

*A Daleckii-Krein formula for the Frechet  
derivative of a generalized matrix function*

Noferini, Vanni

2016

MIMS EPrint: **2016.24**

Manchester Institute for Mathematical Sciences  
School of Mathematics

The University of Manchester

Reports available from: <http://eprints.maths.manchester.ac.uk/>

And by contacting: The MIMS Secretary  
School of Mathematics  
The University of Manchester  
Manchester, M13 9PL, UK

ISSN 1749-9097

# A Daleckii-Kreĭn formula for the Fréchet derivative of a generalized matrix function

Vanni Noferini\*

April 27, 2016

## Abstract

We state and prove an extension of the Daleckii-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 generalized matrix functions are real differentiable but generally not complex differentiable. Finally, we discuss the application of our result 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.

**Keywords:** generalized matrix function, Daleckii-Kreĭn theorem, Gâteaux derivative, Fréchet derivative, condition number

**MSC classification:** 65F60, 15A16, 65F35, 15A12, 15A60

## 1 Introduction

Matrix functions are a central subject in matrix theory and in numerical linear algebra [11, Chapter 9], [17], [19, Chapter 6]. There are several equivalent definitions of matrix functions, based on, for example, the Jordan canonical form, polynomial interpolation, or Cauchy integrals [17].

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

---

\*Department of Mathematical Sciences, University of Essex, Wivenhoe Park, Colchester, UK, CO4 3SQ.  
(vnofer@essex.ac.uk)

basic theory [13]. They forged the name “generalized matrix functions”<sup>1</sup> for their singular value-based definition, and showed that generalized matrix functions satisfy four of the so-called Fantappi  properties [8, 9]. 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 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 [17, Chapter 3]. In particular, a basic result in matrix theory is the Dalecki -Kre n theorem [7], 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 numerical conditioning. Our main result is a “generalized Dalecki -Kre n 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 Dalecki -Kre n 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 Dalecki -Kre n theorem. Section 3 investigates the existence of the 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 consider real matrices first, to later treat separately the technically more involved case of complex matrices. Finally, Section 4 discusses the application of our results to the study of the condition number of generalized matrix functions.

---

<sup>1</sup>In spite of this name, generalized matrix functions do not generally reduce to classical matrix functions when the argument matrix is square [4, 17]. We adhere to the original terminology of [13], following also the more recent paper [4].

## 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) [11] 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. Moreover, the diagonal entries  $S_{ii} = \sigma_i$  are called the singular values of  $A$  and appear in nonincreasing order:  $\sigma_1 \geq \dots \geq \sigma_r > \sigma_{r+1} = \dots = \sigma_\nu = 0$ . 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 orthogonal, 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, 13] the following definition, based on the CSVD, is given. Here and below,  $\mathbb{R}_{\geq} = [0, \infty)$  denotes the set of nonnegative real numbers.

**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  $f : \mathbb{R}_{\geq} \rightarrow \mathbb{R}$  be a scalar function such that  $f(\sigma_i)$  is defined for all  $i = 1, \dots, r$ . The generalized matrix function  $f^\diamond : \mathbb{C}^{m \times n} \rightarrow \mathbb{C}^{m \times n}$  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

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

It is not hard to verify that 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). 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, 13, 17]; below we collect a few results that are useful to us, and we add some new observations of our own.

**Proposition 2.2.** ([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:*

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

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

**Theorem 2.3.** ([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  $X$  [23] and  $f(Y)$  denotes the classical matrix function [17], induced by the same scalar function  $f$ , of the matrix  $Y$ .

Like classical matrix functions, generalized matrix functions can also always be computed by polynomial interpolation. This was already mentioned (without giving details) in [13]. Here we give a more precise statement.

**Proposition 2.4.** *Let  $A \in \mathbb{C}^{m \times n}$  have rank  $r$  and  $f : \mathbb{R}_\geq \rightarrow \mathbb{R}$  be a scalar function such that  $f(\sigma_i)$  is defined for all  $i = 1, 2, \dots, r$ . If  $r = 0$ , let  $p$  be any polynomial. Otherwise, let  $A$  have  $k$  distinct positive singular values*

$$\sigma_1 = \sigma_{i_1} > \sigma_{i_2} > \dots > \sigma_{i_k} = \sigma_r,$$

and let  $p$  be a polynomial satisfying

$$p(\sigma_{i_j}) = f(\sigma_{i_j}), \quad j = 1, \dots, k.$$

Then  $f^\diamond(A) = p^\diamond(A)$ .

*Proof.* Straightforward from Definition 2.1.  $\square$

In the statement of Proposition 2.4, for  $r > 0$  one might make  $p$  unique by restricting its degree to be equal to  $k - 1$ . 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 2.5.** [13] Let first  $p = x^{2k+1}$ . Then the corresponding generalized odd powers of  $A$  are equal to  $p^\diamond(A) = (AA^*)^k A = A(A^*A)^k$  for  $p = x^{2k+1}$ . If  $p = x^{2k}$ , suppose that  $A$  has 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 2.6.** Let  $A \in \mathbb{C}^{m \times n}$  be a rank  $r$  matrix and let  $A = USV^*$  be an SVD. Let  $f : \mathbb{R}_{\geq} \rightarrow \mathbb{R}$  be a scalar function such that  $f(\sigma_i)$  is defined for all  $i = 1, \dots, r$ . Then, we define the scalar function  $f^\diamond(\sigma) : \mathbb{R}_{\geq} \rightarrow \mathbb{R}$  as

$$f^\diamond(\sigma) = \begin{cases} f(\sigma) & \text{if } \sigma > 0; \\ 0 & \text{if } \sigma = 0. \end{cases} \quad (1)$$

