

## SUBSPACE ACCELERATION FOR THE CRAWFORD NUMBER AND RELATED EIGENVALUE OPTIMIZATION PROBLEMS\*

DANIEL KRESSNER<sup>†</sup>, DING LU<sup>‡</sup>, AND BART VANDEREYCKEN<sup>‡</sup>

**Abstract.** This paper is concerned with subspace acceleration techniques for computing the Crawford number, that is, the distance between zero and the numerical range of a matrix  $A$ . Our approach is based on an eigenvalue optimization characterization of the Crawford number. We establish local convergence of order  $1 + \sqrt{2} \approx 2.4$  for an existing subspace method applied to such and other eigenvalue optimization problems involving a Hermitian matrix that depends analytically on one parameter. For the particular case of the Crawford number, we show that the relevant part of the objective function is strongly concave. In turn, this enables us to develop a subspace method that only uses three-dimensional subspaces but still achieves global convergence and a local convergence that is at least quadratic. A number of numerical experiments confirm our theoretical results and reveal that the established convergence orders appear to be tight.

**Key words.** subspace acceleration, eigenvalue optimization, Crawford number, coercivity constant, convergence analysis, complex approximation

**AMS subject classifications.** 30E10, 65F15, 90C26

**DOI.** 10.1137/17M1127545

**1. Introduction.** During the last few years, subspace methods for optimizing eigenvalues of Hermitian matrices or singular values of non-Hermitian matrices have been developed, analyzed, and applied to a variety of problems. A common trait of these methods is that they use the subspace spanned by eigenvectors or singular vectors from previous iterations to build a low-dimensional surrogate model that determines the next iterate. Examples include subspace methods for computing (real) pseudospectral abscissa for linear [15, 18] and nonlinear [19] eigenvalue problems, for computing the  $\mathcal{H}_\infty$  norm of a linear control system [1], as well as for solving Hermitian eigenvalue optimization problems depending on several parameters [13].

In this work, we revisit eigenvalue optimization problems depending on one parameter, with a particular focus on their application to Crawford number computation. The Crawford number [5, 24] of a matrix  $A \in \mathbb{C}^{n \times n}$  is defined as the distance of its numerical range  $\mathcal{F}(A)$  from zero:

$$(1) \quad \gamma(A) = \min\{|z| : z \in \mathcal{F}(A)\}, \quad \mathcal{F}(A) = \{v^H A v : v \in \mathbb{C}^n, \|v\|_2 = 1\}.$$

This quantity plays an important role in the sensitivity analysis of matrix definite pairs [4, 24], and the detection of hyperbolic quadratic eigenvalue problems [7, 11]. Possibly even more importantly,  $\gamma(A)$  and its generalization to bilinear forms on Hilbert spaces play the role of the coercivity constant for (discretized) non-self-adjoint differential/integral operators and their discretization; see [2] for a recent example.

By the Hausdorff–Toeplitz theorem, the numerical range  $\mathcal{F}(A)$  is a compact convex set in  $\mathbb{C}$ . Hence, (1) is a convex optimization problem and a number of algorithms

---

\*Received by the editors April 26, 2017; accepted for publication (in revised form) November 28, 2017; published electronically May 24, 2018.

<http://www.siam.org/journals/simax/39-2/M112754.html>

**Funding:** The work of the second author was supported by the SNSF research project Robust and Fast Solvers for Computing Extremal Quantities of Structured Pseudospectra.

<sup>†</sup>Institute of Mathematics, EPFL, Lausanne CH-1015, Switzerland (daniel.kressner@epfl.ch).

<sup>‡</sup>Department of Mathematics, University of Geneva, Geneva 1211, Switzerland (Ding.Lu@unige.ch, Bart.Vandereycken@unige.ch).

explicitly exploit this property, for example, by approximating the numerical range as a sequence of convex polygons [26]. However, it turns out to be more effective to avoid such geometric considerations in  $\mathbb{C} \cong \mathbb{R} \times \mathbb{R}$  and exploit that (1) is in fact a univariate optimization problem. To see this, we first write

$$A = S + \iota K,$$

where  $S = (A + A^H)/2$  and  $K = (A - A^H)/2\iota$  are both Hermitian matrices. Using the variational characterization of eigenvalues of Hermitian matrices, one then obtains

$$(2) \quad \gamma(A) = \max \left\{ \max_{\theta \in [0, 2\pi]} \lambda_{\min}(S \cos \theta + K \sin \theta), \quad 0 \right\};$$

see, e.g., [4, Thm. 2.1] and [11, eq. (2.8)]. Based on this formula, a level set algorithm [4, 11] and bisection algorithms [3, 8] were developed. While being reliable, these algorithms are only linearly convergent and they typically require a relatively large number of iterations to attain high accuracy.

*Contributions.* The approach considered in this paper consists of applying existing subspace methods, specifically the one from [13], to (2) and, more generally, to univariate eigenvalue optimization problems. We would like to emphasize the following two novel contributions.

Existing convergence analysis [13, 15, 18] of subspace methods for eigenvalue or singular value optimization establishes local superlinear convergence. The restriction to the univariate case enables the use of univariate analytic interpolation results, which in turn allows us to derive in section 2 a much stronger result: local convergence of order  $\sigma = 1 + \sqrt{2} \approx 2.4$ . In contrast to [13], our analysis is sharper since numerical experiments reveal that the obtained value for  $\sigma$  appears to be tight.

Restricting ourselves further to the specific eigenvalue optimization problem (2), we establish in section 3 a strong concavity property for the objective function, a variation of a result from [8]. This allows us to show that subspace methods enjoy a rather useful “bracketing” property when applied to Crawford number computations. Based upon these insights, we develop a “3-vector” subspace method that works with only three-dimensional subspaces but still achieves global convergence and a local convergence that is quadratic or even higher.

**2. Hermitian eigenvalue optimization problem.** In the following, we consider the univariate eigenvalue optimization problem

$$(3) \quad \max_{\theta \in \Omega} \varphi(\theta) \quad \text{with} \quad \varphi(\theta) = \lambda_{\min}(H(\theta)),$$

where  $\Omega \subset \mathbb{R}$  is a closed interval, and  $H(\theta) \in \mathbb{C}^{n \times n}$  is Hermitian and real analytic for  $\theta \in \Omega$ , that is, it admits an analytic extension in an open neighborhood  $\Omega_{\mathbb{C}} \subset \mathbb{C}$ .

Given an orthonormal basis  $V \in \mathbb{C}^{n \times k}$  of a  $k$ -dimensional subspace of  $\mathbb{C}^n$ , the subspace acceleration techniques discussed, e.g., in [13, 22], proceed by considering the reduced objective function

$$\varphi(\theta; V) = \lambda_{\min}(V^H H(\theta) V).$$

Typically  $k \ll n$  and it can therefore be expected that computing the maximum of  $\varphi(\theta; V)$  on  $\Omega$  is cheaper than solving (3). In particular, this is the case when  $H$  admits an affine linear decomposition  $H(\theta) = \sum_{\ell=1}^m f_{\ell}(\theta) H_{\ell}$  with scalar functions  $f_{\ell}$  and Hermitian matrices  $H_{\ell}$  such that  $m \ll n$ . After precomputing  $V^H H_{\ell} V$  for

$\ell = 1, \dots, m$ , each evaluation of  $\varphi(\theta; V)$  only requires the addition of  $m$  matrices and the solution of one  $k \times k$  eigenvalue problem.

The following Lemma 2.1 recalls results from [13, 22] on basic properties of the reduced function  $\varphi(\theta; V)$ ; its proof is included for completeness.

LEMMA 2.1. *Let  $U$  and  $V$  be orthonormal bases of subspaces  $\mathcal{U} \subset \mathcal{V}$  of  $\mathbb{C}^n$ . Then the following properties hold:*

- (a) *Monotonicity:*  $\varphi(\theta; U) \geq \varphi(\theta; V) \geq \varphi(\theta)$ .
- (b) *Interpolation:* Let  $\theta \in \Omega$  and suppose that  $\mathcal{V}$  contains an eigenvector belonging to the eigenvalue  $\lambda_{\min}(H(\theta))$  of  $H(\theta)$ . Then,  $\varphi(\theta; V) = \varphi(\theta)$ . In addition, if  $\lambda_{\min}(H(\theta))$  is a simple eigenvalue, then  $\varphi'(\theta; V) = \varphi'(\theta)$ , where all derivatives are taken with respect to  $\theta$ .

*Proof.* (a) Monotonicity follows directly from the variational characterization of eigenvalues:

$$(4) \quad \varphi(\theta; U) = \min_{u \in \mathcal{U}} \frac{u^H H(\theta) u}{u^H u} \geq \min_{u \in \mathcal{V}} \frac{u^H H(\theta) u}{u^H u} = \varphi(\theta; V).$$

The second inequality is obtained from the first since  $\varphi(\theta) = \varphi(\theta; I_n)$ .

(b) Let  $v(\theta)$  be an eigenvector belonging to  $\lambda_{\min}(H(\theta))$ . Since  $v(\theta) \in \mathcal{V} \setminus \{0\}$  there is  $x \in \mathbb{C}^k$  such that  $v(\theta) = Vx$ . Then, by definition of the eigenvector  $v(\theta)$ , we have

$$\varphi(\theta; V) = \min_{z \in \mathbb{C}^k} \frac{z^H V^H H(\theta) V z}{z^H z} \leq \frac{x^H V^H H(\theta) V x}{x^H x} = \frac{v(\theta)^H H(\theta) v(\theta)}{v(\theta)^H v(\theta)} = \varphi(\theta).$$

Combined with (a), we obtain  $\varphi(\theta; V) = \varphi(\theta)$ . If  $\varphi(\theta)$  is a simple eigenvalue of  $H(\theta)$ , the Cauchy interlacing theorem yields that  $\varphi(\theta; V)$  is also a simple eigenvalue of  $V^H H(\theta) V$  with eigenvector  $x$ . By well-known results on the derivatives of eigenvalues (see, e.g., [16]), we obtain

$$\varphi'(\theta) = \frac{v(\theta)^H H'(\theta) v(\theta)}{v(\theta)^H v(\theta)} = \frac{x^H V^H H'(\theta) V x}{x^H x} = \varphi'(\theta; V). \quad \square$$

By Lemma 2.1,  $\varphi(\theta; V)$  provides an upper bound to  $\varphi(\theta)$  that monotonically decreases as the dimension of  $\mathcal{V}$  increases. In view of the intimate relation between gradients of eigenvalues and eigenvectors, it is a natural idea to extend  $V$  by an eigenvector corresponding to  $\varphi(\theta; V)$ . This idea leads to the basic subspace method from [13], summarized in Algorithm 1. The difficulty of solving the reduced problem in line 4 depends on the application. It turns out to be quite cheap for the Crawford number, as we will see in section 3. Other applications with this property are also listed in [13].

In the exceptional situation that Algorithm 1 stagnates it has found a global maximum.

LEMMA 2.2. *If  $\theta_{k+1} = \theta_j$  for some  $j \leq k$ , then Algorithm 1 terminates and  $\theta_{k+1}, \varphi_{k+1}$  solve (3).*

*Proof.* Without loss of generality, we may assume  $j = k$ . By construction,  $\lambda_k$  is a lower bound and  $\varphi_{k+1}$  is an upper bound for  $\max_{\theta \in \Omega} \varphi(\theta)$ . Now, if  $\theta_k = \theta_{k+1}$ , then

$$\varphi_{k+1} = \lambda_{\min}(V_k^H H(\theta_{k+1}) V_k) = \lambda_{\min}(V_k^H H(\theta_k) V_k) = \lambda_{\min}(H(\theta_k)) = \lambda_k,$$

where we used  $v_k \in \text{span}(V_k)$  in the second-to-last equality. □

---

**Algorithm 1** Subspace method for univariate eigenvalue optimization.

---

**Input:** Real analytic Hermitian matrix-valued function  $H(\theta)$ , initial guess  $\theta_0$ , tolerance  $\text{tol} > 0$ .

