a solution can always be found for $Q = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 & -1 \\ 1 & 1 \end{bmatrix}$. Indeed, imposing that (3.8) is diagonal and focusing only on the $O(t)$ terms we obtain the condition

$$v - u = \frac{1}{2\sigma}.$$

Taking, for example, $u = 0$ and $v = (2\sigma)^{-1}$, we get in particular

$$(U')^* \begin{bmatrix} \sigma & t \\ 0 & \sigma \end{bmatrix} V' = \begin{bmatrix} \sigma + t/2 & 0 \\ 0 & \sigma - t/2 \end{bmatrix} + O(t^2).$$

Define now

$$f_+ = \frac{f(\sigma + t/2) + f(\sigma - t/2)}{2}, \qquad f_- = \frac{f(\sigma + t/2) - f(\sigma - t/2)}{2}.$$

It follows that

$$f^\diamond(A + tE) - f^\diamond(A) = \left( \begin{bmatrix} \gamma & \alpha \\ \beta & \delta \end{bmatrix} \oplus 0_{(m-2) \times (n-2)} \right) + O(t^2),$$

for

$$\alpha = f_- + \frac{t f_+}{2\sigma}, \qquad \beta = f_- - \frac{t f_+}{2\sigma},$$

$$\gamma = f_+ - f(\sigma) - \frac{t f_-}{2\sigma}, \qquad \delta = f_+ - f(\sigma) + \frac{t f_-}{2\sigma}.$$

It now suffices to observe that

$$\lim_{t \to 0} \frac{f_+ - f(\sigma)}{t} = \frac{1}{4}(f'(\sigma) - f'(\sigma)) = 0, \qquad \lim_{t \to 0} \frac{f_-}{t} = \frac{1}{4}(f'(\sigma) + f'(\sigma)) = \frac{f'(\sigma)}{2},$$

$$\lim_{t \to 0} f_- = \frac{1}{2}(f(\sigma) - f(\sigma)) = 0, \qquad \lim_{t \to 0} f_+ = \frac{1}{2}(f(\sigma) + f(\sigma)) = f(\sigma),$$

yielding

$$\lim_{t \to 0} \frac{\alpha}{t} = \frac{\sigma f'(\sigma) + f(\sigma)}{2\sigma}, \qquad \lim_{t \to 0} \frac{\beta}{t} = \frac{\sigma f'(\sigma) - f(\sigma)}{2\sigma}, \qquad \lim_{t \to 0} \frac{\gamma}{t} = \lim_{t \to 0} \frac{\delta}{t} = 0.$$

14

*Subcase 2a-(iii).* It remains to discuss the case $\sigma_i = \sigma_j = 0$. It is immediate to see that

$$f^\diamond(A + tE) - f^\diamond(A) = \begin{bmatrix} 0 & f(t) \\ 0 & 0 \end{bmatrix} \oplus 0_{(m-2)\times(n-2)}.$$

Dividing by $t$ and going to the limit $t \to 0$ yields the statement.

- *Case 2b.* Consider now the case where the unique nonzero element of $E$ lies in the position $(i,j)$ with $m < j \le n$. Again, $A + tE$ is not diagonal, and, similarly to Case 2a, we need first to compute its singular value decomposition. We may assume that $i = 1$, $j = m+1$. Observe that $A + tE = S'(V')^T$ is a singular value decomposition if

$$\begin{bmatrix} \sigma_i & 0 & \ldots & 0 & t & 0 & \ldots & 0 \end{bmatrix} V' = \begin{bmatrix} \sigma & 0 & \ldots & 0 \end{bmatrix},$$

for some $\sigma \ge 0$. As before we can expand $V'$ in powers of $t$. This procedure yields $\sigma = \sigma_i$, $V'_{ii} = 1$ for all $i = 1, \ldots, n$, $V'_{1,n+1} = -t/\sigma_i$, $V'_{n+1,1} = t/\sigma_i$, and $V_{ij} = 0$ in all other cases. Hence, to first order in $t$, there is only one nonzero element in $f^\diamond(A + tE) - f^\diamond(A)$, lying precisely at the position $(i,j)$, and being equal to $tf(\sigma_i)/\sigma_i$.

- *Case 3.* If the unique nonzero element of $E$ is in the position $(i,j)$ with $j < i$, the proof is either analogous to Case 2a if $i \le m$ or to Case 2b if $i > m$. We omit the details.

Now, we turn to the case when the unique nonzero element in $E$ is pure imaginary and equal to $\iota$. Again, there are three cases:

- *Case 4.* Suppose first $\sigma_i \ne 0$. If the unique nonzero element of $E$ is its $i$th diagonal element, then arguing as in Example 3.1 we readily obtain $f^\diamond(A + tE) - f^\diamond(A) = \mathrm{diag}(0, \ldots, 0, \iota t f(\sigma_i)/\sigma_i, 0, \ldots, 0)$, where the nonzero element appears in the $i$th position. Dividing by $t$ and going to the limit, we obtain $\iota \, \mathrm{diag}(0, \ldots, 0, f(\sigma_i)/\sigma_i, 0, \ldots, 0)$.
  Similarly, if $\sigma_i = 0$, we get $f^\diamond(A + tE) - f^\diamond(A) = \mathrm{diag}(0, \ldots, 0, \iota f(t), 0, \ldots, 0)$, and dividing by $t$ and letting $t \to 0$ this yields $\iota \, \mathrm{diag}(0, \ldots, 0, f'(0), 0, \ldots, 0)$ where $f'(0)$ is the right derivative of $f$ at $x = 0$.