The generalized matrix function  $f^\diamond : \mathbb{C}^{m \times n} \rightarrow \mathbb{C}^{m \times n}$  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

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

A third characterization is also possible, as briefly mentioned in [17, Solution to Problem 1.53]. If  $m \geq n$  and  $A = QH$  is a polar decomposition [17, Chapter 8], Definition 2.6 and Proposition 2.2 yield  $f^\diamond(A) = Qf^\diamond(H)$ , where  $f^\diamond(H)$  is the classical matrix function of  $H$  induced by the scalar function (1). If either  $A$  has full rank or  $f(0) = 0$ , then we have the stronger property  $f^\diamond(A) = Qf(H)$ . 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 2.6 makes it manifest that the scalar functions of the form (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 (1), and hence the same observation holds for rank deficient matrices.

**REMARK 2.7.** Suppose that  $f(0) \neq 0$ . If  $A$  does not have full rank, i.e., if  $\text{rank } A < \nu$ , then  $f^\diamond(X)$  is not continuous (let alone differentiable) at  $X = A$ .

*Example 2.8.* For  $t \neq 0$  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,

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

## 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 [22], or in [24] for the finite dimensional case.

Suppose that  $X, Y$  are Banach spaces and let  $f : X \rightarrow Y$ . 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

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

Under these assumptions,  $L_f(x, \cdot)$  is called the Fréchet derivative of  $f$  at  $x$ . It is continuous in  $x$  and, by definition, linear in  $h$ . When it exists, the Fréchet derivative is equal to the Gâteaux derivative, defined as<sup>2</sup>

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

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  $x$  [17, Chapter 3], or (ii)  $f$  is jointly (in  $x$  and  $h$ ) continuously Gâteaux differentiable [12, 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 2.9.* Let  $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ ,  $(x, y) \mapsto \frac{x^3}{x^2+y^2}$  if  $(x, y) \neq (0, 0)$  and  $f(0) = 0$ . Then  $f$  is Gâteaux differentiable for any  $(x, y) \in \mathbb{R}^2$ , with

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

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 diagonal elements of  $D$  (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 [11, 17, 19].

---

<sup>2</sup>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.

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 [17, Chapter 3]. An explicit formula is known for diagonalizable matrices, and was first formulated by Daleckiĭ and Kreĭn. We recall [18] 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 2.10.** [Daleckiĭ–Kreĭn Theorem][7]. *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

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

## 3 Main result

### 3.1 Existence of Gâteaux and Fréchet derivatives

The Daleckiĭ–Kreĭn theorem only applies to diagonalizable (by similarity) square matrices. In this section, we will derive an analogous formula, valid for *any matrix*, either square or not, for the Fréchet derivative of generalized matrix functions.

Let us start by considering the Gâteaux differentiability of generalized matrix functions of real matrices.

**Theorem 3.1.** [Gâteaux differentiability of real generalized matrix functions] *Let  $A \in \mathbb{R}^{m \times n}$  and let  $f : \mathbb{R}_{\geq} \rightarrow \mathbb{R}$  be continuously differentiable on an open set 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  $f^\diamond(X)$ , defined as in Definitions 2.1 or 2.6, is Gâteaux differentiable at  $X = A$ .*

*Proof.* Recall that [5, Theorem 1], [20, Section II.6.2] every real analytic matrix-valued function admits an analytic SVD, so in particular it holds

$$A + tE = U(t)S(t)V(t)^T \tag{2}$$

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

$$G_{f^\diamond}(A, E) = \lim_{t \rightarrow 0} \frac{f^\diamond(A + tE) - f^\diamond(A)}{t} \tag{3}$$



exists provided that  $f^\circ(S(t))$  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  $U(t) = U_0 + tU_1 + O(t^2)$ ,  $V(t) = V_0 + tV_1 + O(t^2)$ ,  $S(t) = S_0 + tS_1 + O(t^2)$ , from (2) we get that  $A = U_0S_0V_0^T$  is an SVD (and without loss of generality we can take it as the SVD that defines  $f^\circ(A)$  in Definition 2.6), that  $U_0U_1^T + U_1U_0^T = 0 = V_0V_1^T + V_1V_0^T$ , that  $S_1$  is diagonal, and that

$$G_{f^\circ}(A, E) = U_1f^\circ(S_0)V_0^T + U_0f^\circ(S_0)V_1^T + U_0 \left. \frac{df^\circ(S(t))}{dt} \right|_{t=0} V_0^T, \quad (4)$$

where by the chain rule

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

□

The computation of the Gâteaux derivative from (4) is, in principle, not impossible employing the sophisticated techniques of [5]; however, this may be very challenging in practice. We will give a much more explicit formula in Theorem 3.7.

If  $A$  is not full rank and  $f(0) \neq 0$ , then by Remark 2.7  $f^\circ(X)$  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 [14, Theorem 3.1], they must be Fréchet differentiable almost everywhere; yet, in principle, there might exist a measure zero subset of  $\mathbb{R}^{m \times n}$  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 [19, Theorem 6.6.14].

