

A Service of

ZBW

Leibniz-Informationszentrum Wirtschaft Leibniz Information Centre for Economics

Mitchell, Tim; Overton, Michael L.

# Article — Published Version On properties of univariate max functions at local maximizers

**Optimization Letters** 

### **Provided in Cooperation with:** Springer Nature

*Suggested Citation:* Mitchell, Tim; Overton, Michael L. (2022) : On properties of univariate max functions at local maximizers, Optimization Letters, ISSN 1862-4480, Springer, Berlin, Heidelberg, Vol. 16, Iss. 9, pp. 2527-2541, https://doi.org/10.1007/s11590-022-01872-y

This Version is available at: https://hdl.handle.net/10419/308732

#### Standard-Nutzungsbedingungen:

Die Dokumente auf EconStor dürfen zu eigenen wissenschaftlichen Zwecken und zum Privatgebrauch gespeichert und kopiert werden.

Sie dürfen die Dokumente nicht für öffentliche oder kommerzielle Zwecke vervielfältigen, öffentlich ausstellen, öffentlich zugänglich machen, vertreiben oder anderweitig nutzen.

Sofern die Verfasser die Dokumente unter Open-Content-Lizenzen (insbesondere CC-Lizenzen) zur Verfügung gestellt haben sollten, gelten abweichend von diesen Nutzungsbedingungen die in der dort genannten Lizenz gewährten Nutzungsrechte.



WWW.ECONSTOR.EU

https://creativecommons.org/licenses/by/4.0/

#### Terms of use:

Documents in EconStor may be saved and copied for your personal and scholarly purposes.

You are not to copy documents for public or commercial purposes, to exhibit the documents publicly, to make them publicly available on the internet, or to distribute or otherwise use the documents in public.

If the documents have been made available under an Open Content Licence (especially Creative Commons Licences), you may exercise further usage rights as specified in the indicated licence.



**ORIGINAL PAPER** 



# On properties of univariate max functions at local maximizers

Tim Mitchell<sup>1</sup> · Michael L. Overton<sup>2</sup>

Received: 17 August 2021 / Accepted: 28 February 2022 / Published online: 28 March 2022 © The Author(s) 2022

#### Abstract

More than three decades ago, Boyd and Balakrishnan established a regularity result for the two-norm of a transfer function at maximizers. Their result extends easily to the statement that the maximum eigenvalue of a univariate real analytic Hermitian matrix family is twice continuously differentiable, with Lipschitz second derivative, at all local maximizers, a property that is useful in several applications that we describe. We also investigate whether this smoothness property extends to max functions more generally. We show that the pointwise maximum of a finite set of q-times continuously differentiable univariate functions must have zero derivative at a maximizer for q = 1, but arbitrarily close to the maximizer, the derivative may not be defined, even when q = 3 and the maximizer is isolated.

**Keywords** Univariate max functions · Eigenvalues of Hermitian matrix families · H-infinity norm · Numerical radius · Optimization of passive systems

#### **1** Introduction

Let  $\mathbb{H}^n$  denote the space of  $n \times n$  complex Hermitian matrices, let  $\mathcal{D} \subseteq \mathbb{R}$  be open, and let  $H : \mathcal{D} \to \mathbb{H}^n$  denote an analytic Hermitian matrix family in one real variable, i.e., for each  $x \in \mathcal{D}$  and each  $i, j \in \{1, ..., n\}$ , there exist coefficients  $a_0, a_1, a_2, ...$ such that the power series  $\sum_{k=0}^{\infty} a_k (t-x)^k$  converges to  $H_{ij}(t) = \overline{H_{ji}(t)}$  for all t in a neighborhood of x. For a generic family H, the eigenvalues of H(t) are simple for all

Tim Mitchell mitchell@mpi-magdeburg.mpg.de
 Michael L. Overton mo1@nyu.edu

<sup>&</sup>lt;sup>1</sup> Max Planck Institute for Dynamics of Complex Technical Systems, Sandtorstr. 1, 39106 Magdeburg, Germany

<sup>&</sup>lt;sup>2</sup> Courant Institute of Mathematical Sciences, New York University, 251 Mercer St., New York, NY 10012, USA

 $t \in \mathcal{D}$ ; often known as the von Neumann–Wigner crossing-avoidance rule [15], this phenomenon is emphasized in [11, section 9.5], where it is also illustrated on the front cover. The reason is simple: the real codimension of the subspace of Hermitian matrices with an eigenvalue of multiplicity *m* is  $m^2 - 1$ , so to obtain a double eigenvalue one would need three parameters generically; when the matrix family is real symmetric, the analogous codimension is  $\frac{m(m+1)}{2} - 1$ , so one would need two parameters generically. When there are no multiple eigenvalues, the ordered eigenvalues of H(t), say,  $\mu_j(t)$ for  $j = 1, \ldots, n$ , are all real analytic functions. Let  $\lambda_{\max} : \mathbb{H}^n \to \mathbb{R}$  and  $\lambda_{\min} : \mathbb{H}^n \to \mathbb{R}$ denote largest and smallest eigenvalue, respectively. In the absence of multiple eigenvalues,  $\lambda_{\max} \circ H$  and  $\lambda_{\min} \circ H$  are both smooth functions of *t*. However, for the nongeneric family H(t) = diag(t, -t), a double eigenvalue occurs at t = 0. By a theorem of Rellich, given in Sect. 4, the eigenvalues can be written as two real analytic functions,  $\mu_1(t) = t$  and  $\mu_2(t) = -t$ , but we must give up the property that these functions are ordered near zero. Consequently, the function  $\lambda_{\max} \circ H$  is not differentiable at its minimizer t = 0.

In contrast, the function  $\lambda_{\max} \circ H$  is *unconditionally*  $C^2$ , i.e., twice continuously differentiable, with Lipschitz second derivative, near all its local *maximizers*, regardless of eigenvalue multiplicity at these maximizers. As we explain below, this observation is a straightforward extension of a well-known result of Boyd and Balakrishnan [2] established more than three decades ago. One purpose of this paper is to bring attention to the more general result, as it is useful in a number of applications. We also investigate whether this smoothness property extends to max functions more generally. We show that the pointwise maximum of a finite set of differentiable univariate functions must have zero derivative at a maximizer. However, arbitrarily close to the maximizer, the derivative may not be defined, even if the functions are three times continuously differentiable and the maximizer is isolated.