- *Case 5a.* Suppose now that the unique nonzero element of $E$ is in the position $(i,j)$ with $i < j \le \nu$. As in Case 2a we can take without loss of generality $i = 1$, $j = 2$, and we need to distinguish three subcases. However, this time we impose that

$$(U')^* \begin{bmatrix} \sigma_i & \iota t \\ 0 & \sigma_j \end{bmatrix} V' \tag{3.9}$$

is real and diagonal and $U', V'$ are unitary.

*Subcase 5a-(i).* Suppose $\sigma_i > \sigma_j$. We can expand $U'$ and $V'$ as in Subcase 2a-(i). Retaining only the $O(t)$ terms and solving for $u, v$ we get

$$u = \frac{\iota \sigma_j}{\sigma_i^2 - \sigma_j^2}, \qquad v = \frac{\iota \sigma_i}{\sigma_i^2 - \sigma_j^2}$$

and $S' = S + O(t^2)$. Computing $f^\diamond(A + tE) = (U' \oplus I_{n-2})f(S')(V' \oplus I_{n-2})^*$ yields

$$f^\diamond(A + tE) - f^\diamond(A) = \iota t \left( \begin{bmatrix} 0 & \alpha \\ -\beta & 0 \end{bmatrix} \oplus 0_{(m-2)\times(n-2)} \right) + O(t^2),$$

15

with

$$\alpha = \frac{\sigma_i f(\sigma_i) - \sigma_j f(\sigma_j)}{\sigma_i^2 - \sigma_j^2}$$

and

$$\beta = \frac{\sigma_j f(\sigma_i) - \sigma_i f(\sigma_j)}{\sigma_i^2 - \sigma_j^2}.$$

*Subcase 5a-(ii).* Suppose now $\sigma_i = \sigma_j > 0$. Similarly to Subcase 2a-(ii) we expand $U' = (I_2 + t \begin{bmatrix} 0 & -u^* \\ u & 0 \end{bmatrix})Q$, $V' = (I_2 + t \begin{bmatrix} 0 & -v^* \\ v & 0 \end{bmatrix})Q$, where $Q$ can be any unitary matrix. This time we impose that (3.9) is real and diagonal and find a solution in $u, v$ for the choice

$$Q = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 & -\iota \\ -\iota & 1 \end{bmatrix}.$$

Specifically, some elementary algebraic manipulations yield the condition

$$u - v = \frac{\iota}{2\sigma}.$$

For example we can take $v = 0$, $u = \iota(2\sigma)^{-1}$. This gives

$$(U')^* \begin{bmatrix} \sigma & it \\ 0 & \sigma \end{bmatrix} V' = \begin{bmatrix} \sigma + t/2 & 0 \\ 0 & \sigma - t/2 \end{bmatrix} + O(t^2);$$

we can then proceed precisely as in Subcase 2a-(ii).

*Subcase 5a-(iii).* Suppose $\sigma_i = \sigma_j = 0$, then it is immediate that

$$f^\diamond(A + tE) - f^\diamond(A) = \begin{bmatrix} 0 & \iota f(t) \\ 0 & 0 \end{bmatrix} \oplus 0_{(m-2) \times (n-2)},$$

and the statement follows dividing by $t$ and going to the limit for $t \to 0$.

- *Case 5b.* If the unique nonzero element of $E$ lies in the position $(i, j)$ with $m < j \leq n$, the procedure is again analogous to Case 2b. Assume that $i = 1$, $j = m + 1$: $A + tE = S'(V')^*$ is a singular value decomposition if

$$\begin{bmatrix} \sigma_i & 0 & \dots & 0 & \iota t & 0 & \dots & 0 \end{bmatrix} V' = \begin{bmatrix} \sigma & 0 & \dots & 0 \end{bmatrix},$$

for some $\sigma \geq 0$. Expanding $V'$ in powers of $t$ yields $\sigma = \sigma_i$, $V'_{ii} = 1$ for all $i = 1, \dots, n$, $V'_{1,n+1} = V'_{n+1,1} = -\iota t/\sigma_i$, and $V_{ij} = 0$ in all other cases. Hence, to first order in $t$, there is only one nonzero element in $f^\diamond(A + tE) - f^\diamond(A)$, lying precisely at the position $(i, j)$, and equal to $\iota t f(\sigma_i)/\sigma_i$.

- *Case 6.* Finally, if the unique nonzero element of $E$ is in the position $(i, j)$ with $j < i$, the proof is either analogous to Case 5a, if $i \leq m$, or to Case 5b, if $i > m$. We omit the details.

□

REMARK 3.9. *The matrices $F$, $G$ and $H$ have a very peculiar structure:*
- *the off-diagonal elements of $H$ and $F$ coincide, i.e., $H - F$ is diagonal;*
- *if $m = n$, $F$, $G$ and $H$ are all symmetric;*
- *$G_{ii} = 0$ for all $i = 1, \dots, \nu$;*

16

- *if $m > n$, then $F = \begin{bmatrix} F_1 \\ ev^T \end{bmatrix}$ with $F_1 = F_1^T \in \mathbb{R}^{n \times n}$, $e = \begin{bmatrix} 1 & \dots & 1 \end{bmatrix}^T \in \mathbb{R}^{m-n}$ and $v \in \mathbb{R}^n$, while $G = \begin{bmatrix} G_1 \\ 0 \end{bmatrix}$ where $G_1 = G_1^T \in \mathbb{R}^{n \times n}$;*

- *if $n > m$, then $F = \begin{bmatrix} F_1 & ve^T \end{bmatrix}$ with $F_1 = F_1^T \in \mathbb{R}^{m \times m}$, $e = \begin{bmatrix} 1 & \dots & 1 \end{bmatrix}^T \in \mathbb{R}^{n-m}$ and $v \in \mathbb{R}^m$, while $G = \begin{bmatrix} G_1 & 0 \end{bmatrix}$ where $G_1 = G_1^T \in \mathbb{R}^{m \times m}$.*

The formula in Theorem 3.8 simplifies considerably if the matrix $A$ is real. Since this is often the case in many applications, we give an explicit version of our main result specialized to real matrices.

