

ON BEST APPROXIMATIONS OF POLYNOMIALS IN MATRICES IN THE MATRIX 2-NORM

JÖRG LIESEN* AND PETR TICHÝ†

Abstract. We show that certain matrix approximation problems in the matrix 2-norm have uniquely defined solutions, despite the lack of strict convexity of the matrix 2-norm. The problems we consider are generalizations of the ideal Arnoldi and ideal GMRES approximation problems introduced by Greenbaum and Trefethen [SIAM J. Sci. Comput., 15 (1994), pp. 359–368]. We also discuss general characterizations of best approximation in normed linear spaces of matrices and show on an example that a known sufficient condition for uniqueness in these characterizations is not necessary.

Key words. matrix approximation problems, polynomials in matrices, matrix functions, matrix 2-norm, GMRES, Arnoldi's method

AMS subject classifications. 41A52, 15A60, 65F35

1. Introduction. Much of the work in approximation theory concerns the approximation of a given function f on some (compact) set Ω in the complex plane by polynomials. Classical results in this area deal with the *best approximation problem*

$$(1.1) \quad \min_{p \in \mathcal{P}_m} \|f - p\|_{\Omega},$$

where $\|g\|_{\Omega} \equiv \max_{z \in \Omega} |g(z)|$, and \mathcal{P}_m denotes the set of polynomials of degree at most m . (Note that since in (1.1) we seek an approximation from a finite dimensional subspace, the minimum is indeed attained by some polynomial $p_* \in \mathcal{P}_m$.)

Scalar approximation problems of the form (1.1) have been studied since the mid 1850s. Accordingly, numerous results on existence and uniqueness of the solution as well as estimates for the value of (1.1) are known. Here we consider a problem that at first sight looks similar, but apparently is much less understood: Let f be a function that is analytic in a neighborhood of the spectrum of a given matrix $A \in \mathbb{C}^{n \times n}$, so that $f(A)$ is well defined, and let $\|\cdot\|$ be a given matrix norm. Consider the *matrix* approximation problem

$$(1.2) \quad \min_{p \in \mathcal{P}_m} \|f(A) - p(A)\|.$$

Does this problem have a unique solution?

An answer to this question of course depends on the norm used in (1.2). If the norm is known to be *strictly convex*, then (1.2) is guaranteed to have a uniquely defined solution as long as the value of (1.2) is positive (see Section 4 at the end of this paper for an informal discussion of strict convexity and matrix norms). A useful matrix norm that is met in many applications is the matrix 2-norm (or spectral norm), which for a given matrix A is equal to the largest singular value of A . This norm is *not* strictly convex, and thus the general result on uniqueness of best approximation

*Institute of Mathematics, Technical University of Berlin, Straße des 17. Juni 136, 10623 Berlin, Germany (liesen@math.tu-berlin.de). The work of this author was supported by the Heisenberg Program of the Deutsche Forschungsgemeinschaft.

†Institute of Computer Science, Academy of Sciences of the Czech Republic, Pod vodárenskou věží 2, 18207 Prague, Czech Republic (tichy@cs.cas.cz). The work of this author was supported by the GAAS grant IAA100300802.

in strictly convex linear spaces does not apply. Hence our question is nontrivial in case of the matrix 2-norm.

It is well known that when the function f is analytic in a neighborhood of the eigenvalues of the matrix $A \in \mathbb{C}^{n \times n}$, then $f(A)$ itself is a complex $n \times n$ matrix. In fact, $f(A) = p_f(A)$, where p_f is a polynomial that depends on the values and possibly the derivatives of f on the spectrum of A . The recent book of Higham [4] gives an extensive overview of definitions, applications, and computational techniques for matrix functions. Our above question now naturally leads to the following mathematical problem: *Let a polynomial b and a nonnegative integer $m < \deg b$ be given. Determine conditions so that the best approximation problem*

$$(1.3) \quad \min_{p \in \mathcal{P}_m} \|b(A) - p(A)\|$$

has a unique solution, where $\|\cdot\|$ is the matrix 2-norm and \mathcal{P}_m denotes the set of polynomials of degree at most m .