#### 2 Properties of max functions at local maximizers

Let  $\mathcal{D} \subset \mathbb{R}$  be open,  $\mathcal{I} = \{1, ..., n\}$ , and  $f_j : \mathcal{D} \to \mathbb{R}$  be continuous for all  $j \in \mathcal{I}$ , and define

$$f_{\max}(t) := \max_{j \in \mathcal{I}} f_j(t).$$
(2.1)

**Lemma 2.1** Let  $x \in D$  be any local maximizer of  $f_{\max}$  with  $f_{\max}(x) = \gamma$  and let  $\mathcal{I}_{\gamma} = \{j \in \mathcal{I} : f_j(x) = \gamma\}$ . Then

(i) for all  $j \in I_{\gamma}$ , x is a local maximizer of  $f_j$  and (ii) for all  $j \in I \setminus I_{\gamma}$ ,  $f_j(x) < \gamma$ .

We omit the proof as it is elementary.

**Theorem 2.1** Let  $x \in D$  be any local maximizer of  $f_{\max}$  with  $f_{\max}(x) = \gamma$ . Suppose that for all  $j \in I$ ,  $f_j$  is differentiable at x. Then  $f_{\max}$  is differentiable at x with  $f'_{\max}(x) = 0$ .

**Proof** Since the functions  $f_j$  are assumed to be continuous, clearly  $f_{\text{max}}$  is also continuous, and without loss of generality, we can assume that  $\gamma = 0$ . Suppose that  $f'_{\text{max}}$  does not exist at x or does not equal zero, i.e., there exists some sequence  $\{\varepsilon_k\}$  with  $\varepsilon_k \to 0$  such that  $\lim_{k\to\infty} \frac{f_{\max}(x+\varepsilon_k)}{\varepsilon_k}$  does not exist or is not zero. Since  $\mathcal{I}$  is finite, there exist a  $j \in \mathcal{I}$  and a subsequence  $\{\varepsilon_{k_\ell}\}$  such that  $f_{\max}(x+\varepsilon_{k_\ell}) = f_j(x+\varepsilon_{k_\ell})$  for all  $k_\ell$ , which implies that  $f'_j$  either does not exist or is not zero at x. However, as  $f_j$  is differentiable and with local maximizer x by Lemma 2.1, it must be that  $f'_j(x) = 0$ ; hence, we have a contradiction.

We now consider adding additional assumptions on the smoothness of the  $f_j$ , writing  $f_j \in C^q$  to mean  $f_j$  is q-times continuously differentiable. Clearly, assuming that the  $f_j$  are  $C^1$  at (or near) a maximizer is not sufficient to obtain that  $f_{\text{max}}$  is twice differentiable at this point. For example, if

$$f_1(t) = \begin{cases} -t^2 & \text{if } t \le 0\\ -3t^2 & \text{if } t > 0 \end{cases} \text{ and } f_2(t) = -2t^2,$$

then the second derivative of  $f_{\text{max}} = \max(f_1, f_2)$  does not exist at the maximizer t = 0, as  $f'_{\text{max}}(t) = -2t$  on the left and -4t on the right, so  $\lim_{t\to 0} \frac{f'_{\text{max}}(t)}{t}$  does not exist at t = 0. In this example,  $f_{\text{max}}$  is continuously differentiable at t = 0, but this does not hold in general, even when assuming that the  $f_j$  are  $C^3$  near a maximizer; see Remark 2.1 below. However, we do have the following result.

**Theorem 2.2** Let  $x \in D$  be any local maximizer of  $f_{\max}$  with  $f_{\max}(x) = \gamma$ . Suppose that for all  $j \in I$ ,  $f_j$  is  $C^3$  near x. Then for all sufficiently small  $|\varepsilon|$ ,

$$f_{\max}(x+\varepsilon) = \gamma + M\varepsilon^2 + O(|\varepsilon|^3), \qquad (2.2)$$

where  $M = \frac{1}{2} \left( \max_{j \in \mathcal{I}_{\gamma}} f_{j}''(x) \right) \leq 0$ . If the  $C^{3}$  assumption is reduced to  $C^{2}$ , then  $f_{\max}(x + \varepsilon) = \gamma + O(\varepsilon^{2})$ .

**Proof** Let  $\gamma = f_{\max}(x)$  and let  $\mathcal{I}_{\gamma} = \{j \in \mathcal{I} : f_j(x) = \gamma\}$ . By Lemma 2.1, we have that *x* is also a local maximizer of  $f_j$  for all  $j \in \mathcal{I}_{\gamma}$  and  $f_j(x) < \gamma$  for all  $j \in \mathcal{I} \setminus \mathcal{I}_{\gamma}$ . Since the  $f_j$  are Lipschitz near *x*,

$$f_{\max}(x+\varepsilon) = \max_{j \in \mathcal{I}_{\gamma}} f_j(x+\varepsilon)$$

holds for all sufficiently small  $|\varepsilon|$ . For each  $j \in \mathcal{I}_{\gamma}$ , by Taylor's Theorem we have that

$$f_j(x+\varepsilon) = f_j(x) + f'_j(x)\varepsilon + \frac{1}{2}f''_j(x)\varepsilon^2 + \frac{1}{6}f'''_j(\tau_j)\varepsilon^3$$
$$= \gamma + \frac{1}{2}f''_j(x)\varepsilon^2 + O(|\varepsilon|^3)$$

for  $\tau_j$  between x and  $x + \varepsilon$ . Taking the maximum of the equation above over all  $j \in \mathcal{I}_{\gamma}$  yields (2.2). The proof for the  $C^2$  case follows analogously.

