

Research Article **A Subspace Iteration Algorithm for Fredholm Valued Functions**

Christian Engström¹ and Luka Grubišić²

¹Department of Mathematics and Mathematical Statistics, Umeå University, MIT-Huset, 90187 Umeå, Sweden ²Department of Mathematics, University of Zagreb, Bijenička 30, 10000 Zagreb, Croatia

Correspondence should be addressed to Luka Grubišić; luka.grubisic@math.hr

Received 2 July 2015; Revised 2 October 2015; Accepted 5 October 2015

Academic Editor: Ruben Specogna

Copyright © 2015 C. Engström and L. Grubišić. This is an open access article distributed under the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

We present an algorithm for approximating an eigensubspace of a spectral component of an analytic Fredholm valued function. Our approach is based on numerical contour integration and the analytic Fredholm theorem. The presented method can be seen as a variant of the FEAST algorithm for infinite dimensional nonlinear eigenvalue problems. Numerical experiments illustrate the performance of the algorithm for polynomial and rational eigenvalue problems.

1. Introduction

In this paper we study analytic operator eigenvalue problems defined on an open connected subset $\Omega \subseteq \mathbb{C}$ in a separable Hilbert space \mathscr{H} . Throughout this paper we assume that S : $\Omega \rightarrow \mathscr{L}(\mathscr{H})$ is an analytic function with values in the space $\mathscr{L}(\mathscr{H})$ of bounded linear operators. A scalar $\lambda \in \Omega$ is called an eigenvalue of S if $S(\lambda)$ is not injective. Hence, the eigenvalue problem is to find $\lambda \in \Omega$ and $u \in \mathscr{H} \setminus \{0\}$ such that

 $S(\lambda) u = 0. \tag{1}$

Such problems are, for example, used to study the dispersion and damping properties of waves [1–3]. Given a closed contour Γ , we would like to approximate all the eigenvalues of *S* inside Γ to the sufficient degree of accuracy. In this paper the numerical method is based on contour integrals of the generalized resolvent S^{-1} . The state-of-the-art results in the contour integration based methods for solving nonlinear matrix eigenvalue problems are presented in [4–6] including the references therein. Results for contour integration based solution methods for Fredholm valued eigenvalue problems can be found in, for example, [7, 8].

The spectrum $\sigma(S)$ of the operator function *S* is the set of all $\lambda \in \Omega$ such that $S(\lambda)$ is not invertible in $\mathscr{L}(\mathscr{H})$ and the resolvent set is defined as the complement $\rho(S) = \Omega \setminus \sigma(S)$. For $F \in \mathscr{L}(\mathscr{H})$ we define the operator norm by the expression $\|F\| = \sqrt{\operatorname{spr}(F^*F)}$, where $\operatorname{spr}(\cdot)$ denotes the spectral radius. We call an operator $F \in \mathscr{L}(\mathscr{H})$ a Fredholm operator if the dimensions of its null space $\operatorname{Ker}(F)$ and of the orthogonal complement of its range $\operatorname{CoKer}(F) = \operatorname{Ran}(F)^{\perp}$ are finite. By $\Phi(\mathscr{H})$ we denote the set of all Fredholm operators on \mathscr{H} and the number $\operatorname{ind}(F) = \dim \operatorname{Ker}(F) - \dim \operatorname{CoKer}(F)$ is called the index of $F \in \Phi(\mathscr{H})$. In what follows we will assume that $S(\lambda) \in \Phi(\mathscr{H})$ for all $\lambda \in \Omega$. If in addition the resolvent set of such *S* is nonempty, the analytic Fredholm theorem, for example, [9, Theorem 1.3.1], implies that the generalized resolvent $z \mapsto S(z)^{-1}$ is finitely meromorphic. This in turn implies that the spectrum $\sigma(S)$ is countable and the geometric multiplicity of λ , that is, dim(Ker $S(\lambda)$), is finite. Moreover, the associated Jordan chains of generalized eigenvectors have finite length bounded by the algebraic multiplicity; see [9].

The results of this paper are a combination of matrix techniques from [6] with the specialization of the results from [7] to Hilbert spaces. In particular, we leverage the technique of block operator matrix representation of Fredholm valued operator functions and prove that the convergence rate for the inexact subspace iteration algorithm depends primarily on the spectral properties of the operator function. Also, we make the case for the problem dependent determination of the number of integration nodes depending on the clustering of eigenvalues towards the contour of integration (cf. [6]), where the use of 16 nodes of Gauss-Legendre integration formula is recommended. To assess the quality of a computed eigenpair, rank one perturbations of the operator function are studied. In particular, we construct a perturbation based on the residual functional and estimate the approximation errors by estimating the norm of the residual by an auxiliary subspace technique. Our algorithm consists of the inexact subspace iteration for the zeroth moment of the resolvent to construct the approximate eigenspace for the eigenvalues contained inside a contour Γ . We then use the moment method of Beyn et al. [4, 7] to extract information on eigenvalues from the computed approximate eigenspace. As the convergence criterion we use a hierarchical residual estimate. In the case in which the convergence criterion is not satisfied the procedure is repeated. This structure of the algorithm will be replicated directly in the structure of the paper and is presented as Algorithm 2.

The paper is organized as follows. In Section 2 we establish a criterion based on the residual norm for assessing approximations of simple eigenvalues. In Section 3 we present inexact subspace iteration algorithm based on contour integration and prove that its convergence rate essentially depends on the properties of the operator function. It is shown that the influence of the integration formula diminishes exponentially with the number of integration nodes. Finally, in Section 4 we present numerical experiments.

2. Notation and Basic Analytic Results

In this section we present the machinery of quasi-matrices from [10-12]. In particular we present basic results on the angles between finite dimensional subspaces of a Hilbert space in terms of quasi-matrix notation. Finally, we will prove an error estimation result for simple eigenvalues of a Fredholm valued function.

A quasi-matrix is a bounded linear operator V from the finite dimensional space \mathbb{C}^r to an (in general) infinite dimensional Hilbert space \mathcal{H} . Then, the product V^*V denotes the Gram matrix:

$$(V^*V)_{ij} = (Ve_i, Ve_j), \quad i, j = 1, 2, \dots, r,$$
 (2)

which depends on the inner product (\cdot, \cdot) of \mathcal{H} . Let P_1, P_2, Q_1 , and Q_2 be orthogonal projections such that $P_1 \oplus P_2 = I$ and $Q_1 \oplus Q_2 = I$, where I is the identity operator on \mathcal{H} . Furthermore, let $\mathcal{H}_i := P_i \mathcal{H}, \ \widehat{\mathcal{H}}_i := Q_i \mathcal{H}, \ i = 1, 2$, and identify \mathcal{H} and the two spaces $\mathcal{H}_1 \oplus \mathcal{H}_2$ and $\widehat{\mathcal{H}}_1 \oplus \widehat{\mathcal{H}}_2$ by isomorphism. Let $B \in \mathcal{L}(\mathcal{H})$ and take

$$u = \begin{bmatrix} u_1 \\ u_2 \end{bmatrix} \in \mathcal{H}_1 \oplus \mathcal{H}_2 := \mathcal{H}.$$
(3)

Then, $Q_i B u = B_{i1} u_1 + B_{i2} u_2$, i = 1, 2, where $B_{ij} : \mathcal{H}_j \to \widehat{\mathcal{H}}_i$ is defined by restricting $B_{ij} := Q_i B P_j$ onto appropriate spaces. Hence, the bounded operator *B* has a block operator matrix representation $B : \mathcal{H}_1 \oplus \mathcal{H}_2 \to \widehat{\mathcal{H}}_1 \oplus \widehat{\mathcal{H}}_2$ in terms of B_{ij} and B u is computed following the rules of matrix algebra:

$$Bu = \begin{bmatrix} B_{11} & B_{12} \\ B_{21} & B_{22} \end{bmatrix} \begin{bmatrix} u_1 \\ u_2 \end{bmatrix} = \begin{bmatrix} B_{11}u_1 + B_{12}u_2 \\ B_{21}u_1 + B_{22}u_2 \end{bmatrix}.$$
 (4)

The multiplication of two block operator matrices also follows the rules of matrix algebra. Represent, for example, $A : \mathcal{H} \to \mathcal{H}$ by $A = \begin{bmatrix} A_1 & A_2 \end{bmatrix}$, where $\mathcal{H} = \mathcal{H}_1 \oplus \mathcal{H}_2$ and $A_1 : \mathcal{H}_1 \to \mathcal{H}$ and $A_2 : \mathcal{H}_2 \to \widetilde{\mathcal{H}}$. Then, the operator $AB : \mathcal{H} \to \mathcal{H}$ has the block representation:

$$AB = \begin{bmatrix} A_1 B_{11} + A_2 B_{21} & A_1 B_{12} + A_2 B_{22} \end{bmatrix}.$$
 (5)

Here we have exemplarily taken the trivial partition of unity $Q_1 = I$ and $Q_0 = 0$ and we assume that both $P_1 \neq 0$ and $P_2 \neq 0$. This illustrates the flexibility we have in choosing the block operator matrix representations for bounded operators. For more details on this construction see a recent book by *C*. Tretter [13]. To make the paper more readable we will use I_r to denote the identity operator on the finite dimensional space \mathbb{C}^r and *I* will denote the identity operator on \mathcal{H} .