**Output:** Approximation solution  $\theta_{k+1}, \varphi_{k+1}$  of eigenvalue optimization problem (3).

- 1: Compute  $\lambda_0 = \lambda_{\min}(H(\theta_0))$  and corresponding normalized eigenvector  $v_0$ .
  - 2: Initialize  $V_0 = v_0$ .
  - 3: **for**  $k = 0, 1, \dots, n - 1$  **do**
  - 4:   Solve  $\theta_{k+1} = \arg \max_{\theta \in \Omega} \lambda_{\min}(V_k^H H(\theta) V_k)$  and set  $\varphi_{k+1} = \lambda_{\min}(V_k^H H(\theta_{k+1}) V_k)$ .
  - 5:   *Stopping criteria:* **if**  $\varphi_{k+1} - \lambda_k < \text{tol} \cdot |\varphi_{k+1}|$  **then** terminate.
  - 6:   Compute the smallest eigenvalue  $\lambda_{k+1}$  with normalized eigenvector  $v_{k+1}$  of  $H(\theta_{k+1})$ .
  - 7:   *Subspace update:*  $V_{k+1} = \text{orth}([V_k, v_{k+1}])$ .
  - 8: **end for**
- 

A consequence of this lemma is that if  $v_{k+1} \in \text{span}(V_k)$ , we have  $V_{k+1} = V_k$  and  $\theta_{k+1} = \theta_k$ . This implies we have solved (3). Let us note that Algorithm 1 converges trivially after at most  $n$  steps, since  $V_{n-1}$  will be a basis of  $\mathbb{C}^n$  and Lemma 2.1 shows that  $\varphi_n$  solves (3) exactly. We would like to emphasize, however, that this observation is merely a curiosity and of little practical value. The use of Algorithm 1 is that it produces very accurate approximations already after a few steps, certainly much less than  $n$  when  $n$  is large. In [13, Thm. 3.1], an infinite-dimensional setting was considered to meaningfully establish a global convergence result, in the sense that every accumulation point of the sequence  $(\theta_k)_{k=0}^\infty$  is in fact a global maximum of the optimization problem (3).

**2.1. Local convergence.** Here and in the following we assume that the eigenvalues corresponding to  $\varphi(\theta_i)$  in Algorithm 1 are simple and hence differentiable. By the interpolation property established in Lemma 2.1(b), the reduced function  $\varphi(\theta; V_k)$  from the  $k$ th iteration of Algorithm 1 then constitutes a (nonpolynomial) Hermite interpolation of  $\varphi(\theta)$ , that is, it interpolates both function values and derivatives at the sampling points:

$$(5) \quad \varphi(\theta_i; V_k) = \varphi(\theta_i) \quad \text{and} \quad \varphi'(\theta_i; V_k) = \varphi'(\theta_i) \quad \text{for } i = 0, 1, \dots, k.$$

From (5) we can therefore expect a highly accurate approximation when the iterates  $\theta_i$  are close to a global maximum  $\theta_*$  of (3). Moreover, this process has a snowball effect: the accuracy of all previous  $\theta_0, \dots, \theta_k$  accumulates to produce the next, even more accurate iterate  $\theta_{k+1}$ . The following theorem quantifies this effect. Its proof is deferred to section 2.2.

**THEOREM 2.3** (error recurrence). *Let  $\theta_*$  solve (3) and assume that  $\varphi(\theta_*) = \lambda_{\min}(H(\theta_*))$  is a simple eigenvalue and  $\varphi''(\theta_*) < 0$ . Suppose that Algorithm 1 starts with  $\theta_0 \in [\theta_* - r, \theta_* + r]$  for sufficiently small  $r > 0$ , and all subsequent iterates  $\theta_1, \theta_2, \theta_3, \dots$  remain in  $[\theta_* - r, \theta_* + r]$ . Then*

$$(6) \quad |\theta_{k+1} - \theta_*| \leq |\theta_k - \theta_*| \quad \forall k \geq 1.$$

*Moreover, there exists a constant  $R > 0$  independent of  $k$  and  $\theta_0$  such that the weighted*

error  $e_\ell = |\theta_\ell - \theta_*|/R$  satisfies

$$(7) \quad e_{k+1} \leq C \cdot (k+1) \cdot e_k \prod_{\ell=0}^{k-1} e_\ell^2, \quad k = 1, 2, \dots,$$

where  $C = 16M/|\varphi''(\theta_*)R^2|$  with  $M = \max_{z \in \mathbb{C}: |z-\theta_*|=R} \|H(z)\|_2$ .

In an infinite-dimensional setting, the global convergence of Algorithm 1 established in [13] implies that the assumption  $\theta_k \in [\theta_* - r, \theta_* + r]$  of Theorem 2.3 is satisfied when starting the analysis at  $k \geq K$  with  $K$  sufficiently large.

We proceed by drawing conclusions on the asymptotic convergence from Theorem 2.3. Because it converges in a finite number of steps, it is meaningless to make any asymptotic statement about Algorithm 1 itself. To avoid this difficulty, we consider the upper bounds produced by the recurrence (7) instead and establish high-order convergence for these bounds.

**THEOREM 2.4.** *Consider a sequence  $(e_k)_{k=0}^\infty$  satisfying the recurrence (7), with  $3Ce_0^2 < 1$  and  $2e_1 \leq 1$ . Then the sequence converges to zero with a local  $R$ -convergence order*

$$(8) \quad \sigma = 1 + \sqrt{2} \approx 2.4142.$$

In other words, there is a sequence  $(\varepsilon_k)_{k=0}^\infty$  such that  $e_k \leq \varepsilon_k$  and  $\lim_{k \rightarrow \infty} |\varepsilon_k|/|\varepsilon_{k-1}|^\sigma < \infty$ .

*Proof.* Introducing the auxiliary terms  $\tilde{e}_\ell = (\ell + 1)e_\ell$ , we get from the error recurrence (7) that

$$\tilde{e}_{k+1} \leq 3C \cdot \tilde{e}_k \prod_{\ell=0}^{k-1} \tilde{e}_\ell^2, \quad k = 1, 2, \dots$$

We now define a sequence  $(\varepsilon_k)_{k=0}^\infty$  by letting  $\varepsilon_k$  satisfy this recurrence with “=” and setting  $\varepsilon_0 = \tilde{e}_0$ ,  $\varepsilon_1 = \tilde{e}_1$ . By induction, it holds that  $\tilde{e}_k \leq \varepsilon_k$ . Taking logarithms, we have

$$(9) \quad \ln \varepsilon_{k+1} = \ln \varepsilon_k + 2 \sum_{\ell=0}^{k-1} \ln \varepsilon_\ell + \ln 3C, \quad k = 1, 2, \dots$$

Setting  $s_k = \frac{1}{2} \ln 3C + \sum_{\ell=0}^k \ln \varepsilon_\ell$ , this becomes equivalent to

$$s_{k+1} = 2s_k + s_{k-1}, \quad s_0 = \ln \varepsilon_0 + \frac{1}{2} \ln 3C < 0, \quad s_1 = s_0 + \ln \varepsilon_1 < 0.$$

By standard techniques for solving difference equations, there exist  $c_1, c_2$  such that

$$s_k = c_1(1 + \sqrt{2})^k + c_2(1 - \sqrt{2})^k.$$

Observe that  $s_k < 0$  for all  $k$ . Since the first term dominates when  $k \rightarrow \infty$ , we have  $c_1 < 0$ . From

$$\frac{\ln \varepsilon_k}{\sigma^k} = \frac{s_k - s_{k-1}}{\sigma^k} \rightarrow \frac{c_1 \sqrt{2}}{\sigma} < 0, \quad \text{as } k \rightarrow \infty,$$

one obtains Q-convergence of  $\varepsilon_k$  with order  $\sigma$ . Since  $e_k \leq \tilde{e}_k \leq \varepsilon_k$  for all  $k$ , this shows the claimed R-convergence of  $e_k$ .  $\square$

*Remark 2.5.* We note that the local convergence analysis in [13] also applies to Algorithm 1. In particular, using [13, Thm. 3.3] one can show superlinear convergence with order at least 1.618. In our setting and analysis, we obtained however a more explicit bound for the error recurrence that lead to the improved order 2.414. As we will see in the numerical experiments, this seems to be the correct order.

**2.2. Proof of Theorem 2.3.** Throughout this section, we suppose that the assumptions of Theorem 2.3 hold. In particular,  $\varphi(\theta_*) = \lambda_{\min}(H(\theta_*))$  is a simple eigenvalue of  $H(\theta_*)$  and  $\varphi''(\theta_*) < 0$  holds. We also recall that  $H(\theta)$  is assumed to admit an analytic extension in an open neighborhood  $\Omega_{\mathbb{C}} \subset \mathbb{C}$  containing  $\Omega$ .

In the following, we measure the distance between a vector  $u \in \mathbb{C}^n$  and the subspace spanned by  $V \in \mathbb{C}^{n \times k}$  as

$$(10) \quad \text{dist}(V, u) = \min \{ \|u - v\|_2 : v \in \text{span}(V) \}.$$

We let

$$(11) \quad v(\theta) = \text{a normalized eigenvector belonging to } \lambda_{\min}(H(\theta)).$$

The lemma below shows that the reduced function  $\varphi(\theta; V)$  can be analytically extended to a region around  $\theta_*$ . Moreover, this region can be chosen uniformly over all subspaces  $V$  containing a sufficiently good approximation of  $v(\theta_*)$ .

**LEMMA 2.6** (analyticity). *Let  $V$  be an orthonormal basis satisfying  $\text{dist}(V, v(\theta_*)) \leq \delta$ . Provided that  $\delta > 0$  is sufficiently small, there exists a constant  $R > 0$  (independent of  $V$ ) such that the following statements hold:*

(a) *The function*

$$(12) \quad \widehat{\varphi}(z; V) = \lambda_*(V^H H(z) V),$$

*where  $\lambda_*$  denotes the eigenvalue closest to  $\varphi(\theta_*; V)$ , is well defined on the disc  $D_{\mathbb{C}}(\theta_*; R) = \{z \in \mathbb{C} : |z - \theta_*| \leq R\}$ .*

(b)  *$\widehat{\varphi}(z; V)$  is complex analytic on  $D_{\mathbb{C}}(\theta_*; R)$ .*

(c)  *$\widehat{\varphi}(\theta; V) = \varphi(\theta; V)$  for  $\theta \in \mathbb{R} \cap D_{\mathbb{C}}(\theta_*; R)$ . Hence,  $\widehat{\varphi}(\cdot; V)$  is the analytic continuation of  $\varphi(\cdot; V)$ .*

*Proof.* (a) and (b) Let  $\lambda_{\min-1}(\cdot)$  denote the second smallest eigenvalue of a Hermitian matrix. By our assumptions,

$$\text{gap} = \lambda_{\min-1}(H(\theta_*)) - \lambda_{\min}(H(\theta_*)) > 0.$$

We claim this induces a nonzero gap for the reduced matrix  $V^H H(\theta_*) V$  provided that  $\delta$  is sufficiently small. To see this, let  $v(\theta_*) = v + d$  with  $v \in \text{span}(V)$  and  $\|d\|_2 = \text{dist}(V, v(\theta_*)) \leq \delta$ . Clearly,  $v$  depends continuously on  $\delta$  around zero and hence  $v^H H(\theta_*) v / v^H v \rightarrow v(\theta_*)^H H(\theta_*) v(\theta_*) = \lambda_{\min}(H(\theta_*))$  as  $\delta \rightarrow 0$ . We can therefore choose  $\delta > 0$  such that

$$\lambda_{\min}(V^H H(\theta_*) V) = \min_{x \in \mathcal{V}} \frac{x^H H(\theta_*) x}{x^H x} \leq \frac{v^H H(\theta_*) v}{v^H v} < \lambda_{\min}(H(\theta_*)) + \text{gap}/2.$$

In turn, using eigenvalue interlacing, the gap for the reduced problem satisfies

$$(13) \quad \begin{aligned} \widetilde{\text{gap}} &= \lambda_{\min-1}(V^H H(\theta_*) V) - \lambda_{\min}(V^H H(\theta_*) V) \\ &> \lambda_{\min-1}(H(\theta_*)) - \lambda_{\min}(H(\theta_*)) - \text{gap}/2 = \text{gap}/2. \end{aligned}$$

In particular,  $\lambda_{\min}(V^H H(\theta_*) V)$  is a simple eigenvalue.

