

A GENERALIZATION OF THE STEEPEST DESCENT METHOD FOR MATRIX FUNCTIONS*

M. AFANASJEW[†] M. EIERMANN[†] O. G. ERNST[†] AND S. GÜTTEL[†]

In memory of Gene Golub

Abstract. We consider the special case of the restarted Arnoldi method for approximating the product of a function of a Hermitian matrix with a vector which results when the restart length is set to one. When applied to the solution of a linear system of equations, this approach coincides with the method of steepest descent. We show that the method is equivalent to an interpolation process in which the node sequence has at most two points of accumulation. This knowledge is used to quantify the asymptotic convergence rate.

Key words. Matrix function, Krylov subspace approximation, restarted Krylov subspace method, restarted Arnoldi/Lanczos method, linear system of equations, steepest descent, polynomial interpolation

AMS subject classifications. 65F10, 65F99, 65M20

1. Introduction. To evaluate the expression $f(A)\mathbf{b}$ for a matrix $A \in \mathbb{C}^{n \times n}$, a vector $\mathbf{b} \in \mathbb{C}^n$ and a function $f : \mathbb{C} \supset D \rightarrow \mathbb{C}$ such that $f(A)$ is defined, approximations based on Krylov subspaces have recently regained new attention, typically for the case when A is large and sparse or structured. In [6] we proposed a technique for restarting the Krylov subspace approximation which permits the calculation to proceed using a fixed number of vectors (and hence storage) in the non-Hermitian case and avoids the additional second Krylov subspace generation phase in the Hermitian case. The method is based on a sequence of standard Arnoldi decompositions

$$AV_j = V_j H_j + \eta_{j+1} \mathbf{v}_{j m+1} \mathbf{e}_m^T, \quad j = 1, 2, \dots, k,$$

with respect to the m -dimensional Krylov subspaces $\mathcal{K}_m(A, \mathbf{v}_{(j-1)m+1})$, where $\mathbf{v}_1 = \mathbf{b}/\|\mathbf{b}\|$. Alternatively, we write

$$A\widehat{V}_k = \widehat{V}_k \widehat{H}_k + \eta_{k+1} \mathbf{v}_{km+1} \mathbf{e}_{km}^T,$$

where $\widehat{V}_k := [V_1 \ V_2 \ \dots \ V_k] \in \mathbb{C}^{n \times km}$,

$$\widehat{H}_k := \begin{bmatrix} H_1 & & & & & \\ E_2 & H_2 & & & & \\ & & \ddots & & & \\ & & & \ddots & & \\ & & & & E_k & H_k \end{bmatrix} \in \mathbb{C}^{km \times km}, \quad E_j := \eta_j \mathbf{e}_1 \mathbf{e}_m^T \in \mathbb{R}^{m \times m}, \quad j = 2, \dots, k.$$

The approximation to $f(A)\mathbf{b}$ associated with this Arnoldi-like decomposition is given by

$$\mathbf{f}_k := \|\mathbf{b}\| \widehat{V}_k f(\widehat{H}_k) \mathbf{e}_1$$

(cf. [6] or [1] for algorithms to compute \mathbf{f}_k) and we refer to this approach as the restarted Arnoldi method with restart length m .

*Received October 16, 2007. Accepted for publication April 22, 2008. Recommended by D. Szyld. This work was supported by the Deutsche Forschungsgemeinschaft.

[†]Institut für Numerische Mathematik und Optimierung, Technische Universität Bergakademie Freiberg, D-09596 Freiberg, Germany ({martin.afanasjew,eiermann,ernst}@math.tu-freiberg.de, guettel@googlemail.com).

The convergence analysis of the sequence $\{f_k\}$ is greatly facilitated by the fact (see, e.g., [6, Theorem 2.4]) that

$$f_k = p_{km-1}(A)b,$$

where $p_{km-1} \in \mathcal{P}_{km-1}$ is the unique polynomial of degree at most $km-1$ which interpolates f at the eigenvalues of \widehat{H}_k (i.e., at the eigenvalues of H_j , $j = 1, 2, \dots, k$) in the Hermite sense. Convergence results for the restarted Arnoldi approximation can be obtained if we are able to answer the following two questions:

1. Where in the complex plane is $\Lambda(\widehat{H}_k)$, the spectrum of \widehat{H}_k , located?
2. For which $\lambda \in \mathbb{C}$ do the interpolation polynomials of f (with nodal set $\Lambda(\widehat{H}_k)$) converge to $f(\lambda)$?

We shall address these issues for the simplest form of this scheme obtained for a restart length of $m = 1$, in which case all Hessenberg matrices H_j are 1×1 and \widehat{H}_k is lower bidiagonal. We refer to this method as the *method of steepest descent for matrix functions*, and we shall present it in greater detail and derive some of its properties in Section 2. In particular, we shall show that, when applied to the function $f(\lambda) = 1/\lambda$, it reduces to the classical *method of steepest descent* for the solution of $Ax = b$, at least if A is Hermitian positive definite.

Although not competitive for the practical solution of systems of linear equations, this method has highly interesting mathematical properties and a remarkable history: More than 100 years after Cauchy [3] introduced it, Forsythe and Motzkin [11] noticed in numerical experiments (see also [8]) that the associated error vectors are asymptotically a linear combination of the eigenvectors belonging to the smallest and largest eigenvalue of A , an observation also made by Stiefel (Stiefel’s cage, see [16]) in the context of relaxation methods. Forsythe and Motzkin also saw that the sequence of error vectors is “asymptotically of period 2”. They were able to prove this statement for problems of dimension $n = 3$ and conjectured that it holds for all n [10]. It was Akaike [2] who first proved this conjecture in 1959. He rephrased the problem in terms of probability distributions and explained the observations of [11, 10] completely. Later, in 1968, Forsythe [9] reconsidered the problem and found a different proof (essentially based on orthogonal polynomials) which generalizes most (but not all) of Akaike’s results from the case of $m = 1$ (method of steepest descent) to the case of $m > 1$ (m -dimensional optimum gradient method).

Drawing on Akaike’s ideas we investigate the first of the two questions mentioned above in Section 3. Under the assumption that A is Hermitian we shall see that, in the case of $m = 1$, the eigenvalues of \widehat{H}_k asymptotically alternate between two values, ρ_1^* and ρ_2^* . Our proofs rely solely on techniques from linear algebra and do not use any concepts from probability theory. We decided to sketch in addition Akaike’s original proof in Section 4 because his techniques are highly interesting and hardly known today: In almost any textbook, the convergence of the method of steepest descent is proven using Kantorovich’s inequality; see, e.g., [7, §70], [12, Theorem 5.35] or [15, §5.3.1]. Such a proof is short and elegant (and also gives the asymptotic rate of convergence, at least in a worst-case-sense) but does not reveal the peculiar way in which the errors tend to zero.

Having answered the first of the above two questions we shall attack the second in Section 5. We have to consider polynomial interpolation processes based, asymptotically, on just two nodes ρ_1^* and ρ_2^* repeated cyclically. We shall use Walsh’s theory [18, Chapter III] on the polynomial interpolation of analytic functions with finite singularities, which we complement by a convergence result for the interpolation of a class of entire functions. Putting the pieces together we shall see that, if A is Hermitian, the method of steepest descent for matrix functions converges (or diverges) geometrically when f has finite singularities, i.e., the error

in step k behaves asymptotically as θ^k , where θ is determined by the eigenvalues of A , the vector \mathbf{b} and the singularities of f . For the function $f(\lambda) = \exp(\tau\lambda)$, the errors behave asymptotically as $\theta^k \tau^k / k!$, where θ depends on the eigenvalues of A and on the vector \mathbf{b} , i.e., we observe superlinear convergence.

Finally, in Section 6 we show why it is so difficult to determine the precise values of the nodes ρ_1^* and ρ_2^* . For a simple example we reveal the complicated relationship between these nodes on the one hand and the eigenvalues of A and the components of the vector \mathbf{b} on the other.

2. Restart length one and the method of steepest descent for matrix functions. We consider a restarted Krylov subspace method for the approximation of $f(A)\mathbf{b}$ with shortest possible restart length, i.e., based on a succession of one-dimensional Krylov subspaces. The restarted Arnoldi method with unit restart length given in Algorithm 1 generates (generally non-orthogonal) bases of the sequence of Krylov spaces $\mathcal{K}_k(A, \mathbf{b})$, $k \leq L$, where L denotes the invariance index of this Krylov sequence. Note that for $m = 1$ restarting and truncating are equivalent and that this algorithm is therefore also an incomplete orthogonalization process with truncation parameter $m = 1$; see [15, §6.4.2].

Algorithm 1: Restarted Arnoldi process with unit restart length.

Given: A, \mathbf{b}
 $\sigma_1 := \|\mathbf{b}\|, \mathbf{v}_1 := \mathbf{b}/\sigma_1$
for $k = 1, 2, \dots$ **do**
 $\mathbf{w} := A\mathbf{v}_k$
 $\rho_k := \mathbf{v}_k^H \mathbf{w}$
 $\mathbf{w} := \mathbf{w} - \rho_k \mathbf{v}_k$
 $\sigma_{k+1} := \|\mathbf{w}\|$
 $\mathbf{v}_{k+1} := \mathbf{w}/\sigma_{k+1}$

Here and in the sequel, $\|\cdot\|$ denotes the Euclidean norm. Obviously, $\sigma_{k+1} = 0$ if and only if \mathbf{v}_k is an eigenvector of A . Since \mathbf{v}_k is a multiple of $(A - \rho_{k-1}I)\mathbf{v}_{k-1}$, this can only happen if already \mathbf{v}_{k-1} is an eigenvector of A and, by induction, if already the initial vector \mathbf{b} is an eigenvector of A . In this case, Algorithm 1 terminates in the first step and $f(A)\mathbf{b} = f(\rho_1)\mathbf{b} = \sigma_1 f(\rho_1)\mathbf{v}_1$. We may therefore assume that $\sigma_k > 0$ for all k .