Let $P_1 \oplus P_2 = I$ be given such that dim $\operatorname{Ran}(P_1) = r$. Let V be a unitary operator such that $V = \begin{bmatrix} V_1 & V_2 \end{bmatrix}$ and $\operatorname{Ran}(P_1) = \operatorname{Ran}(V_1)$ and $\operatorname{Ran}(P_2) = \operatorname{Ran}(V_2)$. Note that in this setting we have $P_1 = V_1V_1^*$. For the quasi-matrix $X : \mathbb{C}^r \to \mathcal{H}$ such that $X^*X = I_r$ we can write

$$X = V_1 C + V_2 S, (6)$$

where $C = V_1^* X$ and $S = V_2^* X$. With this notation we compute

$$I_r = X^* X = C^* C + S^* S,$$
(7)

and so $||S|| \le 1$ and $||C|| \le 1$. Furthermore, note that $Q = XX^*$ is the orthogonal projection onto Ran(X) and so from [14, Theorem 2] it follows that

$$\|P_1 - Q\| = \|V_1 V_1^* - XX^*\| = \left\| \begin{bmatrix} I_r - CC^* & -CS^* \\ -SC^* & -SS^* \end{bmatrix} \right\|$$

= $\|S\|$. (8)

The last identity has been established by spectral calculus using the fact that dimRan(P_1) = r; see [15]. Let ||S|| < 1; then we define the unique number $\theta \in [0, \pi/2)$ such that

$$\sin \theta = \|S\|. \tag{9}$$

We call θ the maximal canonical angle between the spaces $\operatorname{Ran}(X)$ and $\operatorname{Ran}(P_1) = \operatorname{Ran}(V_1)$. Since ||S|| < 1 (7) implies that C^*C is a positive definite matrix, therefore *C* must be invertible. By direct computation using spectral calculus, as in [15], we establish that

$$\cos \theta = \|C\|,$$

$$\tan \theta = \|SC^{-1}\|.$$
(10)

When we are considering several pairs of quasi-matrices V_1 and X we will write $\theta(V_1, X)$ to denote the maximal canonical angle between subspaces $\operatorname{Ran}(V_1)$ and $\operatorname{Ran}(X)$.

This definition can be extended to subspaces of different dimensions using pairs of orthogonal projections and their singular value decomposition [14, 15]. In this case we call the norm of the difference of the two projections the gap distance between subspaces; see [14, 15].

2.1. Application of Abstract Results to Operators Defined by Sesquilinear Forms. The abstract Fredholm analytic theorem is stated for bounded operators between Hilbert spaces. Since we are interested in finite element computations, our problems will always be stated in the variational form which assumes working with unbounded sectorial operators. See [16] for definitions and the terminology relating to unbounded sectorial operators and forms. Below we will formulate the Fredholm analytic theorem in this setting.

Let \mathbf{t}_0 be a densely defined closed symmetric form in \mathcal{H} which is semibounded from below with a strictly positive lower bound; see [16, Section VI.2]. We will call such quadratic forms positive definite forms and use $\text{Dom}(\mathbf{t}_0) \subset \mathcal{H}$ to denote its domain of definition. To simplify the notation we will write $\mathcal{V} := \text{Dom}(\mathbf{t}_0)$ in the rest of the paper. Additionally, let \mathbf{c}_i , i = 1, ..., d, be a sequence of sesquilinear (not necessarily symmetric) forms which are relatively compact [16] with regard to \mathbf{t}_0 . Moreover, assume that $f_i(\cdot)$, i =1, ..., d, is a sequence of scalar analytic functions in Ω . Define the family of sesquilinear forms:

$$\mathbf{t}(z) [\phi, \psi] := \mathbf{t}_0 [\phi, \psi] + \sum_{i=1}^d f_i(z) \, \mathbf{c}_i [\phi, \psi],$$

$$\phi, \psi \in \mathcal{V}, \ z \in \Omega.$$
(11)

In a variationally posed eigenvalue problem we are seeking a scalar $\lambda \in \Omega$ and a vector $\phi \in \mathcal{V} \setminus \{0\}$ such that

$$\mathbf{t}(\lambda)\left[\phi,\psi\right] = 0, \quad \psi \in \mathcal{V}. \tag{12}$$

To this variational formulation we construct a representation by an operator-valued function $T(\cdot)$ and then apply the analytic Fredholm theorem to establish the structure of the spectrum.

Let T_0 , Dom $(T_0^{1/2}) = \mathcal{V}$ be a self-adjoint positive definite operator which represents the form \mathbf{t}_0 in the sense of Kato's second representation theorem; see [16, Theorem VI.2.23]. Let K_i be defined by

$$(K_i u, v) = \mathfrak{c}_i \left[T_0^{-1/2} u, T_0^{-1/2} v \right], \quad u, v \in \mathcal{H}.$$
(13)

The operators K_i , i = 1, ..., d, are obviously compact and for $z \in \Omega \setminus \sigma(T)$ the value of the generalized resolvent is the operator:

$$T^{-1}(z) = T_0^{-1/2} \left(I + \sum_{j=1}^d f_j(z) K_j \right)^{-1} T_0^{-1/2}.$$
 (14)

Here T(z) is the unbounded sectorial operator with domain Dom(T(z)) such that

$$(T(z)\phi,\psi) = \mathbf{t}_{0}[\phi,\psi] + \sum_{i=1}^{d} f_{i}(z) \mathbf{c}_{i}[\phi,\psi],$$

$$\phi \in \text{Dom}(T(z)), \ \psi \in \mathcal{V}.$$
(15)

Obviously the operator-valued function

$$S(z) = I + \sum_{j=1}^{d} f_j(z) K_j, \quad z \in \Omega$$
(16)

satisfies the requirements of the analytic Fredholm theorem and $\sigma(S) = \sigma(T)$. Let us note that we will use

$$(T'(z)\phi,\psi) = \sum_{i=1}^{d} f'_{i}(z) \mathfrak{c}_{i} [\phi,\psi],$$

$$\phi \in \operatorname{Dom} (T'(z)), \ \psi \in \mathcal{V},$$

$$(17)$$

to define a derivative of an operator-valued analytic function. We can now define the notion of a semisimple eigenvalue.

Definition 1. Let T be as in (15) and let $\mu \in \sigma(T)$ be an eigenvalue. The eigenvalue μ is semisimple if for each $\psi \in \text{Ker}(T(\mu)) \setminus \{0\}$ there is a $\psi \in \text{Ker}(T^*(\mu))$ such that

$$\left(T'\left(\mu\right)\phi,\psi\right)\neq0.\tag{18}$$

If dim $\text{Ker}(T(\mu)) = 1$, then μ is called a simple eigenvalue.

Note that in the quasi-matrix notation we will freely write

$$\left(T'\left(\mu\right)\phi,\psi\right) = \psi^*T'\left(\mu\right)\phi. \tag{19}$$

To this end we identify the vectors with a mapping $\psi : \mathbb{C} \to \mathcal{H}$.

With these conventions we state, informally, the generalized argument principle proved by Gohberg and Sigal in [17– 19]. It states that for the closed contour $\Gamma \subset \rho(T)$ the number

$$M(T,\Gamma) = \operatorname{tr}\left(\int_{\Gamma} T(z)^{-1} T'(z) dz\right)$$

= $\operatorname{tr}\left(\int_{\Gamma} T'(z) T(z)^{-1} dz\right)$ (20)

satisfies $M(T, \Gamma) \in \mathbb{N}$ and it equals the total multiplicity of the eigenvalues enclosed by Γ . We also have the following consequence of [9, Theorem 1.3.1].

Proposition 2. Let us assume that we have a variational eigenvalue problem (12) with the operator representation $T : \Omega \to \mathscr{H}$ given by (15). Then the spectrum $\sigma(T)$ consists of a countable collection of eigenvalues with finite multiplicity. Further, let the component of $\sigma(T)$ inside a contour Γ consist only of semisimple eigenvalues λ_i , $i = 1, \ldots, r$, such that $n_{\lambda_i} = \dim \operatorname{Ker}(T(\lambda_i))$. Then there are quasi-matrices $X_i, Y_i : \mathbb{C}^{n_{\lambda_i}} \to \mathscr{H}$, $i = 1, \ldots, r$, such that $\operatorname{Ran}(X_i) \subset \operatorname{Ker}(T(\lambda_i))$, and $\operatorname{Ran}(Y_i) \subset \operatorname{Ker}(T^*(\lambda_i))$, $i = 1, \ldots, r$, and an open once connected neighborhood \mathscr{U} containing λ_i , $i = 1, \ldots, r$ and Γ and an operator-valued function H which is analytic on \mathscr{U} and taking values in $\mathscr{L}(\mathscr{H})$ such that

$$T^{-1}(z) = \sum_{i=1}^{r} \frac{1}{z - \lambda_i} X_i Y_i^* + H(z), \quad z \in \mathcal{U},$$
(21)

and $Y_i^*(T'(\lambda_i)X_i) = I_{n_{\lambda_i}}, i = 1, \ldots, r.$

Proof. For $z \in \Omega \setminus \sigma(T)$ write

$$T^{-1}(z) = T_0^{-1/2} \left(I + \sum_{j=1}^d f_j(z) K_j \right)^{-1} T_0^{-1/2}$$
(22)

and define the Fredholm valued function:

$$S(z) = I + \sum_{j=1}^{d} f_j(z) K_j, \quad z \in \Omega.$$
 (23)

Recall that $\text{Dom}(T_0^{1/2}) = \mathcal{V}$ and that $T_0^{1/2}$ maps \mathcal{V} one to one on \mathcal{H} . Now apply [9, Theorem 1.3.1] on *S*.