We now choose  $R > 0$  sufficiently small such that  $D_{\mathbb{C}}(\theta_*; R)$  is contained in  $\Omega_{\mathbb{C}}$  (where  $H(z)$  is analytic) and, moreover,

$$\|H(z) - H(\theta_*)\|_2 \leq \text{gap}/4 \quad \forall z \in D_{\mathbb{C}}(\theta_*; R).$$

Because  $V$  is an orthonormal basis, this also implies

$$(14) \quad \|V^H H(z) V - V^H H(\theta_*) V\|_2 \leq \text{gap}/4, \quad \forall z \in D_{\mathbb{C}}(\theta_*; R).$$

Note that  $V^H H(\theta_*) V$  is Hermitian while  $V^H H(z) V$  is not. Applying classical eigenvalue perturbation results for *nearly Hermitian* matrices (see, e.g., [12] and [25, Thm. 5.1, Chap. IV, and the discussion below its proof]) implies that there is exactly one eigenvalue of  $V^H H(z) V$  in the disc  $D_{\mathbb{C}}(\lambda_{\min}(V^H H(\theta_*) V); \text{gap}/4)$  for each  $z \in D_{\mathbb{C}}(\theta_*; R)$ , thus showing ((a)), that the eigenvalue function  $\lambda_*(V^H H(z) V)$  is well defined on  $D_{\mathbb{C}}(\theta_*; R)$ . This function is complex analytic due to the analyticity of simple eigenvalues [14]. Since the choice of  $R$  to satisfy (14) is independent of  $V$ , once  $\delta$  is sufficiently small, and we have thus proved ((b)).

((c)) The result follows immediately from the continuity of  $\lambda_{\min}(V^H H(\theta) V)$  for  $\theta \in \mathbb{R}$  and the uniqueness of the eigenvalue  $\lambda_*(V^H H(\theta) V)$  in the interval  $\theta \in [\theta_* - R, \theta_* + R] \subset D_{\mathbb{C}}(\theta_*, R)$ .  $\square$

Trivially, because of  $\varphi(\theta) = \varphi(\theta; I_n)$ , Lemma 2.6 implies that  $\varphi(\theta)$  also has an analytic extension to  $\widehat{\varphi}(z) = \widehat{\varphi}(z; I_n)$  in the complex disc  $D_{\mathbb{C}}(\theta_*; R)$  with the same estimate for  $R$ .

LEMMA 2.7 (analytic approximation). *Let  $\delta$  and  $R$  be as in Lemma 2.6. Then there exists  $0 < r \leq R/2$  such that the following statements holds:*

- (i)  $\text{dist}(v(\theta_*), v(\theta)) \leq \delta$  for all  $\theta \in \mathbb{R}$  with  $|\theta - \theta_*| \leq r$ .
- (ii) Consider any  $k+1$  mutually distinct points  $\{\theta_i\}_{i=0}^k \subset [\theta_* - r, \theta_* + r]$  ordered in decreasing distance to  $\theta_*$ , and suppose that  $v(\theta_0), v(\theta_1), \dots, v(\theta_k) \in \text{span}(V)$ . Then the following hold:

- (a) Both  $\varphi(\theta)$  and  $\varphi(\theta; V)$  are real analytic on  $[\theta_* - R, \theta_* + R]$ .
- (b) With  $M = \max_{z \in \mathbb{C}: |z - \theta_*| = R} \|H(z)\|_2$ ,

$$(15) \quad |\varphi'(\theta_*; V) - \varphi'(\theta_*)| \leq \frac{8M(k+1)}{(R/2)^{2k+2}} |\theta_* - \theta_k| \prod_{i=0}^{k-1} |\theta_* - \theta_i|^2.$$

- (c) For all  $\theta \in \mathbb{R}$  with  $|\theta - \theta_*| \leq r$ ,

$$(16) \quad \varphi''(\theta; V) \leq \varphi''(\theta_*)/2 < 0.$$

*Proof.* (i) This follows immediately from the fact that the eigenvalue  $\lambda_{\min}(H(\theta_*))$  is simple, which implies the eigenvector  $v(\theta)$  can be chosen to be continuous around  $\theta_*$ . Hence,  $v(\theta) \rightarrow v(\theta_*)$  as  $\theta \rightarrow \theta_*$ ; see, e.g., [14].

(ii)(a) This follows from Lemma 2.6, since  $\text{dist}(V, v(\theta_*)) \leq \text{dist}(v(\theta_i), v(\theta_*)) \leq \delta$  according to (i).

(ii)(b) interpolation of  $\varphi(\theta)$  close to  $\theta_*$ . We first recall from Lemma 2.1 that

$$\varphi(\theta_i; V) = \varphi(\theta_i), \quad \varphi'(\theta_i; V) = \varphi'(\theta_i) \quad \text{for } i = 0, \dots, k.$$

Let  $\Pi$  be the permutation such that  $\theta_{\Pi(0)} < \theta_{\Pi(1)} < \dots < \theta_{\Pi(k)}$ . Then the first interpolation condition combined with Rolle's theorem implies that there exist  $\widehat{\theta}_i$  with

$$\theta_{\Pi(i)} < \widehat{\theta}_i < \theta_{\Pi(i+1)} \quad \text{and} \quad \varphi'(\widehat{\theta}_i; V) = \varphi'(\widehat{\theta}_i) \quad \text{for } i = 0, \dots, k-1.$$

This implies  $|\theta_* - \widehat{\theta}_i| < |\theta_* - \theta_{\Pi(i)}|$  if  $\theta_{\Pi(i+1)} \leq \theta_*$  and  $|\theta_* - \widehat{\theta}_i| < |\theta_* - \theta_{\Pi(i+1)}|$  if  $\theta_{\Pi(i)} \geq \theta_*$ . For  $i_*$  with  $\theta_{\Pi(i_*)} \leq \theta_* \leq \theta_{\Pi(i_*+1)}$  the assumed ordering implies  $\Pi(i_*) = k$  or  $\Pi(i_* + 1) = k$ . Without loss of generality, let us suppose that  $\Pi(i_*) = k$ , which in turn implies  $|\theta_* - \widehat{\theta}_{i_*}| < |\theta_* - \theta_{\Pi(i_*+1)}|$ . In summary, we have

$$(17) \quad \prod_{i=0}^{k-1} |\theta_* - \widehat{\theta}_i| = \prod_{i=0}^{i_*-1} |\theta_* - \widehat{\theta}_i| \prod_{i=i_*}^{k-1} |\theta_* - \widehat{\theta}_i| \\ < \prod_{i=0}^{i_*-1} |\theta_* - \theta_{\Pi(i)}| \prod_{i=i_*}^{k-1} |\theta_* - \theta_{\Pi(i+1)}| = \prod_{i=0}^{k-1} |\theta_* - \theta_i|.$$

Since the real analytic function  $\varphi'(\theta; V) - \varphi'(\theta)$  is zero at  $2k + 1$  distinct points  $\theta_i$  and  $\widehat{\theta}_i$  in  $[\theta_* - r, \theta_* + r]$ , standard interpolation error results (see, e.g., [6]) yield

$$(18) \quad |\varphi'(\theta; V) - \varphi'(\theta)| = \left| \frac{\varphi^{(2k+2)}(\xi) - \varphi^{(2k+2)}(\xi; V)}{(2k + 1)!} \prod_{i=0}^k (\theta - \theta_i) \prod_{i=0}^{k-1} (\theta - \widehat{\theta}_i) \right|$$

for all  $\theta \in [\theta_* - r, \theta_* + r]$  with some  $\xi \equiv \xi(\theta) \in [\theta_* - r, \theta_* + r]$ . Next, to bound the derivative term, we use that  $\varphi(\theta; V)$  extends analytically to the complex disc  $D_{\mathbb{C}}(\theta_*; R)$  by Lemma 2.6. Using Cauchy’s integral formula (see, e.g., [21]),

$$(19) \quad |\varphi^{(m)}(\xi; V)| = |\widehat{\varphi}^{(m)}(\xi; V)| = \left| \frac{m!}{2\pi i} \oint_{|z-\theta_*|=R} \frac{\widehat{\varphi}(z; V)}{(z - \xi)^{m+1}} dz \right| \leq \frac{2M \cdot m!}{(R/2)^m},$$

where the last inequality uses  $|\widehat{\varphi}(z; V)| = |\lambda_*(V^H H(z) V)| \leq \|H(z)\|_2 \leq M$ , as well as  $|z - \xi| \geq |z - \theta_*| - |\theta_* - \xi| \geq R - r \geq R/2$ . Setting  $V = I_n$ , we also obtain

$$|\varphi^{(m)}(\xi)| = |\varphi^{(m)}(\xi; I_n)| \leq \frac{2M \cdot m!}{(R/2)^m}.$$

Plugging these two bounds with  $m = 2k + 2$  and (17) into (18) yields (15).

(ii)(c) Because  $\varphi'(\theta; V) - \varphi'(\theta)$  has  $2k + 1$  zeros, Rolle’s theorem implies that  $\varphi''(\theta; V) - \varphi''(\theta)$  has  $2k$  distinct zeros in  $[\theta_* - r, \theta_* + r]$ . Therefore, in analogy to (15) one can show that  $|\varphi''(\theta_*; V) - \varphi''(\theta_*)| \rightarrow 0$  as  $r \rightarrow 0$ . In turn,

$$|\varphi''(\theta_*) - \varphi''(\theta; V)| \leq |\varphi''(\theta; V) - \varphi''(\theta_*; V)| + |\varphi''(\theta_*; V) - \varphi''(\theta_*)| \rightarrow 0.$$

Because of  $\varphi''(\theta_*) < 0$  this shows that (16) is satisfied when  $r > 0$  is chosen sufficiently small.  $\square$

Now we are ready to prove the main result.

*Proof of Theorem 2.3.* Without loss of generality, we may assume that all  $\theta_i$ ,  $i = 0, \dots, k$  are mutually distinct because otherwise  $e_{k+1} = 0$  by Lemma 2.2 and (7) is trivially satisfied.

Let us now assume that  $r$  is sufficiently small such that the statements of Lemma 2.7 and, in turn, Lemma 2.6 hold. Then  $\varphi(\theta; V_k)$  is real analytic on  $[\theta_* - R, \theta_* + R]$ . The optimality condition  $\varphi'(\theta_{k+1}; V_k) = 0$  in the  $k$ th iteration of Algorithm 1 combined with Taylor expansion give

$$0 = \varphi'(\theta_{k+1}; V_k) = \varphi'(\theta_*; V_k) + \varphi''(\xi; V_k)(\theta_{k+1} - \theta_*)$$



for some  $\xi$  with  $|\xi - \theta_*| \leq r$  Therefore,

$$(20) \quad |\theta_{k+1} - \theta_*| = |\varphi'(\theta_*; V_k)|/|\varphi''(\xi; V_k)| \leq 2|\varphi'(\theta_*; V_k)|/|\varphi''(\theta_*)|,$$

where we used (16) in the last inequality. Using the derivative bound (15) and exploiting  $\varphi'(\theta_*) = 0$ , we obtain

$$|\theta_{k+1} - \theta_*| \leq |\theta_k - \theta_*| \cdot \frac{16M(k+1)}{|\varphi''(\theta_*)| \cdot (R/2)^2} \cdot (2r/R)^{2k}.$$

This implies the monotonicity property (6) when choosing  $r$  sufficiently small.