Algorithm 1 generates the Arnoldi-like decomposition

$$(2.1) \quad AV_k = V_{k+1}\tilde{B}_k = V_k B_k + \sigma_{k+1}\mathbf{v}_{k+1}\mathbf{e}_k^T$$

with $V_k := [\mathbf{v}_1 \mathbf{v}_2 \dots \mathbf{v}_k] \in \mathbb{C}^{n \times k}$, the lower bidiagonal matrices

$$\tilde{B}_k := \begin{bmatrix} \rho_1 & & & & & \\ \sigma_2 & \rho_2 & & & & \\ & & \sigma_3 & \ddots & & \\ & & & \ddots & \rho_k & \\ & & & & & \sigma_{k+1} \end{bmatrix} \in \mathbb{C}^{(k+1) \times k}, \quad B_k := [I_k \mathbf{0}] \tilde{B}_k \in \mathbb{C}^{k \times k}$$

and $\mathbf{e}_k \in \mathbb{R}^k$ denoting the k -th unit coordinate vector. The matrices $B_k = B_k(A, \mathbf{b})$ will play a crucial role in our analysis where the following obvious invariance properties will be helpful.

LEMMA 2.1. *For the bidiagonal matrices $B_k = B_k(A, \mathbf{b})$ of (2.1) generated by Algorithm 1 with data A and \mathbf{b} , there holds:*

1.

$$B_k(\tau A, \mathbf{b}) = \begin{bmatrix} \tau\rho_1 & & & & & \\ |\tau|\sigma_2 & \tau\rho_2 & & & & \\ & & \ddots & & & \\ & & & \ddots & & \\ & & & & |\tau|\sigma_k & \tau\rho_k \end{bmatrix}, \quad \tau \neq 0.$$

In particular, for $\tau > 0$, there holds $B_k(\tau A, \mathbf{b}) = \tau B_k(A, \mathbf{b})$.

2. $B_k(A - \tau I, \mathbf{b}) = B_k(A, \mathbf{b}) - \tau I$ for $\tau \in \mathbb{C}$.

3. $B_k(Q^H A Q, Q^H \mathbf{b}) = B_k(A, \mathbf{b})$ for unitary $Q \in \mathbb{C}^{n \times n}$.

Given the Arnoldi-like decomposition (2.1) resulting from the restarted Arnoldi process with restart length one, the approximation of $f(A)\mathbf{b}$ is defined as

$$(2.2) \quad \mathbf{f}_k := \sigma_1 V_k f(B_k) \mathbf{e}_1, \quad k = 1, 2, \dots,$$

with $\mathbf{e}_1 \in \mathbb{R}^k$ denoting the first unit coordinate vector. We state as a first result an explicit representation of these approximants:

LEMMA 2.2. *Let Γ be a Jordan curve which encloses the field of values of A and thereby also $\rho_1, \rho_2, \dots, \rho_k$. Assume that f is analytic in the interior of Γ and extends continuously to Γ . For $r \in \mathbb{N}_0$ and $\ell \in \mathbb{N}$, we denote by*

$$\Delta_\ell^r f := \frac{1}{2\pi i} \int_\Gamma \frac{f(\zeta)}{(\zeta - \rho_\ell)(\zeta - \rho_{\ell+1}) \cdots (\zeta - \rho_{\ell+r})} d\zeta$$

the divided difference of f of order r with respect to the nodes $\rho_\ell, \rho_{\ell+1}, \dots, \rho_{\ell+r}$. Then

$$\mathbf{f}_k = \sum_{r=1}^k \left(\prod_{j=1}^r \sigma_j \right) (\Delta_1^{r-1} f) \mathbf{v}_r = \mathbf{f}_{k-1} + \left(\prod_{j=1}^k \sigma_j \right) (\Delta_1^{k-1} f) \mathbf{v}_k.$$

Proof. A short proof is obtained using a result of Opitz [13]: We have $\mathbf{f}_k = \sigma_1 V_k f(B_k) \mathbf{e}_1$ and Opitz showed that

$$f(B_k) = D \begin{bmatrix} \Delta_1^0 f & & & & \\ \Delta_1^1 f & \Delta_2^0 f & & & \\ \Delta_1^2 f & \Delta_2^1 f & \Delta_3^0 f & & \\ \vdots & \vdots & \vdots & \ddots & \\ \Delta_1^{k-1} f & \Delta_2^{k-2} f & \Delta_3^{k-3} f & \cdots & \Delta_k^0 f \end{bmatrix} D^{-1}$$

with $D := \text{diag} \left(1, \sigma_2, \prod_{j=2}^3 \sigma_j, \dots, \prod_{j=2}^k \sigma_j \right)$, from which the assertion follows immediately. \square

The following convergence result is another immediate consequence of the close connection between \mathbf{f}_k and certain interpolation processes; see [5, Theorem 4.3.1].

THEOREM 2.3. *Let $W(A) := \{\mathbf{v}^H A \mathbf{v} : \|\mathbf{v}\| = 1\}$ denote the field of values of A and let $\delta := \max_{\zeta, \eta \in W(A)} |\zeta - \eta|$ be its diameter. Let f be analytic in (a neighborhood of) $W(A)$ and let $\rho > 0$ be maximal such that f can be continued analytically to $W_\rho := \{\lambda \in \mathbb{C} : \min_{\zeta \in W(A)} |\lambda - \zeta| < \rho\}$. (If f is entire, we set $\rho = \infty$.)*

If $\rho > \delta$ then $\lim_{k \rightarrow \infty} \mathbf{f}_k = f(A)\mathbf{b}$ and this convergence is at least linear.

Proof. We choose $0 < \tau < \rho$ and a Jordan curve Γ such that $\tau \leq \min_{\lambda \in W(A)} |\zeta - \lambda| < \rho$ for every $\zeta \in \Gamma$. Hermite's representation of the interpolation error

$$(2.3) \quad f(\lambda) - p_{k-1}(\lambda) = \frac{1}{2\pi i} \int_{\Gamma} \frac{(\lambda - \rho_1)(\lambda - \rho_2) \cdots (\lambda - \rho_k)}{(\zeta - \rho_1)(\zeta - \rho_2) \cdots (\zeta - \rho_k)} \frac{f(\zeta)}{\zeta - \lambda} d\zeta$$

(see, e.g., [5, Theorem 3.6.1]) gives, for $\lambda \in W(A)$,

$$|f(\lambda) - p_{k-1}(\lambda)| \leq C_1 \left[\frac{\delta}{\tau} \right]^k,$$

with the constant $C_1 = \text{length}(\Gamma) \max_{\zeta \in \Gamma} |f(\zeta)| / [2\pi \min_{\zeta \in \Gamma, \lambda \in W(A)} |\zeta - \lambda|]$. The assertion follows from a result of Crouzeix [4], who showed that

$$\|f(A) - p_{k-1}(A)\| \leq C_2 \max_{\lambda \in W(A)} |f(\lambda) - p_{k-1}(\lambda)|,$$

with a constant $C_2 \leq 12$. \square

Note that Theorem 2.3 holds for Arnoldi approximations of arbitrary restart length (and also for its unrestarted variant). Note further that we always have superlinear convergence if f is an entire function; see also [6, Theorem 4.2].

We conclude this section by considering the specific function $f(\lambda) = 1/\lambda$. For a non-singular matrix A , computing $f(A)\mathbf{b}$ is nothing but solving the linear system $A\mathbf{x} = \mathbf{b}$. It is known (cf. [6, §4.1.1]) that the Arnoldi method with restart length $m = 1$ is equivalent to FOM(1) (restarted full orthogonalization method with restart length 1; see [15, §6.4.1]) as well as to IOM(1) (incomplete orthogonalization method with truncation parameter $m = 1$; see [15, §6.4.2]). If we choose $\mathbf{f}_0 = \mathbf{0}$ as the initial approximation and express the approximations \mathbf{f}_k in terms of the residual vectors $\mathbf{r}_k := \mathbf{b} - A\mathbf{f}_k$, there holds

$$\mathbf{f}_k = \mathbf{f}_{k-1} + (\sigma_1 \sigma_2 \cdots \sigma_k) (\Delta_1^{k-1} f) \mathbf{v}_k = \mathbf{f}_{k-1} + \alpha_k \mathbf{r}_{k-1},$$

where

$$\alpha_k = \frac{1}{\rho_k} = \frac{1}{\mathbf{v}_k^H A \mathbf{v}_k} = \frac{\mathbf{r}_{k-1}^H \mathbf{r}_{k-1}}{\mathbf{r}_{k-1}^H A \mathbf{r}_{k-1}},$$

which is known as the method of steepest descent, at least if A is Hermitian positive definite.

3. Asymptotics of B_k in the Hermitian case. The aim of this section is to show how the entries of the bidiagonal matrix B_k in (2.1) behave for large k .

We first consider a very special situation.