2.2. Estimating Eigenvalues inside a Contour. To count the semisimple eigenvalues inside a contour we will use the approach of [20]. We will limit our consideration on the case of semisimple eigenvalues and rank one perturbations. First we will present results for Fredholm valued operator functions and then formulate the result for operators defined by sesquilinear forms.

Lemma 3. Let $S : \Omega \to \Phi(\mathcal{H})$ be an analytic function and let E be a bounded operator such that dim Ran(E) = 1. Assume that $\Gamma \subset \rho(S)$ is a simple closed contour such that the component of $\sigma(S)$ inside Γ consists solely of a simple eigenvalue λ . If $S(z) + \tau E$ is invertible for all $z \in \Gamma$ and all $\tau \in [0, 1]$, then T + E has a simple eigenvalue $\tilde{\lambda}$ inside Γ and

$$\left|\lambda - \widetilde{\lambda}\right| \le C \left\|E\right\|,\tag{24}$$

where C essentially depends on $\max_{z \in \Gamma} ||S(z)^{-1}|_{\operatorname{Ran}(E)}||$ and $\max_{z \in \Gamma} ||S(z)^{-*}|_{\operatorname{Ran}(E^*)}||$ and the length of the the integration curve Γ .

Proof. Recall that (S(z) + E)' = S'(z) and define the function:

$$f(\tau) = \frac{1}{2\pi i} \operatorname{tr}\left(\int_{\Gamma} S'(z) \left(S(z) + \tau E\right)^{-1} dz\right).$$
(25)

By the generalized argument principle—see [17-19]—the value of $f(\tau)$ equals the total multiplicity of the eigenvalues inside Γ . In particular, by Proposition 2 we have that there are vectors x and y such that $(S'(\lambda)x, y) = y^*S'(\lambda)x = 1$ and

$$f(0) = \operatorname{tr}\left(S'(\lambda)xy^*\right) = \operatorname{tr}\left(y^*S'(\lambda)x\right) = 1, \quad (26)$$

where we used the circularity of the trace. By the assumptions of the theorem S(z) is invertible for all $z \in \Gamma$ and E is a rank one bounded operator. Therefore, there are vectors v and u so that $E = uv^*$ and using Sherman-Morrison formula (see [21–23] and the references therein), we write

$$(S(z) + \tau E)^{-1} = S(z)^{-1} - \frac{\tau}{1 + v^* S^{-1}(z) u} S(z)^{-1} (uv^*) S(z)^{-1},$$
(27)

and so it follows that

$$(S(z) + \tau_1 E)^{-1} - (S(z) + \tau_2 E)^{-1}$$

$$= \frac{\tau_2 - \tau_1}{1 + v^* S^{-1}(z) u} S(z)^{-1} (uv^*) S(z)^{-1}$$

$$= \frac{\tau_2 - \tau_1}{1 + v^* S^{-1}(z) u} (S(z)^{-1} u) (S(z)^{-*} v)^* .$$
(28)

And in particular f is a smooth function. Since f, due to Rouche's theorem, conclude that f takes values only in natural numbers, it must be constant for all $\tau \in [0, 1]$. Let us denote this eigenvalue of $\tilde{S} = S + E$ by $\tilde{\lambda}$. Define the operators

$$A^{(q)} = \int_{\Gamma} z^{q} S^{-1}(z) dz,$$

$$\widetilde{A}^{(q)} = \int_{\Gamma} z^{q} \widetilde{S}^{-1}(z) dz,$$

$$q = 0, 1.$$
(29)

Proposition 2 implies $A^{(1)} = \lambda x y^*$, $\widetilde{A}^{(1)} = \widetilde{\lambda} \widetilde{x} \widetilde{y}^*$, where $S(\lambda)x = 0$, $S(\lambda)^* y = 0$, $\widetilde{S}(\widetilde{\lambda})\widetilde{x} = 0$, and $\widetilde{S}(\lambda)^* \widetilde{y} = 0$ and so there exists a vector ϕ such that $\phi^* A^{(0)} \phi \neq 0$ and $\phi^* \widetilde{A}^{(0)} \phi \neq 0$. We now compute

$$\lambda = \frac{\phi^* \left(\int_{\Gamma} z S^{-1}(z) dz \right) \phi}{\phi^* \left(\int_{\Gamma} S^{-1}(z) dz \right) \phi},$$

$$\tilde{\lambda} = \frac{\phi^* \left(\int_{\Gamma} z \tilde{S}^{-1}(z) dz \right) \phi}{\phi^* \left(\int_{\Gamma} \tilde{S}^{-1}(z) dz \right) \phi}$$
(30)

and so the second assertion follows by the following computation:

$$\begin{split} \left|\lambda - \widetilde{\lambda}\right| &= \left|\frac{\phi^*\left(\int_{\Gamma} zS^{-1}\left(z\right)dz\right)\phi}{\phi^*\left(\int_{\Gamma} S^{-1}\left(z\right)dz\right)\phi} - \frac{\phi^*\left(\int_{\Gamma} z\widetilde{S}^{-1}\left(z\right)dz\right)\phi}{\phi^*\left(\int_{\Gamma} \widetilde{S}^{-1}\left(z\right)dz\right)\phi}\right| \\ &\leq \frac{1}{\left|\phi^*\left(\int_{\Gamma} S^{-1}\left(z\right)dz\right)\phi\right|} \left[\left|\phi^*\left(\int_{\Gamma} zS^{-1}\left(z\right)dz\right)\phi - \phi^*\left(\int_{\Gamma} z\widetilde{S}^{-1}\left(z\right)dz\right)\phi\right| \\ &+ \left|\frac{\phi^*\left(\int_{\Gamma} z\widetilde{S}^{-1}\left(z\right)dz\right)\phi}{\phi^*\left(\int_{\Gamma} \widetilde{S}^{-1}\left(z\right)dz\right)\phi}\right| \left|\phi^*\left(\int_{\Gamma} \widetilde{S}^{-1}\left(z\right)dz\right)\phi - \phi^*\left(\int_{\Gamma} S^{-1}\left(z\right)dz\right)\phi\right| \right] \\ &\leq \frac{1}{\left|\phi^*\left(\int_{\Gamma} S^{-1}\left(z\right)dz\right)\phi\right|} \left[\left|\phi^*\left(\int_{\Gamma} zS^{-1}\left(z\right) - z\widetilde{S}^{-1}\left(z\right)dz\right)\phi\right| + \left|\widetilde{\lambda}\right|\left|\phi^*\left(\int_{\Gamma} \widetilde{S}^{-1}\left(z\right) - S^{-1}\left(z\right)dz\right)\phi\right|\right] \end{split}$$

$$= \frac{1}{\left|\phi^{*}\left(\int_{\Gamma} S^{-1}(z) dz\right)\phi\right|} \left[\left|\phi^{*}\left(\int_{\Gamma} \frac{z}{1 + \nu^{*} S^{-1}(z) u} S(z)^{-1} ES(z)^{-1} dz\right)\phi\right| + \left|\tilde{\lambda}\right| \left|\phi^{*}\left(\int_{\Gamma} \frac{1}{1 + \nu^{*} S^{-1}(z) u} S(z)^{-1} ES(z)^{-1} dz\right)\phi\right| \right] \le C \|E\|.$$
(31)

The claim on the constant *C* follows from the observation that $\operatorname{Ran}(E) = \operatorname{span}\{u\}$ and $\operatorname{Ran}(E^*) = \operatorname{span}\{v\}$.

With this result we can now formulate the main result of this section. It will be used to assess the quality of a given approximation, regardless of its origin.

Proposition 4. Let $\mathbf{t}(\cdot)$ denote the family of sesquilinear forms (11) with operator representation $T(\cdot) : \Omega \to \mathcal{H}$ given by (15). Assume that a contour $\Gamma \subset \rho(T)$ encloses only a simple eigenvalue λ and $\mu \in \mathbb{C}$ but no other points of $\sigma(T)$. Let $u \in \mathcal{V}$, with $\|u\| = 1$. Then there is $\lambda \in \sigma(T)$ such that

$$\left|\lambda - \mu\right| \le C \sup_{\phi \in \mathscr{V} \setminus \{0\}} \frac{\left|\mathbf{t}\left(\mu\right) \left[u, \phi\right]\right|}{\sqrt{\mathbf{t}_{0}\left[\phi, \phi\right]}}.$$
(32)

The constant C does not depend on μ *and u but on contour* Γ *and the restriction of* $||T^{-1}(\cdot)||$ *and* $||T^{-*}(\cdot)||$ *to* Γ .