The error recurrence (7) follows immediately from (20), combined with monotonicity (6) and the derivative bound (15). (For notational simplicity, we replaced  $R/2$  by  $R$ .)  $\square$

**3. Crawford number computation.** We will now apply and specialize the developments from the previous section to the computation of the Crawford number, which has the form of an eigenvalue optimization problem characterization (2). We apply the subspace method, Algorithm 1, to maximize

$$(21) \quad \varphi(\theta) = \lambda_{\min}(H(\theta)) \quad \text{with} \quad H(\theta) = S \cos \theta + K \sin \theta$$

with the Hermitian matrices  $S = (A + A^H)/2$  and  $K = (A - A^H)/2i$ . This function is clearly  $2\pi$  periodic and continuous. When maximizing  $\varphi(\theta)$ , we can therefore restrict its domain to any interval  $[\theta_0, \theta_0 + 2\pi]$ .

**3.1. Properties of the objective function.** The result of the following lemma implies that, assuming  $\gamma(A) > 0$ , the objective function  $\varphi$  is quasi-concave (unimodal) on the open set  $\{\theta : \varphi(\theta) > 0\}$  and after restricting  $\theta$  to an appropriately chosen interval of length  $2\pi$ . A similar result<sup>1</sup> was also shown in [8, sec. 4.3] but we prove in addition that  $\varphi$  is strongly concave.

LEMMA 3.1 (strong concavity). *Let  $\gamma(A) > 0$ . Then there exists  $\theta_0 \in \mathbb{R}$  such that*

$$(22) \quad \{\theta : \varphi(\theta) > 0\} \cap [\theta_0, \theta_0 + 2\pi]$$

*is an open, nonempty interval  $(\ell, u)$  of length at most  $\pi$ . Moreover,  $\varphi(\theta)$  is strongly concave on any closed subinterval of  $(\ell, u)$ .*

*Proof.* Let  $\theta_*$  be a maximizer of  $\varphi(\theta)$ . By replacing  $A$  with  $e^{-i\theta_*}A$ , we may assume without loss of generality that  $\theta_* = 0$ , and in turn  $\varphi(0) = \gamma(A) > 0$ . We will prove the result for  $\theta_0 = -\pi$ .

We have

$$(23) \quad \varphi(\theta) = \min_{\|v\|_2=1} v^H H(\theta)v = \min_{\|v\|_2=1} v^H S v \cdot \cos \theta + v^H K v \cdot \sin \theta = \min_{x+iy \in \mathcal{F}(A)} x \cos \theta + y \sin \theta.$$

Writing  $z = x + iy = |z|(\cos \text{Arg } z + i \sin \text{Arg } z)$  with the argument  $\text{Arg } z \in (-\pi, \pi]$ , we obtain

$$(24) \quad \varphi(\theta) = \min_{z \in \mathcal{F}(A)} |z| \cdot \cos(\theta - \text{Arg } z).$$

<sup>1</sup>The matrix  $K = (A - A^H)/2i$  needs to be invertible in order to apply Theorem 4.1 in [8] and thus conclude quasi-concavity of  $\varphi$ . However, as one of the referees pointed out, one can apply the result to a rotated matrix  $e^{-i\theta}A$  and reach the same conclusion. Since  $\gamma(A) > 0$ , there exists a  $\theta$  such that the corresponding matrix  $K$  becomes invertible.

Note that

$$0 < \varphi(0) = \min_{z \in \mathcal{F}(A)} |z| \cos(\operatorname{Arg} z)$$

implies that  $\operatorname{Arg} z \in (-\pi/2, \pi/2)$  and, in turn,  $(\operatorname{Arg} z - \pi/2, \operatorname{Arg} z + \pi/2) \subset (-\pi, \pi)$  for every  $z \in \mathcal{F}(A)$ . This allows us to write

$$\begin{aligned} \{\theta: \varphi(\theta) > 0\} \cap [-\pi, \pi] &= \bigcap_{z \in \mathcal{F}(A)} \{\theta: \cos(\theta - \operatorname{Arg} z) > 0\} \cap [-\pi, \pi] \\ (25) \qquad \qquad \qquad &= \bigcap_{z \in \mathcal{F}(A)} (\operatorname{Arg} z - \pi/2, \operatorname{Arg} z + \pi/2). \end{aligned}$$

For the last equality, we have used that  $\cos(\theta - \operatorname{Arg} z) > 0$  is equivalent to  $\theta - \operatorname{Arg} z \in (-\pi/2, \pi/2)$ . The right-hand side of (25) shows that this set is an interval of length at most  $\pi$ . The left-hand side, together with  $\varphi(0) > 0$  and the continuity of  $\varphi(\theta)$ , shows that this interval is nonempty and open. This proves the first part of the lemma.

To show strong concavity, let  $[\widehat{\ell}, \widehat{u}] \subset (\ell, u)$ . Because  $\varphi(\theta)$  is continuous there exists  $\delta > 0$  such that

$$\varphi(\theta) \geq \delta \quad \text{for } \theta \in [\widehat{\ell}, \widehat{u}].$$

Denoting  $\varphi_z(\theta) = |z| \cdot \cos(\theta - \operatorname{Arg} z)$ , it follows from (24) that for every  $\tilde{z} \in \mathcal{F}(A)$

$$\varphi''_{\tilde{z}}(\theta) = -\varphi_{\tilde{z}}(\theta) \leq -\min_{z \in \mathcal{F}(A)} \varphi_z(\theta) = -\varphi(\theta) \leq -\delta \quad \text{for } \theta \in [\widehat{\ell}, \widehat{u}].$$

Hence, the function  $\varphi(\theta)$  is the minimum of strongly concave functions  $\varphi_z(\theta)$  with uniform bound on the second derivative on a closed interval. This proves strong concavity of  $\varphi(\theta)$  on  $[\widehat{\ell}, \widehat{u}]$ .  $\square$

The strong concavity property implies, in particular, that the global maximizer  $\theta_*$  of  $\varphi(\theta)$  within any interval of length  $2\pi$  when  $\gamma(A) > 0$ . Moreover, if  $\varphi(\theta)$  is twice differentiable at  $\theta_*$ , we also have

$$(26) \qquad \qquad \qquad \varphi''(\theta_*) < 0.$$

This will for example be the case when  $\varphi(\theta_*) = \lambda_{\min}(H(\theta_*))$  is a simple eigenvalue.

*Remark 3.2.* Being continuous and  $2\pi$ -periodic, the function  $\varphi(\theta)$  cannot be strongly concave on an interval of length  $2\pi$ . On the other hand, Figure 1 (see also [4, Fig. 2]) demonstrates that  $\varphi(\theta)$  may well be concave on an interval larger than guaranteed by Lemma 3.1. In addition, one observes that the negative part of  $\varphi(\theta)$  is not necessarily convex.

As illustrated in Figure 1, the function value  $\varphi(\theta)$  admits a geometric interpretation in terms of the numerical range  $\mathcal{F}(A)$ . Recalling (23), we have

$$\varphi(\theta) = \min_{x+iy \in \mathcal{F}(A)} x \cos \theta + y \sin \theta.$$

This corresponds to the *signed distance* of the numerical range to zero along the direction  $\cos \theta + i \sin \theta$  (since the right-hand side corresponds to the length of the orthogonal projection of points in the numerical range  $\mathcal{F}(A)$  onto the direction vector). It also implies that

$$(27) \qquad \qquad \qquad \{x + iy: x \cos \theta + y \sin \theta - \varphi(\theta) = 0\}$$

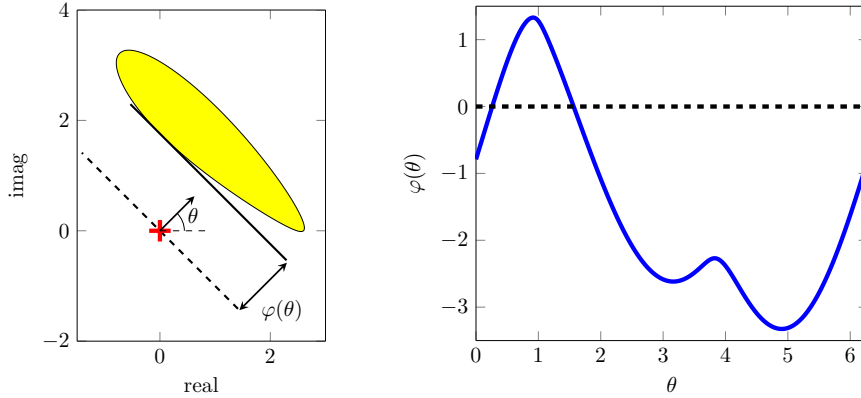


FIG. 1. Left: Numerical range of a random 4-by-4 matrix (left). Right: Corresponding function  $\varphi(\theta)$  on the interval  $[0, 2\pi]$ .

defines the supporting line for the convex set  $\mathcal{F}(A) = \mathcal{F}(S + \iota K)$  with normal  $\cos \theta + \iota \sin \theta$ . Its supporting point is given by

$$(28) \quad p(\theta) = v(\theta)^H (S + \iota K) v(\theta),$$

since the minimum of (23) is attained at an eigenvector  $v(\theta)$ , as defined in (11). We will use these geometric properties in the following section to get a better understanding of the subspace method.

**3.2. Subspace method applied to Crawford number computation.** As mentioned before, the full subspace method for computing the Crawford number simply consists of applying Algorithm 1 to the function  $\varphi$  from (21). The objective function of the reduced maximization problem takes the form

$$\varphi(\theta; V_k) = \lambda_{\min}(V_k^H H(\theta) V_k) = \lambda_{\min}(S_k \cos \theta + K_k \sin \theta),$$

where  $S_k = V_k^H S V_k = (A_k + A_k^H)/2$  and  $K_k = V_k^H K V_k = (A_k - A_k^H)/2\iota$  with  $A_k = V_k^H A V_k$ . This shows that the reduced problem amounts to computing the Crawford number of  $A_k$ :

$$(29) \quad \gamma(A_k) = \max \left\{ \max_{\theta \in [0, 2\pi]} \varphi(\theta; V_k), 0 \right\}.$$

Because  $\gamma(A_k) \geq \gamma(A)$ , one can immediately stop the algorithm and return  $\gamma(A) = 0$  once it detects  $\gamma(A_k) = 0$ .

Since  $A_k$  is expected to be a small matrix, the computation of  $\gamma(A_k)$ , together with the optimal value  $\theta_{k+1}$  in line 4 of Algorithm 1, is not expensive and can be performed by many algorithms; see, e.g., [11, 8, 26, 3]. Our implementation uses a straightforward adaptation of the criss-cross type algorithm from [20], which we observed to work very well in our numerical experiments.

**3.2.1. Properties of the reduction objective function and convergence.**

Because of (29), the strong concavity property established in Lemma 3.1 also holds for the reduced objective function  $\varphi(\theta; V_k)$ . In particular, the global maximizer of  $\varphi(\theta; V_k)$  is unique (within one period). This property allows us to gain additional insights into the behavior of the subspace method.