**Remark 2.1** Even with the  $C^3$  assumption,  $f_{\text{max}}$  is not necessarily continuously differentiable at maximizers. For example, consider  $f_1(t) = t^8(\sin(\frac{1}{t}) - 1)$  and  $f_2(t) = t^8(\sin(\frac{1}{2t}) - 1)$ , with  $f_1(0) = f_2(0) = 0$ , where  $f_1$  and  $f_2$  are  $C^3$  but not  $C^4$  at the maximizer t = 0. Although  $f_{\text{max}}$  is differentiable at 0 by Theorem 2.1, it is easy to see that it is not  $C^1$  there. However, in this case, 0 is not an isolated maximizer of  $f_{\text{max}}$ . In contrast, in Sect. 3, we construct a counterexample where the  $f_j$  are  $C^3$  functions, and for which  $f_{\text{max}}$  has an isolated maximizer, yet  $f_{\text{max}}$  is not  $C^1$  there. It seems that this counterexample can be extended to apply to  $C^q$  functions for any  $q \ge 1$ . The key point of both of these counterexamples is not that the  $f_j$  are insufficiently smooth *per se*, but that the  $f_j$  cross each other infinitely many times near maximizers.

In light of Remark 2.1, we now make a much stronger assumption.

**Theorem 2.3** Given a maximizer x of  $f_{\max}$ , suppose there exist  $j_1, j_2 \in \mathcal{I}$ , possibly equal, such that, for all sufficiently small  $\varepsilon > 0$ ,  $f_{\max}(x - \varepsilon) = f_{j_1}(x - \varepsilon)$  and  $f_{\max}(x + \varepsilon) = f_{j_2}(x + \varepsilon)$ , with  $f_{j_1}$  and  $f_{j_2}$  both  $C^3$  near x. Then  $f_{\max}$  is twice continuously differentiable, with Lipschitz second derivative, near x.

**Proof** It is clear that  $f_{j_1}(x) = f_{j_2}(x) = \gamma$  and  $f'_{j_1}(x) = f'_{j_2}(x) = 0$ . By Theorem 2.2, both  $f''_{j_1}(x)$  and  $f''_{j_2}(x)$  are equal to M, so  $f_{\max}$  is locally described by two  $C^3$  pieces whose function values and first and second derivatives agree at x. Hence,  $f_{\max}$  is  $C^2$  with Lipschitz second derivative near x.

*Remark 2.2* The assumptions of Theorem 2.3 hold if the  $f_j$  are real analytic [10, Corollary 1.2.7].

In particular, the  $f_j$  are real analytic functions if they are eigenvalues of a univariate real analytic Hermitian matrix function, as we discuss in Sect. 4. First, we present the  $C^3$  counterexample mentioned above.

# 3 An example with $C^3$ functions $f_j$ and an isolated maximizer for which $f_{max}$ is not continuously differentiable at x

Let  $l_k = \frac{1}{2^k}$ , and  $f_1 : [-1, 1] \to \mathbb{R}$  be defined by

$$f_1(t) = \begin{cases} p_k(t) & \text{if } t \in [l_{k+1}, l_k] \text{ for } k = 0, 1, 2, \dots \\ -t^2 & \text{if } t \in [-1, 0] \end{cases}$$

where  $p_k$  is a (piece of a) degree-nine polynomial chosen such that at

- 1.  $l_{k+1}$  (the left endpoint),  $p_k$  and  $-t^2$  agree up to and including their respective third derivatives,
- 2.  $l_k$  (the right endpoint),  $p_k$  and  $-t^2$  agree up to and including their respective third derivatives,



**Fig. 1** Plots of  $f_1$  and  $f_2$  with  $s_k = 2$ ; their  $-t^2$  parts are shown in solid, while their  $p_k$  parts are shown in dotted for k even and dash-dot for k odd

3.  $t_k = \frac{1}{2}(l_{k+1} + l_k)$  (the midpoint),  $p_k$  and  $-t^2$  agree, but the first derivative of  $p_k$  is  $s_k \neq 1$  times the value of the first derivative of  $-t^2$ .

For any k, the degree-nine polynomial  $p_k$  is uniquely determined by the ten algebraic constraints given above. If we choose  $s_k = 1$ , then  $p_k$  is simply  $-t^2$ . However, by choosing  $s_k > 1$  but sufficiently close to 1, then  $p_k$  must be strictly decreasing between its endpoints  $l_{k+1}$  and  $l_k$  and cross  $-t^2$  at  $t_k$ . If this is done for all k, it follows that t = 0 must be an isolated maximizer of  $f_1$ . See Fig. 1a for a plot of  $f_1$  with  $s_k = 2$ for all k; the choice  $s_k = 2$  is not close to 1 but was chosen to make the features of  $f_1$ easily seen.

Now define  $f_2(t) = f_1(-t)$ , i.e., the graph of  $f_2$  is a reflection of the graph of  $f_1$  across the vertical line t = 0. Figure 1b shows  $f_1$  and  $f_2$  plotted together, again with  $s_k = 2$ , showing how they cross at every  $t_k$ . Recall that by our construction, their respective first three derivatives match at each  $l_k$ , but their first derivatives do not match at any  $t_k$ . Figure 2 shows plots of the first three derivatives of  $f_1$  for two different sequences  $\{s_k\}$  respectively defined by  $s_k = 1 + 2^{-k}$  and  $s_k = 1 + 2^{-2k}$ . The rightmost plots in Fig. 2 indicate that the first choice for sequence  $\{s_k\}$  does not converge to 1 fast enough for  $f_1'''$  to exist and be continuous at t = 0, but that the second sequence does. In fact, for this latter choice of sequence, we have the following pair of theorems respectively proving that  $f_1$  is indeed  $C^3$  with t = 0 being an isolated maximizer. We defer the proofs to Appendix A as they are a bit technical, and in Appendix B, we discuss why  $s_k = 1 + 2^{-k}$  does not converge to 1 sufficiently fast for  $f_1'''(0)$  to exist.

**Theorem 3.1** For  $f_1$  defined in (3.1), if  $s_k = 1 + 2^{-2k}$ , then  $f_1$  is  $C^3$  on its domain [-1, 1].

**Theorem 3.2** For  $f_1$  defined in (3.1), if  $s_k = 1 + 2^{-2k}$ , then t = 0 is an isolated maximizer of  $f_1$ , as well as an isolated maximizer of  $f_{max} = \max(f_1, f_2)$ .