Proof. We will now construct a particular operator $E_{\mu,u}$ which will be used to assess the quality of the pair μ , u. Define the sesquilinear form $e_{\mu,u} : \mathcal{V} \times \mathcal{V} \to \mathbb{C}$ by the formula

$$\mathbf{e}_{\mu,u}\left[\psi,\phi\right] = -\mathbf{t}\left(\mu\right)\left[u,\phi\right]\left(\psi,u\right), \quad \psi,\phi\in\mathcal{V}. \tag{33}$$

It is obviously relatively compact with respect to \mathbf{t}_0 and so we can define the Fredholm valued operator function $\Omega \ni z \rightarrow \tilde{T}(z)$, where $\tilde{T}(z)$ is the operator defined by the form

$$\widetilde{\mathbf{t}}(\lambda) \left[\psi, \phi \right] = \mathbf{t}(\lambda) \left[\psi, \phi \right] + \mathbf{e}_{\mu, \mu} \left[\psi, \phi \right]. \tag{34}$$

Further,

$$\left(\widetilde{T}(\mu) u, \phi\right) = \mathbf{t}(\mu) [u, \phi] + \mathbf{e}_{\mu, u} [u, \phi]$$
$$= \mathbf{t}(\mu) [u, \phi] - \mathbf{t}(\mu) [u, \phi] (u, u)$$
$$= \mathbf{t}(\mu) [u, \phi] - \mathbf{t}(\mu) [u, \phi] = 0,$$
$$\phi \in \mathcal{V},$$
$$(35)$$

and so $\mu \in \sigma(\tilde{T})$. We construct the operator $E_{\mu,\mu}$ as the operator defined by

$$\left(E_{\mu,u}\phi,\psi\right) = \mathfrak{e}_{\mu,u}\left[T_0^{-1/2}\phi,T_0^{-1/2}\psi\right], \quad \phi,\psi\in\mathscr{H}.$$
 (36)

Now recall that (22) and (35) imply

$$\widetilde{T}^{-1}(z) = T_0^{-1/2} \left(I + \sum_{j=1}^d f_j(z) K_j + E_{\mu,\mu} \right)^{-1} T_0^{-1/2}.$$
 (37)

By construction the operator function $\widetilde{S} : \Omega \to \mathscr{L}(\mathscr{H})$

$$\widetilde{S}(z) = I + \sum_{j=1}^{d} f_{j}(z) K_{j} + E_{\mu,u} = S(z) + E_{\mu,u}$$
(38)

satisfies the assumption of Lemma 3. Finally, note that

$$\sup_{\phi \in \mathscr{V} \setminus \{0\}} \frac{\left| \mathbf{e}_{\mu,u} \left(\mu \right) \left[u, \phi \right] \right|}{\sqrt{\mathbf{t}_0 \left[\phi, \phi \right]}} = \sup_{\phi \in \mathscr{V} \setminus \{0\}} \frac{\left| \mathbf{t} \left(\mu \right) \left[u, \phi \right] \right|}{\sqrt{\mathbf{t}_0 \left[\phi, \phi \right]}}$$
(39)
$$=: \left\| \mathbf{t}_{\mu,u} \left(\mu \right) \left[u, \cdot \right] \right\|_{\mathbf{t}_0, -1},$$

and so from (36) it follows that

$$\left\| E_{\mu,u} \right\| = \left\| \mathbf{t}_{\mu,u} \left(\mu \right) [u, \cdot] \right\|_{\mathbf{t}_{0}, -1}.$$
(40)

Now recall the definition of *C* from Lemma 3 to conclude the proof. \Box

2.3. A Sketch for a Practical Algorithm for Error Estimation. Proposition 4 will be the basis for practical error estimation. Let $\mathcal{V} \supset \mathcal{Q}_h$, h > 0 be a family of finite dimensional spaces such that the orthogonal projections Q_h onto \mathcal{Q}_h converge strongly to *I* as $h \rightarrow \infty$. Let further $\mathcal{Q}_{h_1} \subset \mathcal{Q}_{h_2}$ for $h_1 \le h_2$. Assume that $\phi \in \mathcal{Q}_{h_1}$ and $\mu \in \mathbb{C}$ are given. For h_2 , $h_2 > h_1$ define

$$\left\|\mathbf{t}\left(\mu\right)\left[u,\cdot\right]\right\|_{\mathcal{Q}_{h_{2}},\mathbf{t}_{0},-1} = \sup_{\phi\in\mathcal{Q}_{h_{2}}\setminus\{0\}}\frac{\left|\mathbf{t}\left(\mu\right)\left[u,\phi\right]\right|}{\sqrt{\mathbf{t}_{0}\left[\phi,\phi\right]}}.$$
(41)

. . . .

Obviously it holds, recalling the definition from (39), that

$$\| \mathbf{t}(\mu) [u, \cdot] \|_{\mathcal{Q}_{h_2}, \mathbf{t}_0, -1} \le \| \mathbf{t}_{\mu, u}(\mu) [u, \cdot] \|_{\mathbf{t}_0, -1}.$$
(42)

However, in the case in which \mathcal{Q}_h , h > 0 satisfies the standard saturation property, for example, $||(I-Q_{h_2})\phi|| \le q||(I-Q_{h_1})\phi||$ for fixed q, 0 < q < 1 and some ϕ , we can also prove

$$\|\mathbf{t}(\mu)[u,\cdot]\|_{\hat{\mathcal{Q}}_{h_{2}},\mathbf{t}_{0},-1} \leq \|\mathbf{t}(\mu)[u,\cdot]\|_{\mathbf{t}_{0},-1}$$

$$\leq C \|\mathbf{t}(\mu)[u,\cdot]\|_{\hat{\mathcal{Q}}_{h_{2}},\mathbf{t}_{0},-1}.$$
(43)

The constant *C* essentially depends on *q* which in turn depends on $h_2 - h_1$ but not on the magnitude of h_1 ; for more details see [24, 25].

3. Contour Integration Based Subspace Iteration

In the following section we will present the inexact subspace iteration algorithm based on contour integration and prove basic convergence results using quasi-matrix notation. We consider the spectral transform functions from [5, 6, 26, 27] in the context of eigenvalues of operator-valued analytic functions. Let $T(\cdot)$ be given and let $\Gamma \subset \Omega$ be a closed curve which encloses either a set of *r* simple eigenvalues or a single semisimple eigenvalue whose multiplicity is *r*.

Let us consider the operators

$$A^{(q)} = \frac{1}{2\pi i} \int_{\Gamma} z^{q} T(z)^{-1} dz,$$

$$B^{(q)} = \frac{1}{2\pi i} \int_{\Gamma} z^{q} T(z)^{-1} T'(z) dz,$$

$$q = 0, 1,$$
(44)

and their approximations:

$$A_{N}^{(q)} = \sum_{k=1}^{N} \omega_{k} z_{k}^{q} T(z_{k})^{-1} ,$$

$$B_{N}^{(q)} = \sum_{k=1}^{N} \omega_{k} z_{k}^{q} T(z_{k})^{-1} T'(z_{k}) ,$$

$$q = 0, 1.$$
(45)

Here ω_i and $z_i \in \Gamma$, i = 1, ..., N, are integration weights and integration nodes. Based on Theorem 4.7 in [4] we establish, for ω_i and $z_i \in \Gamma$ defined by the *N* node trapeze integration formula for the contour integral (44), the following estimates:

$$\begin{split} \left\| A^{(q)} - A_N^{(q)} \right\| &\leq C_1 d \left(T \right)^{-\kappa} e^{-C_2 N d(T)}, \\ \left\| B^{(q)} - B_N^{(q)} \right\| &\leq C_1 d \left(T \right)^{-\kappa} e^{-C_2 N d(T)}. \end{split}$$
(46)

Here Γ is simple closed contour in Ω such that $\sigma(T) \cap \Gamma = \emptyset$, $d(T) = \min_{\lambda \in \sigma(T)} \operatorname{dist}(\lambda, \Gamma)$, and κ is the maximum order of poles for the inverse of *T*. The constants in (46) are in general different and depend on the maximum of the integrand on the contour Γ . For more details see [4, 7]. Subsequently we conclude that $A_N^{(q)} \to A^{(q)}$ and $B_N^{(q)} \to B^{(q)}$, q = 0, 1 in the norm resolvent sense.

3.1. Extracting Information on Eigenvalues Enclosed by a Contour. Based on Proposition 2 we see that operators $A^{(0)}$ and $B^{(0)}$ are finite rank operators such that $\operatorname{Ran}(A^{(0)}) = \operatorname{Ran}(B^{(0)})$ is the space spanned by linearly independent eigenvectors associated with semisimple eigenvalues which are enclosed by Γ . Rather than providing a technical proof of these claims, which will be a subject of subsequent reports, we will present numerical evidence on two judiciously chosen examples. Further, we see that based on [7] we can establish the following technical result for the operators $A^{(1)}$ and $B^{(1)}$.

Proposition 5. Let *T* be the operator-valued function from (15) and let Γ be the contour which encloses solely the *r*, counting according to algebraic multiplicity, semisimple eigenvalues of *T*. Define the matrices $A^{(q)}$ and $B^{(q)}$, q = 0, 1, as in (44) and let $Q : \mathbb{C}^r \to \widetilde{\mathcal{H}}, Q^*Q = I$, be a quasi-matrix such that $Q^*A^{(0)}Q$ and $Q^*B^{(0)}Q$ are invertible. Then the eigenvalues of the matrix pairs $(Q^*A^{(1)}Q, Q^*A^{(0)}Q)$ and $(Q^*B^{(1)}Q, Q^*B^{(0)}Q)$ are precisely the eigenvalues λ_i , $i = 1, \ldots, r$, where we count according to multiplicity. Furthermore, if $a_i, b_i \in \mathbb{C}^r$ satisfy $Q^*A^{(1)}Qa_i = \lambda_iQ^*A^{(0)}Qb_i$ and $Q^*B^{(1)}Qa_i = \lambda_iQ^*B^{(0)}Qb_i$, then Qa_i and Qb_i are eigenvectors of *T* associated with λ_i , $i = 1, \ldots, r$.