LEMMA 3.1. *For a Hermitian matrix $A \in \mathbb{C}^{n \times n}$, assume that \mathbf{b} and therefore \mathbf{v}_1 are linear combinations of two (orthonormal) eigenvectors of A ,*

$$\mathbf{v}_1 = \frac{1}{\sqrt{1 + |\gamma|^2}} \mathbf{z}_1 + \frac{\gamma}{\sqrt{1 + |\gamma|^2}} \mathbf{z}_2,$$

where $A\mathbf{z}_j = \lambda_j \mathbf{z}_j$ ($j = 1, 2$), $\lambda_1 < \lambda_2$ and $\|\mathbf{z}_1\| = \|\mathbf{z}_2\| = 1$. Then, for $k = 1, 2, \dots$, there holds

$$\begin{aligned} \mathbf{v}_{2k-1} &= \mathbf{v}_1 = \frac{1}{\sqrt{1 + |\gamma|^2}} \mathbf{z}_1 + \frac{\gamma}{\sqrt{1 + |\gamma|^2}} \mathbf{z}_2, \\ \mathbf{v}_{2k} &= \mathbf{v}_2 = -\frac{|\gamma|}{\sqrt{1 + |\gamma|^2}} \mathbf{z}_1 + \frac{\gamma}{|\gamma| \sqrt{1 + |\gamma|^2}} \mathbf{z}_2. \end{aligned}$$

Proof. A straightforward calculation shows

$$A\mathbf{v}_1 - \rho_1\mathbf{v}_1 = \frac{\lambda_1 - \lambda_2}{(1 + |\gamma|^2)^{3/2}} (|\gamma|^2\mathbf{z}_1 - \gamma\mathbf{z}_2)$$

and

$$\mathbf{v}_2 = \frac{A\mathbf{v}_1 - \rho_1\mathbf{v}_1}{\|A\mathbf{v}_1 - \rho_1\mathbf{v}_1\|} = \frac{-|\gamma|}{\sqrt{1 + |\gamma|^2}} \mathbf{z}_1 + \frac{\gamma}{|\gamma|\sqrt{1 + |\gamma|^2}} \mathbf{z}_2.$$

By the same token,

$$A\mathbf{v}_2 - \rho_2\mathbf{v}_2 = \frac{|\gamma|(\lambda_1 - \lambda_2)}{(1 + |\gamma|^2)^{3/2}} (\mathbf{z}_1 + \gamma\mathbf{z}_2),$$

and therefore

$$\mathbf{v}_3 = \frac{A\mathbf{v}_2 - \rho_2\mathbf{v}_2}{\|A\mathbf{v}_2 - \rho_2\mathbf{v}_2\|} = \frac{1}{\sqrt{1 + |\gamma|^2}} (\mathbf{z}_1 + \gamma\mathbf{z}_2) = \mathbf{v}_1. \quad \square$$

Another elementary calculation leads to the following result.

COROLLARY 3.2. *Under the assumptions of Lemma 3.1, the entries ρ_k and σ_{k+1} ($k = 1, 2, \dots$) of the bidiagonal matrices B_k are given by*

$$\begin{aligned} \rho_{2k-1} &= \theta\lambda_1 + (1 - \theta)\lambda_2, \\ \rho_{2k} &= (1 - \theta)\lambda_1 + \theta\lambda_2, \text{ and} \\ \sigma_{k+1} &= \sqrt{\theta(1 - \theta)} (\lambda_2 - \lambda_1), \end{aligned}$$

with $\theta := 1/(1 + |\gamma|^2)$.

In an asymptotic sense Corollary 3.2 covers the general case if A is Hermitian, which we shall assume throughout the remainder of this section.

THEOREM 3.3. *If A is Hermitian with extremal eigenvalues λ_{\min} and λ_{\max} and if the vector \mathbf{b} has nonzero components in the associated eigenvectors, then there is a real number $\theta \in (0, 1)$, which depends on the spectrum of A and on \mathbf{b} , such that the entries ρ_k and σ_{k+1} ($k = 1, 2, \dots$) of the bidiagonal matrices B_k in (2.1) satisfy*

$$\begin{aligned} \lim_{k \rightarrow \infty} \rho_{2k-1} &= \theta\lambda_{\min} + (1 - \theta)\lambda_{\max} =: \rho_1^*, \\ \lim_{k \rightarrow \infty} \rho_{2k} &= (1 - \theta)\lambda_{\min} + \theta\lambda_{\max} =: \rho_2^*, \\ \lim_{k \rightarrow \infty} \sigma_{k+1} &= \sqrt{\theta(1 - \theta)} (\lambda_{\max} - \lambda_{\min}) =: \sigma^*. \end{aligned}$$

The proof of this result will be broken down into the following three lemmas. It simplifies if we assume that A has only simple eigenvalues,

$$\lambda_1 < \lambda_2 < \dots < \lambda_n, \quad n \geq 2,$$

otherwise we replace A by $A|_{\mathcal{K}_L(A, \mathbf{b})}$. By $\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_n$ we denote corresponding normalized eigenvectors: $A\mathbf{z}_j = \lambda_j\mathbf{z}_j$, $\|\mathbf{z}_j\| = 1$. We also assume, again without loss of generality, that the vector \mathbf{b} and therefore \mathbf{v}_1 have nonzero components in all eigenvectors: $\mathbf{z}_j^H \mathbf{b} \neq 0$ for $j = 1, 2, \dots, n$. Next, we may assume that A is diagonal (otherwise we

replace A by $Q^H A Q$ and \mathbf{b} by $Q^H \mathbf{b}$, where $Q = [z_1, z_2, \dots, z_n]$; cf. Lemma 2.1). Finally, we assume that $\mathbf{b} = [b_1, b_2, \dots, b_n]^T$ is real. (If not, we replace \mathbf{b} by $Q^H \mathbf{b}$, where $Q = \text{diag}(b_1/|b_1|, b_2/|b_2|, \dots, b_n/|b_n|)$ is a diagonal unitary matrix. Note that $Q^H A Q = A$ if A is diagonal.)

In summary, for Hermitian A we may assume that A is a real diagonal matrix with pairwise distinct diagonal entries and that \mathbf{b} is a real vector with nonzero entries.

LEMMA 3.4. *The sequence $\{\sigma_{k+1}\}_{k \in \mathbb{N}}$ of the subdiagonal entries of B_k is bounded and nondecreasing and thus convergent. Moreover, $\sigma_{k+1} = \sigma_{k+2}$ if and only if \mathbf{v}_k and \mathbf{v}_{k+2} are linearly dependent.*

Proof. Boundedness of the sequence $\{\sigma_{k+1}\}_{k \in \mathbb{N}}$ follows easily via

$$0 \leq \sigma_{k+1} = \|(A - \rho_k I) \mathbf{v}_k\| \leq \|A - \rho_k I\| \leq \|A\| + |\rho_k| \leq 2\|A\|.$$

Monotonicity is shown as follows:

$$\begin{aligned} \sigma_{k+1} &= \|(A - \rho_k I) \mathbf{v}_k\| && \text{since } \|\mathbf{v}_{k+1}\| = 1 \\ &= \|\mathbf{v}_{k+1}\| \|(A - \rho_k I) \mathbf{v}_k\| && \text{since } \sigma_{k+1} \mathbf{v}_{k+1} = (A - \rho_k I) \mathbf{v}_k \\ &= |\mathbf{v}_{k+1}^H (A - \rho_k I) \mathbf{v}_k| && \text{since } A \text{ is Hermitian} \\ &= |\mathbf{v}_{k+1}^H (A - \rho_k I)^H \mathbf{v}_k| && \\ &= |\mathbf{v}_{k+1}^H (A - \rho_{k+1} I + (\rho_{k+1} - \rho_k) I)^H \mathbf{v}_k| && \\ &= |\mathbf{v}_{k+1}^H (A - \rho_{k+1} I)^H \mathbf{v}_k + (\rho_{k+1} - \rho_k) \mathbf{v}_{k+1}^H \mathbf{v}_k| && \\ &= |\mathbf{v}_{k+1}^H (A - \rho_{k+1} I)^H \mathbf{v}_k| && \text{since } \mathbf{v}_k \perp \mathbf{v}_{k+1} \\ &\leq \|(A - \rho_{k+1} I) \mathbf{v}_{k+1}\| \|\mathbf{v}_k\| && \text{by the Cauchy-Schwarz inequality} \\ &= \|(A - \rho_{k+1} I) \mathbf{v}_{k+1}\| && \text{since } \|\mathbf{v}_k\| = 1 \\ &= \sigma_{k+2}. \end{aligned}$$

Equality holds if and only if \mathbf{v}_k and $(A - \rho_{k+1} I) \mathbf{v}_{k+1} = \sigma_{k+2} \mathbf{v}_{k+2}$ are linearly dependent. \square

LEMMA 3.5. *Every accumulation point of the vector sequence $\{\mathbf{v}_k\}_{k \in \mathbb{N}}$ generated by Algorithm 1 is contained in $\text{span}\{z_1, z_n\}$, i.e., in the linear hull of the eigenvectors of A associated with its extremal eigenvalues.*

Proof. By the compactness of the unit sphere in \mathbb{C}^n , the sequence of unit vectors $\{\mathbf{v}_k\}_{k \in \mathbb{N}}$ must have at least one point of accumulation. Each such accumulation point is the limit of a subsequence $\{\mathbf{v}_{k_\nu}\}$, for which, by Lemma 3.4, the associated sequence $\{\sigma_{k_\nu+1}\}$ converges, and we denote its limit by σ^* . We conclude that for each accumulation point \mathbf{u}_1 there holds $\sigma_1 = \|A \mathbf{u}_1 - (\mathbf{u}_1^H A \mathbf{u}_1) \mathbf{u}_1\| = \sigma^*$. Furthermore, one step of Algorithm 1 starting with an accumulation point \mathbf{u}_1 as the initial vector yields another accumulation point \mathbf{u}_2 , and therefore also $\sigma_2 = \|A \mathbf{u}_2 - (\mathbf{u}_2^H A \mathbf{u}_2) \mathbf{u}_2\| = \sigma^*$. Two steps of Algorithm 1 with initial vector \mathbf{u}_1 thus result in the decomposition

$$A [\mathbf{u}_1 \ \mathbf{u}_2] = [\mathbf{u}_1 \ \mathbf{u}_2 \ \mathbf{u}_3] \begin{bmatrix} \mathbf{u}_1^H A \mathbf{u}_1 & 0 \\ \sigma^* & \mathbf{u}_2^H A \mathbf{u}_2 \\ 0 & \sigma^* \end{bmatrix},$$

and from the fact that $\sigma_2 = \sigma_3 = \sigma^*$ we conclude using Lemma 3.4 that \mathbf{u}_1 and \mathbf{u}_3 must be linearly dependent, i.e., $\mathbf{u}_3 = \kappa \mathbf{u}_1$ for some κ . We thus obtain

$$A [\mathbf{u}_1 \ \mathbf{u}_2] = [\mathbf{u}_1 \ \mathbf{u}_2] A_2 \quad \text{with } A_2 := \begin{bmatrix} \mathbf{u}_1^H A \mathbf{u}_1 & \kappa \sigma^* \\ \sigma^* & \mathbf{u}_2^H A \mathbf{u}_2 \end{bmatrix},$$

which means that $\text{span}\{\mathbf{u}_1, \mathbf{u}_2\}$ is an A -invariant subspace, which in turn is only possible if $\text{span}\{\mathbf{u}_1, \mathbf{u}_2\} = \text{span}\{z_\ell, z_m\}$ for some $1 \leq \ell < m \leq n$.