**Theorem 3.2.** *Let  $A \in \mathbb{R}^{m \times n}$  and let  $f : \mathbb{R}_\geq \rightarrow \mathbb{R}$  be continuously differentiable on an open set 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

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

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

$$G_{f^\circ}(A, E) = G_{q^\circ}(A, E) \quad \forall E \in \mathbb{R}^{m \times n}.$$

*Proof.* It is immediate from (4). 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  $f^\circ(S_0) = q^\circ(S_0)$  and that  $(f')^\circ(S_0) = (q')^\circ(S_0)$ . □

If  $f^\circ$  is the generalized matrix power induced by  $f(x) = 1 = x^0$ , then  $f^\circ(X)$  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 3.3.** *Let  $f(x) = x^h$ ,  $h = 1, 2, 3, \dots$ , and let  $A \in \mathbb{R}^{m \times n}$ . Then the generalized power matrix  $f^\circ(X)$ , defined as in Definitions 2.1 or 2.6, is Fréchet differentiable at  $X = A$ .*

*Proof.* By Remark 2.5, if  $h = 2k + 1$  is odd then  $f^\circ(A)$  is manifestly Fréchet differentiable: indeed, each entry of  $(AA^T)^k A$  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 2.5,  $f^\circ(A) = (A^T A)^k U_r V_r^T$ , where  $U_r$  and  $V_r$  are the partial isometries defining a CSVD of  $A = U_r S_r V_r^T$ . Suppose first that  $r = m = n$ , i.e.,  $A$  is square and nonsingular. Observe that

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

where  $H = V_r S_r V_r^T$  is the Hermitian factor of any polar decomposition of  $A$  [17, Chapter 8]. That  $H$  is Gâteaux differentiable can be argued as in the proof of Theorem 3.1 via the existence of the analytic SVD (2). Indeed, observe that if  $A(t) = A + tE = U(t)S(t)V(t)^T$  is an SVD then  $H(A(t)) = \sqrt{V(t)S(t)^T S(t)V(t)^T}$ . Note also that  $H$  is a Lipschitz continuous function of  $A$  [17, Theorem 8.9]. Therefore, if we can show that the Gâteaux derivative is linear in  $E$ , it will follow that  $H$  is a Fréchet differentiable function of  $A$ . Let  $G_H(A, E)$  be the Gâteaux derivative of  $H$  as a function of  $A$ , applied to the direction  $E$ . Then it holds [16, Proof of Theorem 2.5]

$$H G_H(A, E) + G_H(A, E) H = A^T E + E^T A,$$

displaying linearity in  $E$  (note that, by assumption,  $H$  is symmetric positive definite). 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  $A^T$  and invoking item (i) in Proposition 2.2. Let  $A = QH$ ,  $Q = U_r V_r^T$ ,  $H = V_r S_r V_r^T$ , be a polar decomposition. When  $A$  does not have full rank, it is not any more true that the Hermitian factor  $H$  is Fréchet differentiable at  $A$ . However, we will argue that  $B = AH$  is, implying that

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

is differentiable as well. Observe that  $B = AH = f^\circ(A)$  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 3.1 it suffices to show that its Gâteaux derivative is linear in the perturbation  $E$ . Note that by item (ii) in Proposition 2.2 the generalized matrix function  $B$  is differentiable at  $A$  if and only if it is differentiable at  $S$ , where  $A = USV^T$  is an SVD. Therefore, with no loss of generality, we may take  $A$  to be of the form

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

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  $A(t) = U(t)S(t)V(t)^T$  is an analytic SVD (2). Partition

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

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

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  $f^\circ$  is the generalized matrix function induced by  $f(x) = x^2$ , one has

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

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

$$B(t) = B + t \begin{bmatrix} U_{11}^{(0)}(S_r^2)'(0)(V_{11}^{(0)})^T + U_{11}^{(0)}S_r^2(V_{11}^{(1)})^T + U_{11}^{(1)}S_r^2(V_{11}^{(0)})^T & U_{11}^{(0)}S_r^2V_{21}^T \\ U_{21}S_r^2(V_{11}^{(0)})^T & 0 \end{bmatrix} + O(t^2). \quad (5)$$

We deduce that the Gâteaux derivative of  $B$  at  $A$ , applied to the perturbation  $E$ , has the form

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

and by (5) it is immediate to check that  $X \in \mathbb{R}^{r \times r}$  is precisely  $G_B(S_r, E_{11})$ , i.e., the Gâteaux derivative of the generalized matrix function  $f^\circ$  induced by  $f(x) = x^2$  (but seen as a function defined on  $\mathbb{R}^{r \times r}$  rather than on  $\mathbb{R}^{m \times n}$  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 Fréchet derivative, and therefore it is linear in  $E_{11}$ , and hence, in  $E$ .

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

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

and

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

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 in  $E$ , and this concludes the proof.  $\square$

The next theorem is the main result of this subsection.