When searching the literature we found a number of results on general characterizations of best approximation in normed linear spaces of matrices, e.g. in [6, 8, 14, 15], but just a few papers related to our specific problem. In particular, Greenbaum and Trefethen consider in [3] the two approximation problems

$$(1.4) \quad \min_{p \in \mathcal{P}_m} \|A^{m+1} - p(A)\|,$$

$$(1.5) \quad \min_{p \in \mathcal{P}_m} \|I - Ap(A)\|.$$

They state that both (1.4) and (1.5) (for nonsingular A) have a unique minimizer^a. The problem (1.4) is equal to (1.3) with $b(A) = A^{m+1}$. Because of its relation to the convergence of the Arnoldi method [1] for approximating the eigenvalues of A , the uniquely defined monic polynomial $z^{m+1} - p_*$ that solves (1.4) is called the $(m+1)$ st *ideal Arnoldi polynomial* of A . In a paper that is mostly concerned with algorithmic and computational results, Toh and Trefethen [12] call this polynomial the $(m+1)$ st *Chebyshev polynomial* of A . The reason for this terminology is the following: When the matrix A is *normal*, i.e. unitarily diagonalizable, problem (1.4) becomes a scalar approximation problem of the form (1.1) with $f(z) = z^{m+1}$ and Ω being the set of eigenvalues of A . The resulting monic polynomial is the $(m+1)$ st Chebyshev polynomial on the (discrete) set of eigenvalues of A . In this sense, the matrix approximation problem (1.3) we study here can be considered a generalization of the classical scalar approximation problem (1.1). Some further results on Chebyshev polynomials of matrices are given in [10], and [13, Chapter 29].

The quantity (1.5) can be used for bounding the relative residual norm in the GMRES method [7]; for details see, e.g., [9, 11]. Therefore, the uniquely defined polynomial $1 - zp_*$ that solves (1.5) is called the $(m+1)$ st *ideal GMRES polynomial* of A .

In this paper we show that, despite the lack of strict convexity of the matrix 2-norm, the approximation problem (1.3) as well as a certain related problem that generalizes (1.5) have a unique minimizer. Furthermore, we discuss some of the above mentioned general characterizations of best approximations with respect to the 2-norm in linear spaces of matrices. On the example of a Jordan block we show that a

^aThe statement of uniqueness is true but the proof given in [3], which was later repeated in [13, Chapter 29], contains a small error at the very end. After the error was spotted by Oliver Ernst, it was fixed by Anne Greenbaum in 2005, but the correction was not published.

sufficient condition for uniqueness of the best approximation obtained by Ziętak [14] does not hold. We are not aware that such an example for a nonnormal matrix has been given before.

2. Uniqueness results. Let $\ell \geq 0$ and $m \geq 0$ be given integers, and consider a given polynomial b of the form

$$b = \sum_{j=0}^{\ell+m+1} \beta_j z^j \in \mathcal{P}_{\ell+m+1}.$$

Let us rewrite the approximation problem (1.3) in a more convenient equivalent form:

$$\begin{aligned} \min_{p \in \mathcal{P}_m} \|b(A) - p(A)\| &= \min_{p \in \mathcal{P}_m} \|b(A) - (p(A) + \sum_{j=0}^m \beta_j A^j)\| \\ &= \min_{p \in \mathcal{P}_m} \left\| \sum_{j=m+1}^{\ell+m+1} \beta_j A^j - p(A) \right\| \\ (2.1) \quad &= \min_{p \in \mathcal{P}_m} \|A^{m+1} \sum_{j=0}^{\ell} \beta_{j+m+1} A^j - p(A)\|. \end{aligned}$$

The polynomials in (2.1) are of the form $z^{m+1}g + h$, where the polynomial $g \in \mathcal{P}_\ell$ is given, and $h \in \mathcal{P}_m$ is sought. Hence (1.3) is equivalent to the problem

$$(2.2) \quad \min_{h \in \mathcal{P}_m} \|A^{m+1}g(A) + h(A)\|,$$

where $g \in \mathcal{P}_\ell$ is a given polynomial, or,

$$(2.3) \quad \min_{p \in \mathcal{G}_{\ell,m}^{(g)}} \|p(A)\|, \text{ where } \mathcal{G}_{\ell,m}^{(g)} \equiv \{z^{m+1}g + h : g \in \mathcal{P}_\ell \text{ is given, } h \in \mathcal{P}_m\}.$$

With $\ell = 0$ and $g = 1$, (2.3) reduces to (1.4).

Similarly, we may consider the approximation problem

$$(2.4) \quad \min_{p \in \mathcal{H}_{\ell,m}^{(h)}} \|p(A)\|, \text{ where } \mathcal{H}_{\ell,m}^{(h)} \equiv \{z^{m+1}g + h : h \in \mathcal{P}_m \text{ is given, } g \in \mathcal{P}_\ell\}.$$

Setting $m = 0$ and $h = 1$ in (2.4), we retrieve a problem of the form (1.5).