Theorem 2.1 shows that  $f_{\text{max}} = \max(f_1, f_2)$  is differentiable at t = 0 with  $f'_{\text{max}}(0) = 0$ . However, even though  $f_1$  and  $f_2$  are  $\mathcal{C}^3$  and t = 0 is an isolated



**Fig. 2** Plots of the first three derivatives of  $f_1$  for two different sequences  $\{s_k\}$ ; their  $-t^2$  parts are shown in solid, while their  $p_k$  parts are shown in dotted for k even and dash-dot for k odd

maximizer of  $f_{\text{max}}$  with the choice of  $s_k = 1 + 2^{-2k}$ , by construction, we have that (i)  $t_k \to 0$  as  $k \to \infty$ , and (ii)  $f_{\text{max}}$  is nondifferentiable at every  $t_k$ . Hence, although  $f_{\text{max}}$  is differentiable at t = 0, it is not  $C^1$  at this point. Plots of  $f_{\text{max}}$  and its first and second derivatives are shown in Fig. 3, where we see the discontinuities in  $f'_{\text{max}}$  for all  $t_k$  and  $-t_k$ .

**Remark 3.1** For any  $q \ge 1$ , it seems that the same argument extends to show that  $f_{\text{max}}$  is not necessarily  $C^1$  at t = 0 when defined by functions  $f_j$  that are  $C^q$ , using polynomials  $p_k$  of degree 2q + 3. From computational investigations for  $q \in \{1, 2, 3, 4, 5\}$ , we conjecture that  $s_k = 1+2^{-(k+1)}$  for q = 1 and  $s_k = 1+2^{-(q-1)k}$  for  $q \ge 2$  are suitable choices in general to obtain that  $f_1$  is  $C^q$  with t = 0 being an isolated maximizer. It is not clear how to extend such an argument to the  $C^{\infty}$  case.

#### 4 Smoothness of eigenvalue extrema and applications

We will need the following well-known theorem, whose history is discussed in Kato's treatise [9, pp. XI–XIII]; specifically, see Eq. (2.2) on p. XII and Theorem II-6.1 on p. 139 of [9].



Fig. 3 Plots of  $f_{\text{max}}$  and its first and second derivatives; their  $-t^2$  parts are shown in solid, while their  $p_k$  parts are shown in dash-dot

**Theorem 4.1** [*Rellich*] Let  $H : \mathcal{D} \to \mathbb{H}^n$  be an analytic Hermitian matrix family in one real variable. Let  $x \in \mathcal{D}$  be given, and let H(x) have eigenvalues  $\tilde{\mu}_j \in \mathbb{R}$ , j = 1, ..., n, not necessarily distinct. Then, for sufficiently small  $|\varepsilon|$ , the eigenvalues of  $H(x + \varepsilon)$  can be expressed as convergent power series

$$\mu_j(\varepsilon) = \tilde{\mu}_j + \tilde{\mu}_j^{(1)}\varepsilon + \tilde{\mu}_j^{(2)}\varepsilon^2 + \cdots, \quad j = 1, \dots, n.$$

$$(4.1)$$

We now apply Theorems 4.1 and 2.3 to obtain smoothness results for eigenvalue extrema of univariate real analytic Hermitian matrix families, as well as analogous results for singular value extrema. Subsequently, we discuss how these results are useful in several important applications.

**Theorem 4.2** Let  $H : \mathcal{D} \to \mathbb{H}^n$  be an analytic Hermitian matrix family in one real variable on an open domain  $\mathcal{D} \subseteq \mathbb{R}$ , and let  $\lambda_{\max} : \mathbb{H}^n \to \mathbb{R}$  denote largest eigenvalue. Then  $\lambda_{\max} \circ H$  is  $\mathcal{C}^2$  with Lipschitz second derivative near all of its local maximizers.

**Proof** Let  $x \in D$  be any local maximizer of  $\lambda_{\max} \circ H$ , with H(x) having eigenvalues  $\tilde{\mu}_j$ . By Theorem 4.1, in a neighborhood of x, the eigenvalues of  $H(x + \varepsilon)$  can be expressed as  $\mu_j(\varepsilon)$ , j = 1, ..., n, where the  $\mu_j(\varepsilon)$  are locally given by the power series (4.1). Since  $\lambda_{\max}(H(x + \varepsilon)) = \max_{j \in \{1,...,n\}} \mu_j(\varepsilon)$  with all the  $\mu_j$  analytic, we can apply Theorem 2.3 to these functions, as mentioned in Remark 2.2, thus completing the proof.

*Remark 4.1* The proof of Theorem 4.2 is essentially the same as the proof given by Boyd and Balakrishnan [2], presented differently and in a more general context.

**Corollary 4.1** Let  $H : \mathcal{D} \to \mathbb{H}^n$  be an analytic Hermitian matrix family in one real variable on an open domain  $\mathcal{D} \subseteq \mathbb{R}$ . Then:

- (i)  $\lambda_{\min} \circ H$  is  $C^2$  near all of its local minimizers, where  $\lambda_{\min}$  denotes algebraically smallest eigenvalue;
- (ii)  $\rho \circ H$  is  $C^2$  near all of its local maximizers, where  $\rho$  denotes spectral radius  $(\max(\lambda_{\max}, -\lambda_{\min}));$

(iii)  $\rho_{in} \circ H$  is  $C^2$  near all of its local minimizers at which the minimal value is nonzero, where  $\rho_{in}$  denotes inner spectral radius (0 if H is singular,  $\rho(H^{-1})^{-1}$ otherwise).

Furthermore, in each case the second derivative is Lipschitz near the relevant maximizers/minimizers.

**Proof** Statements (i) and (ii) follow from applying Theorem 4.2 to -H and diag(H, -H), respectively. For (iii), apply (ii) to  $\rho \circ H^{-1}$  and take the reciprocal.

**Corollary 4.2** Let  $A : \mathcal{D} \to \mathbb{C}^{m \times n}$  be an analytic matrix family in one real variable on an open domain  $\mathcal{D} \subseteq \mathbb{R}$ , let  $\sigma_{\max}$  denote largest singular value, and let  $\sigma_{\min}$  denote smallest singular value, noting that the latter is nonzero if and only if the matrix has full rank. Then:

- (i)  $\sigma_{\max} \circ A$  is  $C^2$  near all of its local maximizers, and (ii)  $\sigma_{\min} \circ A$  is  $C^2$  near all of its local minimizers at which the minimal value is nonzero.