**Theorem 3.4. [Fréchet differentiability of real generalized matrix functions]** *Let  $A \in \mathbb{R}^{m \times n}$  and let  $f : \mathbb{R}_{\geq} \rightarrow \mathbb{R}$  be continuously differentiable on an open set 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  $f^\circ(X)$ , defined as in Definitions 2.1 or 2.6, is Fréchet differentiable at  $X = A$ .*

*Proof.* If  $A$  has full rank then the statement follows by Theorem 2.3 and by standard results on the differentiability of standard matrix functions [17, Chapter 3] and pseudoinverses of real full rank matrices [10]. To deal with the case of a rank deficient  $A$ , we may suppose that  $f(0) = 0$ . By Theorem 3.2, 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 3.3.  $\square$

REMARK 3.5. Theorem 3.4 implies that the Gâteaux and Fréchet derivatives coincide, are continuous in  $A$ , and are linear in  $E$ . These facts have very useful practical consequences, as Gâteaux derivatives are easy to compute, and because we may obtain derivatives at a matrix  $A$  as limits of derivatives at a sequence converging to  $A$ . More in detail, if we have a basis  $E_i$ ,  $i = 1, \dots, mn$ , of  $\mathbb{R}^{m \times n}$ , and if we can compute the  $mn$  Gâteaux derivatives  $G_{f^\circ}(A, E_i)$ , then for  $E = \sum_i \alpha_i E_i$  we can obtain  $L_{f^\circ}(A, E) = \sum_i \alpha_i G_{f^\circ}(A, E_i)$ . Moreover, if we have a converging sequence  $A_n \rightarrow A$ , and we can compute  $L_{f^\circ}(A_n, E)$ , then we can obtain  $L_{f^\circ}(A, E) = \lim_{n \rightarrow \infty} L_{f^\circ}(A_n, E)$ . These properties are crucial for the proof of Theorem 3.7.

The study of the differentiability of generalized matrix functions is subtler if we turn to complex matrices. In the real case, our proof of Theorem 3.1 is based on [5, Theorem 1], which in turn applies the analysis of [20, Section II.6.2] to the matrix  $\begin{bmatrix} 0 & A + tE \\ A^T + tE^T & 0 \end{bmatrix}$ . The theory of [20] applies to any matrix-valued function, possibly complex, which is Hermitian for any value of the *real* parameter  $t$ . Thus, if  $A \in \mathbb{C}^{m \times n}$  one may adapt the argument by starting from the matrix  $\begin{bmatrix} 0 & A + tE \\ A^* + tE^* & 0 \end{bmatrix}$ . Hence, by slightly modifying the proofs of Theorems 3.1, 3.2 and 3.4, it is not hard to show that complex generalized matrix functions are real differentiable under the same assumptions on  $f$ . On the other hand, if  $t$  is allowed to be *complex* then we cannot invoke the results of [20, Section II.6.2] any more. Indeed, not only the argument fails, but the conclusion does not hold: generalized matrix functions are generally<sup>3</sup> not complex differentiable (neither in the Gâteaux nor in the Fréchet sense), not even in the scalar case.

*Example 3.6.* Let us compute  $f^\circ(\rho + z) - f^\circ(\rho)$  for  $0 \neq \rho \in \mathbb{C}$  and  $z = \epsilon e^{i\zeta} \in \mathbb{C}$ . By item (ii) in Proposition 2.2, without loss of generality we may take  $\rho$  real and positive. Defining

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

we get

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

---

<sup>3</sup>That is, except a few trivial exceptions, e.g., linear functions.

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

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

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  $\epsilon \rightarrow 0^+$  with  $\zeta = \text{const.}$  in  $z = \epsilon e^{i\zeta}$  yields different results of the limit

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

depending on  $\zeta$ .

In summary, 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 will study the Fréchet derivative of complex generalized matrix functions in this context.

## 3.2 Real matrices

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

### Theorem 3.7. [Daleckii-Kreĭn Theorem for real generalized matrix functions]

Let  $A \in \mathbb{R}^{m \times n}$  have an SVD  $A = USV^T$ , 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  $f : \mathbb{R}_{\geq} \rightarrow \mathbb{R}$  be continuously differentiable on an open set containing the 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 Fréchet derivative at  $X = A$  of the generalized matrix function  $f^\diamond(X)$ , applied to the perturbation  $E$ , is

$$L_{f^\diamond}(A, E) = U \left( F \circ \widehat{E} + G \circ \Upsilon(\widehat{E}) \right) V^T, \quad (6)$$

where

- the symbol  $\circ$  denotes the Schur product;
- $\widehat{E} = U^T E V$ ;
- the linear operator  $\Upsilon$  is the following generalization of the transposition operator: for any  $X \in \mathbb{R}^{m \times n}$ ,  $\Upsilon(X) \in \mathbb{R}^{m \times n}$  and

$$\begin{aligned} & \text{if } m = n, & \Upsilon(X) &= X^T; \\ & \text{if } m > n \text{ and } X = \begin{bmatrix} X_1 \\ X_2 \end{bmatrix}, X_1 \in \mathbb{R}^{n \times n}, & \Upsilon(X) &= \begin{bmatrix} X_1^T \\ X_2 \end{bmatrix}; \\ & \text{if } m < n \text{ and } X = \begin{bmatrix} X_1 & X_2 \end{bmatrix}, X_1 \in \mathbb{R}^{m \times m}, & \Upsilon(X) &= \begin{bmatrix} X_1^T & X_2 \end{bmatrix}; \end{aligned}$$