LEMMA 3.3. *Let  $V$  be an orthonormal basis of a subspace  $\mathcal{V}$  of  $\mathbb{C}^n$ . Assume that  $\gamma(A) > 0$  and let  $\theta_\ell, \theta_u$  be such that  $\theta_\ell < \theta_u < \theta_\ell + 2\pi$ . If  $v(\theta_\ell), v(\theta_u) \in \mathcal{V}$ , then*

$$\arg \max_{\theta \in [\theta_\ell, \theta_\ell + 2\pi]} \varphi(\theta; V) \in (\theta_\ell, \theta_u) \quad \text{if and only if} \quad \arg \max_{\theta \in [\theta_\ell, \theta_\ell + 2\pi]} \varphi(\theta) \in (\theta_\ell, \theta_u).$$

*Proof.* Suppose that

$$\theta_V = \arg \max_{\theta \in [\theta_\ell, \theta_\ell + 2\pi]} \varphi(\theta; V) \in (\theta_\ell, \theta_u) \quad \text{and} \quad \theta_* = \arg \max_{\theta \in [\theta_\ell, \theta_\ell + 2\pi]} \varphi(\theta) \notin (\theta_\ell, \theta_u),$$

which implies  $\theta_\ell < \theta_V < \theta_u < \theta_* \leq \theta_\ell + 2\pi$ . Letting  $\beta = \max\{\varphi(\theta_\ell), \varphi(\theta_u), 0\}$ , Lemmas 2.1 and 3.1 give

$$\begin{aligned} \varphi(\theta_\ell) = \varphi(\theta_\ell; V) &\leq \beta, & \varphi(\theta_V; V) &> \beta, \\ \varphi(\theta_u) = \varphi(\theta_u; V) &\leq \beta, & \varphi(\theta_*; V) &\geq \varphi(\theta_*) > \beta. \end{aligned}$$

In particular, this implies that  $\varphi(\theta; V)$  has two different local, strictly positive maxima within an interval of length  $2\pi$ . This contradicts the strong concavity property from Lemma 3.1. Hence,  $\theta_V \in (\theta_\ell, \theta_u)$  implies  $\theta_* \in (\theta_\ell, \theta_u)$ . The proof of the other direction proceeds analogously.  $\square$

The particular structure of  $\varphi$  allows us, for the special case of the Crawford number, to remove some of the assumptions of Theorem 2.3 on the local convergence of the subspace method.

THEOREM 3.4 (convergence of subspace method). *Consider  $A \in \mathbb{C}^{n \times n}$  with  $\gamma(A) > 0$ . Let  $\theta_*$  be the maximizer of (21) and assume that  $\lambda_{\min}(H(\theta_*))$  is a simple eigenvalue. If the initial guess  $\theta_0$  is sufficiently close to  $\theta_*$ , then the iterates  $\theta_1, \theta_2, \theta_3, \dots$  produced by Algorithm 1 satisfy the monotonicity property (6) and the error recurrence (7).*

*Proof.* Lemma 3.1 shows that the assumption  $\varphi''(\theta_*) < 0$  of Theorem 2.3 is satisfied. It remains to show that the other assumption—all iterates remain close to  $\theta_*$ —is also satisfied.

Without loss of generality, we may assume that  $\theta_0 < \theta_* = 0$ ,  $\varphi(\theta_0) > 0$ . Considering eigenvectors  $v(\theta_0), v(\theta_*)$  with  $\|v(\theta_0)\|_2 = \|v(\theta_*)\|_2 = 1$ , this implies

$$(30) \quad \gamma(A) = v(\theta_*)^H H(\theta_*) v(\theta_*) = v(\theta_*)^H S v(\theta_*), \quad 0 = v(\theta_*)^H H'(\theta_*) v(\theta_*) = v(\theta_*)^H K v(\theta_*).$$

The reduced function during the first iteration takes the form

$$(31) \quad \begin{aligned} \varphi(\theta; v(\theta_0)) &= \lambda_{\min}(v(\theta_0)^H H(\theta) v(\theta_0)) \\ &= \alpha \cos \theta + \beta \sin \theta = |z_0| \cos(\theta - \text{Arg } z_0) \end{aligned}$$

with  $\alpha = v(\theta_0)^H S v(\theta_0)$ ,  $\beta = v(\theta_0)^H K v(\theta_0)$ , and  $z_0 = \alpha + i\beta$ . Since  $\varphi(\theta_0; v(\theta_0)) = \varphi(\theta_0)$ , we obtain the super-level set

$$\Omega_0 = \{\theta: \varphi(\theta; v(\theta_0)) \geq \varphi(\theta_0)\} = \{\theta: \cos(\theta - \text{Arg } z_0) \geq \cos(\theta_0 - \text{Arg } z_0)\}.$$

Since  $\varphi(\theta_0) > 0$ , we get  $\cos(\theta_0 - \text{Arg } z_0) > 0$ . In addition, by strong concavity  $\varphi'(\theta_0) > 0$  and thus Lemma 2.1 implies  $\sin(\theta_0 - \text{Arg } z_0) < 0$ . Hence, by elementary trigonometry  $\theta_0 - \text{Arg } z_0 \in (-\pi/2, 0)$  and so we obtain up to a  $2\pi$  period that  $\Omega_0 = [\theta_0, -\theta_0 + 2 \text{Arg } z_0]$ . Because of the monotonicity property in Lemma 2.1, the iterates of Algorithm 1 satisfy for  $k \geq 2$  the inequalities

$$\varphi(\theta_k; v(\theta_0)) \geq \varphi(\theta_k; V_{k-1}) \geq \varphi(\theta_*) \geq \varphi(\theta_0),$$

and, hence, they stay in the superlevel set:

$$\{\theta_0, \theta_1, \dots\} \subset \Omega_0.$$

We now show that the length of  $\Omega_0$  becomes arbitrarily small when choosing  $\theta_0$  sufficiently close to 0. First, from the interpolation conditions in Lemma 2.1 it follows that

$$\begin{aligned} \varphi(\theta_0) &= \varphi(\theta_0; v(\theta_0)) = \alpha \cos \theta_0 + \beta \sin \theta_0, \\ \varphi'(\theta_0) &= \varphi'(\theta_0; v(\theta_0)) = -\alpha \sin \theta_0 + \beta \cos \theta_0. \end{aligned}$$

By solving this linear equation in  $\alpha \equiv \alpha(\theta_0)$ ,  $\beta \equiv \beta(\theta_0)$ , we obtain

$$\alpha(\theta_0) = \varphi(\theta_0) \cos \theta_0 - \varphi'(\theta_0) \sin \theta_0, \quad \beta(\theta_0) = \varphi(\theta_0) \sin \theta_0 + \varphi'(\theta_0) \cos \theta_0.$$

Using Taylor expansions of  $\alpha(\theta_0)$ ,  $\beta(\theta_0)$  around 0, and exploiting  $\varphi'(0) = 0$ , we have

$$(32) \quad \frac{\beta}{\alpha} = \theta_0 + \frac{\varphi''(0)}{\varphi(0)}\theta_0 + O(\theta_0^2)$$

and therefore

$$|\Omega_0| = 2|\operatorname{Arg} z_0 - \theta_0| = 2\left|\arctan \frac{\beta}{\alpha} - \theta_0\right| = 2\left|\frac{\varphi''(0)}{\varphi(0)}\theta_0\right| + O(\theta_0^2).$$

Since  $\varphi(0) = \gamma(A) > 0$  and  $\varphi''(0)$  are constants independent of  $\theta_0$ ,  $|\Omega_0| \rightarrow 0$  as  $\theta_0 \rightarrow 0$ .

In summary, for every  $r > 0$ , we can attain

$$\{\theta_0, \theta_1, \dots\} \subset \Omega_0 \subset [-r, r],$$

by choosing  $\theta_0 < 0$  sufficiently close to 0, which completes the proof. □

We end this section with a geometric interpretation of the subspace method in terms of the numerical range; see Figure 2. Recalling (1), we have the reduced Crawford number

$$(33) \quad \gamma(A_k) = \min\{|z|: z \in \mathcal{F}(A_k)\}, \quad A_k = S_k + \iota K_k.$$

It follows immediately from definition (1) that the numerical range of the reduced matrix  $A_k$  is a subset of the original one:

$$\mathcal{F}(A_k) = \mathcal{F}(V_k^H A V_k) \subset \mathcal{F}(A).$$

The boundaries of  $\mathcal{F}(A)$  and  $\mathcal{F}(A_k)$  touch each other at the *supporting points*

$$p(\theta_i) = v_i^H (S + \iota K) v_i, \quad i = 1, \dots, k,$$

because the eigenvectors  $v_1, \dots, v_k$  in Algorithm 1 are contained in  $\operatorname{span}(V_k)$ . The two boundaries also share the supporting lines

$$\{x + \iota y: x \cos \theta_i + y \sin \theta_i - \varphi(\theta_i) = 0\}, \quad i = 1, \dots, k,$$

due to the interpolation  $\varphi(\theta_i; V_k) = \varphi(\theta_i)$  together with (27)–(28). In the particular case of a polygonal numerical range (which is equivalent to  $A$  being a normal matrix), the reduced numerical range interpolates the edge of the polygon exactly if the two vertices belong to the supporting points, as illustrated in the right plot of Figure 2.

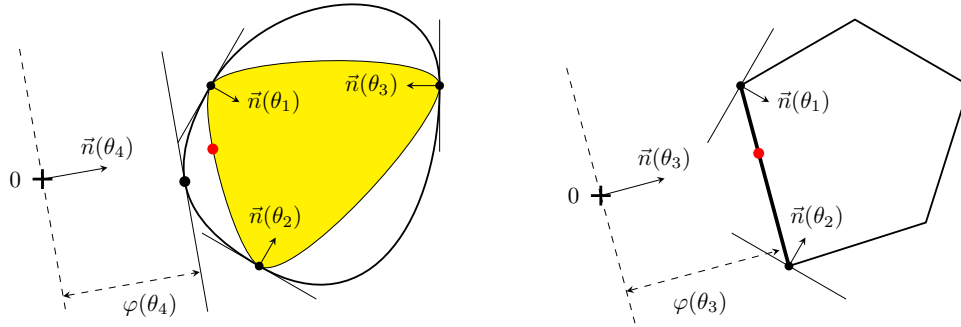


FIG. 2. Geometric illustration of the subspace method, Algorithm 1, for Crawford number computation. Left: The thick black line and the yellow region correspond to the full and reduced numerical range, respectively. The black dots are the supporting points, and the red dot is the closest point in the yellow region to 0. Right: Example of a polygonal numerical range for which the reduced numerical range reduces to a line on the boundary (the thick black line with the red dot).

**3.3. The 3-vector subspace method.** The concavity property of  $\varphi(\theta)$  leading to the bracketing property in Lemma 3.3 allows us to develop a variant of the subspace method that only uses 3 vectors to generate the projection subspace but still enjoys favorable convergence properties.

Given three sampling points  $0 < \ell_k < \theta_k < u_k \leq 2\pi$  such that the global maximum  $\theta_*$  is contained in  $[\ell_k, u_k]$ , the  $k$ th iteration of our 3-vector subspace method proceeds by first determining an orthonormal basis  $V_k$  of the corresponding three eigenvectors  $v(\ell_k), v(\theta_k), v(u_k)$ . It then computes

$$\hat{\theta} = \arg \max_{\theta \in [0, 2\pi]} \lambda_{\min}(V_k^H H(\theta) V_k).$$

By Lemma 3.3,  $\hat{\theta} \in [\ell_k, u_k]$  and, moreover,

$$\theta_* \in \begin{cases} [\ell_k, \theta_k] & \text{if } \hat{\theta} \in [\ell_k, \theta_k], \\ [\theta_k, u_k] & \text{otherwise.} \end{cases}$$

Therefore, we set  $[\ell_{k+1}, \theta_{k+1}, u_{k+1}] = [\ell_k, \hat{\theta}, \theta_k]$  if  $\hat{\theta} \in [\ell_k, \theta_k]$  and  $[\ell_{k+1}, \theta_{k+1}, u_{k+1}] = [\theta_k, \hat{\theta}, u_k]$  otherwise. In turn, we obtain a shorter interval  $[\ell_{k+1}, u_{k+1}]$  that still contains  $\theta_*$ . This allows us to repeat the process until convergence is achieved, leading to Algorithm 2.