COROLLARY 3.10. [**Daleckiĭ-Kreĭn Theorem for real generalized matrix functions**]
*Let $A = USV^T$ be an SVD of $A \in \mathbb{R}^{m \times n}$, where $U \in \mathbb{R}^{m \times m}$, $V \in \mathbb{R}^{n \times n}$, $S \in \mathbb{R}^{m \times n}$, and $S_{ii} =: \sigma_i$, $i = 1, \dots, \nu$, are the singular values of $A$. Let $f : \mathcal{S} \to \mathbb{R}$ be differentiable on an open subset of $\mathcal{S}$ containing the positive singular values of $A$. Moreover, if $A$ is rank deficient, i.e., if $\operatorname{rank} A < \nu$, suppose further that $f(0) = 0$ and that $f$ is right differentiable at $0$. Then the Fréchet derivative at $X = A$ of the generalized matrix function $f^\diamond(X)$, applied to the perturbation $E$, is*

$$L_{f^\diamond}(A, E) = U\left(F \circ \widehat{E} + G \circ \Upsilon(\widehat{E})\right) V^T, \tag{3.10}$$

*where $\circ, \widehat{E}$ and $\Upsilon$ are defined as in the statement of Theorem 3.8, and $F, G \in \mathbb{R}^{m \times n}$ are defined as in (3.6) and (3.7).*

EXAMPLE 3.11. *Take $f(x) = e^x$ and $A = \begin{bmatrix} 2 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix}$ so that an SVD is $A = USV^T$ with $U = I_2$, $V = I_3$, $S = A$. Since $A$ is full rank, $f^\diamond(X)$ is differentiable at $X = A$ in spite of the fact that $f(0) \neq 0$. Moreover, Corollary 3.10 holds with*

$$F = \begin{bmatrix} e^2 & \frac{e(2e-1)}{3} & \frac{e^2}{2} \\ \frac{e(2e-1)}{3} & e & e \end{bmatrix}, \qquad G = \begin{bmatrix} 0 & \frac{e(e-2)}{3} & 0 \\ \frac{e(e-2)}{3} & 0 & 0 \end{bmatrix}.$$

*Taking for example $E = \begin{bmatrix} 1 & 3 & 0 \\ 0 & -1 & 1 \end{bmatrix}$ we obtain*

$$L_{f^\diamond}(A, E) = \begin{bmatrix} e^2 & e(2e-1) & 0 \\ e(e-2) & -e & e \end{bmatrix}.$$

Further simplifications to Theorem 3.8 are possible by specializing $f$. We give one of the many possible examples: if $f$ is the constant function 1 then, if $A$ is full rank, $f^\diamond(A)$ is the unitary factor in the polar decomposition of $A$ [18, Chapter 8]. An implicit formula for the Fréchet derivative of the latter appeared in [22, Theorem 2.1], which collected it from [17, Proof of Theorem 2.5]. However, the formula in [17, 22] is a theoretical result that was not proposed for the computation of the Fréchet derivative, and in fact, can lead to numerical instabilities if implemented as given. Below, we specialize Theorem 3.8 to obtain an explicit formula, equivalent to [22, Theorem 2.1], which can be used [5] to devise an efficient and stable SVD-based algorithm for the computation of the Fréchet derivative of the unitary factor in a polar decomposition.

COROLLARY 3.12. *For any full rank matrix $X \in \mathbb{C}^{m \times n}$, $m \geq n$, let $X = QH$ be the polar decomposition of $X$ and consider the unitary factor $Q(X)$ as a function*

of $X$. Suppose morever that $A \in \mathbb{C}^{m \times n}$ is full rank and that $A = USV^*$ is an SVD. Then, the real Fréchet derivative of $Q(X)$ at $X = A$, applied to the perturbation $E$, is

$$L_Q(A, E) = U \left( F \circ \Re(\widehat{E}) + \iota H \circ \Im(\widehat{E}) + G \circ \Upsilon(\widehat{E}) \right) V^*, \qquad (3.11)$$

with the same notation as in Theorem 3.8 and where:
- $F \in \mathbb{R}^{m \times n}$ is defined as follows:

$$F_{ij} = \begin{cases} (\sigma_i + \sigma_j)^{-1} & \text{if } i \neq j \text{ and } i \leq n; \\ (\sigma_j)^{-1} & \text{if } i > n, \text{ and } \sigma_j \neq 0; \\ 0 & \text{otherwise}; \end{cases} \qquad (3.12)$$

- $G \in \mathbb{R}^{m \times n}$ is such that $G_{ij} = -F_{ij}$ for $i \leq n$ and $G_{ij} = 0$ for $i > n$;
- and $H \in \mathbb{R}^{m \times n}$ is defined as $H_{ii} = \sigma_i^{-1}$ and $H_{ij} = F_{ij}$ for $i \neq j$.

In particular, if $A \in \mathbb{R}^{n \times n}$ is square and invertible, then the real Fréchet derivative of $U(X)$ at $X = A$, applied to the perturbation $E$, is

$$L_Q(A, E) = U \left( F \circ (\widehat{E} - \widehat{E}^*) + \iota(H - F) \circ \Im(\widehat{E}) \right) V^T. \qquad (3.13)$$

**4. Application to conditioning.** In this section we apply the theory developed so far to the analysis of the conditioning of generalized matrix functions. To some extent, part of the analysis that we will be deriving may also be inferred starting from the Lipschitz continuity of generalized matrix functions, proved in [1] (assuming $f(0) = 0$); however, there is no explicit conditioning analysis there, and our treatment includes the case of $f(0) \neq 0$. Since the real case is the most relevant for the applications [4], and to keep the paper within a reasonable length, we focus on generalized matrix functions of real matrices and only allow real perturbations. We emphasize, however, that an analogous analysis can be performed for generalized matrix functions of complex matrices, starting from Theorem 3.8 rather than its specialization to real matrices, i.e., Corollary 3.10.