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

$$F_{ij} = \begin{cases} \frac{\sigma_i f(\sigma_i) - \sigma_j f(\sigma_j)}{\sigma_i^2 - \sigma_j^2} & \text{if } i \neq j, \max(i, j) \leq \nu, \text{ and } \sigma_i \neq \sigma_j; \\ \frac{\sigma_i f'(\sigma_i) + f(\sigma_i)}{2\sigma_i} & \text{if } i \neq j, \max(i, j) \leq \nu, \text{ and } \sigma_i = \sigma_j \neq 0; \\ \frac{f(\sigma_j)}{\sigma_j} & \text{if } i > n, \text{ and } \sigma_j \neq 0; \\ \frac{f(\sigma_i)}{\sigma_i} & \text{if } j > m, \text{ and } \sigma_i \neq 0; \\ f'(\sigma_i) & \text{otherwise;} \end{cases} \quad (7)$$

$$G_{ij} = \begin{cases} \frac{\sigma_j f(\sigma_i) - \sigma_i f(\sigma_j)}{\sigma_i^2 - \sigma_j^2} & \text{if } i \neq j, i, j \leq \nu, \text{ and } \sigma_i \neq \sigma_j; \\ \frac{\sigma_i f'(\sigma_i) - f(\sigma_i)}{2\sigma_i} & \text{if } i \neq j, i, j \leq \nu, \text{ and } \sigma_i = \sigma_j \neq 0; \\ 0 & \text{otherwise.} \end{cases} \quad (8)$$

*Proof.* The formulae for  $\sigma_i = \sigma_j \neq 0$  can be obtained from those valid for  $\sigma_i \neq \sigma_j$  by setting  $\sigma_i = \sigma_j + h$  and by taking the limit  $h \rightarrow 0$ . Similarly, under the assumption that  $f(0) = 0$ , we can obtain the formulae for the case  $\sigma_i = \sigma_j = 0$  from those valid for  $\sigma_i = \sigma_j \neq 0$ , by going to the limit  $\sigma_i \rightarrow 0$ . Note that this argument is valid because the Fréchet derivative, when it exists, is continuous in  $A$ : see also Remark 3.5.

Moreover, item (ii) in Proposition 2.2 yields

$$f^\circ(A + tE) = Uf(S + tU^T E V)V^T \simeq f^\circ(A) + tUL_f(A, U^T E V)V^T,$$

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

The strategy of the proof is to first prove the result when  $E$  is zero except for one element, equal to 1. Using the fact that the Fréchet derivative is equal to the Gâteaux derivative, we will compute  $L_f(A, E)$  for such an  $E$  as the limit at the right hand side of (3). The result for a general  $E$  will then follow by linearity. We now examine a few separate cases according to the exact position of the unique nonzero element of  $E$ .

- *Case 1.* If the unique nonzero element of  $E$  is its  $i$ th diagonal element, then  $f^\circ(A + tE) - f^\circ(A) = \text{diag}(0, \dots, 0, f(\sigma_i + t) - f(\sigma_i), 0, \dots, 0)$ , where the nonzero element appears in the  $i$ th position. Dividing by  $t$  and going to the limit  $t \rightarrow 0$  we obtain  $\text{diag}(0, \dots, 0, f'(\sigma_i), 0, \dots, 0)$ , 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  $f^\circ(A + tE)$ . Let us take, without loss of generality (modulo applying a permutation equivalence),  $i = 1, j = 2$ . Then,  $A + tE = (U' \oplus I_{m-2})S'(V' \oplus I_{n-2})^T$  is a singular value decomposition if

$$(U')^T \begin{bmatrix} \sigma_i & t \\ 0 & \sigma_j \end{bmatrix} V' \quad (9)$$

is diagonal and  $U', V'$  are orthogonal. 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 orthogonal matrices is the subspace of skew-symmetric matrices. Imposing that (9) 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

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

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

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

with

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

and

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

- *Case 2b.* Consider now the case where the unique nonzero element of  $E$  lies in the position  $(i, j)$  with  $m < j \leq n$ . 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

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

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

- *Case 3.* Finally, 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  $i \leq m$  or to Case 2b. if  $i > m$ . We omit the details.

□

*Example 3.8.* Take  $f(x) = e^x$  and  $A = \begin{bmatrix} 2 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix}$ . Since  $A$  is full rank,  $f^\circ(X)$  is differentiable at  $X = A$  in spite of the fact that  $f(0) \neq 0$ . Moreover, Theorem 3.7 holds with

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

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

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

REMARK 3.9. We conclude this subsection with some observations on the forms of the matrices  $F$  and  $G$ :

- $F, G \in \mathbb{R}^{m \times n}$ ;
- if  $m = n$ ,  $F$  and  $G$  are symmetric;
- $G_{ii} = 0$  for all  $i = 1, \dots, \nu$ ;
- if  $m > n$ , then  $F = \begin{bmatrix} F_1 \\ ev^T \end{bmatrix}$  with  $F_1 = F_1^T \in \mathbb{R}^{n \times n}$ ,  $e = [1 \ \dots \ 1]^T \in \mathbb{R}^{m-n}$  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  $F = [F_1 \ ve^T]$  with  $F_1 = F_1^T \in \mathbb{R}^{m \times m}$ ,  $e = [1 \ \dots \ 1]^T \in \mathbb{R}^{n-m}$  and  $v \in \mathbb{R}^m$ , while  $G = [G_1 \ 0]$  where  $G_1 = G_1^T \in \mathbb{R}^{m \times m}$ .

### 3.3 Complex matrices

As we will see, the statement for complex matrices differs in two details from its analogue for real matrices: transposition is replaced by conjugate transposition and there is an additional term depending only on the imaginary part of the (rotated) perturbation. This is coherent with the discussion in Subsection 3.1, since the resulting expression is real linear, but *not* complex linear, in the elements of the perturbation.