Note that in contrast to the (full) subspace method, Algorithm 2 will not produce the exact solution after  $n - 1$  steps. The following theorem shows that it still converges globally and attains fast local convergence. For simplicity and generality, the initial search interval in Algorithm 2 equals  $[0, 2\pi]$ . Smaller intervals can, however, be obtained using the procedure in [3, section 7] but at the expense of computing all the eigenvalues of the pair  $(S, K)$  with  $S = (A + A^H)/2$  and  $K = (A - A^H)/2i$ .

**THEOREM 3.5** (convergence of 3-vector method). *Let  $\gamma(A) > 0$ . Then the iterates  $\theta_1, \theta_2, \theta_3, \dots$  produced by Algorithm 2 are globally convergent to  $\theta_*$ . The convergence is locally at least quadratic, provided that  $\lambda_{\min}(H(\theta_*))$  is a simple eigenvalue. Moreover, if the sequence is also locally alternating around  $\theta_*$ , i.e.,*

$$(\theta_{k+1} - \theta_*)(\theta_k - \theta_*) < 0 \quad \forall k > p$$

with some  $p \in \mathbb{N}$ , then the local convergence order is improved to  $\sigma \approx 2.26953$ .

---

**Algorithm 2** 3-vector subspace method.

---

**Input:** Matrix  $A \in \mathbb{C}^{n \times n}$ , initial guess  $\theta_0 \in (0, 2\pi)$ , tolerance  $\text{tol} > 0$ .

**Output:** Approximation of Crawford number  $\gamma(A)$ .

- 1: Initialize:  $\ell_0 = 0$ ,  $u_0 = 2\pi$ , and  $V_0 = \text{orth}([v(\ell_0), v(\theta_0), v(u_0)])$ .
  - 2: **for**  $k = 0, 1, 2, \dots$  **do**
  - 3:    $\theta_{k+1} = \arg \max_{[0, 2\pi]} \lambda_{\min}(V_k^H H(\theta) V_k)$  with  $\varphi_{k+1} = \lambda_{\min}(V_k^H H(\theta_{k+1}) V_k)$ .
  - 4:   **if**  $\varphi_{k+1} \leq 0$ , **then** return  $\gamma(A) = 0$ .
  - 5:   *Stopping criteria:* **if**  $k > 0$  **and**  $\varphi_{k+1} - \lambda_k < \text{tol} \cdot |\varphi_{k+1}|$ , **then** return  $\gamma(A) = \varphi_{k+1}$ .
  - 6:   *Interval update:* **if**  $\theta_{k+1} \in [\ell_k, \theta_k]$ , **then**  $(\ell_{k+1}, u_{k+1}) = (\ell_k, \theta_k)$ , **otherwise**  $(\ell_{k+1}, u_{k+1}) = (\theta_k, u_k)$ .
  - 7:   Compute smallest eigenvalue  $\lambda_{k+1}$  and corresponding eigenvector  $v(\theta_{k+1})$  of  $H(\theta_{k+1})$ .
  - 8:   *Subspace update:*  $V_{k+1} = \text{orth}([v(\ell_{k+1}), v(\theta_{k+1}), v(u_{k+1})])$ .
  - 9: **end for**
- 

*Proof.* Since  $(\ell_k)_{k=0}^\infty$  is monotonically increasing, it converges to a limiting point  $\ell$ . Similarly,  $(u_k)_{k=0}^\infty$  monotonically decreases to a limiting point  $u$ . As  $(\ell_k)$  and  $(u_k)$  contain all iterates  $\theta_k$  for  $k = 0, 1, 2, \dots$ , it follows that  $(\theta_k)_{k=0}^\infty$  has at most two limiting points  $\ell$  and  $u$ . Let  $(\theta_{k_i})$ ,  $i = 0, 1, \dots$ , be a converging subsequence and assume, without loss of generality, that it converges to  $\ell$ .

Since  $\text{span}(V_{k_i-1})$  contains  $v(\ell_{k_i-1})$ , the interpolation and monotonicity properties from Lemma 2.1 imply

$$(36) \quad \varphi(\ell_{k_i-1}; V_{k_i-1}) = \varphi(\ell_{k_i-1}) \leq \varphi(\theta_*) \leq \max_{\theta \in [0, 2\pi]} \varphi(\theta; V_{k_i-1}) = \varphi(\theta_{k_i}; V_{k_i-1})$$

and

$$\begin{aligned} & |\varphi(\ell_{k_i-1}; V_{k_i-1}) - \varphi(\theta_{k_i}; V_{k_i-1})| \\ &= |\lambda_{\min}(V_{k_i-1}^H H(\ell_{k_i-1}) V_{k_i-1}) - \lambda_{\min}(V_{k_i-1}^H H(\theta_{k_i}) V_{k_i-1})| \\ &\leq \|V_{k_i-1}^H (H(\ell_{k_i-1}) - H(\theta_{k_i})) V_{k_i-1}\|_2. \end{aligned}$$

The right-hand side converges to zero due to the continuity of  $H(\theta)$ , since both subsequences  $(\theta_{k_i})_{i=0}^\infty$  and  $(\ell_{k_i-1})_{i=1}^\infty$  converge to the same limiting point  $\ell$ . Therefore, the inequalities in (36) become equalities as  $k_i \rightarrow \infty$ , and it holds that

$$\lim_{k_i \rightarrow \infty} \varphi(\ell_{k_i-1}; V_{k_i-1}) = \lim_{k_i \rightarrow \infty} \varphi(\ell_{k_i-1}) = \varphi(\ell) = \varphi(\theta_*),$$

using the continuity of  $\varphi(\theta)$ . This shows that any converging subsequence of  $(\theta_k)_{k=0}^\infty$  converges to  $\theta_*$ , and hence  $\theta_k \rightarrow \theta_*$ .

It remains to discuss the local convergence order. According to the updating strategy in lines 6 and 8,  $\text{span}(V_{k+1})$  contains the *two* eigenvectors  $v(\theta_{k+1})$  and  $v(\theta_k)$ . Therefore,  $\theta_{k+2}$  can be viewed as being produced by the subspace method, Algorithm 1, with a subspace containing  $v(\theta_k)$  and  $v(\theta_{k+1})$ . Note that we may assume  $\theta_k \neq \theta_{k+1}$  because otherwise (36) implies  $\theta_k = \theta_{k+1} = \theta_*$  and Algorithm 2 returns the exact solution. Due to the local monotonicity behavior of the subspace method (see (6)), it holds that  $|\theta_{k+2} - \theta_*| \leq |\theta_{k+1} - \theta_*|$  provided that  $\theta_k$  and  $\theta_{k+1}$  are suffi-

ciently close to  $\theta_*$ . Moreover, using the local error recurrence (7), we obtain

$$e_{k+2} \leq 2C \cdot e_{k+1}e_k^2.$$

Then, following the proof of Theorem 2.4, we obtain a convergence order of  $\sigma = 2$ , which is the root of  $\sigma - \sigma^2 - 2 = 0$  that has magnitude larger than 1.

If the sequence is locally alternating, then the subspace updating scheme in line 6 yields  $\{\ell_{k+1}, m_{k+1}, u_{k+1}\} = \{\theta_{k+1}, \theta_k, \theta_{k-1}\}$  for  $k > p$ . Therefore, the subspace contains the *three* vectors  $v(\theta_{k-1})$ ,  $v(\theta_k)$  and  $v(\theta_{k+1})$ . For the reason mentioned above, we may assume that  $\theta_{k-1}$ ,  $\theta_k$ ,  $\theta_{k+1}$  are mutually distinct. In turn, the error recurrence (6) improves to

$$e_{k+2} \leq C \cdot e_{k+1}e_k^2e_{k-1}^2.$$

This leads to a convergence order of about  $\sigma = 2.26953$ , determined by the root of the cubic equation  $\sigma^3 - \sigma^2 - 2\sigma - 2 = 0$  with magnitude larger than 1.  $\square$

Although the 3-vector subspace method enjoys global convergence, the function values  $\varphi_k$  produced by Algorithm 2 are not guaranteed to monotonically decrease because the subspace monotonicity  $\text{span}(V_{k+1}) \subset \text{span}(V_k)$  is lost.

Numerically, we usually observed that the alternating property, required for faster local convergence in Theorem 3.5, holds once  $\theta_k$  is sufficiently close to  $\theta_*$ . In the following, we provide some theoretical evidence, by showing that this property holds in the first step. For this purpose, we first establish the following technical result.

**LEMMA 3.6.** *Let  $\gamma(A) > 0$  and let  $\theta_*$  be the maximizer of (21). Denote  $v(\theta)$  as in (11). Assume that  $\lambda_{\min}(H(\theta_*))$  is a simple eigenvalue. Then the following statements hold:*

- (a)  $\varphi''(\theta_*) \leq -\gamma(A)$ .
- (b)  $\varphi''(\theta_*) = -\gamma(A)$  if and only if  $v(\theta_*)$  is an eigenvector of  $H'(\theta_*)$ .
- (c) If  $v(\theta_*)$  is an eigenvector of  $H'(\theta_*)$ , then there is  $r > 0$  such that  $v(\theta) = v(\theta_*)$  for all  $|\theta - \theta_*| \leq r$ .

*Proof.* Consider the spectral decomposition  $H(\theta_*) = \sum_{i=1}^n \lambda_i v_i v_i^H$  with  $\lambda_1 = \lambda_{\min}(H(\theta_*))$  and  $v_1 = v(\theta_*)$ . Using existing results on the second derivative of eigenvalues [16] and  $v_1^H H''(\theta_*) v_1 = -v_1^H H(\theta_*) v_1$ , we obtain

$$(37) \quad \varphi''(\theta_*) = v_1^H H''(\theta_*) v_1 + 2 \sum_{i=2}^n \frac{|v_i^H H'(\theta_*) v_1|^2}{\lambda_1 - \lambda_i} = -\gamma(A) - 2 \sum_{i=2}^n \frac{|v_i^H H'(\theta_*) v_1|^2}{\lambda_i - \lambda_1}.$$

This proves (a).

Concerning part (b), (37) shows that  $\varphi''(\theta_*) = -\gamma(A)$  if and only if  $|v_i^H H'(\theta_*) v_1|^2 = 0$  for  $i = 2, \dots, n$ . Since  $\{v_1, \dots, v_n\}$  is an orthonormal basis, the later condition is equivalent to  $H'(\theta_*) v_1 \in \text{span}\{v_1\}$ , that is,  $v_1$  is an eigenvector of  $H'(\theta_*)$ .

To prove part (c), we assume  $\theta_* = 0$  without loss of generality (by replacing  $A$  with  $e^{-i\theta_*} A$ ). Because of

$$H(\theta) = S \cos \theta + K \sin \theta = \sin(\theta)H'(0) + \cos(\theta)H(0),$$

and  $H'(0)v(0) = \alpha v(0)$  for some (fixed)  $\alpha \in \mathbb{R}$ , we have  $H(\theta)v(0) = \tilde{\lambda}v(0)$  with  $\tilde{\lambda} = \sin(\theta)\alpha + \cos(\theta)\lambda_1$ . Hence  $v(0)$  is an eigenvector of  $H(\theta)$  with eigenvalue  $\tilde{\lambda}$ . For sufficiently small  $|\theta|$ , this remains the smallest eigenvalue of  $H(\theta)$  because of the simplicity assumption.  $\square$



The alternating property for the first iterate now follows.

LEMMA 3.7. *Let  $\gamma(A) > 0$  and assume that  $\lambda_{\min}(H(\theta_*))$  is a simple eigenvalue. Then the first iterate of Algorithm 2 satisfies*

$$(38) \quad (\theta_1 - \theta_*)(\theta_0 - \theta_*) \leq 0$$

for  $\theta_0$  sufficiently close to  $\theta_*$ .