We note in passing that $A_2 = [\mathbf{u}_1 \ \mathbf{u}_2]^H A [\mathbf{u}_1 \ \mathbf{u}_2]$ must be Hermitian—in fact, real and symmetric—and that consequently $\kappa = 1$ and $\mathbf{u}_3 = \mathbf{u}_1$; cf. Lemma 3.1.

Expanding the vectors $\mathbf{v}_k = \sum_{j=1}^n \gamma_{j,k} \mathbf{z}_j$ generated by Algorithm 1 in the orthonormal eigenbasis of A , we note that, by our assumption that the initial vector not be deficient in any eigenvector, there holds $\gamma_{j,1} \neq 0$ for all $j = 1, 2, \dots, n$. In addition, since $\rho_k \in (\lambda_1, \lambda_n)$ and $\gamma_{j,k+1} = \gamma_{j,k}(\lambda_j - \rho_k)/\sigma_{k+1}$, we see that $\gamma_{1,k}$ and $\gamma_{n,k}$ are both nonzero for all k . For the interior eigenvalues $\{\lambda_j\}_{j=2}^{n-1}$ it may well happen that $\rho_{k_0} = \lambda_j$ for some k_0 (cf. Section 6 for examples), whereupon subsequent vectors of the sequence $\{\mathbf{v}_k\}$ will be deficient in \mathbf{z}_j , i.e., $\gamma_{j,k} = 0$ for all $k > k_0$. It follows that, for the sequence considered above starting with an accumulation point \mathbf{u}_1 , $\gamma_{m,k}$ and $\gamma_{\ell,k}$ must also be nonzero for all k .

Assume now that $m < n$ and consider a subsequence $\{\mathbf{v}_{k_\nu}\}$ converging to \mathbf{u}_1 (without loss of generality, $\mathbf{v}_{k_1} = \mathbf{v}_1$). For $\nu \rightarrow \infty$ the Rayleigh quotients ρ_{k_ν} , being continuous functions of the vectors \mathbf{v}_{k_ν} , then converge to a limit contained in $(\lambda_\ell, \lambda_m)$. Consequently, $\lambda_n - \rho_{k_\nu} > \lambda_m - \rho_{k_\nu} > 0$ for all sufficiently large ν . Since $\mathbf{z}_n^H \mathbf{u}_1 = 0$ by assumption, we further have

$$0 = \lim_{\nu \rightarrow \infty} \left| \frac{\gamma_{n,k_\nu}}{\gamma_{m,k_\nu}} \right| = \left| \frac{\gamma_{n,1}}{\gamma_{m,1}} \right| \lim_{\nu \rightarrow \infty} \prod_{\eta=1}^{\nu} \frac{|\lambda_n - \rho_{k_\eta}|}{|\lambda_m - \rho_{k_\eta}|}.$$

But this is impossible since none of the factors on the right-hand side is zero and $|\lambda_n - \rho_{k_\nu}|/|\lambda_m - \rho_{k_\nu}| > 1$ for all sufficiently large ν . In a similar way, the assumption $1 < \ell$ is also found to lead to a contradiction. \square

LEMMA 3.6. *For the vector sequence $\{\mathbf{v}_k\}_{k \in \mathbb{N}}$ of Algorithm 1 there exist nonzero real numbers α and β , $\alpha^2 + \beta^2 = 1$, which depend on the spectrum of A and on \mathbf{b} , such that*

$$\lim_{k \rightarrow \infty} \mathbf{v}_{2k-1} = \alpha \mathbf{z}_1 + \beta \mathbf{z}_n \quad \text{and} \quad \lim_{k \rightarrow \infty} \mathbf{v}_{2k} = \text{sign}(\alpha\beta)[- \beta \mathbf{z}_1 + \alpha \mathbf{z}_n],$$

where $\text{sign}(\lambda)$ denotes the sign of the real number λ .

Proof. We count the candidates for accumulation points \mathbf{u} of the sequence $\{\mathbf{v}_k\}$. By Lemma 3.5, $\mathbf{u} \in \text{span}\{\mathbf{z}_1, \mathbf{z}_n\}$ and, since $\|\mathbf{u}\| = 1$, every accumulation point can be written as $\mathbf{u} = \alpha \mathbf{z}_1 + \beta \mathbf{z}_n$ with $\alpha^2 + \beta^2 = 1$. For every vector of this form, there holds

$$\|A\mathbf{u} - (\mathbf{u}^H A \mathbf{u})\mathbf{u}\|^2 = \alpha^2 \beta^2 (\lambda_n - \lambda_1)^2 = \alpha^2 (1 - \alpha^2) (\lambda_n - \lambda_1)^2.$$

Since \mathbf{u} is an accumulation point of the sequence $\{\mathbf{v}_k\}$, we have, as in the proof of Lemma 3.4, $\|A\mathbf{u} - (\mathbf{u}^H A \mathbf{u})\mathbf{u}\| = \sigma^*$, i.e.,

$$\alpha^2 (1 - \alpha^2) = \left(\frac{\sigma^*}{\lambda_n - \lambda_1} \right)^2.$$

This equation has at most four solutions α which shows that there are at most eight points of accumulation.

Assume now that \mathbf{v}_k is sufficiently close to such an accumulation point $\mathbf{u} = \mathbf{u}_1 = \alpha \mathbf{z}_1 + \beta \mathbf{z}_n$. Since all operations in Algorithm 1 are continuous, \mathbf{v}_{k+1} , for k sufficiently large, will be arbitrarily close to

$$\mathbf{u}_2 = \frac{A - (\mathbf{u}_1^H A \mathbf{u}_1)\mathbf{u}_1}{\|A - (\mathbf{u}_1^H A \mathbf{u}_1)\mathbf{u}_1\|} = \text{sign}(\alpha\beta)[- \beta \mathbf{z}_1 + \alpha \mathbf{z}_n]$$

(which is also an accumulation point of $\{\mathbf{v}_k\}$ different from \mathbf{u}_1 since $\alpha\beta \neq 0$). Moreover, \mathbf{v}_{k+2} must then be close to

$$\mathbf{u}_3 = \frac{A - (\mathbf{u}_2^H A \mathbf{u}_2)\mathbf{u}_2}{\|A - (\mathbf{u}_2^H A \mathbf{u}_2)\mathbf{u}_2\|} = \alpha \mathbf{z}_1 + \beta \mathbf{z}_n = \mathbf{u}_1.$$

Since we already know there are only finitely many accumulation points of $\{v_k\}$, we conclude that the sequence $\{v_k\}$ must asymptotically alternate between u_1 and u_2 . \square

The assertion of Theorem 3.3 now follows by elementary calculations, e.g.,

$$\begin{aligned} \lim_{k \rightarrow \infty} \rho_{2k-1} &= \lim_{k \rightarrow \infty} v_{2k-1}^H A v_{2k-1} = u_1^H A u_1 = (\alpha z_1 + \beta z_n)^H A (\alpha z_1 + \beta z_n) \\ &= \alpha^2 \lambda_{\min} + \beta^2 \lambda_{\max} = \theta \lambda_{\min} + (1 - \theta) \lambda_{\max}, \end{aligned}$$

where $\theta := \alpha^2$.

4. Akaike's probability theory setting. Theorem 3.3, the main result of the previous section, is implicitly contained in Akaike's paper [2] from 1959. His proof is based on the analysis of a transformation of probability measures: As is well-known (see, e.g., [17]) a Hermitian matrix $A \in \mathbb{C}^{n \times n}$ and any vector $v \in \mathbb{C}^n$ of unit length give rise to a probability measure μ on \mathbb{R} , assigning to any set $M \subset \mathbb{R}$ the measure

$$(4.1) \quad \mu(M) := \int_M w(\lambda) d\lambda, \quad w(\lambda) := \sum_{j=1}^n \omega_j^2 \delta(\lambda - \lambda_j),$$

where δ denotes the Dirac δ -function, $\lambda_1 < \lambda_2 < \dots < \lambda_n$ are the eigenvalues of A (we assume again without loss of generality that A has n simple eigenvalues) with corresponding eigenvectors z_j , $\|z_j\| = 1$, $j = 1, 2, \dots, n$, and where the weights are given by $\omega_j^2 = |z_j^H v|^2$. For a fixed matrix A , this correspondence between unit vectors v and probability measures μ supported on $\Lambda(A)$ is essentially one-to-one (if we do not distinguish between vectors $v = [v_1, v_2, \dots, v_n]^T$ and $w = [w_1, w_2, \dots, w_n]^T$ with $|v_j| = |w_j|$ for all $j = 1, 2, \dots, n$). In this way, each basis vector v_k generated by the restarted Arnoldi process with unit restart length (Algorithm 1) induces a probability measure μ_k whose support is a subset of $\Lambda(A)$.