Proof. We will prove the statement for the matrix pair $(Q^*B^{(1)}Q, Q^*B^{(0)}Q)$ and note that the proof for the matrix pair $(Q^*A^{(1)}Q, Q^*A^{(0)}Q)$ is equivalent. Based on Proposition 2 we see that there are quasi-matrices $X : \mathbb{C}^r \to \widetilde{\mathcal{H}}$ and $Y : \mathbb{C}^r \to \widetilde{\mathcal{H}}$ and a neighborhood \mathscr{U} of Γ such that

$$T(z)^{-1}$$

$$= X \begin{bmatrix} (z - \lambda_1)^{-1} & 0 \cdots & 0 \\ 0 & (z - \lambda_2)^{-1} & \cdots & 0 \\ \cdots & \cdots & \cdots & \cdots \\ 0 & \cdots & 0 & (z - \lambda_r)^{-1} \end{bmatrix} Y^* \quad (47)$$

+ $H(z), \quad z \in \mathcal{U}$

Here H is an analytic operator-valued function. Now we compute

$$Q^* B^{(1)} Q = (Q^* X) \begin{bmatrix} \lambda_1 & 0 \cdots & 0 \\ 0 & \lambda_2 & \cdots & 0 \\ \cdots & \cdots & \cdots & \cdots \\ 0 & \cdots & 0 & \lambda_r \end{bmatrix} (Y^* Q)$$

$$Q^* B^{(0)} Q = (Q^* X) (Y^* Q).$$
(48)

Based on the assumptions of the theorem we conclude that (Q^*X) and (Y^*Q) have to be invertible and so the conclusion readily follows.

Before we proceed, note that norm resolvent convergence of $A_N^{(q)}$ and $B_N^{(q)}$ implies the convergence of spectra and associated spaces to those of $A^{(q)}$ and $B^{(q)}$, q = 0, 1. In particular, this means that $A_N^{(0)}$ and $B_N^{(0)}$, for N large enough, will have two well separated components in the spectrum and so we are motivated to use subspace iteration on the operator level to improve the quality of the approximate eigenvalues. Note that for a vector v the vectors $A_N^{(0)}v$ and $B_N^{(0)}v$ can be computed using formula (45) without ever forming a representation for the underlying operator.

3.2. Inexact Subspace Iteration Based on Quadrature Formula. We have the following algorithm for a generic bounded Pick a random quasi-matrix $Q_1 = [q_1 \ q_2 \ \cdots \ q_r], Q_1 : \mathbb{C}^r \to L^2(\Omega)$ Set n = 1 **repeat** Compute $Y_{n+1} = B_N Q_n$ Compute an orthogonal quasi-matrix $Q_{n+1}, Q_{n+1}^* Q_{n+1} = I$ using Gram-Schmidt so that $\operatorname{Ran}(Q_{n+1}) = \operatorname{Ran}(Y_{n+1})$. Compute criterion = $||Q_{n+1} - Q_n(Q_n^*Q_{n+1})||$ n := n + 1 **until** criterion < tol **return** Q_{n+1}



$$\begin{array}{l} \textbf{Require } Q_1: \mathbb{C}^r \to \widetilde{\mathscr{H}}, Q_1^*Q_1 = I_r, \texttt{tol}_1 \texttt{ and } N. \\ \text{Compute } Q \texttt{ as output from Algorithm 1 with tolerance } \texttt{tol}_1 \texttt{ and } Q_1 \texttt{ as starting value.} \\ \text{Compute } Z_1 = B_N^{(1)}Q \texttt{ and } Z_0 = B_N^{(0)}Q. \\ \text{Form } W^{(q)} = Q^*Z_q, q = 0, 1 \\ \text{Compute } v_i \texttt{ and } \lambda_i \texttt{ so that } W^{(1)}u_i = \lambda_i W^{(0)}u_i, i = 1, \ldots, r \\ \text{Compute eigenvectors } v_i = Q_{n+1}u_i, i = 1, \ldots, r. \\ \text{Compute } \|\texttt{t}(\lambda_i)[v_i,\cdot]\|_{\mathscr{Q}_{h_2}, t_0, -1} \\ \text{Form the index set } \{i_1, \ldots, i_{r_2}\} = \{i: \|\texttt{t}(\lambda_i)[v_i,\cdot]\|_{\mathscr{Q}_{h_2}, t_0, -1} \leq \texttt{tol}_2 \} \\ \text{Set } V = \begin{bmatrix} v_{i_1} & \cdots & v_{i_r_2} \end{bmatrix} \\ \text{return } V: \mathbb{C}^{r_2} \to \widetilde{\mathscr{H}} \texttt{ and } \lambda_{i_j}, j = 1, \ldots, r_2 \end{array}$$

ALGORITHM 2: Approximating eigenvectors and eigenvalues with tolerance tol₂.

operator *B* of type (44) which has a component of $r, r \in \mathbb{N}$ dominant eigenvalues and whose action on a vector can be represented by a formula (45).

For Algorithm 1 we present standard convergence results in Theorems 6 and 8. First we consider a special case in which B is a Hermitian (bounded self-adjoint) operator. We will see that a possible presence of singularities in the operator function will not pollute the convergence rate.

Theorem 6. Let B be a bounded Hermitian operator such that $\sigma(B) = \{\lambda_1, \ldots, \lambda_r\} \cup \Sigma$, where λ_i , $i = 1, \ldots, r$, are eigenvalues of finite multiplicity and $|\lambda_1| \ge \cdots \ge |\lambda_r| > \tau > 0$, and $\tau > |\lambda|$, $\lambda \in \Sigma$. Then for the sequence $V^{(n+1)} := BV^{(n)}$, where $V^{(0)}$: $\mathbb{C}^r \to \mathcal{H}$, $(V^{(0)})^*V^{(0)} = I_r$, we have

$$\tan\theta\left(V^{(n+1)},U\right) \le \frac{\tau^n}{\left|\lambda_r\right|^n}\tan\theta\left(V^{(0)},U\right),\qquad(49)$$

where $U : \mathbb{C}^r \to \mathcal{H}, U^*U = I$, is such that $U = [u_1 \ u_2 \ \cdots \ u_r]$ and $Bu_i = \lambda_i u_i, i = 1, \dots, r$.

Proof. Let $P_1 = UU^*$ and represent $U : \mathbb{C}^r \to \mathcal{H}$ with respect to the partition of unity $P_1 \oplus P_2 = I$ by $U = [U_1 \ U_2]$, with $\operatorname{Ran}(P_1) = \operatorname{Ran}(U_1)$. Since *B* is self-adjoint it can be represented by the block operator matrix:

$$B = \begin{bmatrix} B_{11} & 0\\ 0 & B_{22} \end{bmatrix}.$$
 (50)

Let us represent, with respect to the same partition of unity, the quasi-matrix $V^{(0)}$ as

$$V^{(0)} = \begin{bmatrix} V_1^{(0)} \\ V_2^{(0)} \end{bmatrix}.$$
 (51)

Without reducing the level of generality we may assume that $V_1^{(0)}$ is invertible, since otherwise $\text{Ran}(V^{(0)})$ would contain eigenvectors of *B*. This corresponds to the trivial situation and will not be further discussed.

Since $(V^{(0)})^* V^{(0)} = I_r$ it follows from (6) that

$$\tan\theta\left(U_{1},V^{(0)}\right) = \left\|V_{2}^{(0)}\left(V_{1}^{(0)}\right)^{-1}\right\|.$$
(52)

The assumption of the separation of the spectra of *B* implies that B_{11} must also be invertible, and so

$$V^{(n)} := B^{n} V^{(0)} = \begin{bmatrix} B_{11}^{n} V_{1}^{(0)} \\ B_{22}^{k} V_{2}^{(0)} \end{bmatrix} = \begin{bmatrix} I \\ F^{(n)} \end{bmatrix} B_{11}^{n} V_{1}^{(0)}$$
(53)

for $F^{(n)} = B_{22}^n V_2^{(0)} (V_1^{(0)})^{-1} B_{11}^{-n}$ which implies that

$$V^{(n)} \left(V_1^{(0)}\right)^{-1} B_{11}^{-n} = \begin{bmatrix} I\\F^{(n)} \end{bmatrix} = U_1 + U_2 F^{(n)}.$$
(54)

And so from the optimality of the canonical angles (see [15]), it follows that

$$\tan\theta\left(V^{(n)},U\right) \le \left\|F^{(n)}\right\|.$$
(55)

Consequently, it follows that

$$\tan \theta \left(V^{(n)}, U \right) \leq \left\| F^{(n)} \right\| \leq \left\| B_{22}^{n} \right\| \left\| F^{(0)} \right\| \left\| B_{11}^{-n} \right\|$$

$$\leq \frac{\tau^{n}}{\left| \lambda_{r} \right|^{n}} \tan \theta \left(V^{(0)}, U_{1} \right).$$
(56)

Now the conclusion of the theorem readily follows.