**Theorem 3.10.** [Daleckiĭ-Kreĭn Theorem for complex generalized matrix functions]

Let  $A \in \mathbb{C}^{m \times n}$  have an SVD  $A = USV^*$ , where  $U \in \mathbb{C}^{m \times m}$ ,  $V \in \mathbb{C}^{n \times n}$ ,  $S \in \mathbb{R}^{m \times n}$ , and  $S_{ii} =: \sigma_i$ ,  $i = 1, \dots, \nu$ , are the singular values of  $A$ . Let  $f : \mathbb{R}_{\geq} \rightarrow \mathbb{R}$  be continuously differentiable on an open set containing the 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, if we see  $\mathbb{C}^{m \times n}$  as a  $2mn$ -dimensional real vector space, the Fréchet derivative



at  $X = A$  of the generalized matrix function  $f^\diamond(X)$ , applied to the complex perturbation  $E$ , is

$$L_{f^\diamond}(A, E) = U \left( F \circ \widehat{E} + iH \circ \Im \widehat{E} + G \circ \Upsilon(\widehat{E}) \right) V^*, \quad (10)$$

where

- the symbol  $\circ$  denotes the Schur product;
- $\widehat{E} = U^* E V$ , 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

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

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

$$\text{if } m < n \text{ and } X = [X_1 \ X_2], X_1 \in \mathbb{C}^{m \times m}, \quad \Upsilon(X) = [X_1^* \ X_2];$$

- $F, G \in \mathbb{R}^{m \times n}$  are defined as in (7) and in (8), respectively;
- and  $H \in \mathbb{R}^{m \times n}$  is diagonal with  $H_{ii} = f(\sigma_i)/\sigma_i - F_{ii}$  if  $\sigma_i \neq 0$  and  $H_{ii} = f'(0) - F_{ii}$  if  $\sigma_i = 0$ .

*Proof.* The proof is very similar to the one of Theorem 3.7, except that we need to consider separately the cases of a perturbation  $E$  with only one real nonzero element and of a perturbation  $E$  with only one pure imaginary nonzero element. The former case is precisely the same as in the proof of Theorem 3.7, so we omit the details. We now argue on the latter case, i.e.,  $E$  pure imaginary with only one nonzero element equal to  $i$ .

- *Case 1.* If the unique nonzero element of  $E$  is its  $i$ th diagonal element, then arguing as in Example 3.6 we readily obtain  $f^\diamond(A+tE) - f^\diamond(A) = \text{diag}(0, \dots, 0, itf(\sigma_i)/\sigma_i, 0, \dots, 0)$ , where the nonzero element appears in the  $i$ th position. Dividing by  $t$  and going to the limit, we obtain  $i \text{diag}(0, \dots, 0, f(\sigma_i)/\sigma_i, 0, \dots, 0)$ .
- *Case 2a.* Suppose now that the unique nonzero element of  $E$  is in the position  $(i, j)$  with  $i < j \leq \nu$ . As in Theorem 3.7 we can take without loss of generality  $i = 1, j = 2$ . We look for a singular value decomposition  $A + tE = (U' \oplus I_{m-2})S'(V' \oplus I_{n-2})^*$ , and therefore we impose that

$$(U')^* \begin{bmatrix} \sigma_i & it \\ 0 & \sigma_j \end{bmatrix} V' \quad (11)$$

is real and diagonal and  $U', V'$  are unitary. We expand  $U' = I_2 + t \begin{bmatrix} 0 & u^* \\ -u & 0 \end{bmatrix} + O(t^2)$  and  $V' = I_2 + t \begin{bmatrix} 0 & v^* \\ -v & 0 \end{bmatrix} + O(t^2)$ . Retaining only the  $O(t)$  terms, assuming  $\sigma_i \neq \sigma_j$ , and solving for  $u, v$  we get

$$u = \frac{i\sigma_j}{\sigma_i^2 - \sigma_j^2}, \quad v = \frac{i\sigma_i}{\sigma_i^2 - \sigma_j^2}$$

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

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

with

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

and

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

- *Case 2b.* If the unique nonzero element of  $E$  lies in the position  $(i, j)$  with  $m < j \leq n$ , the procedure is again analogous to Theorem 3.7. Assume that  $i = 1, j = m + 1$ :  $A + tE = S'(V')^*$  is a singular value decomposition if

$$\begin{bmatrix} \sigma_i & 0 & \dots & 0 & it & 0 & \dots & 0 \end{bmatrix} V' = \begin{bmatrix} \sigma & 0 & \dots & 0 \end{bmatrix},$$

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

- *Case 3.* Finally, 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  $i \leq m$ , or to Case 2b., if  $i > m$ . We omit the details.

□

## 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  $f(0) \neq 0$ . To keep the paper within a reasonable

length, and since the real case is the most relevant for the applications [4], we focus on generalized matrix functions of real matrices and only allow real perturbations.