The absolute conditioning of a generalized matrix function can be defined as

$$\text{cond } f^\diamond(A) = \lim_{t \to 0} \sup_{\|E\| \leq 1} \frac{\|f^\diamond(A + tE) - f^\diamond(A)\|}{|t| \|E\|}. \qquad (4.1)$$

There are two cases. If $f(0) \neq 0$ and $A$ is rank deficient, then clearly $\text{cond } f^\diamond(A) = \infty$, as $f^\diamond(X)$ is not continuous at $X = A$. More interestingly, it may happen that either $f(0) = 0$ or $f(0) \neq 0$ but $A$ is full rank. Then, $f^\diamond(X)$ is differentiable at $X = A$, and $\|f^\diamond(A + tE) - f^\diamond(A)\| = \|L_{f^\diamond}(A, tE) + O(t^2)\| = |t| \|L_{f^\diamond}(A, E)\| + O(t^2)$. If we specialize to any unitarily invariant norm, and if $A = USV^T$ is an SVD, it is immediate that $\|L_{f^\diamond}(A, E)\| = \|L_{f^\diamond}(S, \widehat{E})\|$ having defined $\widehat{E} = U^T EV$. For example, the Frobenius norm is unitarily invariant, and this choice leads to the condition number

$$\text{cond}_F f^\diamond(A) = \|K_{f^\diamond}(S)\|_2,$$

18

where $K_{f^\diamond}(X)$ is the Kronecker form of the Fréchet derivative [18] of $f^\diamond$ at $X \in \mathbb{R}^{m \times n}$. To define $K_{f^\diamond}(X)$ it is convenient to introduce the vec operator [16]

$$\text{vec} : \mathbb{R}^{m \times n} \to \mathbb{R}^{mn}, X = \begin{bmatrix} \mathbf{x_1} & \ldots & \mathbf{x_n} \end{bmatrix} \mapsto \text{vec}(X) = \begin{bmatrix} \mathbf{x_1} \\ \vdots \\ \mathbf{x_n} \end{bmatrix}.$$

Then, $K_{f^\diamond}(X)$ is the unique matrix such that, for any $E \in \mathbb{R}^{m \times n}$, $\text{vec}(L_{f^\diamond}(X, E)) = K_{f^\diamond}(X) \text{vec}(E)$.

Let us now consider the linear map $\Upsilon$, defined in the statement of Corollary 3.10. Via the vec operator, it can be represented by the unique matrix $P \in \mathbb{R}^{mn \times mn}$ satisfying

$$P \text{vec}(A) = \text{vec}(\Upsilon(A)) \qquad \forall A \in \mathbb{R}^{m \times n}. \tag{4.2}$$

Note that in the special case $m = n$ we recover the well-studied vec-permutation operator [16].

LEMMA 4.1. *The matrix $P$ defined in (4.2) is a permutation matrix, and it is symmetric, orthogonal, and involutory. Moreover, it has precisely $mn + \nu(1 - \nu)/2$ eigenvalues equal to $+1$ and $\nu(\nu - 1)/2$ eigenvalues equal to $-1$.*

*Proof.* Since $\text{vec}(A)$ and $\text{vec}(\Upsilon(A))$ always contain the same elements, although possibly in a different order, we see that $P$ must be a permutation matrix, and hence, orthogonal: $PP^T = I_{mn}$. Moreover, from the fact that $\Upsilon$ is involutory, i.e., $\Upsilon(\Upsilon(A)) \equiv A$, we deduce that $P$ is also involutory: $P^2 = I_{mn}$. Therefore $P$ is also symmetric, $P = P^T$.

Any symmetric orthogonal matrix must have all semisimple eigenvalues equal to $\pm 1$. Suppose for simplicity $m \geq n$ (the proof for $m < n$ is analogous). Consider the two subspaces $\mathcal{V}_1 = \{X = \begin{bmatrix} X_1 \\ X_2 \end{bmatrix} \in \mathbb{R}^{m \times n} | X_1 = X_1^T \in \mathbb{R}^{n \times n}\}$ and $\mathcal{V}_2 = \{X = \begin{bmatrix} X_1 \\ 0 \end{bmatrix} \in \mathbb{R}^{m \times n} | X_1 = -X_1^T \in \mathbb{R}^{n \times n}\}$. Observe that $X \in \mathcal{V}_1 \Rightarrow \Upsilon(X) = X$, that $X \in \mathcal{V}_2 \Rightarrow \Upsilon(X) = -X$, and that $\mathbb{R}^{m \times n}$ is equal to the direct sum $\mathcal{V}_1 \oplus \mathcal{V}_2$. Noting that $m \geq n$ implies $n = \nu$, this concludes the proof. $\square$

For any vector $v \in \mathbb{R}^n$, we define $\text{diag}(v) \in \mathbb{R}^{n \times n}$ to be the diagonal matrix such that $(\text{diag}(v))_{ii} = v_i$. We then have the following corollary of Corollary 3.10.