Remark 7. Note that from (8) and (9) it is clear that $\theta(V^{(n)}, U)$ does not depend on the quasi-matrices $V^{(n)}$ and U, but rather on the orthogonal projections onto $\operatorname{Ran}(V^{(n)})$ and $\operatorname{Ran}(U) = \operatorname{Ran}(P_1)$. In Theorem 6 we have stated the convergence result for the angle $\theta(V^{(n)}, U)$. Here we have ignored the choice of a basis of $\operatorname{Ran}(V^{(n)})$ which is a very important part of Algorithm 1. Note that in practical computations the choice of the basis is crucial to achieve an efficient implementation; see [28].

In the case in which the operator B is not Hermitian but has clearly separated invariant subspace associated with the finite collection, counting according to the algebraic multiplicity of dominant eigenvalues we have the theorem below.

Theorem 8. Let *B* be a bounded operator such that its spectrum has two disjoint components Σ_1 and Σ_2 so that $\sigma(B) = \Sigma_1 \cup \Sigma_2$. Let further $\Sigma_1 = \{\lambda_1, \lambda_2, ..., \lambda_r\}$, where we have counted $\lambda_i \in \sigma(B)$ according to the algebraic multiplicity of an eigenvalue. Further let a partition of unity $P_1 \oplus P_2 = I$ exist such that $r = \dim(\operatorname{Ran}(P_1))$ and

$$B = \begin{bmatrix} B_{11} & B_{12} \\ 0 & B_{22} \end{bmatrix},$$
 (57)

and let further $sep(B_{11}, B_{22}) > 1$ and $||B_{22}|| ||B_{11}^{-1}|| < 1$. Then for a quasi-matrix $V : \mathbb{C}^r \to \mathcal{H}$ we have

$$\tan \theta \left(B^{n}V, P_{1} \right) \leq \left(\left\| B_{22} \right\| \left\| B_{11}^{-1} \right\| \right)^{n} (1 + \|G\|) \tan \theta \left(V, P_{1} \right).$$
(58)

Here G is the solution of the Sylvester equation $GB_{22} - B_{11}G = B_{12}$ and $sep(B_{11}, B_{22}) = \inf_{\|G\|=1} \|GB_{22} - B_{11}G\|.$

Proof. For *B* such that $\sigma(B) = \Sigma_1 \cup \Sigma_2$ with $\Sigma_1 = \{\lambda_1, \lambda_2, ..., \lambda_r\}$ satisfying the assumptions of the theorem choose P_1 as the orthogonal projection onto the maximal invariant subspace associated with $\{\lambda_1, \lambda_2, ..., \lambda_r\}$ and $P_2 = I - P_1$; then *B* has the block operator representation as claimed in the theorem. We note that

$$B = \begin{bmatrix} B_{11} & B_{12} \\ 0 & B_{22} \end{bmatrix} = \begin{bmatrix} I & -G \\ 0 & I \end{bmatrix} \begin{bmatrix} B_{11} & 0 \\ 0 & B_{22} \end{bmatrix} \begin{bmatrix} I & G \\ 0 & I \end{bmatrix}, \quad (59)$$

where $GB_{22} - B_{11}G = B_{12}$. Operator *G*, such that $GB_{22} - B_{11}G = B_{12}$ exists and is unique, for example, [29, Lemma A.1], providing the spectra of B_{11} and B_{22} , are disjoint. The rest of the proof follows analogously as in Theorem 6.

Remark 9. In [29, Lemma A.1] there is also an estimate of $sep(B_{11}, B_{22})$ in terms of the pseudospectra distances.

This result implies that the convergence rate for the subspace iteration essentially depends on the properties of the operator and not on the subspace which is iterated. Further, the dependence on the number of integration nodes N diminishes exponentially. To extract eigenvalues and eigenvectors we use the method from [7] and apply Proposition 5 on the quasi-matrix Q_{n+1} which is returned by Algorithm 1. More precisely we have the following. For simplicity we write $B = B^{(0)}$.

Theorem 10. Let the conditions of Theorem 8 be satisfied and also assume the same notation conventions. For a quasi-matrix $V : \mathbb{C}^r \to \mathcal{H}$, we set $Q_1 := V$ and apply Algorithm 1 with fixed $N \in \mathbb{N}$. Then

$$\sin \theta \left(\operatorname{Ran} \left(Q_{n} \right), P_{1} \right)$$

$$\leq \left(\left\| B_{22} \right\| \left\| B_{11}^{-1} \right\| + C_{3} d \left(T \right)^{-\kappa} e^{-C_{2} N d(T)} \right)^{n} + \left(1 + \left\| G \right\| + C_{4} d \left(T \right)^{-\kappa} e^{-C_{2} N d(T)} \right) \tan \theta \left(Q_{1}, P_{1} \right)$$

$$+ C_{5} d \left(T \right)^{-\kappa} e^{-C_{2} N d(T)}.$$
(60)

Proof. Recall that

$$\|B - B_N\| \le C_1 d(T)^{-\kappa} e^{-C_2 N d(T)}.$$
 (61)

From standard results on norm convergence of bounded operators (see [16]) we conclude that

$$\sin \theta (P_1, (P_N)_1) = ||P_1 - (P_N)_1|| \\\leq C_5 d (T)^{-\kappa} e^{-C_2 N d(T)}.$$
(62)

Here $(P_N)_1 + (P_N)_2 = I$ is the partition of unity in which B_N has a block triangular form:

$$B_N = \begin{bmatrix} (B_N)_{11} & (B_N)_{12} \\ 0 & (B_N)_{22} \end{bmatrix}.$$
 (63)

Assumption on the spectral separation from Theorem 8 implies that B_{11} is invertible and obviously

$$\|B_{11} - (B_N)_{11}\| \le C_1 d(T)^{-\kappa} e^{-C_2 N d(T)}$$
(64)

$$\|B_{22} - (B_N)_{22}\| \le C_1 d(T)^{-\kappa} e^{-C_2 N d(T)}.$$
 (65)

Norm convergence (64) implies that for N large enough $(B_N)_{11}$ must also be invertible; for details see [30]. For such N we have

$$\begin{split} \left\| B_{11}^{-1} - (B_N)_{11}^{-1} \right\| &\leq \left\| (B_{11})^{-1} (B_{11} - (B_N)_{11}) (B_N)_{11}^{-1} \right\| \\ &\leq \left\| (B_{11})^{-1} \right\| \left\| (B_N)_{11}^{-1} \right\| \left\| B_{11} - (B_N)_{11} \right\| \quad (66) \\ &\leq C_3 d \left(T \right)^{-\kappa} e^{-C_2 N d(T)}. \end{split}$$

Triangle inequality implies that

$$\| (B_N)_{22} \| \le \| B_{22} \| + \| (B_N)_{22} - B_{22} \|,$$

$$\| (B_N)_{11}^{-1} \| \le \| B_{11}^{-1} \| + \| (B_N)_{11}^{-1} - B_{11}^{-1} \|$$
(67)

and so

$$\|(B_N)_{22}\| \|(B_N)_{11}^{-1}\| \le \|B_{22}\| \|B_{11}^{-1}\| + C_4 d(T)^{-\kappa} e^{-C_2 N d(T)}.$$
(68)

Let us now assume that N is large enough so that

$$||B_{22}|| ||B_{11}^{-1}|| + C_4 d(T)^{-\kappa} e^{-C_2 N d(T)} < 1.$$
(69)

Then we can apply Theorem 8 to B_N and conclude that

$$\sin \theta \left(B_{N}^{n} V, \left(P_{N} \right)_{1} \right) \leq \tan \theta \left(B_{N}^{n} V, \left(P_{N} \right)_{1} \right)$$
$$\leq \left(\left\| \left(B_{N} \right)_{22} \right\| \left\| \left(B_{N} \right)_{11}^{-1} \right\| \right)^{n} \left(1 + \left\| G_{N} \right\| \right)$$
$$\cdot \tan \theta \left(V, \left(P_{N} \right)_{1} \right).$$
(70)

Equivalently since $sep((B_N)_{22}, (B_N)_{11}) \rightarrow sep(B_{22}, B_{11})$ and $(B_N)_{12} \rightarrow B_{12}$ in norm we conclude that

$$\|G_N\| \le \|G\| + C_5 d(T)^{-\kappa} e^{-C_2 N d(T)}.$$
(71)

We now apply the triangle inequality for the sine of the maximal canonical angle to conclude the proof. \Box

Using Theorem 10 we will define the *effective convergence rate* of the inexact contour based subspace iteration. Recall that

$$B = \frac{1}{2\pi i} \int_{\Gamma} T(z)^{-1} T'(z) dz,$$

$$B_{N} = \sum_{k=1}^{N} \omega_{k} z_{k}^{q} T(z_{k})^{-1} T'(z_{k}).$$
(72)

Starting from (63) we define the *effective convergence rate* for the inexact subspace iteration for *B* as

$$\eta(\Gamma, N) := \left\| (B_N)_{22} \right\| \left\| (B_N)_{11}^{-1} \right\|.$$
(73)

Remark 11. We will project B_N on carefully constructed finite dimensional spaces, for example, finite element spaces, to experimentally estimate $\eta(\Gamma, N)$ in the following section. We will not further elaborate on this procedure but solely present the results of experiments for illustration purposes.