*Proof.* Without loss of generality, assume that  $\theta_0 < \theta_* = 0$ . The first iterate  $\theta_1$  is determined as the maximizer of  $\varphi(\theta; v(\theta_0))$ .

If  $v(\theta_*)$  is an eigenvector of  $H'(\theta_*)$ , then Lemma 3.6(c) implies  $v(\theta_0) = v(\theta_*)$  (provided that  $\theta_0$  is sufficiently close to  $\theta_*$ ). Thus,  $\theta_1$  is already the exact maximizer and (38) is trivially satisfied. Now, suppose that  $v(\theta_*)$  is not an eigenvector of  $H'(\theta_*)$  and, hence,  $\gamma(A) + \varphi''(\theta_*) < 0$  holds by Lemma 3.6(a) and (b). Recalling from (31) that  $\varphi(\theta; v(\theta_0)) = \alpha \cos \theta + \beta \sin \theta$  with  $\alpha = v(\theta_0)^H S v(\theta_0)$  and  $\beta = v(\theta_0)^H K v(\theta_0)$ , we obtain from  $0 = \varphi'(\theta_1; v(\theta_0))$  that  $\tan \theta_1 = \beta/\alpha$ . Using the expansion (32) for  $\beta/\alpha$ , it thus follows that

$$\tan \theta_1 = \theta_0 + \frac{\varphi''(0)}{\varphi'(0)}\theta_0 + O(\theta_0^2) = \frac{\gamma(A) + \varphi''(0)}{\gamma(A)}\theta_0 + O(\theta_0^2).$$

This implies  $\theta_1 > 0$  for  $\theta_0 < 0$  sufficiently close to 0 and shows (38). □

**4. Numerical experiments.** All algorithms discussed in this paper have been implemented in MATLAB.<sup>2</sup> The primary purpose of the first experiment is to confirm our theoretical results and to demonstrate that the derived local convergence orders seem to be tight. We also compare with another algorithm [26] for Crawford number computation. In the end, we apply our subspace algorithms to the discretized boundary integral operators from [2] and show the potential of the algorithms for large-scale problems.

In finite precision, the user-supplied tolerance  $\text{tol}$  in Algorithms 1 and 2 should not be taken too strict. Provided that eigenvalues are computed in a backward stable way, it is advisable to terminate when  $\varphi_{k+1} - \lambda_k < \|A\|_2 \varepsilon_m$  with  $\varepsilon_m$  the machine precision.

*Example 1.* Our convergence analysis implies for  $k$  sufficiently large that

$$|\varphi_k - \varphi_*| \approx c |\varphi_{k-1} - \varphi_*|^\sigma$$

with  $\sigma \approx 2.4142$  for the (full) subspace method (Algorithm 1) and  $\sigma \approx 2.2695$  for the 3-vector variant (Algorithm 2). In turn,

$$\ln e_k \approx \sigma \ln e_{k-1} + \ln c \quad \text{with} \quad e_k = |\varphi_k - \varphi_*|.$$

In the following, we numerically verify the value of the slope  $\sigma$ . Since it is rather difficult to test high convergence orders in double precision, we perform the computation with 620 decimal digits using the Advanpix MP toolbox.<sup>3</sup> We substitute the “exact” value for  $\varphi_*$  by the result of a computation with 700 decimal digits.

*Scenario 1:*  $\lambda_{\min}(H(\theta_*))$  is a simple eigenvalue. Figure 3 plots  $\ln e_k$  vs.  $\ln e_{k-1}$  for Algorithms 1 and 2 applied to the  $120 \times 120$  matrix  $A = F + \iota M - cI_n$ , where  $F$  is the Fiedler matrix,  $M$  is the Moler matrix, and  $c = 4000 - 4000i$ ; see also [26]. The numerical range of  $A$  is shown in the left plot of Figure 3. Figure 4 shows the

<sup>2</sup>This is available at [www.unige.ch/~dlu](http://www.unige.ch/~dlu).

<sup>3</sup>This is available from [www.advanpix.com](http://www.advanpix.com).

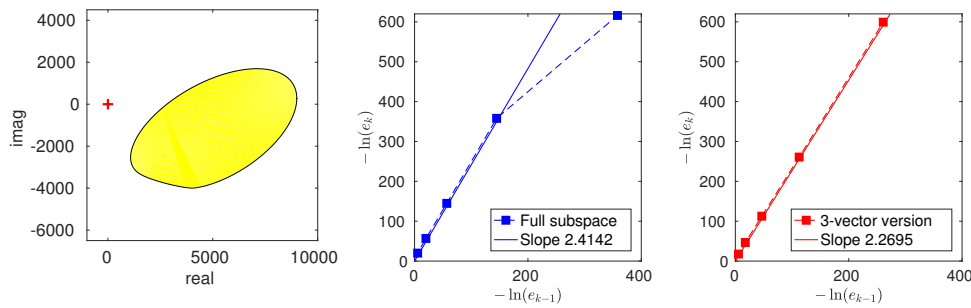


FIG. 3. Numerical range (left) and observed order of convergence of Algorithm 1 (middle) and Algorithm 2 (right) for a shifted sum of the Fiedler and Moler matrices.

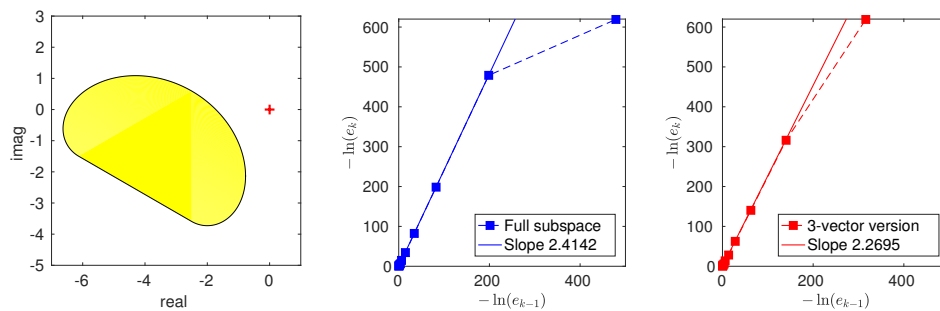


FIG. 4. Numerical range (left) and observed order of convergence of Algorithm 1 (middle) and Algorithm 2 (right) for a rotated Grcar matrix.

results obtained for  $A = G \cdot \exp(\pi i/3) - (4 + 2i) \cdot I_n$ , where  $G$  is the Grcar matrix of size  $n = 120$ .

For both matrices, we observe that the derived convergence orders from Theorems 3.4 and 3.5 seem to be tight.

*Scenario 2:  $\lambda_{\min}(H(\theta_*))$  is a multiple eigenvalue.* Figure 5 shows the convergence of the subspace methods for the tridiagonal matrix

$$(39) \quad A = \text{tridiag} \begin{pmatrix} \iota & \iota & \dots & \iota \\ 1 & 1 & a_3 & \dots & a_n \\ \iota & \iota & \dots & \iota \end{pmatrix} + 0.5\iota \cdot I_n \quad \text{with} \quad a_j = 2 + \frac{j}{n}$$

with  $n = 120$ . Visually, the point of the numerical range closest to zero is on the real line. Therefore the eigenvalue relevant for the convergence analysis is the smallest eigenvalue 1 of  $(A + A^H)/2$ . As this eigenvalue has multiplicity 2, our local convergence results do not apply. Indeed, the convergence plots reveal that the convergence orders established in Theorems 3.4 and 3.5 are not attained. On the other hand, one still obtains fast local convergence, seemingly quadratic convergence for the full subspace method, and superlinear convergence for the 3-vector variant. Further experiments show that the observed order of convergence remains the same when setting  $a_3 = \dots = a_d = 1$ , despite the increase of the eigenvalue multiplicity to  $d$ . We have also constructed matrices  $A$  with  $\lambda_{\min}(H(\theta_*))$  a multiple eigenvalue for which the convergence orders were less than the one in Figure 5. They always seem to be at least linear.

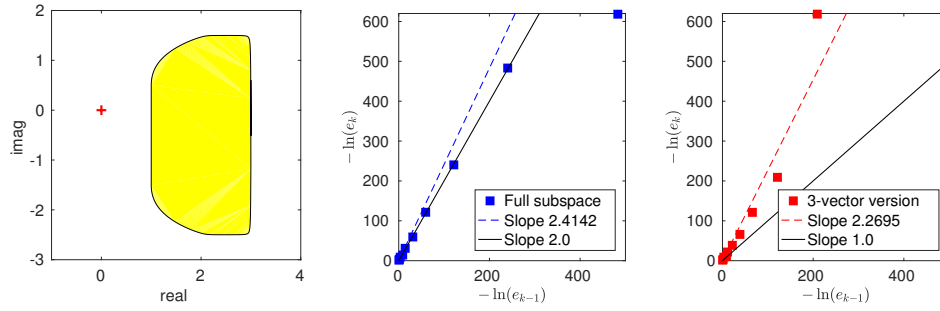


FIG. 5. Numerical range (left) and observed order of convergence of Algorithm 1 (middle) and Algorithm 2 (right) for the tridiagonal matrix (39).

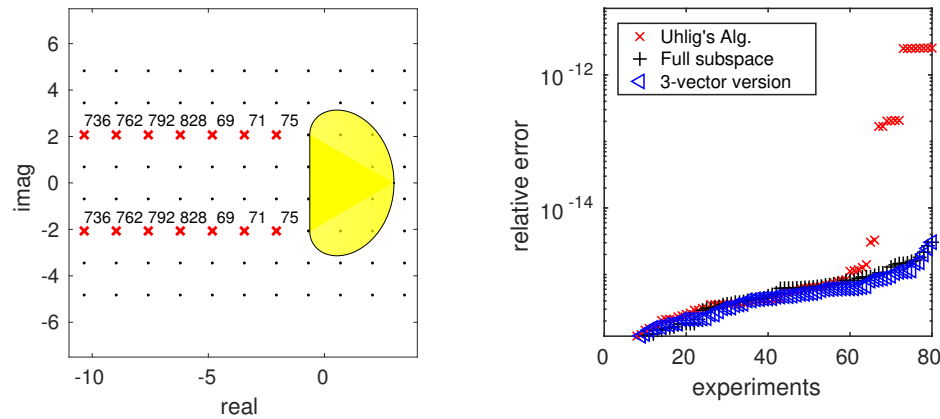


FIG. 6. Left: Black dots and red crosses denote shifts used in the experiment. Red crosses denote shifts for which Uhlig’s algorithm takes more than 60 iterations to converge, with the number of required iterations marked beside. Right: Relative error of the computed Crawford number, sorted in increasing order.

*Example 2.* In this example, we compare the performance of the subspace methods with Uhlig’s algorithm<sup>4</sup> [26]. As a test matrix  $A$ , we use the Gcar matrix of size  $n = 120$  mentioned above, but this time the matrix is not rotated. Although the smallest eigenvalue is not multiple, the small eigenvalues of the symmetric part of  $A$  form a cluster. This corresponds to a visually flat vertical portion of the boundary of the numerical range; see Figure 6. In this flat portion, the angle has little impact on the distance which could make it difficult to determine the optimal angle. In our experiment, we compute the Crawford number of the shifted matrices

$$A - (x_0 + y_0i)I_n$$

with  $x_0 + y_0i$  marked as black dots in Figure 6. One can equivalently view these shifts as zero and compute its distance to the numerical range.

In all experiments, we used the tolerance  $\text{tol} = 10^{-13}$  to terminate the algorithms. The iteration numbers shown in Figure 6 show that Uhlig’s algorithm sometimes faces

<sup>4</sup>The source code is available from <http://www.auburn.edu/~uhligfd>.