Furthermore, in each case the second derivative is Lipschitz near the relevant maximizers/minimizers.

**Proof** If  $m \geq n$ , consider the real analytic Hermitian matrix family  $H : \mathcal{D} \to \mathbb{H}^n$ defined by

$$H(t) = A(t)^* A(t) = (\Re A(t) - i \Im A(t))^{\mathsf{T}} (\Re A(t) + i \Im A(t)),$$

whose eigenvalues are the squares of the singular values of A(t). Then (i) and (ii), respectively, follow from applying Corollary 4.1 (ii) and (iii), respectively, to H(t), and then taking the square root. If n > m, set  $H(t) = A(t)A(t)^*$  instead. П

Corollary 4.2 (i) is the regularity result that Boyd and Balakrishnan established in [2]. For Corollary 4.2 (ii), note that the assumption that the minimal value of  $\sigma_{\min} \circ A$ is nonzero is necessary; e.g.,  $\sigma_{\min}(t)$  is nonsmooth at its minimizer t = 0.

#### 4.1 The $\mathcal{H}_{\infty}$ norm

This application was the original motivation for Boyd and Balakrishnan's work. Let  $A \in \mathbb{C}^{n \times n}, B \in \mathbb{C}^{n \times m}, C \in \mathbb{C}^{p \times n}$ , and  $D \in \mathbb{C}^{p \times m}$  and consider the linear timeinvariant system with input and output:

$$\dot{x} = Ax + Bu, \tag{4.2a}$$

$$y = Cx + Du. \tag{4.2b}$$

Assume that A is asymptotically stable, i.e., its eigenvalues are all in the open left half-plane. An important quantity in control systems engineering and model-order reduction is the  $\mathcal{H}_{\infty}$  norm of (4.2), which measures the sensitivity of the system to perturbation and can be computed by solving the following optimization problem:

$$\max_{\omega \in \mathbb{R}} \sigma_{\max}(G(\mathrm{i}\omega)), \tag{4.3}$$

where  $G(\lambda) = C(\lambda I - A)^{-1}B + D$  is the transfer matrix associated with (4.2). Even though there is only one real variable, finding the global maximum of this function is nontrivial.

By extending Byer's breakthrough result on computing the distance to instability [5], Boyd et al. [3] developed a globally convergent bisection method to solve (4.3) to arbitrary accuracy. Shortly thereafter, a much faster algorithm, based on computing level sets of  $\sigma_{max}(G(i\omega))$ , was independently proposed in [2] and [4], with Boyd and Balakrishnan showing that this iteration converges quadratically [2, Theorem 5.1]. As part of their work, they showed that, with respect to the real variable  $\omega$ ,  $\sigma_{max}(G(i\omega))$  is  $C^2$  with Lipschitz second derivative near any of its local maximizers [2, pp. 2–3]. Subsequently, this smoothness property has been leveraged to further accelerate computation of the  $\mathcal{H}_{\infty}$  norm [1,6].

#### 4.2 The numerical radius

Now consider the numerical radius of a matrix  $A \in \mathbb{C}^{n \times n}$ :

$$r(A) = \max\{|z| : z \in W(A)\},$$
(4.4)

where  $W(A) = \{v^*Av : v \in \mathbb{C}^n, \|v\|_2 = 1\}$  is the field of values (numerical range) of *A*. Following [8, Ch. 1], the numerical radius can be computed by solving either

$$r(A) = \max_{\theta \in [0,2\pi)} \lambda_{\max}(H(\theta)) \quad \text{or} \quad r(A) = \max_{\theta \in [0,\pi)} \rho(H(\theta)), \tag{4.5}$$

where  $H(\theta) = \frac{1}{2} \left( e^{i\theta} A + e^{-i\theta} A^* \right)$ .

In [13], Mengi and the second author proposed the first globally convergent method guaranteed to compute r(A) to arbitrary accuracy. This was done by employing a level-set technique that converges to a global maximizer of  $\lambda_{\max} \circ H$ , similar to the aforementioned method of [2,4] for the  $\mathcal{H}_{\infty}$  norm, and observing, but not proving, quadratic convergence of the method. Quadratic convergence was later proved by Gürbüzbalaban in his PhD thesis [7, Lemma 3.4.2], following the proof used in [2], showing that  $\lambda_{\max} \circ H$  is  $\mathcal{C}^2$  near maximizers.

#### 4.3 Optimization of passive systems

Let  $\mathcal{M} = \{A, B, C, D\}$  denote the system (4.2), but now with m = p and the associated transfer function *G* being minimal and proper [16]. Mehrmann and Van Dooren [12] have recently shown that another important problem is to compute the maximal value  $\Xi \in \mathbb{R}$  such that for all  $\xi < \Xi$ , the related system  $\mathcal{M}_{\xi} = \{A_{\xi}, B, C, D_{\xi}\}$  is

strictly passive<sup>1</sup>, where  $A_{\xi} = A + \frac{\xi}{2}I_n$  and  $D_{\xi} = D - \frac{\xi}{2}I_m$ . Letting  $G_{\xi}$  be the transfer matrix associated with  $\mathcal{M}_{\xi}$ , by [12, Theorem 5.1], the quantity  $\Xi$  is the unique root of

$$\gamma(\xi) := \min_{\omega \in \mathbb{R}} \lambda_{\min} \left( G_{\xi}(i\omega)^* + G_{\xi}(i\omega) \right) = 0.$$
(4.6)

Note that in contrast to the univariate optimization problems discussed previously, computing  $\Xi$  is a problem in two real parameters, namely,  $\xi$  and  $\omega$ . In [12, section 5], Mehrmann and Van Dooren introduced both a bisection algorithm to compute  $\Xi$ , and an apparently faster "improved iteration" whose exact convergence properties were not established. However, using the fact that  $\lambda_{\min}$  in (4.6) is  $C^2$  with Lipschitz second derivative near all its minimizers, as well as some other tools, the first author and Van Dooren have since established a rate-of-convergence result for this "improved iteration" and also presented a much faster and more numerically reliable algorithm to compute  $\Xi$  with quadratic convergence [14].