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

$$\text{cond } f^\circ(A) = \limsup_{t \rightarrow 0} \sup_{E \neq 0} \frac{\|f^\circ(A + tE) - f^\circ(A)\|}{t\|E\|}. \quad (12)$$

There are two cases. If  $f(0) \neq 0$  and  $A$  is rank deficient, then clearly  $\text{cond } f^\circ(A) = \infty$ , as  $f^\circ(X)$  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,  $f^\circ(X)$  is differentiable at  $X = A$ , and  $\|f^\circ(A + tE) - f^\circ(A)\| = \|L_{f^\circ}(A, tE) + O(t^2)\| = |t|\|L_{f^\circ}(A, E)\| + O(t^2)$ . If we specialize to any unitarily invariant norm, and if  $A = USV^T$  is an SVD, it is immediate that  $\|L_{f^\circ}(A, E)\| = \|L_{f^\circ}(S, \hat{E})\|$  having defined  $\hat{E} = U^T E V$ . For example, the Frobenius norm is unitarily invariant, and this choice leads to the condition number

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

where  $K_{f^\circ}(X)$  is the Kronecker form of the Fréchet derivative [17] of  $f^\circ$  at  $X \in \mathbb{R}^{m \times n}$ . To define  $K_{f^\circ}(X)$  it is convenient to introduce the vec operator [15]

$$\text{vec} : \mathbb{C}^{m \times n} \rightarrow \mathbb{C}^{mn}, X = \begin{bmatrix} \mathbf{x}_1 & \dots & \mathbf{x}_n \end{bmatrix} \mapsto \text{vec}(X) = \begin{bmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_n \end{bmatrix}.$$

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

Let us now consider the linear map  $\Upsilon$ , defined in the statement of Theorem 3.7. Via the vec operator, it can be represented by the unique matrix  $P \in \mathbb{R}^{mn \times mn}$  satisfying

$$P \text{vec}(A) = \text{vec}(\Upsilon(A)) \quad \forall A \in \mathbb{R}^{m \times n}. \quad (13)$$

Note that in the special case  $m = n$  we recover the well-studied vec-permutation operator [15].

**Lemma 4.1.** *The matrix  $P$  defined in (13) 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

$$\mathcal{V}_1 = \left\{ X = \begin{bmatrix} X_1 \\ X_2 \end{bmatrix} \in \mathbb{R}^{m \times n} \mid X_1 = X_1^T \in \mathbb{R}^{n \times n} \right\} \text{ and } \mathcal{V}_2 = \left\{ X = \begin{bmatrix} X_1 \\ 0 \end{bmatrix} \in \mathbb{R}^{m \times n} \mid X_1 = -X_1^T \in \mathbb{R}^{n \times n} \right\}$$

$\mathbb{R}^{n \times n}$ . 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

**Corollary 4.2.** *Let  $A \in \mathbb{R}^{m \times n}$ , and suppose  $A = USV^T$  is an SVD. Let  $f : \mathbb{R}_{\geq} \rightarrow \mathbb{R}$  be differentiable on an open set containing the 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 generalized matrix function  $f^\diamond(X)$  is*

$$K_{f^\diamond}(A) = (V \otimes U)(\Phi + \Gamma P)(V^T \otimes U^T), \quad (14)$$

where  $\Phi = \text{diag}(\text{vec}(F))$ ,  $\Gamma = \text{diag}(\text{vec}(G))$ ,  $F$  and  $G$  are defined as in Theorem 3.7 and  $P$  is the matrix defined by (13).

*Proof.* It is a corollary of Theorem 3.7. Indeed, applying the  $\text{vec}$  operator to (6), 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  $\text{vec}(\hat{E}) = \text{vec}(U^T E V) = (V^T \otimes U^T) \text{vec}(E)$ .  $\square$

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

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

Similarly, it is easy to check that

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

Applying the  $\text{vec}$  operator to each of these equations yields the statement.  $\square$

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

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

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

**Theorem 4.4.** *It holds*

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

where  $a = \max_i |F_{ii}|$ ,  $b = \max_{j \leq n \leq i} |F_{ij}|$ ,  $c = \max_{i < j} |F_{ij} + G_{ij}|$  and  $d = \max_{i < j} |F_{ij} - G_{ij}|$ , and  $F$  and  $G$  are defined as in Theorem 3.7.

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 Theorem 3.7 by a direct computation.  $\square$

## 4.2 On the conditioning of real generalized matrix functions

If all the singular values of the matrix  $A$  are known then Theorems 3.7, 4.4, and 4.5 can be combined to compute  $\|K_{f^\diamond}(S)\|_2$ , and hence  $\text{cond}_F f^\diamond(A)$ . In this subsection, we give some upper bounds for  $\text{cond}_F f^\diamond(A)$  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 [11, 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*

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

*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 3.7, 4.4, and 4.5.

More precisely, in the notation of Theorem 4.4, if  $A$  has full rank then  $a \leq K$ ,  $b \leq M\sigma_r^{-1}$ ,  $c \leq K$ ,  $d \leq M\sigma_r^{-1}$ . If  $A$  is rank deficient, then  $a \leq K$ ,  $b \leq \max\{K, M\sigma_r^{-1}\}$ ,  $c \leq K$ ,  $d \leq \max\{K, M\sigma_r^{-1}\}$ .  $\square$

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*

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

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

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 function 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  $A$  [17, 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 [21, Theorems 2.2 and 2.3] that the absolute condition number is

$$\text{cond}_F f^\diamond(A) = \frac{2}{\sigma_r + \sigma_{r-1}},$$

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  $1/\sigma_r$ , which is slack for  $\sigma_r < \sigma_{r-1}$ , specializing Theorem 4.4 to  $f = 1$  gives  $a = c = 0$  and  $d = 2/(\sigma_r + \sigma_{r-1})$ . 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 [6] the absolute condition number is  $1/\sigma_r$ , 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,  $\text{cond}_F \exp^\diamond(A) \leq \exp(\|A\|_2) \max\{1, \sigma_r^{-1}\}$ . 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