The Lebesgue integral associated with (4.1) is given by

$$\int_{\mathbb{R}} f(\lambda) w(\lambda) d\lambda = \sum_{j=1}^n \omega_j^2 f(\lambda_j)$$

for any function f defined on $\Lambda(A)$. In particular, the mean of μ ,

$$\rho_\mu := \int_{\mathbb{R}} \lambda w(\lambda) d\lambda = \sum_{j=1}^n \omega_j^2 \lambda_j = v^H A v,$$

is the Rayleigh quotient of v and A , and the variance of μ is given by

$$\sigma_\mu^2 := \int_{\mathbb{R}} (\lambda - \rho_\mu)^2 w(\lambda) d\lambda = \sum_{j=1}^n \omega_j^2 (\lambda_j - \rho_\mu)^2 = \|(A - \rho_\mu I)v\|^2,$$

the squared norm of the vector $Av - (v^H Av)v$. We now see that the (nonlinear) vector transformation

$$v_{k+1} = T v_k, \quad \text{where } T v := \frac{Av - (v^H Av)v}{\|Av - (v^H Av)v\|},$$

underlying Algorithm 1 can be rephrased as a transformation of probability measures, $\mu_{k+1} = T \mu_k$, where

$$(4.2) \quad (T \mu)(M) := \int_M \left[\frac{\lambda - \rho_\mu}{\sigma_\mu} \right]^2 w(\lambda) d\lambda, \quad \text{if } \mu(M) = \int_M w(\lambda) d\lambda.$$

As above, we assume that \mathbf{v}_1 and thus \mathbf{v}_k , $k \geq 1$, is not an eigenvector of A , which implies that the support of μ_k consists of more than one point and therefore $\sigma_k = \sigma_{\mu_k} > 0$. We also remark that the transformation (4.2) $\mu \mapsto T\mu$ is not only well-defined for probability measures with finite support but for any probability measure whose first and second moments exist.

The crucial points in the proof of Theorem 3.3 were to show that the subdiagonal entries σ_k of B_k , which we have now identified as the standard deviations of μ_k , form a nondecreasing sequence (see Lemma 3.4) and that $\sigma_{k+1} = \sigma_k$ can only hold if \mathbf{v}_k is a linear combination of two eigenvectors of A ; see Lemma 3.5. Akaike based his proof on explicit formulas for the mean and variance of the transformed measure $T\mu$:

$$\rho_{T\mu} = \rho_\mu + \frac{1}{\sigma_\mu^2} \int_{\mathbb{R}} (\lambda - \rho_\mu)^3 w(\lambda) d\lambda$$

(cf. [2, Lemma 1]) and

$$\sigma_{T\mu}^2 = \sigma_\mu^2 + \frac{1}{\sigma_\mu^4} \det(M_3), \quad \text{where } M_3 := \left[\int_{\mathbb{R}} (\lambda - \rho_\mu)^{k+j} w(\lambda) d\lambda \right]_{0 \leq k, j \leq 2}$$

is the (3×3) -moment matrix associated with $(f, g) = \int_{\mathbb{R}} f(\lambda) \overline{g(\lambda)} w(\lambda) d\lambda$; cf. [2, Lemma 2]. Since M_3 is positive semidefinite it follows that $\sigma_{T\mu}^2 \geq \sigma_\mu^2$, with equality holding if and only if M_3 is singular, which can only happen if the support of μ consists of two points or less.

5. Convergence for functions of Hermitian matrices. As mentioned previously our convergence analysis is based on the close connection between Krylov subspace methods for approximating $f(A)\mathbf{b}$ and polynomial interpolation; see, e.g., [6, Theorem 2.4]. For the vectors \mathbf{f}_k of (2.2), we have

$$\mathbf{f}_k = \sigma_1 V_k f(B_k) \mathbf{e}_1 = p_{k-1}(A) \mathbf{b},$$

where $p_{k-1} \in \mathcal{P}_{k-1}$ interpolates f in the Hermite sense at the Rayleigh quotients $\rho_j = \mathbf{v}_j^H A \mathbf{v}_j$ ($j = 1, 2, \dots, k$). If A is Hermitian there holds (see Theorem 3.3)

$$(5.1) \quad \lim_{k \rightarrow \infty} \rho_{2k-1} = \rho_1^* \quad \text{and} \quad \lim_{k \rightarrow \infty} \rho_{2k} = \rho_2^*,$$

with two numbers ρ_1^* and ρ_2^* both contained in the convex hull of $\Lambda(A)$. In other words, asymptotically, the restarted Arnoldi approximation of $f(A)\mathbf{b}$ with unit restart length is equivalent to interpolating f at just the two nodes ρ_1^* and ρ_2^* with increasing multiplicity. Interpolation processes of such simple nature are well understood. To formulate the convergence results we need additional terminology: For $\delta \geq 0$, we define the curves

$$(5.2) \quad \Gamma_\delta := \{\lambda \in \mathbb{C} : |\lambda - \rho_1^*| |\lambda - \rho_2^*| = \delta^2\},$$

known as *lemniscates*[†] with foci ρ_1^* and ρ_2^* . If $\rho_1^* = \rho_2^*$ these reduce to concentric circles of radius δ . Otherwise, if $0 < \delta < \delta_0 := |\rho_1^* - \rho_2^*|/2$, Γ_δ consists of two mutually exterior analytic Jordan curves. When $\delta = \delta_0$, we obtain what is known as the *Bernoulli lemniscate*, for which these curves touch at $(\rho_1^* + \rho_2^*)/2$, whereas for $\delta > \delta_0$ the lemniscate is a single analytic Jordan curve. Obviously, its interior $\text{int } \Gamma_\delta$ contains ρ_1^* and ρ_2^* for every $\delta > 0$, $\Gamma_\gamma \subset \text{int } \Gamma_\delta$ for $0 \leq \gamma < \delta$, and every $\lambda \in \mathbb{C}$ is located on exactly one Γ_δ .

[†]Lemniscates of polynomials of degree 2 are also known as *Ovals of Cassini*.

We first assume that f is analytic in (an open set containing) ρ_1^* and ρ_2^* but not entire, i.e., that it has finite singularities in the complex plane. There exists thus a unique $\delta_f > 0$ such that f is analytic in $\text{int } \Gamma_{\delta_f}$ but not in $\text{int } \Gamma_\delta$ for any $\delta > \delta_f$,

$$(5.3) \quad \delta_f := \max \{ \delta : f \text{ is analytic in } \text{int } \Gamma_\delta \}.$$

THEOREM 5.1 (Walsh [18, Theorems 3.4 and 3.6]). *Let the sequence ρ_1, ρ_2, \dots be asymptotic to the sequence $\rho_1^*, \rho_2^*, \rho_1^*, \rho_2^*, \dots$ in the sense of (5.1). Assume that f is defined in all nodes ρ_1, ρ_2, \dots and let $p_{k-1} \in \mathcal{P}_{k-1}$ be the polynomial which interpolates f at $\rho_1, \rho_2, \dots, \rho_k$. Then $\lim_{k \rightarrow \infty} p_{k-1} = f$ uniformly on compact subsets of $\text{int } \Gamma_{\delta_f}$. More precisely, there holds*

$$\limsup_{k \rightarrow \infty} |f(\lambda) - p_{k-1}(\lambda)|^{1/k} \leq \frac{\delta}{\delta_f} \quad \text{for } \lambda \in \overline{\text{int } \Gamma_\delta}.$$

For $\lambda \notin \overline{\text{int } \Gamma_{\delta_f}}$ the sequence $\{p_{k-1}(\lambda)\}_{k \geq 1}$ diverges (unless λ is one of the nodes ρ_j).

It remains to investigate the convergence of the interpolation polynomials if f is an entire function. We here concentrate on $f(\lambda) = \exp(\tau\lambda)$, $\tau \neq 0$, which among entire functions is of the most practical interest. We remark, however, that the following argument applies to all entire functions of order 1 and type $|\tau|$ and can be generalized to arbitrary entire functions. We further note that the following theorem could be easily deduced from more general results of Winiarski [19] or Rice [14], but we prefer to present an elementary and self-contained proof.

THEOREM 5.2. *Let the sequence ρ_1, ρ_2, \dots satisfy the assumptions of Theorem 5.1 and let $p_{k-1} \in \mathcal{P}_{k-1}$ be the polynomial which interpolates $f(\lambda) = \exp(\tau\lambda)$ at $\rho_1, \rho_2, \dots, \rho_k$. Then $\{p_{k-1}\}$ converges to f uniformly on compact subsets of \mathbb{C} . More precisely, there holds*

$$\limsup_{k \rightarrow \infty} \left[k |f(\lambda) - p_{k-1}(\lambda)|^{1/k} \right] \leq \delta |\tau| e \quad \text{for } \lambda \in \overline{\text{int } \Gamma_\delta},$$

where $e = \exp(1)$.