#### **5** Concluding remarks

We have shown that the maximum eigenvalue of a univariate real analytic Hermitian matrix family is unconditionally  $C^2$  near all its maximizers, with Lipschitz second derivative. Although the result is well known in the context of the maximum singular value of a transfer function, its generality and simplicity have apparently not been fully appreciated. We believe that this result and its corollaries may be useful in many applications, some of which were summarized in this paper. We also investigated whether this smoothness property extends to max functions more generally, showing that the pointwise maximum of a finite set of *q*-times continuously differentiable univariate functions must have zero derivative at a maximizer for q = 1, but arbitrarily close to the maximizer, the derivative may not be defined, even when q = 3 and the maximizer is isolated.

All figures and the symbolically computed coefficients of  $p_k$  given in Appendices A and B were generated by MATLAB codes that can be downloaded from https://doi.org/10.5281/zenodo.5831694.

Acknowledgements The second author was supported in part by U.S. National Science Foundation Grant DMS-2012250.

Funding Open Access funding enabled and organized by Projekt DEAL.

**Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit http://creativecommons.org/licenses/by/4.0/.

<sup>&</sup>lt;sup>1</sup> A strictly passive system is one whose stored energy is decreasing; for more a formal treatment, see [12].

#### A Proofs of Theorems 3.1 and 3.2

Lemma A.1 For  $f_1$  defined in (3.1), if  $s_k = 1 + 2^{-2k}$ , then the coefficients of the polynomial  $p_k(t) = \sum_{j=0}^{9} c_j t^j$  are:  $c_j = \begin{cases} z_j 2^{(j-4)k} - 1 & \text{if } j = 2 \\ z_j 2^{(j-4)k} & \text{otherwise} \end{cases}$  with  $z_8 = 663552, z_4 = 2585088, z_0 = 4608, z_7 = -1966080, z_3 = -1210368, z_6 = 3354624, z_2 = 359424, z_7 = 359424$ 

**Proof** The coefficients were computed symbolically in MATLAB by solving the linear system defined by the generalized Vandermonde matrix and right-hand side determining each  $p_k$  in (3.1). These formulas were also verified by comparing with numerical computations.

**Proof** (of Theorem 3.1) Function  $f_1$  defined in (3.1) is clearly  $C^3$  near any nonzero t, since our construction ensures that the first three derivatives of  $p_k$  and  $p_{k+1}$  match where they meet. We must show that it is also  $C^3$  at t = 0. First note that for the coefficients given in Lemma A.1, we can replace their dependency on k with a dependency on t by using  $k = -\lceil \log_2 t \rceil$ . Thus,  $f_1$  can be written as follows:

$$f_1(t) = \begin{cases} \sum_{j=0}^9 \tilde{c}_j t^j & \text{if } t > 0\\ -t^2 & \text{if } t \in [-1, 0] \end{cases}$$
(A.1)

where  $\tilde{c}_i$  is obtained by replacing k in  $c_i$  with  $-\lceil \log_2 t \rceil$ .

We begin by looking at the first derivative. For  $f'_1$  to exist and be continuous at t = 0,

$$f_1'(0) = \lim_{\varepsilon \to 0^+} \frac{f_1(0+\varepsilon) - f_1(0)}{\varepsilon} = \lim_{\varepsilon \to 0^+} \frac{f_1(\varepsilon)}{\varepsilon} = 0$$
(A.2)

must hold, i.e., the derivative from the right (over the  $p_k$  pieces) must match the derivative from the left (over the  $-t^2$  piece). To show that (A.2) holds, we show that each term in the sum in (A.1) divided by t goes to zero as  $t \to 0^+$ , i.e., that  $\lim_{t\to 0^+} \tilde{c}_j t^{j-1} = 0$  for  $j \in \{0, 1, \ldots, 9\}$ . It is obvious that this holds for j = 4 since  $c_4 = \tilde{c}_4 = z_4$  is a fixed number. To show the highest-order term (j = 9) vanishes as  $t \to 0^+$ , we can make use of the fact that  $0 < 2^{-\lceil \log_2 t \rceil} \le 2^{-\log_2 t} = t^{-1}$  holds for all t > 0, i.e.,

$$\lim_{t \to 0^+} \left| z_9 2^{-5 \lceil \log_2 t \rceil} t^8 \right| \le \lim_{t \to 0^+} \left| z_9 (t^{-1})^5 t^8 \right| = \lim_{t \to 0^+} \left| z_9 t^3 \right| = 0.$$

Similar arguments show that  $\lim_{t\to 0^+} \tilde{c}_j t^{j-1} = 0$  holds for  $j \in \{5, 6, 7, 8\}$ . Using the fact that  $0 < 2^{\lceil \log_2 t \rceil} \le 2^{1+\log_2 t} = 2t$  for all t > 0, for j = 3, we have that

$$\lim_{t \to 0^+} |z_3 2^{\lceil \log_2 t \rceil} t^2| \le \lim_{t \to 0^+} |z_3 2 t^3| = 0,$$

Deringer

while for j = 2 and j = 1 we respectively have that

$$\lim_{t \to 0^+} |z_2 2^{2\lceil \log_2 t \rceil} - 1| t \le \lim_{t \to 0^+} |z_2 (2t)^2 - 1| t = \lim_{t \to 0^+} |z_2 (4t^3 - t)| = 0.$$

and

$$\lim_{t \to 0^+} \left| z_1 2^{3 \lceil \log_2 t \rceil} \right| \le \lim_{t \to 0^+} \left| z_1 (2t)^3 \right| = 0.$$

Finally, for j = 0, we have that

$$\lim_{t \to 0^+} \frac{z_0 2^{4\lceil \log_2 t \rceil}}{t} \le \lim_{t \to 0^+} \frac{z_0 (2t)^4}{t} = 0.$$

Hence, we have shown that  $f_1$  is at least  $C^1$  on its domain.