3.3. Auxiliary Subspace Error Estimates and the Convergence Criterion. In this paper we will use the inequality $\|\mathbf{t}(\mu)[u,\cdot]\|_{\bar{\mathcal{Q}}_{h_2},t_0,-1} \leq \|\mathbf{t}(\mu)[u,\cdot]\|_{t_0,-1}$ to filter the eigenpairs which have been extracted using Proposition 5 from the subspace returned by Algorithm 1. More precisely, eigenpairs for which $\|\mathbf{t}(\mu)[u,\cdot]\|_{\bar{\mathcal{Q}}_{h_2},t_0,-1}$ is too large will be discarded. We will indicate the importance of this step of the algorithm in the experiments section. 9

In the case in which the number of eigenvalues returned by Algorithm 2 is not satisfactory, the procedure can be extended as an innerouter iteration scheme by setting $Q_1 = V$ and repeating the procedure. One further possibility to improve performance is to modify Algorithm 1 so that at the beginning of the repeat loop we remove all those directions from Q_1 which are identified by Algorithm 2 as converged and then proceed as in the original versions of Algorithm 1.

4. Numerical Experiments

For the finite element discretizations, we use the space of piecewise linear and continuous finite elements on a given subdivision δ_h of an interval [a, b]. We denote this space by $\mathcal{V}_h(a, b)$. For our computations we set the tolerance tol so as to balance the error when solving the source problem by finite element approximation with the error in the integration. For the contour Γ we chose a circle and as integration nodes and weights we use the trapeze formula from [7] and Gauss-Legendre nodes and weights from [5].

To estimate the approximation errors in our experiments we will use an auxiliary subspace error estimation technique. To this end we use $\mathcal{Q}_h(a, b)$ to denote the space of piecewise quadratic and continuous function on the same subdivision of the interval [a, b] which was used to define the space $\mathcal{V}_h(a, b)$. To estimate the error $E_{\mu,u}$ from (39) for the function $u \in \mathcal{V}_h(a, b)$ and $\mu \in \mathbb{C}$ we use the formula; see [24] for further information:

$$\left\|\mathbf{t}\left(\boldsymbol{\mu}\right)\left[\boldsymbol{u},\cdot\right]\right\|_{\mathcal{Q}_{h}\left(\boldsymbol{a},\boldsymbol{b}\right),\mathbf{t}_{0},-1} = \sup_{\boldsymbol{\phi}\in\mathcal{Q}_{h}\left(\boldsymbol{a},\boldsymbol{b}\right)\setminus\{0\}} \frac{\left|\mathbf{t}\left(\boldsymbol{\mu}\right)\left[\boldsymbol{u},\boldsymbol{\phi}\right]\right|}{\sqrt{\mathbf{t}_{0}\left[\boldsymbol{\phi},\boldsymbol{\phi}\right]}}.$$
 (74)

Example 12. We study the quadratic eigenvalue problem:

$$\partial_{xx}u = \lambda \left(\gamma \partial_{xx}u + \delta u\right) + \lambda^2 u$$

(75)
$$u \left(-2\right) = u \left(2\right) = 0$$

for $\gamma = 0.05$ and $\delta = 0.3$. As a benchmark we use the eigenvalue computed with a spectral discretization using the chebfun system; see [10]. The timings are 2.5 seconds for chebfun and 1.6 seconds for subspace iteration. We used the trapeze formula to define the operator B_N with N = 25 and achieved the residuals as small as 1e-12. The results are presented in Figure 1. From Figure 1 we see that the spectrum of (75) is well separated, and so the sufficient separation of the spectra of $(B_N)_{11}$ and $(B_N)_{22}$ can be achieved with few integration nodes N: for example, N = 25. Figure 2(a) depicts the dependence of the effective convergence rate for the inexact subspace iteration $\eta(\Gamma, N) = ||(B_N)_{22}|||(B_N)_{11}^{-1}||$ on the number of integration nodes N for the trapeze formula. The contour was a circle which enclosed the five eigenvalues marked by crosses.

Remark 13. Note that in [5, 6] it was shown that Gauss-Legendre and trapeze integration require similar number of integration nodes to be applicable in inexact subspace iteration algorithm. A slight experimental advantage was noticed in favor of Gauss-Legendre, but both approaches



FIGURE 1: Eigenvalues for the quadratic eigenvalue problem.



 $\bigcirc \eta(\Gamma, N)$

(a) Quadratic eigenvalue problem with well-separated spectrum from Example 12 and Γ as in Figure 1

(b) Rational eigenvalue problem with tightly clustered eigenvalues from Example 14 and Γ as in Figure 4

FIGURE 2: The dependence of the effective convergence rate $\eta(\Gamma, N)$ as defined in (73) on the number of integration nodes N. This illustrates the negative effect of the distance to Γ on the convergence rate.

can be used in competitive algorithms depending on the particular context of the application.

Example 14. We will now consider the following rational eigenvalue problem:

$$\frac{d^{2}}{dx^{2}}u + \frac{\lambda}{1-\lambda}\chi_{|x| \le 9/10} \cdot u + \frac{\lambda}{3-\lambda}\chi_{|x| \le 9/10} \cdot u + \lambda\left(\chi_{|x| \le 9/10}u + 10\chi_{|x| > 9/10}u\right) = 0$$
(76)

$$u \in H^1(-1, 1)$$
 with periodic b.c.,

where

$$\chi_{|x| \le 9/10} : x \longmapsto \begin{cases} 0: & x \notin \left[-\frac{9}{10}, \frac{9}{10} \right] \\ & & \\ 1: & x \in \left[-\frac{9}{10}, \frac{9}{10} \right] \end{cases}$$
(77)

is the characteristic function of the interval $\left[-9/10, 9/10\right]$.

This eigenvalue problem has a singularity at $\lambda = 1$ and λ = 3. The eigenvalues below the singularity λ = 1 and those in $\langle 1, 3 \rangle$ are discrete (finite multiplicity) and they accumulate at $\lambda = 1$ and at $\lambda = 3$ from below [31].



FIGURE 3: Eigenvectors for eigenvalues in [0, 1).



FIGURE 4: Eigenvectors for eigenvalues in $\langle 1, 3 \rangle$.

The results in Figures 3 and 4 are presented solely for benchmark purposes. They were computed using trapeze formula with 4500 nodes and validated using the approach from [8]. Looking at Figures 3 and 4 we see that in general there are two classes of eigenfunctions: those that are confined to the interval [-9/10, 9/10] and those that are not. Further discussion of the qualitative properties of eigenfunctions is outside the scope of this paper. We will address the modeling aspects elsewhere.

As a first experiment in Table 1 we present the empirical study of the convergence of the residual measures and the relative eigenvalue error computed from the sequence $V^{(n+1)} = B_N V^{(n)}, V_0 : \mathbb{C}^2 \to \mathcal{V}$, where N = 16 and ω_i and z_i are Gauss-Legendre integration nodes and weights from [5]. We have used the contour which encloses only two well-separated eigenvalues, two left most crosses, in Figure 3.

TABLE 1: Convergence history for the rational eigenvalue problemExample 14 and the two leftmost eigenvalues from Figure 3.

Iteration number <i>n</i>	Maximal	Maximal residual
	eigenvalue error	Cottillate
1	2.0880e - 008	1.2607e - 009
2	4.9292e - 009	1.3579e - 010
3	4.9282e - 009	1.3580e - 010

The convergence rate for the well-separated eigenvalues and the standard choice of the integration nodes from [6] seem satisfactory. However, the eigenvalues of (76) cluster towards z = 1. In Figure 2(b) we see that many more integration nodes are necessary to achieve a similar estimate of the effective convergence rate $\eta(\Gamma, N)$ compared to the



FIGURE 5: We display $||Q_{n+1} - Q_n(Q_n^*Q_{n+1})||$ for n = 1, ..., 10.

quadratic eigenvalue problem in Figure 2(a). To this end we will compare below the measured convergence rates for the criterion in Algorithm 1 on Examples (76) and (75).

Remark 15. Results presented in Figure 5 show that, unlike what was suggested in [6], the number of integration nodes might have to be adaptively adjusted for some contours Γ . Also we emphasize that the residual criterion

$$\left\| \mathbf{t} \left(\boldsymbol{\mu} \right) \left[\boldsymbol{u}, \cdot \right] \right\|_{\mathcal{Q}_{\boldsymbol{\mu}}(\boldsymbol{a}, \boldsymbol{b}), \mathbf{t}_{0}, -1} \le \mathtt{tol}_{2} \tag{78}$$

in Algorithm 2 is particularly important in the case of the clustering of eigenvalues towards Γ like in Figures 3 and 4. In comparison with benchmark results from Figures 3 and 4 three spurious eigenvalues, which would otherwise have been declared as converged, were discarded because their residuals were larger than a threshold. With this in mind we justify the application of the modified algorithm from [7]. The algorithm in [7] had residual error control but controlled the accuracy solely by increasing the number of integration nodes (e.g., we used 4500 nodes for benchmark results in Figures 3 and 4). In our modification we combine subspace iteration with contour integration with problem adapted number of nodes (e.g., with N = 16 and 3 steps of inexact subspace iteration acceleration we reach the same accuracy (see Table 1) as a priori predicted for N = 4500 based on [8]). Finally, we note that performance issues are outside the scope of this paper. As is indicated in [6] the algorithm offers large potential to leverage parallel processing. This will be the topic of further research.