COROLLARY 4.2. *Let $A \in \mathbb{R}^{m \times n}$, and suppose $A = USV^T$ is an SVD. Let $f : \mathcal{S} \to \mathbb{R}$ be differentiable on an open subset of $\mathcal{S}$ containing the positive singular values of $A$. Moreover, if $A$ is rank deficient, i.e., if $\text{rank}\, A < \nu$, suppose further that $f(0) = 0$ and that $f$ is right differentiable at $0$. Then the Kronecker form of the Fréchet derivative at $X = A$ of the real generalized matrix function $f^\diamond(X)$ is*

$$K_{f^\diamond}(A) = (V \otimes U)(\Phi + \Gamma P)(V^T \otimes U^T), \tag{4.3}$$

*where $\Phi = \text{diag}(\text{vec}(F))$, $\Gamma = \text{diag}(\text{vec}(G))$, $F$ and $G$ are defined as in Corollary 3.10 and $P$ is the matrix defined by (4.2).*

*Proof.* Applying the vec operator to (3.10), and using the properties $\text{vec}(ABC) = (C^T \otimes A) \text{vec}(B)$ and $\text{vec}(A \circ B) = \text{diag}(\text{vec}(A)) \text{vec}(B)$, we obtain

$$\text{vec}(L_{f^\diamond}(A, E)) = (V \otimes U) (\Phi + \Gamma P) \text{vec}(\hat{E}.)$$

The statement follows noting that $\text{vec}(\hat{E}) = \text{vec}(U^T E V) = (V^T \otimes U^T)\,\text{vec}(E)$. □

Slightly different formulae for $K_{f^\diamond}(A)$ may be deduced by the following lemma.

LEMMA 4.3. *In the notation of Corollary 4.2, $\Gamma P = P\Gamma$ and $\Phi P = P\Phi$.*

*Proof.* The structure of the matrix $G$ and the definition of $\Upsilon$ readily yield the property

$$G \circ \Upsilon(X) = \Upsilon(G \circ X) \qquad \forall\, X \in \mathbb{R}^{m \times n}.$$

Similarly, it is easy to check that

$$F \circ \Upsilon(X) = \Upsilon(F \circ X) \qquad \forall\, X \in \mathbb{R}^{m \times n}.$$

Applying the vec operator to each of these equations yields the statement. □

Lemma 4.1 and Lemma 4.3 imply that $K_{f^\diamond}(A)$ is symmetric: indeed, $(\Gamma P)^T = P^T \Gamma^T = P\Gamma = \Gamma P$. Moreover, taking $U = I_m$ and $V = I_n$ in Corollary 4.2 it is immediate that $K_{f^\diamond}(S) = \Phi + P\Gamma$. As by Lemma 4.1 $P$ is orthogonal, this immediately yields the bound $\text{cond}_F f^\diamond(A) = \|K_{f^\diamond}(S)\|_2 \leq \max|F_{ij}| + \max|G_{ij}|$. It is easy to improve the latter estimate by diagonalizing $K_f(S)$. The next subsection is devoted to this goal.

**4.1. The eigenvalues of the Kronecker form of the Fréchet derivative.** For simplicity of exposition, in this subsection we will assume $m \geq n$. The results, however, do not change if $m < n$, except that in certain formulae the roles of the pairs $(i, m)$ and $(j, n)$ must be exchanged. Observe first that, due to the zero structure of $G$ and to the symmetry of $P$, a simple simultaneous permutation $Q$ of rows and columns leads to the block diagonalization

$$QK_f(S)Q^T = \overset{\nu}{\underset{i=1}{\bigoplus}} F_{ii} \oplus \overset{n}{\underset{j=1}{\bigoplus}}\, \overset{m}{\underset{i=n+1}{\bigoplus}} F_{ij} \oplus \overset{n-1}{\underset{i=1}{\bigoplus}}\, \overset{n}{\underset{j=i+1}{\bigoplus}} \begin{bmatrix} F_{ij} & G_{ij} \\ G_{ij} & F_{ij} \end{bmatrix}.$$

Each $2 \times 2$ block has eigenvalues $F_{ij} \pm G_{ij}$, and therefore we have the following theorem.

THEOREM 4.4. *For the condition number of a real generalized matrix function, it holds*

$$\text{cond}_F f^\diamond(A) = \max\{a, b, c, d\},$$

*where $a = \max_i |F_{ii}|$, $b = \max_{j \leq n \leq i} |F_{ij}|$, $c = \max_{i<j} |F_{ij} + G_{ij}|$, $d = \max_{i<j} |F_{ij} - G_{ij}|$, and $F$ and $G$ are defined as in Corollary 3.10.*

In order to estimate the values of $a, b, c, d$, it is useful to give an explicit expression for the eigenvalues of $K_{f^\diamond}(S)$.

THEOREM 4.5. *It holds*

$$F_{ij} + G_{ij} = \begin{cases} \frac{f(\sigma_i) - f(\sigma_j)}{\sigma_i - \sigma_j} & \text{if } \sigma_i \neq \sigma_j; \\ f'(\sigma_i) & \text{if } \sigma_i = \sigma_j \end{cases}$$

*and*

$$F_{ij} - G_{ij} = \begin{cases} \frac{f(\sigma_i) + f(\sigma_j)}{\sigma_i + \sigma_j} & \text{if } \sigma_i \neq \sigma_j; \\ \frac{f(\sigma_i)}{\sigma_i} & \text{if } \sigma_i = \sigma_j \neq 0; \\ f'(0) & \text{if } \sigma_i = \sigma_j = 0. \end{cases}$$

*Proof.* It follows from Corollary 3.10 by a direct computation. □

20

**4.2. On the conditioning of real generalized matrix functions.** If all the singular values of the matrix $A$ are known then Theorems 3.10, 4.4, and 4.5 can be combined to compute $\|K_{f^\diamond}(S)\|_2$, and hence $\mathrm{cond}_F\, f^\diamond(A)$. In this subsection, we give some upper bounds for $\mathrm{cond}_F\, f^\diamond(A)$ that only require the knowledge of the function $f$ and of the largest and smallest nonzero singular values of $A$, $\sigma_1 = \|A\|_2$ and $\sigma_r$. In practice, these estimates may be useful: for very large matrices it is expensive to compute a full SVD, but algorithms exist to cheaply compute the extremal singular values only, e.g., Lanczos-based methods [12, Chapter 10].