Analogously, for  $f_1''$  to exist and be continuous at t = 0,

$$f_1''(0) = \lim_{\varepsilon \to 0^+} \frac{f_1'(0+\varepsilon) - f_1'(0)}{\varepsilon} = \lim_{\varepsilon \to 0^+} \frac{f_1'(\varepsilon)}{\varepsilon} = -2$$
(A.3)

must hold. We have that

$$f_1'(t) = \begin{cases} \sum_{j=1}^9 j \tilde{c}_j t^{j-1} & \text{if } t > 0\\ -2t & \text{if } t \in [-1, 0] \end{cases}$$
(A.4)

and so we consider  $\lim_{t\to 0^+} j\tilde{c}_j t^{j-2}$  for  $j \in \{1, \ldots, 9\}$ , i.e., the limit of each term in the sum in (A.4) divided by *t*. We show that for all but j = 2, these values goes to zero, while the j = 2 value goes to -2 as  $t \to 0^+$ . For j = 9, we have that

$$\lim_{t \to 0^+} |9z_9 2^{-5\lceil \log_2 t \rceil} t^7| \le \lim_{t \to 0^+} |9z_9 (t^{-1})^5 t^7| = 0,$$

with similar arguments showing that  $j \in \{5, 6, 7, 8\}$  values also diminish to zero. For j = 4, we simply have  $\lim_{t\to 0^+} 4z_4t^2 = 0$ . For j = 3,

$$\lim_{t \to 0^+} |3z_3 2^{\lceil \log_2 t \rceil} t| \le \lim_{t \to 0^+} |3z_3 (2t)t| = 0.$$

For j = 2, we have that

$$\lim_{t \to 0^+} 2(z_2 2^{2\lceil \log_2 t \rceil} - 1) = \lim_{t \to 0^+} 2z_2 (2^{\lceil \log_2 t \rceil})^2 - 2 = -2.$$

Lastly, for j = 1, we have that

$$\lim_{t \to 0^+} \frac{\left| z_1 2^{3 \lceil \log_2 t \rceil} \right|}{t} \le \lim_{t \to 0^+} \frac{\left| z_1 (2t)^3 \right|}{t} = 0,$$

Deringer

and so we have now shown that  $f_2$  is at least  $C^2$  on its domain.

Finally, for  $f_1'''$  to exist and be continuous at t = 0,

$$f_1'''(0) = \lim_{\epsilon \to 0^+} \frac{f_1''(0+\epsilon) - f_1''(0)}{\epsilon} = \lim_{\epsilon \to 0^+} \frac{f_1''(\epsilon) + 2}{\epsilon} = 0$$
(A.5)

must hold. We have that

$$f_1''(t) = \begin{cases} \sum_{j=2}^9 j(j-1)\tilde{c}_j t^{j-2} & \text{if } t > 0\\ -2 & \text{if } t \in [-1,0] \end{cases}$$
(A.6)

and so we consider  $\lim_{t\to 0^+} j(j-1)\tilde{c}_j t^{j-3}$  for  $j \in \{2, \ldots, 9\}$ , i.e., the limit of each term in the sum in (A.6) divided by *t*. For  $j \in \{5, 6, 7, 8, 9\}$ , we again have similar arguments showing that the corresponding values vanish, so we just show the j = 9 case, which follows because

$$\lim_{t \to 0^+} \left| 72z_9 2^{-5 \lceil \log_2 t \rceil} t^6 \right| \le \lim_{t \to 0^+} \left| 72z_9 (t^{-1})^5 t^6 \right| = 0.$$

Again, it is clear that the value for j = 4 vanishes. For j = 3, we have that

$$\lim_{t \to 0^+} |6z_3 2^{\lceil \log_2 t \rceil}| \le \lim_{t \to 0^+} |6z_3(2t)| = 0.$$

Finally, for j = 2, we can rewrite (A.5) as follows, making use of these aforementioned limits which vanish and replacing  $\varepsilon$  by t, to obtain a limit only involving the j = 2 term:

$$f_1'''(0) = \lim_{t \to 0^+} \frac{f_1''(t) + 2}{t} = \lim_{t \to 0^+} \frac{2(z_2 2^{2\lceil \log_2 t \rceil} - 1) + 2}{t}$$
$$= \lim_{t \to 0^+} \frac{2z_2 (2^{\lceil \log_2 t \rceil})^2}{t} \le \lim_{t \to 0^+} \frac{2z_2 (2t)^2}{t} = 0.$$

Thus,  $f_1$  is indeed  $C^3$  on its domain.

**Proof** (of Theorem 3.2) Since  $l_k$  is a power of two, we can rewrite the derivative of  $p_k$ , i.e.,  $p'_k(t) = \sum_{j=1}^9 jc_j t^{j-1}$ , as a function of  $\zeta \in [1, 2]$ :

$$\begin{split} \tilde{p}'_k(\zeta) &= \sum_{j=1}^9 j c_j (l_{k+1}\zeta)^{j-1} = \sum_{j=1}^9 \frac{j c_j}{2^{(k+1)(j-1)}} \zeta^{j-1} \\ &= \frac{2(z_2 2^{-2k} - 1)}{2^{k+1}} \zeta + \sum_{\substack{j=1\\j \neq 2}}^9 \frac{j z_j 2^{(j-4)k}}{2^{(k+1)(j-1)}} \zeta^{j-1} \end{split}$$

$$= \frac{z_2 - 2^{2k}}{2^{3k}} \zeta + \sum_{\substack{j=1\\j \neq 2}}^9 \frac{j z_j 2^{1-j}}{2^{3k}} \zeta^{j-1}$$
$$= \frac{1}{2^{3k}} \left( (z_2 - 2^{2k}) \zeta + \sum_{\substack{j=1\\j \neq 2}}^9 \tilde{z}_j \zeta^{j-1} \right),$$

where  $\tilde{z}_j = jz_j 2^{1-j}$ . From Lemma A.1, we see that  $z_2 - 2^{2k} < 0$  for all  $k \ge 10$ , while for any k, we have that  $\tilde{z}_j < 0$  for  $j \in \{1, 3, 5, 7, 9\}$  and  $\tilde{z}_j > 0$  for  $j \in \{4, 6, 8\}$ . Since  $\zeta \in [1, 2]$ , an upper bound for  $\tilde{p}'_k$  can be obtained by evaluating its negative terms at  $\zeta = 1$  and its positive terms at  $\zeta = 2$ , i.e., for all  $k \ge 10$  and any  $\zeta \in [1, 2]$ , we have that

$$\tilde{p}'_k(\zeta) \le \frac{1}{2^{3k}} \bigg( (z_2 - 2^{2k}) + \sum_{j \in \{1,3,5,7,9\}} \tilde{z}_j + \sum_{j \in \{4,6,8\}} \tilde{z}_j 2^{j-1} \bigg).$$

For  $k \ge 13$ , the upper bound on the derivative is negative. Thus, for  $k \ge 13$ ,  $\tilde{p}'_k(\zeta) < 0$  for any  $\zeta \in [1, 2]$ , so  $p_k$  must be decreasing. Consequently, the t = 0 maximizer of  $f_1$  is isolated. Finally, it immediately follows that the t = 0 maximizer of  $f_{\text{max}} = \max(f_1, f_2)$  is also isolated.