4.1. Stability of Convergence Rate of Inexact Subspace Iteration. From Figure 2 we see that the convergence rate for the subspace iteration essentially depends on $||B_{22}|| ||B_{11}^{-1}|| \approx \eta(\Gamma, N)$. We will now see that this convergence rate is relatively robust



FIGURE 6: Eigenvalues for the eigenvalue problem (79).

to perturbations as long as the distance from the curve Γ and $\sigma(T)$ is not too small. To this end we consider the problem

ı

$$\partial_{xx}u = \lambda \left(\gamma \partial_{xx}u + \delta u\right) + \lambda^2 u + \frac{1}{4} \exp\left(-\lambda\right)$$

$$u(-2) = u(2) = 0$$
(79)

as a perturbation of (75). Its eigenvalues are presented on Figure 6. We see that both the eigenvalues of (75) and (79) are equally well separated from Γ and so the convergence rates, as measured by the decay rate of $||Q_{n+1} - Q_n(Q_n^*Q_{n+1})||$ in Algorithm 1, are similar. However, when comparing the results from Figure 7(c) with those from Figure 7(a) we see that we need more iterations to achieve a comparable reduction in the convergence criterion for Algorithm 1 in the



FIGURE 7: We display $||Q_{n+1} - Q_n(Q_n^*Q_{n+1})||$ for n = 1, ..., 10. For Γ we have used the contours as in Figures 1, 6, and 5(a), respectively.

case of eigenvalue clustering towards Γ . On the other hand, Figure 7(b) indicates that the nonlinear perturbation did not change the convergence rate significantly. This justifies the consideration of the adaptivity both in the choice of integration nodes and in representing the operators.

5. Conclusion

We have shown that the subspace iteration nonlinear eigensolver based on spectral transformation of the analytic Fredholm valued function converges at rates that depend primarily on the problem and not on the discretization. Further, we have seen that the distance from the contour and the spectrum does limit the accuracy of the approximation based on numerical integration of the resolvent. Also, residual norm is a reliable estimator of the approximation error even in the presence of poles. In comparison, the convergence in the case of the quadratic eigenvalue problem which had a more pronounced spectral gap than the rational problem was much faster and more robust. This suggests that a subspace iteration algorithm that increases the number of integration nodes based on a posteriori computable criterion is promising for nonlinear eigenvalue problems.

Conflict of Interests

The authors declare that there is no conflict of interests regarding the publication of this paper.

Acknowledgments

Christian Engström gratefully acknowledges the support of the Swedish Research Council under Grant no. 621-2012-3863. The work of Luka Grubišić has in part been supported

References

- C. Effenberger, D. Kressner, and C. Engström, "Linearization techniques for band structure calculations in absorbing photonic crystals," *International Journal for Numerical Methods in Engineering*, vol. 89, no. 2, pp. 180–191, 2012.
- [2] K. F. Gurski, R. Kollár, and R. L. Pego, "Slow damping of internal waves in a stably stratified fluid," *Proceedings A: Mathematical, Physical and Engineering Sciences*, vol. 460, no. 2044, pp. 977– 994, 2004.
- [3] P. J. Schmid and D. S. Henningson, Stability and Transition in Shear Flows, vol. 142 of Applied Mathematical Sciences, Springer, New York, NY, USA, 2001.
- [4] W.-J. Beyn, "An integral method for solving nonlinear eigenvalue problems," *Linear Algebra and Its Applications*, vol. 436, no. 10, pp. 3839–3863, 2012.
- [5] E. Polizzi, "Density-matrix-based algorithm for solving eigenvalue problems," *Physical Review B: Condensed Matter and Materials Physics*, vol. 79, no. 11, Article ID 115112, 2009.
- [6] P. T. P. Tang and E. Polizzi, "Feast as a subspace iteration eigensolver accelerated by approximate spectral projection," *SIAM Journal on Matrix Analysis and Applications*, vol. 35, no. 2, pp. 354–390, 2014.
- [7] W.-J. Beyn, Y. Latushkin, and J. Rottmann-Matthes, "Finding eigenvalues of holomorphic Fredholm operator pencils using boundary value problems and contour integrals," *Integral Equations and Operator Theory*, vol. 78, no. 2, pp. 155–211, 2014.
- [8] L. Grubišić and A. Grbić, "Discrete perturbation estimates for eigenpairs of fredholm operator-valued functions," *Applied Mathematics and Computation*, vol. 267, pp. 632–647, 2015.
- [9] R. Mennicken and M. Möller, Non-Self-Adjoint Boundary Eigenvalue Problems, vol. 192 of North-Holland Mathematics Studies, North-Holland Publishing Company, Amsterdam, The Netherlands, 2003.
- [10] R. B. Platte and L. N. Trefethen, "Chebfun: a new kind of numerical computing," in *Progress in Industrial Mathematics at ECMI* 2008, vol. 15 of *Mathematics in Industry*, pp. 69–87, Springer, Berlin, Germany, 2010.
- [11] A. Townsend and L. N. Trefethen, "Continuous analogues of matrix factorizations," *Proceedings of the Royal Society A: Mathematical, Physical and Engineering Sciences*, vol. 471, no. 2173, 2015.
- [12] L. N. Trefethen, "Householder triangularization of a quasimatrix," *IMA Journal of Numerical Analysis*, vol. 30, no. 4, pp. 887– 897, 2010.
- [13] C. Tretter, Spectral Theory of Block Operator Matrices and Applications, Imperial College Press, London, UK, 2008.
- [14] P. R. Halmos, "Two subspaces," *Transactions of the American Mathematical Society*, vol. 144, pp. 381–389, 1969.
- [15] C. Davis and W. M. Kahan, "The rotation of eigenvectors by a perturbation. III," *SIAM Journal on Numerical Analysis*, vol. 7, no. 1, pp. 1–46, 1970.

- [16] T. Kato, Perturbation Theory for Linear Operators, Classics in Mathematics, Springer, Berlin, Germany, 1995, Reprint of the 1980 edition.
- [17] I. C. Gohberg and E. I. Sigal, "Global factorization of a meromorphic operator-function and some of its applications," *Matematicheskie Issledovaniya*, vol. 6, no. 1, pp. 63–82, 1971.
- [18] I. C. Gohberg and E. I. Sigal, "An operator generalization of the logarithmic residue theorem and the theorem of Rouché," *Matematicheskii Sbornik*, vol. 84, no. 126, pp. 607–629, 1971.
- [19] I. Gohberg and E. I. Sigal, "The root multiplicity of the product of meromorphic operator functions," *Istoriko-Matematicheskie Issledovaniya*, vol. 6, pp. 33–50, 1971.
- [20] D. Bindel and H. Amanda, "Localization theorems for nonlinear eigenvalue problems," *SIAM Journal on Matrix Analysis and Applications*, vol. 34, no. 4, pp. 1728–1749, 2013.
- [21] W. W. Hager, "Updating the inverse of a matrix," *SIAM Review*, vol. 31, no. 2, pp. 221–239, 1989.
- [22] R. Harte, Invertibility and Singularity for Bounded Linear Operators, vol. 109 of Monographs and Textbooks in Pure and Applied Mathematics, Marcel Dekker, New York, NY, USA, 1988.
- [23] J. Sherman and W. J. Morrison, "Adjustment of an inverse matrix corresponding to a change in one element of a given matrix," *The Annals of Mathematical Statistics*, vol. 21, no. 1, pp. 124–127, 1950.
- [24] R. E. Bank, "Hierarchical bases and the finite element method," Acta Numerica, vol. 5, pp. 1–43, 1996.
- [25] J. S. Ovall, "Function, gradient, and Hessian recovery using quadratic edge-bump functions," *SIAM Journal on Numerical Analysis*, vol. 45, no. 3, pp. 1064–1080, 2007.
- [26] L. Krämer, E. Di Napoli, M. Galgon, B. Lang, and P. Bientinesi, "Dissecting the FEAST algorithm for generalized eigenproblems," *Journal of Computational and Applied Mathematics*, vol. 244, no. 1, pp. 1–9, 2013.
- [27] L. Krämer, Integration based solvers for standard and generalized Hermitian eigenvalue problems [Ph.D. thesis], Bergische Universität, Wuppertal, Germany, 2014.
- [28] G. W. Stewart, "Simultaneous iteration for computing invariant subspaces of non-hermitian matrices," *Numerische Mathematik*, vol. 25, no. 2, pp. 123–136, 1976.
- [29] C. Beattie, "Galerkin eigenvector approximations," *Mathematics of Computation*, vol. 69, no. 232, pp. 1409–1434, 2000.
- [30] T. Kato, Perturbation Theory for Linear Operators, vol. 132 of Grundlehren der Mathematischen Wissenschaften, Band, Springer, Berlin, Germany, 2nd edition, 1976.
- [31] C. Engström, H. Langer, and C. Tretter, "Non-linear eigenvalue problemsand applications to photonic crystals," http://arxiv .org/abs/1507.06381.



The Scientific World Journal





Decision Sciences







Journal of Probability and Statistics



Hindawi Submit your manuscripts at http://www.hindawi.com



(0,1),

International Journal of Differential Equations





International Journal of Combinatorics





Mathematical Problems in Engineering



Abstract and Applied Analysis



Discrete Dynamics in Nature and Society







Function Spaces



International Journal of Stochastic Analysis



Journal of Optimization