Proof. We first interpolate $f(\lambda) = \exp(\tau\lambda)$ at the nodes ρ_1^* and ρ_2^* repeated cyclically, i.e., at $\rho_1^*, \rho_2^*, \rho_1^*, \rho_2^*, \rho_1^*, \dots$. By $p_{k-1}^* \in \mathcal{P}_{k-1}$ we denote the polynomial which interpolates f at the first k points of this sequence, and by $q_k^* \in \mathcal{P}_k$ the corresponding nodal polynomial. For $\lambda \in \overline{\text{int } \Gamma_\delta}$, Hermite's error formula (2.3) implies

$$f(\lambda) - p_{k-1}^*(\lambda) = \frac{1}{2\pi i} \int_{\Gamma_\eta} \frac{q_k^*(\lambda)}{q_k^*(\zeta)} \frac{\exp(\tau\zeta)}{\zeta - \lambda} d\zeta = \sum_{j=0}^{\infty} \frac{\tau^j}{j!} \frac{1}{2\pi i} \int_{\Gamma_\eta} \frac{q_k^*(\lambda)}{q_k^*(\zeta)} \frac{\zeta^j}{\zeta - \lambda} d\zeta,$$

where $\eta > \delta$. Note that $\int_{\Gamma_\eta} \frac{q_k^*(\lambda)}{q_k^*(\zeta)} \frac{\zeta^j}{\zeta - \lambda} d\zeta$ is the interpolation error for the function $g(\lambda) = \lambda^j$ which vanishes for $j = 0, 1, \dots, k-1$. Hence,

$$\begin{aligned} f(\lambda) - p_{k-1}^*(\lambda) &= \sum_{j=k}^{\infty} \frac{\tau^j}{j!} \frac{1}{2\pi i} \int_{\Gamma_\delta} \frac{q_k^*(\lambda)}{q_k^*(\zeta)} \frac{\zeta^j}{\zeta - \lambda} d\zeta \\ &= q_k^*(\lambda) \frac{\tau^k}{k!} \sum_{j=0}^{\infty} \frac{\tau^j k!}{(k+j)!} \frac{1}{2\pi i} \int_{\Gamma_\delta} \frac{1}{q_k^*(\zeta)} \frac{\zeta^{k+j}}{\zeta - \lambda} d\zeta \end{aligned}$$

and therefore,

$$|f(\lambda) - p_{k-1}^*(\lambda)| \leq |q_k^*(\lambda)| \frac{|\tau|^k}{k!} \sum_{j=k}^{\infty} \frac{|\tau|^j}{j!} \frac{1}{2\pi} \frac{\text{length}(\Gamma_\eta)}{\text{dist}(\Gamma_\delta, \Gamma_\eta)} \left[\max_{\zeta \in \Gamma_\eta} |\zeta| \right]^j \left[\max_{\zeta \in \Gamma_\eta} \frac{|\zeta|^k}{|q_k^*(\zeta)|} \right],$$

where we used $k!/(k+j)! \leq 1/j!$. Assume that k is even, then $q_k^*(\lambda) = [(\lambda - \rho_1^*)(\lambda - \rho_2^*)]^{k/2}$ and thus $|q_k^*(\lambda)| \leq \delta^k$. We further set $C_1 := \max_{\zeta \in \Gamma_\eta} |\zeta| \sim \eta$ (for $\eta \rightarrow \infty$), $C_2 := \max_{\zeta \in \Gamma_\eta} \frac{|\zeta|^2}{|\zeta - \rho_1^*| |\zeta - \rho_2^*|} \sim 1$ (for $\eta \rightarrow \infty$) and $C_3 := \frac{1}{2\pi} \frac{\text{length}(\Gamma_\eta)}{\text{dist}(\Gamma_\delta, \Gamma_\eta)} \sim 1$ (for $\eta \rightarrow \infty$). Now,

$$k! |f(\lambda) - p_{k-1}^*(\lambda)| \leq \delta^k |\tau|^k C_2^{k/2} C_3 \exp(|\tau| C_1).$$

Using Stirling's formula, $k! \sim \sqrt{2\pi k} (k/e)^k$ (for $k \rightarrow \infty$), and taking the k -th root we obtain

$$(5.4) \quad \limsup_{k \rightarrow \infty} \left[k |f(\lambda) - p_{k-1}^*(\lambda)|^{1/k} \right] \leq \delta |\tau| e \sqrt{C_2},$$

which is valid for every $\eta > \delta$. Since $C_2 \rightarrow 1$ for $\eta \rightarrow \infty$ we arrived at the desired conclusion, at least if the two nodes ρ_1^* and ρ_2^* are cyclically repeated. A minor modification shows that this inequality holds also for odd k .

It remains to show that (5.4) is valid if we interpolate $f(\lambda) = \exp(\tau\lambda)$ in nodes $\rho_1, \rho_2, \rho_3, \dots$ satisfying (5.1). We use again a result of Walsh [18, §3.5], who proved that

$$\lim_{k \rightarrow \infty} |(\lambda - \rho_1)(\lambda - \rho_2) \cdots (\lambda - \rho_k)|^{1/k} = \lim_{k \rightarrow \infty} |q_k^*(\lambda)|^{1/k}$$

uniformly on any compact set that does not contain one of the nodes $\rho_1, \rho_2, \rho_3, \dots$, which, together with (2.3), completes the proof. \square

Now all that remains is to translate the preceding interpolation results to the matrix setting. Introducing the quantity

$$\delta_A := \inf\{\delta : \Lambda(A) \subseteq \text{int } \Gamma_\delta\} = \max\{|(\lambda - \rho_1^*)(\lambda - \rho_2^*)|^{1/2} : \lambda \in \Lambda(A)\},$$

we are now in position to formulate our main result.

THEOREM 5.3. *Let A be Hermitian and let f denote a function which is analytic in a neighborhood of the spectral interval $[\lambda_{\min}, \lambda_{\max}]$ of A . For the approximants \mathbf{f}_k generated by the Arnoldi method with unit restart length and initial vector \mathbf{b} , there holds: If f possesses finite singularities, then*

$$\limsup_{k \rightarrow \infty} \|f(A)\mathbf{b} - \mathbf{f}_k\|^{1/k} \leq \frac{\delta_A}{\delta_f},$$

where δ_f is defined by (5.3).

If $f(\lambda) = \exp(\tau\lambda)$, $\tau \neq 0$, then

$$\limsup_{k \rightarrow \infty} \left[k \|f(A)\mathbf{b} - \mathbf{f}_k\|^{1/k} \right] \leq \delta_A |\tau| e.$$

Proof. Since A is Hermitian, i.e., normal, there holds

$$\|f(A)\mathbf{b} - \mathbf{f}_k\| \leq \max_{\lambda \in \Lambda(A)} |f(\lambda) - p_{k-1}(\lambda)| \|\mathbf{b}\|.$$

Now Theorems 5.1 and 5.2 imply the desired conclusion. \square

We next derive a necessary and sufficient condition for the convergence of the method of steepest descent. As before, we expand the k -th basis vector \mathbf{v}_k generated by the Arnoldi method with unit restart length in the orthonormal eigenbasis of A as $\mathbf{v}_k = \sum_{j=1}^n \gamma_{j,k} \mathbf{z}_j$. As noted already in the proof of Lemma 3.5, it is possible that at some index k_0 in Algorithm 1

the Rayleigh quotient ρ_{k_0} coincides with an eigenvalue λ_{j_0} ($2 \leq j_0 \leq n-1$). In this case $\gamma_{j_0, k+1} = \gamma_{j_0, k}(\lambda_j - \rho_k)/\sigma_{k+1} = 0$ for all $k > k_0$. But since

$$f(A)\mathbf{b} - \mathbf{f}_k = \tilde{f}(A)\mathbf{v}_{k+1} = \sum_{j=1}^n \tilde{f}(\lambda_j)\gamma_{j, k+1}\mathbf{z}_j, \quad \text{for some 'restart function' } \tilde{f}$$

(cf. [6, Theorem 2.6]), there follows $\mathbf{z}_{j_0}^H(f(A)\mathbf{b} - \mathbf{f}_k) = 0$ for all $k > k_0$ or, in other words, \mathbf{f}_k has no error component in the direction of \mathbf{z}_{j_0} .

Consider now an eigenvalue λ_{j_0} ($2 \leq j_0 \leq n-1$) with $\lambda_{j_0} \neq \rho_k$ for all k . The sequence

$$\left| \frac{\gamma_{j_0, k+2}}{\gamma_{n, k+2}} \right| = \left| \frac{\gamma_{j_0, k}}{\gamma_{n, k}} \right| \left| \frac{(\lambda_{j_0} - \rho_k)(\lambda_{j_0} - \rho_{k+1})}{(\lambda_n - \rho_k)(\lambda_n - \rho_{k+1})} \right|$$

tends to 0 for $k \rightarrow \infty$ (see Lemma 3.6), the second factor of the right-hand side tends to $|(\lambda_{j_0} - \rho_1^*)(\lambda_{j_0} - \rho_2^*)|/|(\lambda_n - \rho_1^*)(\lambda_n - \rho_2^*)|$. Consequently, we have

$$|(\lambda_{j_0} - \rho_1^*)(\lambda_{j_0} - \rho_2^*)| < |(\lambda_n - \rho_1^*)(\lambda_n - \rho_2^*)|,$$

i.e., the lemniscate Γ_{δ^*} , with

$$\delta^* := |(\lambda_n - \rho_1^*)(\lambda_n - \rho_2^*)|^{1/2} = |(\lambda_1 - \rho_1^*)(\lambda_1 - \rho_2^*)|^{1/2},$$

which passes through the extremal eigenvalues of A , contains all other eigenvalues in its interior (at least those which are relevant for the convergence of the steepest descent method).

THEOREM 5.4. *Denote by Γ_{δ^*} the lemniscate of the family (5.2) which passes through λ_{\min} and λ_{\max} . Then the method of steepest descent for computing $f(A)\mathbf{b}$ converges if and only if Γ_{δ^*} and its interior contain no singularity of f .*

We conclude this section by an obvious consequence.

COROLLARY 5.5. *Let f be a function analytic in $[\lambda_{\min}, \lambda_{\max}]$ which has no singularities in $\mathbb{C} \setminus \mathbb{R}$. Then the method of steepest descent for computing $f(A)\mathbf{b}$ converges. The convergence is at least linear with convergence factor*

$$\theta = \frac{\lambda_{\max} - \lambda_{\min}}{|\zeta_0 - \lambda_{\max}| + |\zeta_0 - \lambda_{\min}|},$$

where ζ_0 is a singularity of f closest to $[\lambda_{\min}, \lambda_{\max}]$.