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

In practice, quoting Nick Higham [17, 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  $\text{cond}_F f^\diamond(A)$  by the formula

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

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

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

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

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

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  $\mu := \sqrt{f(\|A\|_2)^2 + f(\sigma_r)^2}$ , the relative condition number of the generalized matrix function  $f^\diamond(A)$  is bounded above by*

$$\sqrt{\nu}\|A\|_2 \max\left\{\frac{K}{\mu}, \frac{1}{\sigma_r}\right\}.$$

*In the notation and under the assumptions of Theorem 4.7, the relative condition number of the generalized matrix function  $f^\diamond(A)$  is bounded above by*

$$\frac{\sqrt{\nu}\|A\|_2 K}{\mu}.$$

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

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

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

$$\frac{1}{\epsilon\sqrt{5}} + O(1).$$

Figure 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

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



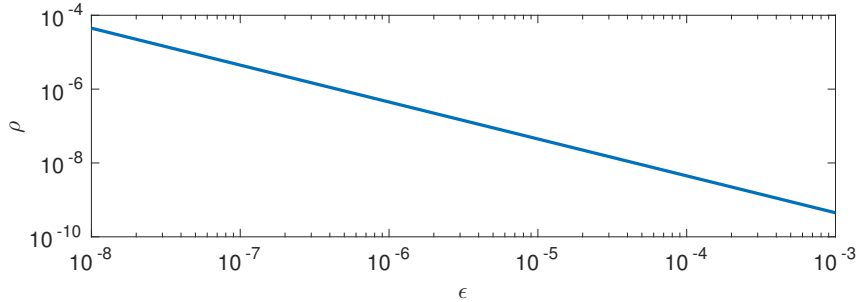


Figure 1: Computed relative error for Example 4.12

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 of the eigenvector matrix  $Z$  in Theorem 2.10; on the other hand, since  $U$  and  $V$  in Theorem 3.7 are orthogonal, the information on the conditioning of generalized matrix functions is directly encoded in the Daleckiĭ-Kreĭn formula developed in this paper.

An exception to this generally optimistic situation is when  $f(0) \neq 0$  and  $f^\circ(A)$  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  $f^\circ(A)$ .

## 5 Acknowledgements

I would like to thank Kate Fenu and Yuji Nakatsukasa for useful discussions on generalized matrix functions and on the continuity of the singular value decomposition. I am also grateful to Michele Benzi for pointing me to reference [8].

## References

- [1] F. ANDERSSON, M. CARLSSON AND K.-M. PERFECT, *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] 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.
- [6] 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.
- [7] J. L. DALECKIĀ AND S. G. KREĀN, *Integration and differentiation of functions of Hermitian operators and applications to the theory of perturbations*, Amer. Math. Soc. Transl., Series 2, 47 (1965), pp. 1–30.
- [8] L. FANTAPPIÉ, *Le calcul des matrices*, C. R. Ac. des Sc. Paris 186 (1928), pp. 619–621.
- [9] L. FANTAPPIÉ, *Sulle funzioni di una matrice*, An. Acad. Brasil. Cienc. 26 (1954), pp. 25–33.
- [10] 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.
- [11] G. H. GOLUB AND C. VAN LOAN, *Matrix Computations*, 4th edition, John Hopkins University Press, Baltimore, MD, United States, 2012.
- [12] R. S. HAMILTON, *The inverse function theorem of Nash and Moser*, Bulletin of the AMS 7(1) (1982), pp. 65–222.
- [13] J. B. HAWKINS AND A. BEN–ISRAEL, *On generalized matrix functions*, Linear and Multilinear Algebra 1(2) (1973), pp. 163–171.
- [14] J. HEINONEN, *Lectures on Lipschitz analysis*, Rep. Univ. Jyväskylä Dept. Math. Stat. 100 (2005), 1–77.
- [15] 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.
- [16] N. J. HIGHAM, *Computing the Polar Decomposition — with Applications*, SIAM J. Sci. Stat. Comput. 7(1) (1986), pp. 1160–1174.
- [17] N. J. HIGHAM, *Functions of Matrices: Theory and Computation*, SIAM, Philadelphia, PA, United States, 2008.

- [18] R. A. HORN AND C. R. JOHNSON, *Matrix Analysis*, 2nd edition, *Cambridge University Press*, New York, NY, United States, 2013.
- [19] R. A. HORN AND C. R. JOHNSON, *Topics in Matrix Analysis*, *Cambridge University Press*, New York, NY, United States, 1991.
- [20] T. KATO, *Perturbation Theory for Linear Operators*, *Springer-Verlag*, New York, NY, United States, 1966.
- [21] 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.
- [22] A. J. KURDILA AND M. ZABARANKIN, *Convex Functional Analysis*, *Birkhäuser-Verlag*, Basel, Switzerland, 2005.
- [23] R. PENROSE, *A generalized inverse for matrices*, *Proc. Cambridge Philos. Soc.* 51 (1955), pp. 406–413.
- [24] W. RUDIN, *Principles of Mathematical Analysis*, *McGraw-Hill*, New York, NY, United States, 1976.