THEOREM 4.6. *Let $A \in \mathbb{R}^{m \times n}$ have full rank, and let $\sigma_r$ be the smallest singular value of $A$. Denote by $\mathcal{I}$ the interval $[\sigma_r, \|A\|_2]$, and set $M = \max_{x \in \mathcal{I}} |f(x)|$. Suppose moreover that $f$ is continuously differentiable, and hence locally Lipschitz continuous, on $\mathcal{I}$, with Lipschitz constant $K$. Then, it holds*

$$\mathrm{cond}_F\, f^\diamond(A) \leq \max\{K, M\sigma_r^{-1}\}.$$

*Proof.* Let $x > y \in \mathcal{I}$ be singular values of $A$. We can bound $|\frac{f(x)-f(y)}{x-y}| \leq K$, $|f'(y)| \leq K$, $|\frac{f(x)+f(y)}{x+y}| \leq \frac{M}{\sigma_r}$, $|\frac{f(x)}{x}| \leq \frac{M}{\sigma_r}$. The statement then follows from Theorems 3.10, 4.4, and 4.5.

More precisely, in the notation of Theorem 4.4, $a \leq K$, $b \leq M\sigma_r^{-1}$, $c \leq K$, $d \leq M\sigma_r^{-1}$. □

If we further assume that $f(0) = 0$, a stronger result can be derived. It could also be obtained as a consequence of [1, Theorem 1.1], proved with a different approach. Here, we give our own proof.

THEOREM 4.7. *Let $A \in \mathbb{R}^{m \times n}$. Denote by $\mathcal{I}$ the interval $[0, \|A\|_2]$. Suppose moreover that $f(0) = 0$ and that $f$ is continuously differentiable, and hence locally Lipschitz continuous, on $\mathcal{I}$, with Lipschitz constant $K$. Then, it holds*

$$\mathrm{cond}_F\, f^\diamond(A) \leq K.$$

*Proof.* This time, for any $x, y \in \mathcal{I}$ we can bound $|\frac{f(x)-f(y)}{x-y}| \leq K$, $|f'(x)| \leq K$, $|\frac{f(x)+f(y)}{x+y}| \leq \frac{|f(x)|+|f(y)|}{x+y} \leq \frac{Kx+Ky}{x+y} = K$, $|\frac{f(x)}{x}| = \frac{|f(x)-f(0)|}{|x-0|} \leq K$. □

Variations on the theme of Theorems 4.6 and 4.7 can be obtained assuming that more singular values are known. Intuitively, since generalized matrix functions are computed via the SVD, one expects that they should be better conditioned, in the sense of being closer to having the same conditioning of scalar functions, than classical matrix functions. When $f(0) = 0$ and $f$ is Lipschitz, this is essentially the case, as shown by Theorem 4.7: in this scenario, generalized matrix functions are "as well conditioned as their scalar counterparts". If $f(0) \neq 0$, Theorems 4.4 and 4.5 show that the absolute condition number for the generalized matrix function is controlled by the maximum of the absolute values of the functions $f'(x)$ and $f(x)/x$, both evaluated at the singular values of $A$. Note that the norm of the derivative is the absolute condition number of the scalar function $f$. It may happen that a generalized matrix funcion is worse conditioned than its scalar counterpart applied to each singular value individually *only if* it happens that $\max_i |f(\sigma_i)/\sigma_i| \gg \max_i |f'(\sigma_i)|$.

EXAMPLE 4.8. *Let $A = \begin{bmatrix} 2 & 0 \\ 0 & 1 \end{bmatrix}$ and let $f(x) = M(-2x^3 + 9x^2 - 12x + 6)$ for some arbitrary $M > 0$. Observe that $f^\diamond(A) = MA$.*

Then the eigenvalues of $K_f(A)$ are equal to $f'(1) = 0, f'(2) = 0, \frac{f(2)+f(1)}{3} = M$, and $f(2) - f(1) = M$. Hence, the absolute condition number of $f^\diamond(A)$ is $M$, to be compared with the absolute condition number of $f(x)$ at the individual singular values, which is 0 for both of them.

By specializing to a fixed generalized matrix function $f^\diamond$ stronger results may be obtained. We give a couple of examples.

EXAMPLE 4.9. *Letting* $f(x) = 1$, *computing* $f^\diamond(A)$ *for a full rank matrix* $A$ *corresponds to the computation of the orthogonal polar factor in a polar decomposition of* $A$ *[18, Chapter 8]. (If* $A$ *is rank deficient, then this is no longer true, and the orthogonal factor in a polar decomposition is not unique).*

*If* $m = n$ *and* $A \in \mathbb{R}^{n \times n}$ *is invertible, Kenney and Laub showed [22, Theorems 2.2 and 2.3] that the absolute condition number is*

$$\operatorname{cond}_F f^\diamond(A) = \frac{2}{\sigma_n + \sigma_{n-1}},$$

*where* $\sigma_r$ *and* $\sigma_{r-1}$ *are, respectively, the smallest and second smallest singular values. Although Theorem 4.6 only gives a bound of* $1/\sigma_n$, *which is slack for* $\sigma_n < \sigma_{n-1}$, *specializing Theorem 4.4 to* $f = 1$ *gives* $a = c = 0$ *and* $d = 2/(\sigma_n + \sigma_{n-1})$. *We thus recover the result by Kenney and Laub as a special case of our analysis. If* $m > n$ *and* $A$ *has full rank, then [7] the absolute condition number is* $1/\sigma_n$, *and hence the upper bound of Theorem 4.6 is tight.*