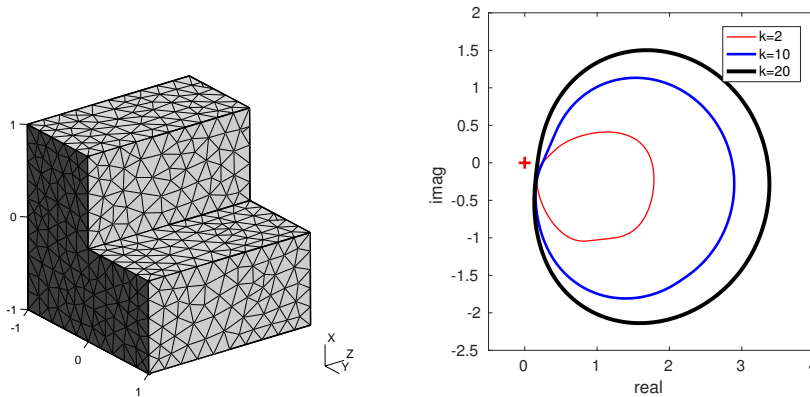


FIG. 7. *Left: L-shaped obstacle discretized with mesh size  $h = 0.2$ . Right: Numerical range of the discretized boundary integral operator  $A_k^h$  for wave numbers  $k = 2, 10, 20$  and  $h = 0.2$ .*

convergence difficulties. In the worst case it requires more than 800 iterations—and so also 800 eigenvalue computations of the matrix  $H(\theta)$ —until convergence. In contrast, both Algorithms 1 and 2 converge with a maximum number of 9 iterations and an average of 5.5. The error plot in Figure 6 reveals that the subspace methods are as accurate as Uhlig’s algorithm.<sup>5</sup>

*Example 3.* We now consider the Crawford number computation problem arising from the study of coercivity constants of boundary integral operators in acoustic scattering [2]. The integral operators of interest are defined by

$$(40) \quad a_k(u, v) = \int_{\Gamma} B_k u(y) \cdot \overline{v(y)} \, ds(y) \quad \text{with} \quad B_k = I + K_k - ikS_k,$$

where  $k > 0$  is the wave number,  $\Gamma$  is the boundary of a sound-soft bounded obstacle in  $\mathbb{R}^3$ ,  $u(x), v(x) \in L^2(\Gamma)$ ,  $I$  is the identity operator, and  $K_k$  and  $S_k$  are defined by

$$K_k u(x) = 2 \int_{\Gamma} \frac{\partial \Phi(x, y)}{\partial n(x)} u(y) \, ds(y), \quad S_k u(x) = 2 \int_{\Gamma} \Phi(x, y) u(y) \, ds(y), \quad x \in \Gamma.$$

Here,  $\Phi(x, y) = e^{ik|x-y|}/(2\pi|x-y|)$  for  $x, y \in \mathbb{R}^3$ ,  $x \neq y$ , and  $n(x)$  is the outpoint unit normal at  $\Gamma$ . To estimate the coercivity constant  $\gamma$  of the boundary integral operator (40) (i.e., the largest  $\gamma$  with  $\gamma\|u\|^2 \leq |a_k(u, u)|$  for all  $u \in L^2(\Gamma)$ ), we discretize it to a matrix  $A_k^h$  and compute its Crawford number. In order to achieve a good estimation, it is preferable to use a small mesh size  $h$  in discretization. This in general leads to a large dense matrix  $A_k^h$  for which the eigenvalue computation of  $H(\theta)$  is expensive. So it is crucial to compute the Crawford number using as few eigenvalue evaluations as possible.

In our experiment, we consider the three-dimensional time-harmonic acoustic scattering from an L-shaped sound-soft bounded obstacle  $\Gamma \subset \mathbb{R}^3$  displayed in Figure 7. We discretize the boundary integral operators (40) by the Galerkin boundary element library BEM++,<sup>6</sup> with triangular mesh (generated by gmsh<sup>7</sup> with Delaunay’s

<sup>5</sup>For the relative error computation, we compute the “accurate” solution by bisection methods with tolerance equal to the machine precision.

<sup>6</sup>The software is available from <http://www.bempp.org/>.

<sup>7</sup>The software is available from <http://gmsh.info/>.

TABLE 1  
 Computation results for Example 3 with mesh size  $h = 0.02$ .

wave number $k$	memory $A_k^h$		Crawford number $\gamma(A_k^h)$	its.	timing (h)
2	46 GB	Alg. 1	1.556145884392413e-01	5	2.5
		Alg. 2	1.556145884392399e-01	5	2.5
		Uhlig's	1.556145884392416e-01	11	4.4
10	54 GB	Alg. 1	1.880394259192281e-01	7	6.1
		Alg. 2	1.880394259192268e-01	7	6.1
		Uhlig's	1.880394259192323e-01	30	29.9
20	60 GB	Alg. 1	1.777716873410842e-01	8	8.1
		Alg. 2	1.777716873410808e-01	10	10.8
		Uhlig's	1.777716873410810e-01	24	24.9

algorithm) and piecewise linear basis functions. For a detailed description of the Galerkin discretization of the boundary integral operator, we refer to [23]. In order to speed up the assembly of the boundary integral operator, BEM++ stores the coefficient matrix  $A_k^h$  as a hierarchical matrix [9]. Since matrix vector and transpose matrix vector products can be efficiently evaluated for hierarchical matrices, we compute the smallest eigenvalue of  $H(\theta)$  using the MATLAB function `eigs` [17].

We summarize in Table 1 the experimental results for three wave numbers  $k = 2, 10, 20$ . The boundary integral operators were discretized with mesh size  $h = 0.02$  to matrices  $A_k^h$  of size  $n = 396\,162$ . From the timing statistics, it is clear that the subspace acceleration in Algorithms 1 and 2 has significantly reduced the computational cost. As in the previous experiment, all three algorithms terminate with tolerances  $\text{tol} = 10^{-13}$ .

**5. Conclusions.** We have analyzed the convergence of the subspace method for univariate eigenvalue optimization. The obtained convergence order not only improves upon existing results but it also appears to be tight. For the special case of the Crawford number, we have established novel properties of the objective function. This has resulted in a three-dimensional subspace method that is proven to enjoy favorable convergence properties and seems to work very well for the examples considered in this paper.

The developments in this paper offer several possibilities for future research. In particular, this concerns the proper treatment of large-scale problems and computing variations of the Crawford number, such as the inner numerical radius [4]. In addition, the subspace acceleration and its bracketing property may also prove useful for the shrinking problem in [10].

**Acknowledgments.** The authors thank the referees for their valuable suggestions. In addition, we appreciated the valuable feedback and suggestions of Emre Mengi on eigenvalue optimization.

#### REFERENCES

- [1] N. ALIYEV, P. BENNER, E. MENGI, AND M. VOIGT, *Large-Scale Computation of  $\mathcal{H}_\infty$ -Infinity Norms by a Greedy Subspace Method*, Technical report, Max Planck Institute for Dynamics of Complex Technical Systems, 2016.
- [2] T. BETCKE AND E. A. SPENCE, *Numerical estimation of coercivity constants for boundary integral operators in acoustic scattering*, SIAM J. Numer. Anal., 49 (2011), pp. 1572–1601.

- [3] S. BORA AND R. SRIVASTAVA, *Distance problems for Hermitian matrix pencils with eigenvalues of definite type*, SIAM J. Matrix Anal. Appl., 37 (2016), pp. 53–85.
- [4] S. H. CHENG AND N. J. HIGHAM, *The nearest definite pair for the Hermitian generalized eigenvalue problem*, Linear Algebra Appl., 302 (1999), pp. 63–76.
- [5] C. R. CRAWFORD, *A stable generalized eigenvalue problem*, SIAM J. Numer. Anal., 13 (1976), pp. 854–860.
- [6] C. DE BOOR, *A Practical Guide to Splines*, Vol. 27, Springer, New York, 1978.
- [7] C.-H. GUO, N. J. HIGHAM, AND F. TISSEUR, *Detecting and solving hyperbolic quadratic eigenvalue problems*, SIAM J. Matrix Anal. Appl., 30 (2009), pp. 1593–1613.
- [8] C.-H. GUO, N. J. HIGHAM, AND F. TISSEUR, *An improved arc algorithm for detecting definite hermitian pairs*, SIAM J. Matrix Anal. Appl., 31 (2009), pp. 1131–1151.
- [9] W. HACKBUSCH, *Hierarchical Matrices: Algorithms and Analysis*, Springer, New York, 2015.
- [10] N. J. HIGHAM, N. STRABIC, AND V. SEGO, *Restoring definiteness via shrinking, with an application to correlation matrices with a fixed block*, SIAM Rev., 58 (2016), pp. 245–263.
- [11] N. J. HIGHAM, F. TISSEUR, AND P. M. VAN DOOREN, *Detecting a definite Hermitian pair and a hyperbolic or elliptic quadratic eigenvalue problem, and associated nearness problems*, Linear Algebra Appl., 351 (2002), pp. 455–474.
- [12] W. KAHAN, *Spectra of nearly Hermitian matrices*, Proc. Amer. Math. Soc., 48 (1975), pp. 11–17.
- [13] F. KANGAL, K. MEERBERGEN, E. MENGI, AND W. MICHIELS, *A Subspace Method for Large Scale Eigenvalue Optimization*, preprint, arXiv:1508.04214, 2015.
- [14] T. KATO, *Perturbation Theory for Linear Operators*, Vol. 132, Springer, New York, 1995.
- [15] D. KRESSNER AND B. VANDEREYCKEN, *Subspace methods for computing the pseudospectral abscissa and the stability radius*, SIAM J. Matrix Anal. Appl., 35 (2014), pp. 292–313.
- [16] P. LANCASTER, *On eigenvalues of matrices dependent on a parameter*, Numer. Math., 6 (1964), pp. 377–387.
- [17] R. B. LEHOUCQ, D. C. SORENSEN, AND C. YANG, *ARPACK Users' Guide: Solution of Large-Scale Eigenvalue Problems with Implicitly Restarted Arnoldi Methods*, SIAM, Philadelphia, 1997.
- [18] D. LU AND B. VANDEREYCKEN, *Criss-cross type algorithms for computing the real pseudospectral abscissa*, SIAM J. Matrix Anal. Appl., 38 (2017), pp. 891–923.
- [19] K. MEERBERGEN, W. MICHIELS, R. VAN BEEUMEN, AND E. MENGI, *Computation of pseudospectral abscissa for large-scale nonlinear eigenvalue problems*, IMA J. Numer. Anal., 37 (2017), pp. 1831–1863.
- [20] E. MENGI AND M. L. OVERTON, *Algorithms for the computation of the pseudospectral radius and the numerical radius of a matrix*, IMA J. Numer. Anal., 25 (2005), pp. 648–669.
- [21] W. RUDIN, *Real and Complex Analysis*, Tata McGraw-Hill Education, New York, 1987.
- [22] P. SIRKOVIĆ AND D. KRESSNER, *Subspace acceleration for large-scale parameter-dependent Hermitian eigenproblems*, SIAM J. Matrix Anal. Appl., 37 (2016), pp. 695–718.
- [23] W. ŚMIGAJ, T. BETCKE, S. ARRIDGE, J. PHILLIPS, AND M. SCHWEIGER, *Solving boundary integral problems with BEM++*, ACM Trans. Math. Software, 41 (2015), p. 6.
- [24] G. W. STEWART, *Perturbation bounds for the definite generalized eigenvalue problem*, Linear Algebra Appl., 23 (1979), pp. 69–85.
- [25] G. W. STEWART AND J. SUN, *Matrix Perturbation Theory*, Academic Press, New York, 1990.
- [26] F. UHLIG, *On computing the generalized Crawford number of a matrix*, Linear Algebra Appl., 438 (2013), pp. 1923–1935.