Proof. Convergence follows from Theorem 5.4. Denoting the foci of the lemniscates Γ_{δ} (5.2) by $\rho_1 = \frac{1}{2}(\lambda_{\min} + \lambda_{\max}) - \gamma$ and $\rho_2 = \frac{1}{2}(\lambda_{\min} + \lambda_{\max}) + \gamma$, $\gamma \in [0, \frac{1}{2}(\lambda_{\max} - \lambda_{\min})]$, the convergence factor is given by

$$\left[\frac{|\lambda_{\max} - \rho_1^*||\lambda_{\max} - \rho_2^*|}{|\zeta_0 - \rho_1^*||\zeta_0 - \rho_2^*|} \right]^{1/2} = \left[\frac{(\lambda_{\max} - \lambda_{\min})^2 - 4\gamma^2}{(|\zeta_0 - \lambda_{\max}| + |\zeta_0 - \lambda_{\min}|)^2 - \gamma^2} \right]^{1/2}$$

which is a monotonically decreasing function of γ , i.e., it attains its maximal value for $\gamma = 0$. \square

Functions satisfying the assumptions of this corollary, such as e.g., $f(\lambda) = \log(\lambda)$, $f(\lambda) = \sqrt{\lambda}$ etc., play important roles in applications. Among them is also $f(\lambda) = 1/\lambda$ and, if we assume that A is positive (or negative) definite, then we regain the well-known result that the classical method of steepest descent converges with a convergence factor which is not greater than

$$\frac{\lambda_{\max} - \lambda_{\min}}{|\lambda_{\max}| + |\lambda_{\min}|} = \frac{\kappa - 1}{\kappa + 1},$$

where $\kappa = \lambda_{\max}/\lambda_{\min}$ is the condition number of A .

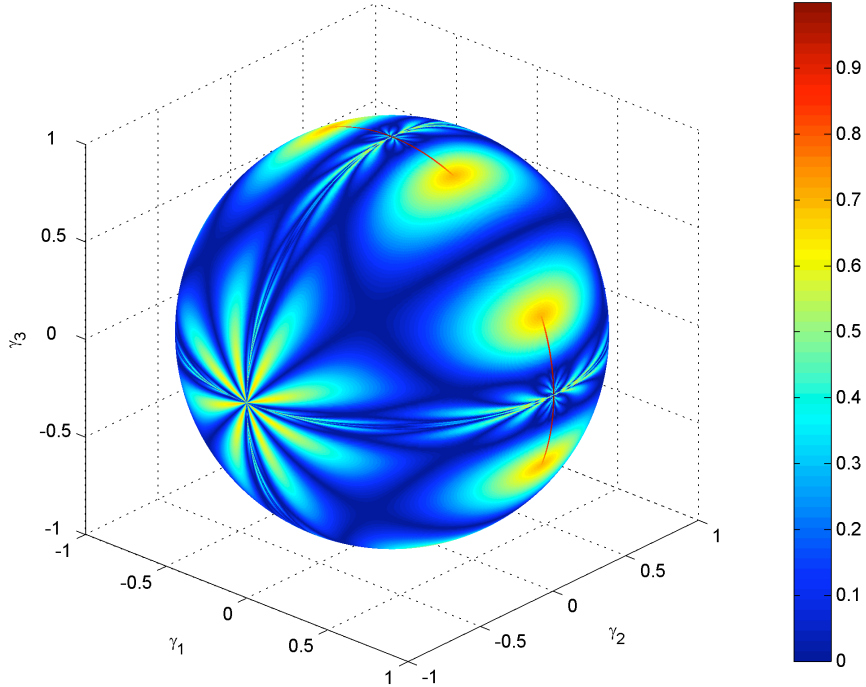


FIG. 6.1. The function $[\gamma_1, \gamma_2, \gamma_3] \mapsto \rho$.

6. Location of the foci. We have not been able to determine the exact location of the foci ρ_1^* and ρ_2^* . Of course, by Theorem 3.3 they are contained in the open interval $(\lambda_{\min}, \lambda_{\max})$ and lie symmetric to $\frac{1}{2}(\lambda_{\min} + \lambda_{\max})$. If $\frac{1}{2}(\lambda_{\min} + \lambda_{\max})$ is an eigenvalue of A (and if v_1 has a nonzero component in the corresponding eigenvector) then

$$|\rho_1^* - \rho_2^*| < \frac{1}{2}\sqrt{2}(\lambda_{\max} - \lambda_{\min})$$

because otherwise the lemniscate passing through λ_{\min} and λ_{\max} would not contain $\frac{1}{2}(\lambda_{\min} + \lambda_{\max})$ in its interior.

More precise information is only available in very special situations: Assume that $\Lambda(A)$ is symmetric with respect to $\frac{1}{2}(\lambda_{\min} + \lambda_{\max})$,

$$|\lambda_j - \frac{1}{2}(\lambda_{\min} + \lambda_{\max})| = |\lambda_{n+1-j} - \frac{1}{2}(\lambda_{\min} + \lambda_{\max})|$$

for $j = 1, 2, \dots, n/2$ if n is even and for $j = 1, 2, \dots, (n-1)/2$ if n is odd. In the latter case this means that $\lambda_{(n+1)/2} = \frac{1}{2}(\lambda_{\min} + \lambda_{\max})$. In addition, we require that the coefficients of $v_1 = \sum_{j=1}^n \gamma_{j,1} z_j$ are symmetric as well:

$$\gamma_{j,1} = \pm \gamma_{n+1-j,1}, \quad j = 1, 2, \dots, \lfloor n/2 \rfloor.$$

It is then easy to see that $\rho_k = \frac{1}{2}(\lambda_{\min} + \lambda_{\max})$ for every k and thus, $\rho_1^* = \rho_2^* = \frac{1}{2}(\lambda_{\min} + \lambda_{\max})$.

As a case study, we consider the fixed matrix $A = \text{diag}(-1, 0, 1)$ together with a real vector $v_1 = [\gamma_1, \gamma_2, \gamma_3]^T$, $\|v_1\| = 1$, $\gamma_1 \gamma_3 \neq 0$. The restarted Arnoldi process with unit

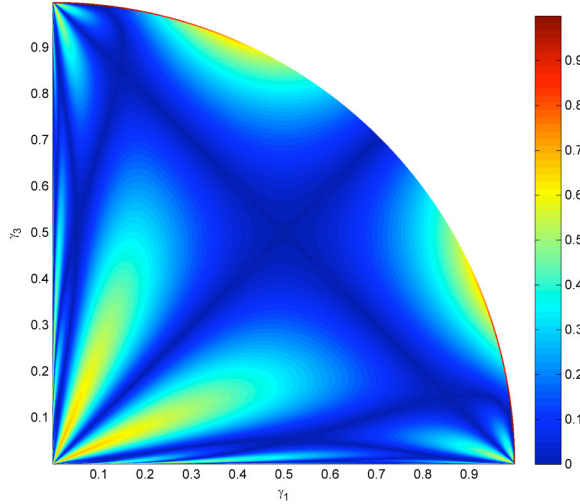


FIG. 6.2. The function $[\gamma_1, \gamma_3] \mapsto \rho$.

restart length (Algorithm 1) then generates a sequence $\{\mathbf{v}_k = [\gamma_{1,k}, \gamma_{2,k}, \gamma_{3,k}]^T\}_{k \geq 1}$ of unit vectors as follows,

$$(6.1) \quad \begin{bmatrix} \gamma_{1,k+1} \\ \gamma_{2,k+1} \\ \gamma_{3,k+1} \end{bmatrix} = T \begin{bmatrix} \gamma_{1,k} \\ \gamma_{2,k} \\ \gamma_{3,k} \end{bmatrix} := \frac{1}{\sigma_{k+1}} \begin{bmatrix} \gamma_{1,k}(-1 - \rho_k) \\ -\gamma_{2,k}\rho_k \\ \gamma_{3,k}(1 - \rho_k) \end{bmatrix},$$

with $\rho_k = -\gamma_{1,k}^2 + \gamma_{3,k}^2$, $\sigma_{k+1} = \sqrt{(\gamma_{1,k}^2 + \gamma_{3,k}^2) - (\gamma_{1,k}^2 - \gamma_{3,k}^2)^2}$ and the initial vector $[\gamma_{1,1}, \gamma_{2,1}, \gamma_{3,1}]^T := [\gamma_1, \gamma_2, \gamma_3]^T$. We know from Lemma 3.6 that

$$\lim_{k \rightarrow \infty} \mathbf{v}_{2k-1} = \begin{bmatrix} \alpha \\ 0 \\ \beta \end{bmatrix} \quad \text{and} \quad \lim_{k \rightarrow \infty} \mathbf{v}_{2k} = \text{sign}(\alpha\beta) \begin{bmatrix} -\beta \\ 0 \\ \alpha \end{bmatrix},$$

for some nonzero real numbers α, β , $\alpha^2 + \beta^2 = 1$, and that consequently (cf. Theorem 3.3)

$$\rho_1^* = \lim_{k \rightarrow \infty} \mathbf{v}_{2k-1}^T A \mathbf{v}_{2k-1} = -\alpha^2 + \beta^2 \quad \text{and} \quad \rho_2^* = \lim_{k \rightarrow \infty} \mathbf{v}_{2k}^T A \mathbf{v}_{2k} = \alpha^2 - \beta^2 = -\rho_1^*.$$

Denoting by $\rho = |\rho_1^*| = |\rho_2^*|$ the common modulus of these two nodes we are interested in the mapping $\rho = \rho(\gamma_1, \gamma_2, \gamma_3)$ which is defined on the unit sphere in \mathbb{R}^3 with the exception of the great circles $\gamma_1 = 0$ and $\gamma_3 = 0$. Figure 6.1 illustrates this function.