EXAMPLE 4.10. *Let* $f(x) = \exp(x)$ *and let us consider* $f^\diamond(A)$ *for a full rank matrix* $A \in \mathbb{R}^{m \times n}$. *Then,* $M = K = \exp(\|A\|_2)$, *and hence,* $\operatorname{cond}_F \exp^\diamond(A) \leq \exp(\|A\|_2) \max\{1, \sigma_r^{-1}\}$. *Again, a careful examination of the explicit expressions of Theorem 4.5 can improve the general bound of Theorem 4.6. In particular,* $\exp(x)$ *is convex and increasing, while* $\exp(x)/x$ *is convex and has a minimum at* $x = 1$. *Therefore, we conclude that*

$$\operatorname{cond}_F \exp^\diamond(A) = \max\{\exp(\sigma_r)/\sigma_r, \exp(\|A\|_2)/\|A\|_2, \exp(\|A\|_2)\}.$$

In practice, quoting Nick Higham [18, p. 56], "it is the relative condition number that is of interest, but it is more convenient to state results for the absolute condition number". In the Frobenius norm, the relative condition number for the generalized matrix function $f^\diamond(A)$ is given in terms of the absolute condition number $\operatorname{cond}_F f^\diamond(A)$ by the formula

$$\operatorname{rcond}_F f^\diamond(A) = \operatorname{cond}_F f^\diamond(A) \cdot \frac{\|A\|_F}{\|f^\diamond(A)\|_F}.$$

Suppose that we have an upper bound for the absolute condition number, say, $\operatorname{cond}_F f^\diamond(A) \leq \beta$. Using $\|A\|_F \leq \sqrt{\nu}\|A\|_2$, we then see that an upper bound for the relative condition number is

$$\operatorname{rcond}_F f^\diamond(A) \leq \frac{\beta\sqrt{\nu}\|A\|_2}{\|f^\diamond(A)\|_F}.$$

For a general $f$, calculating $\|f^\diamond(A)\|_F$, or its lower bound $\|f^\diamond(A)\|_2$, might be nontrivial without computing $f^\diamond(A)$ explicitly or knowing the full singular spectrum of $A$. In the spirit of this subsection, we provide a lower bound assuming that only the largest and smallest nonzero singular values of $A$ are known. Observe that

$$\|f^\diamond(A)\|_F \geq \mu := \sqrt{f(\|A\|_2)^2 + f(\sigma_r)^2}.$$

22

Moreover, it is easy to see that in the statement and proof of Theorem 4.6 we could replace $M$ by $\|f^\diamond(A)\|_2$ (the reason for not having done so is that the latter may be more difficult to compute in practice). Hence, we obtain the following corollary.

COROLLARY 4.11. *In the notation and under the assumptions of Theorem 4.6, setting $\mu := \sqrt{f(\|A\|_2)^2 + f(\sigma_r)^2}$, it holds*

$$\mathrm{rcond}_F f^\diamond(A) \leq \sqrt{\nu}\|A\|_2 \max\{\frac{K}{\mu}, \frac{1}{\sigma_r}\}.$$

*In the notation and under the assumptions of Theorem 4.7, it holds*

$$\mathrm{rcond}_F f^\diamond(A) \leq \frac{\sqrt{\nu}\|A\|_2 K}{\mu}.$$

Example 4.8 showed that the absolute conditioning of generalized matrix functions can be much higher than that of the scalar functions they are induced by. However, the relative condition number for that example is 1. Can generalized matrix functions be much worse conditioned, in the *relative* sense, than their scalar counterparts? Using Corollary 4.11, one may expect trouble if $f(0) \neq 0$ and $A$ is numerically near to being rank deficient. We illustrate this with a concrete example.

EXAMPLE 4.12. *For some $0 < \epsilon \ll 1$, let*

$$A = \begin{bmatrix} \epsilon & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix} \qquad \text{and} \qquad f(x) = 1 + (x - \epsilon)^2.$$

*It is immediate to check that*

$$f^\diamond(A) = \begin{bmatrix} f(\epsilon) & 0 & 0 \\ 0 & f(1) & 0 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 2 - 2\epsilon + \epsilon^2 & 0 \end{bmatrix}.$$

*The relative condition numbers of the scalar function $f$ at the singular values of $A$ are, respectively, $0$ at $x = \epsilon$ and $1 + O(\epsilon^2)$ at $x = 1$. However, the relative condition number of $f^\diamond(A)$ is*

$$\mathrm{rcond}_F f^\diamond(A) = \frac{1}{\epsilon\sqrt{5}} + O(1).$$

*Figure 4.1 plots the relative error*

$$\rho = \frac{\|f^\diamond(A + E) - f^\diamond(A)\|_F \|A\|_F}{\|f^\diamond(A)\|_F},$$

*computed with MATLAB version R2015b, against the parameter $\epsilon$ for the perturbation*

$$E = \begin{bmatrix} 0 & 0 & 10^{-15} \\ 0 & 0 & 0 \end{bmatrix}.$$

To summarize, unlike classical matrix functions, real generalized matrix functions induced by Lipschitz continuous functions satisfying $f(0) = 0$ – two conditions that are commonly met in practical applications, see [4] – are never numerically dodgier than the scalar functions they are induced by. Informally speaking, this is because the Jordan decomposition is not numerically tame, but the SVD is. Indeed, classical functions of non-normal matrices may encounter issues due to the ill conditioning
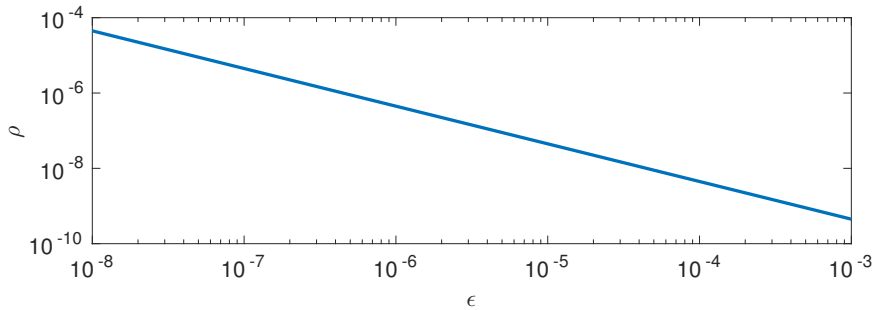
FIG. 4.1. *Computed relative error for Example 4.12*

of the eigenvector matrix $Z$ in Theorem 2.11; on the other hand, since $U$ and $V$ in Corollary 3.10 are orthogonal, the information on the conditioning of generalized matrix functions is directly encoded in the Daleckiĭ-Kreĭn type formula developed in this paper.