## B Why $s_k = 1 + 2^{-k}$ is insufficient to make (3.1) a $C^3$ function

For  $s_k = 1 + 2^{-k}$ , symbolic computation shows that the coefficients of  $p_k(t) = \sum_{j=0}^{9} c_j t^j$  are:

$$c_j = \begin{cases} z_j 2^{(j-3)k} - 1 & \text{if } j = 2\\ z_j 2^{(j-3)k} & \text{otherwise} \end{cases}$$

where the integers  $z_j$  remain the same as given in Lemma A.1. To see if (A.5) still holds for this new choice of  $s_k$  we look at  $\lim_{t\to 0^+} j(j-1)\tilde{c}_j t^{j-3}$  for  $j \in \{2, \ldots, 9\}$ . However, now none of the individual limits vanish. For example, for j = 9, we have that

$$\lim_{t \to 0^+} \left| 72z_9 2^{-6\lceil \log_2 t \rceil} t^6 \right| \ge \lim_{t \to 0^+} \left| 72z_9 (2^{-1}t^{-1})^6 t^6 \right| = \frac{9}{8} |z_9| \ne 0,$$

where we have used the fact that  $0 < 2^{-1}t^{-1} = 2^{-1-\log_2 t} \le 2^{-\lceil \log_2 t \rceil}$ ; similarly, the limits for  $j \in \{4, 5, 6, 7, 8\}$  do not vanish either. For j = 3, we simply have that

$$\lim_{t \to 0^+} 6z_3 = 6z_3 \neq 0.$$

Finally, even if all of the terms considered above were to vanish and we substitute in the value for j = 2 into (A.5), we nevertheless would end up attaining another limit that does not vanish:

$$\lim_{t \to 0^+} \frac{2z_2 2^{\lceil \log_2 t \rceil}}{t} \ge \lim_{t \to 0^+} \frac{2z_2(t)}{t} = 2z_2 \neq 0.$$

The only remaining way that (A.5) could hold is if all of these non-vanishing terms cancel, but from our experiments (see Fig. 2a), we know this is not the case.

#### References

- 1. Benner, P., Mitchell, T.: Faster and more accurate computation of the  $\mathcal{H}_{\infty}$  norm via optimization. SIAM J. Sci. Comput. **40**(5), A3609–A3635 (2018). https://doi.org/10.1137/17M1137966
- 2. Boyd, S., Balakrishnan, V.: A regularity result for the singular values of a transfer matrix and a quadratically convergent algorithm for computing its  $L_{\infty}$ -norm. Syst. Control Lett. **15**(1), 1–7 (1990). https://doi.org/10.1016/0167-6911(90)90037-U
- Boyd, S., Balakrishnan, V., Kabamba, P.: A bisection method for computing the H<sub>∞</sub> norm of a transfer matrix and related problems. Math. Control Signals Syst. 2, 207–219 (1989)
- Bruinsma, N.A., Steinbuch, M.: A fast algorithm to compute the H<sub>∞</sub>-norm of a transfer function matrix. Syst. Control Lett. 14(4), 287–293 (1990)
- Byers, R.: A bisection method for measuring the distance of a stable matrix to unstable matrices. SIAM J. Sci. Stat. Comput. 9, 875–881 (1988). https://doi.org/10.1137/0909059
- Genin, Y., Van Dooren, P., Vermaut, V.: Convergence of the calculation of H<sub>∞</sub>-norms and related questions. In: A. Beghi, L. Finesso, G. Picci (eds.) Mathematical Theory of Networks and Systems, 13 ed., Proceedings of the MTNS-98 Symposium, Padova, pp. 629–632 (1998)
- Gürbüzbalaban, M.: Theory and methods for problems arising in robust stability, optimization and quantization. Ph.D. thesis, New York University, New York, NY, USA (2012). https://mert-g.org/wpcontent/uploads/2018/06/Mert-Thesis.pdf
- 8. Horn, R.A., Johnson, C.R.: Topics in Matrix Analysis. Cambridge University Press, Cambridge (1991)
- Kato, T.: A Short Introduction to Perturbation Theory for Linear Operators. Springer, New York (1982). https://doi.org/10.1007/978-1-4612-5700-4
- Krantz, S.G., Parks, H.R.: A Primer of Real Analytic Functions. Birkhäuser Advanced Texts, 2nd edn. Birkhäuser, Boston, MA (2002). https://doi.org/10.1007/978-0-8176-8134-0
- 11. Lax, P.D.: Linear Algebra and its Applications, 2nd edn. Wiley, Hoboken, NJ (2007)
- Mehrmann, V., Van Dooren, P.M.: Optimal robustness of port-Hamiltonian systems. SIAM J. Matrix Anal. Appl. 41(1), 134–151 (2020). https://doi.org/10.1137/19M1259092
- Mengi, E., Overton, M.L.: Algorithms for the computation of the pseudospectral radius and the numerical radius of a matrix. IMA J. Numer. Anal. 25(4), 648–669 (2005). https://doi.org/10.1093/imanum/ dri012
- Mitchell, T., Van Dooren, P.: Root-max problems, hybrid expansion-contraction, and quadratically convergent optimization of passive systems. ArXiv arXiv:2109.00974 (2021)
- 15. von Neumann, J., Wigner, E.P.: Über merkwürdige diskrete Eigenwerte. Phys. Z. 40, 467–470 (1929)
- Zhou, K., Doyle, J.C., Glover, K.: Robust and Optimal Control. Prentice-Hall, Upper Saddle River, NJ (1996). https://doi.org/10.1007/978-1-4471-6257-5

Publisher's Note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.