We first observe certain symmetries: Obviously the eight vectors $[\pm\gamma_1, \pm\gamma_2, \pm\gamma_3]^T$ lead to the same value of ρ . Moreover, we have $\rho(\gamma_1, \gamma_2, \gamma_3) = \rho(\gamma_3, \gamma_2, \gamma_1)$; see (6.1). The great circle $\gamma_2 = 0$ is of special interest: If we select $\mathbf{v}_1 = [\gamma_1, 0, \sqrt{1 - \gamma_1^2}]^T$ as the starting vector of the iteration (6.1), then $\mathbf{v}_{2k-1} = \mathbf{v}_1$ and $\mathbf{v}_{2k} = \mathbf{v}_2 = [-\sqrt{1 - \gamma_1^2}, 0, \gamma_1]^T$ for every k ; cf. Lemma 3.1. A simple computation yields

$$\rho = \rho(\gamma_1, 0, \sqrt{1 - \gamma_1^2}) = |\mathbf{v}_1^T A \mathbf{v}_1| = |1 - 2\gamma_1^2|.$$

Therefore, for a suitable choice of γ_1 , the function ρ attains every value in $[0, 1)$. Values of ρ contained in $(\sqrt{2}/2, 1)$ are attained if we select \mathbf{v}_1 on the ‘red subarcs’ of the great

circle $\gamma_2 = 0$; see Figure 6.1. Note that $\rho = \rho(\gamma_1, \gamma_2, \gamma_3) \in [0, \sqrt{2}/2)$ whenever $\gamma_2 \neq 0$. Consequently, ρ is discontinuous at every point of those arcs.

We next determine combinations of γ_1, γ_2 and γ_3 which lead to the value $\rho = 0$: If we start the iteration with $\mathbf{v}_1 = \sqrt{2}/2 [\pm 1, 0, \pm 1]^T$ then, for the Rayleigh quotients $\rho_k = \mathbf{v}_k^T A \mathbf{v}_k$, there holds $\rho_k = 0$ for all $k > 0$. We set $S_0 := \{\sqrt{2}/2 [\pm 1, 0, \pm 1]^T\}$. Now we define inductively the sets

$$S_\ell := \{\mathbf{v} : T\mathbf{v} \in S_{\ell-1}\}, \quad \ell = 1, 2, \dots,$$

and note that, for starting vectors $\mathbf{v}_1 \in S_\ell$, there holds $\rho_k = 0$ for all $k > \ell$.

To illustrate these sets in a more convenient way, we eliminate $\gamma_{2,1} = (1 - \gamma_{1,1}^2 - \gamma_{3,1}^2)^{1/2}$ from the transformation T defined in (6.1) and consider ρ as a function of the two variables γ_1 and γ_3 ; see Figure 6.2. For symmetry reasons we can restrict our attention to $0 < \gamma_1, \gamma_3 < 1$. The intersection of the sets S_ℓ and this restricted domain will be denoted by R_ℓ . We have

$$\begin{aligned} R_0 &= \{[\gamma_1, \gamma_3]^T : \gamma_1 = \gamma_3 = \sqrt{2}/2\}, \\ R_1 &= \{[\gamma_1, \gamma_3]^T : \gamma_3 = \gamma_1\}, \\ R_2 &= \{[\gamma_1, \gamma_3]^T : \gamma_3 = 1 - \gamma_1\}, \\ R_3 &= \{[\gamma_1, \gamma_3]^T : p(\gamma_1, \gamma_3) = 0\}, \\ &\vdots \\ &\vdots \end{aligned}$$

where $p(\gamma_1, \gamma_3) = \gamma_1^6 + \gamma_3^6 - \gamma_1^4 - \gamma_3^4 + 2\gamma_1^2\gamma_3^2 - \gamma_1^4\gamma_3^2 - \gamma_1^2\gamma_3^4 + 2\gamma_1^5\gamma_3 + 2\gamma_1\gamma_3^5 - 4\gamma_1^3\gamma_3^3 - 2\gamma_1\gamma_3$. Figure 6.3 shows these sets $R_\ell, \ell = 0, 1, \dots, 5$, where R_4 and R_5 were computed numerically using Newton's method.

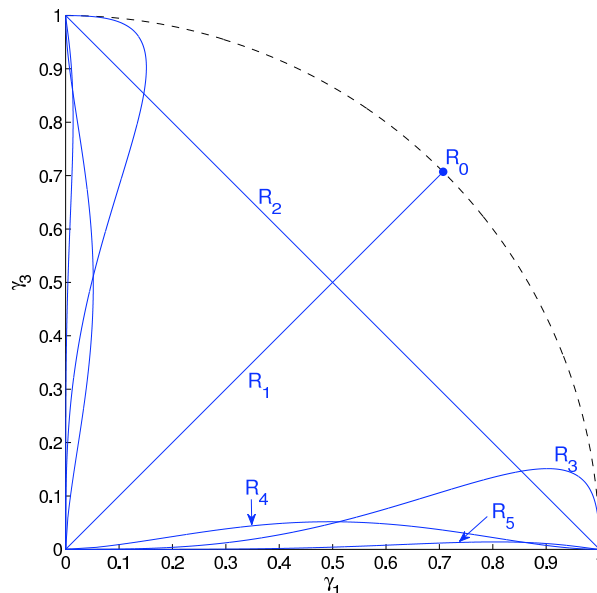


FIG. 6.3. The sets $R_\ell, \ell = 0, 1, \dots, 5$.

In summary, determining the foci ρ_1^* and ρ_2^* requires the analytic evaluation of the function $\rho = \rho(\gamma_1, \gamma_2, \gamma_3)$ which, even in the simple example considered here, appears intractable.

7. Conclusion. We have given a convergence analysis of the restarted Arnoldi approximation for functions of Hermitian matrices in the case when the restart length is one. The analysis is based on an earlier result of Akaike given in a probability theory setting, which we have translated into the terminology of linear algebra, and results of Walsh on the convergence of interpolation polynomials. In particular, we have shown that the restarted Arnoldi method exhibits, asymptotically, a two-periodic behavior. Moreover, we have characterized the asymptotic behavior of the entries of the associated Hessenberg matrix. The precise location of the asymptotic interpolation nodes is a complicated task, as was illustrated for a simple example. These results may be viewed as a first step towards understanding the asymptotic behavior of the restarted Arnoldi process.

Acknowledgments. We thank Ken Hayami of the National Institute of Informatics, Tokyo, for valuable comments, suggestions, and information.

REFERENCES

- [1] M. AFANASJEW, M. EIERMANN, O. G. ERNST, AND S. GÜTTEL, *Implementation of a restarted Krylov subspace method for the evaluation of matrix functions*, Linear Algebra Appl., (to appear).
- [2] H. AKAIKE, *On a successive transformation of probability distribution and its application to the analysis of the optimum gradient method*, Ann. Inst. Statist. Math. Tokio, 11 (1959), pp. 1–16.
- [3] A. CAUCHY, *Méthode générale pour la résolution des systèmes d'équations simultanées*, Comp. Rend. Sci. Paris, 25 (1847), pp. 536–538.
- [4] M. CROUZEIX, *Numerical range and functional calculus in Hilbert space*, J. Funct. Anal., 244 (2007), pp. 668–690.
- [5] P. J. DAVIS, *Interpolation and Approximation*, Dover Publications, Inc., New York, NY, 1975.
- [6] M. EIERMANN AND O. G. ERNST, *A restarted Krylov subspace method for the evaluation of matrix functions*, SIAM J. Numer. Anal., 44 (2006), pp. 2481–2504.
- [7] D. K. FADDEJEW AND W. N. FADDEJEW, *Numerische Methoden der Linearen Algebra*, vol. 10 of Mathematik für Naturwissenschaft und Technik, VEB Deutscher Verlag der Wissenschaften, Berlin, 3. ed., 1973.
- [8] A. I. FORSYTHE AND G. E. FORSYTHE, *Punched-card experiments with accelerated gradient methods for linear equations*, National Bureau of Standards, Appl. Math. Ser., 39 (1954), pp. 55–69.
- [9] G. E. FORSYTHE, *On the asymptotic directions of the s -dimensional optimum gradient method*, Numer. Math., 11 (1968), pp. 57–76.
- [10] G. E. FORSYTHE AND T. S. MOTZKIN, *Acceleration of the optimum gradient method (Abstract)*, Bull. Amer. Math. Soc., 57 (1951), pp. 304–305.
- [11] ———, *Asymptotic properties of the optimum gradient method (Abstract)*, Bull. Amer. Math. Soc., 57 (1951), p. 183.
- [12] G. MEURANT, *Computer Solution of Large Linear Systems*, vol. 28 of Studies in Mathematics and its Applications, Elsevier, Amsterdam, 1999.
- [13] G. OPITZ, *Steigungsmatrizen*, Z. Angew. Math. Mech., 44 (1964), pp. T52–T54.
- [14] J. R. RICE, *The degree of convergence for entire functions*, Duke Math. J., 38 (1971), pp. 429–440.
- [15] Y. SAAD, *Iterative Methods for Sparse Linear Systems*, SIAM, Philadelphia, PA, 2nd ed., 2003.
- [16] E. STIEFEL, *Über einige Methoden der Relaxationsrechnung*, Z. Angew. Math. Phys., 3 (1952), pp. 1–33.
- [17] ———, *Relaxationsmethoden bester Strategie zur Lösung linearer Gleichungssysteme*, Comment. Math. Helv., 29 (1955), pp. 157–179.
- [18] J. L. WALSH, *Interpolation and Approximation by Rational Functions in the Complex Domain*, vol. XX of American Mathematical Society Colloquium Publications, AMS, Providence, RI, 5th ed., 1969.
- [19] T. WINIARSKI, *Approximation and interpolation of entire functions*, Ann. Polon. Math., XXIII (1970), pp. 259–273.