An exception to this generally optimistic situation is when $f(0) \neq 0$ and $f^\diamond(A)$ is computed for some rank-deficient, or near-rank deficient, matrix $A$. In this scenario, one is trying to evaluate numerically a function at, or close to, a point of discontinuity, and the closer $A$ is to having some zero singular values, the harsher potential challenges are to be expected for the numerical computation of $f^\diamond(A)$.

REFERENCES

[1] F. ANDERSSON, M. CARLSSON AND K.-M. PERFEKT, *Operator-Lipschitz estimates for the singular value functional calculus*, Proc. AMS, 144(5), (2016), pp. 1867–1875.

[2] B. ANDREWS AND C. HOPPER, *The Ricci Flow in Riemannian Geometry. A Complete Proof of the Differentiable 1/4-Pinching Sphere Theorem*, Lecture Notes in Mathematics, Vol. 2011, (2011), Edited by J.-M. Morel, F. Takens, B. Teissier, P.K. Maini, Springer.

[3] F. ARRIGO AND M. BENZI, *Edge modification criteria for enhancing the communicability of digraphs*, SIAM J. Matrix Anal. Appl. 37(1) (2016), pp. 443–468.

[4] F. ARRIGO, M. BENZI AND C. FENU, *Computation of generalized matrix functions*, To appear in SIAM J. Matrix Anal. Appl.

[5] B. ARSLAN, V. NOFERINI AND F. TISSEUR, *The structured condition number of a differentiable map between matrix manifolds, with applications*, In preparation, 2016.

[6] A. BUNSE-GERSTNER, R. BYERS, V. MEHRMANN AND N. K. NICHOLS, *Numerical computation of an analytic singular value decomposition of a matrix valued function*, Numer. Math. 60 (1991), pp. 1–39.

[7] F. CHAITIN–CHATELIN AND S. GRATTON, *On the condition numbers associated with the polar factorization of matrix*, Numer. Linear Algebra Appl. 7 (2000), pp. 337–354.

[8] J. L. DALECKIĬ AND S. G. KREĬN, *Integration and differentiation of functions of Hermitian operators and applications to the theory of perturbatitions*, Amer. Math. Soc. Transl., Series 2, 47 (1965), pp. 1–30.

[9] L. FANTAPPIÉ, *Le calcul des matrices*, C. R. Ac. des Sc. Paris 186 (1928), pp. 619–621.

[10] L. FANTAPPIÉ, *Sulle funzioni di una matrice*, An. Acad. Brasil. Cienc. 26 (1954), pp. 25–33.

[11] G. H. GOLUB AND V. PEREYRA, *The differentiation of pseudo-inverses and nonlinear least square problems whose variables separate*, SIAM J. Matrix Anal. Appl. 10(2) (1973),

pp. 413–432.

[12] G. H. Golub and C. Van Loan, Matrix Computations, 4th edition, *John Hopkins University Press*, Baltimore, MD, United States, 2012.

[13] R. S. Hamilton, *The inverse function theorem of Nash and Moser*, Bulletin of the AMS 7(1) (1982), pp. 65–222.

[14] J. B. Hawkins and A. Ben–Israel, *On generalized matrix functions*, Linear and Multilinear Algebra 1(2) (1973), pp. 163–171.

[15] J. Heinonen, *Lectures on Lipschitz analysis*, Rep. Univ. Jyväskylä Dept. Math. Stat. 100 (2005), 1–77.

[16] H. V. Henderson and S. R. Searle, *The vec-permutation matrix, the vec operator and Kronecker products: a review*, Linear and Multilinear Algebra 9(4) (1980/81), pp. 271–288.

[17] N. J. Higham, *Computing the Polar Decomposition — with Applications*, SIAM J. Sci. Stat. Comput. 7(1) (1986), pp. 1160–1174.

[18] N. J. Higham, Functions of Matrices: Theory and Computation, *SIAM*, Philadelphia, PA, United States, 2008.

[19] R. A. Horn and C. R. Johnson, Matrix Analysis, 2nd edition, *Cambridge University Press*, New York, NY, United States, 2013.

[20] R. A. Horn and C. R. Johnson, Topics in Matrix Analysis, *Cambridge University Press*, New York, NY, United States, 1991.

[21] T. Kato, Perturbation Theory for Linear Operators, *Springer-Verlag*, New York, NY, United States, 1966.

[22] C. S. Kenney and A. J. Laub, *Polar decomposition and matrix sign function condition estimates*, SIAM J. Sci. Statist. Comput. 12(3) (1991), pp. 488–504.

[23] A. J. Kurdila and M. Zabarankin, Convex Functional Analysis, *Birkhäuser-Verlag*, Basel, Switzerland, 2005.

[24] R. Penrose, *A generalized inverse for matrices*, Proc. Cambridge Philos. Soc. 51 (1955), pp. 406–413.

[25] W. Rudin, Principles of Mathematical Analysis, *McGraw-Hill*, New York, NY, United States, 1976.

[26] K. Tapp, Matrix Groups for Undergraduates, *American Mathematical Society*, United States, 2005.