The problems (2.3) and (2.4) are trivial for $g = 0$ and $h = 0$, respectively. Both cases are unconstrained minimizations problems, and it is easily seen that the resulting minimum value is zero. In the following we will therefore exclude the cases $g = 0$ in (2.3) and $h = 0$ in (2.4). Under this assumption, both $\mathcal{G}_{\ell,m}^{(g)}$ and $\mathcal{H}_{\ell,m}^{(h)}$ are subsets of $\mathcal{P}_{\ell+m+1}$ where certain coefficients are fixed. In case of $\mathcal{G}_{\ell,m}^{(g)}$ these are the coefficients at the $\ell + 1$ largest powers of z , namely $z^{m+1}, \dots, z^{\ell+m+1}$. For $\mathcal{H}_{\ell,m}^{(h)}$ these are the coefficients at the $m + 1$ smallest powers of z , namely $1, \dots, z^m$.

We start with conditions so that the values of (2.3) and (2.4) are positive for all given nonzero polynomials $g \in \mathcal{P}_\ell$ and $h \in \mathcal{P}_m$, respectively.

LEMMA 2.1. *Consider the approximation problems (2.3) and (2.4), where $\ell \geq 0$ and $m \geq 0$ are given integers. Denote by $d(A)$ the degree of the minimal polynomial of the given matrix $A \in \mathbb{C}^{n \times n}$. Then the following two assertions are equivalent:*

- (1) $\min_{p \in \mathcal{G}_{\ell,m}^{(g)}} \|p(A)\| > 0$ for all nonzero polynomials $g \in \mathcal{P}_\ell$.
 (2) $m + \ell + 1 < d(A)$.

If A is nonsingular, the two assertions are equivalent with:

- (3) $\min_{p \in \mathcal{H}_{\ell,m}^{(h)}} \|p(A)\| > 0$ for all nonzero polynomials $h \in \mathcal{P}_m$.

Proof. (1) \Rightarrow (2): We suppose that $m + \ell + 1 \geq d(A)$ and show that (1) fails to hold. Denote the minimal polynomial of A by Ψ_A . If $m + 1 \leq d(A) \leq \ell + m + 1$, then there exist uniquely determined polynomials $\hat{g} \in \mathcal{P}_\ell$, $\hat{g} \neq 0$, and $\hat{h} \in \mathcal{P}_m$, so that $z^{m+1} \cdot \hat{g} + \hat{h} = \Psi_A$. Hence $\min_{p \in \mathcal{G}_{\ell,m}^{(g)}} \|p(A)\| = 0$ for $g = \hat{g}$. If $0 \leq d(A) \leq m$, let \hat{g} be any nonzero polynomial of degree at most ℓ . By the division theorem for polynomials^b, there exist uniquely defined polynomials q and $h \in \mathcal{P}_{m-1}$, so that $z^{m+1} \cdot \hat{g} = q \cdot \Psi_A - h$, or, equivalently, $z^{m+1} \cdot \hat{g} + h = q \cdot \Psi_A$. Hence $A^{m+1} \hat{g}(A) + h(A) = 0$, which means that $\min_{p \in \mathcal{G}_{\ell,m}^{(g)}} \|p(A)\| = 0$ for the nonzero polynomial $g = \hat{g} \in \mathcal{P}_\ell$.

(2) \Rightarrow (1): If $m + \ell + 1 < d(A)$, then $\mathcal{G}_{\ell,m}^{(g)} \subset \mathcal{P}_{m+\ell+1}$ implies $\min_{p \in \mathcal{G}_{\ell,m}^{(g)}} \|p(A)\| > 0$ for every nonzero polynomial $g \in \mathcal{P}_\ell$.

(2) \Rightarrow (3): If $m + \ell + 1 < d(A)$, then $\mathcal{H}_{\ell,m}^{(h)} \subset \mathcal{P}_{m+\ell+1}$ implies $\min_{p \in \mathcal{H}_{\ell,m}^{(h)}} \|p(A)\| > 0$ for every nonzero polynomial $h \in \mathcal{P}_m$.

(3) \Rightarrow (2): For this implication we use that A is nonsingular. Suppose that (2) does not hold, i.e. that $0 \leq d(A) \leq m + \ell + 1$. Then there exist uniquely defined polynomials $\hat{g} \in \mathcal{P}_\ell$ and $\hat{h} \in \mathcal{P}_m$, such that $z^{m+1} \cdot \hat{g} + \hat{h} = \Psi_A$. Since A is assumed to be nonsingular, we must have $\hat{h} \neq 0$. Consequently, $\min_{p \in \mathcal{H}_{\ell,m}^{(h)}} \|p(A)\| = 0$ for the nonzero polynomial $h = \hat{h} \in \mathcal{P}_m$. \square

In the following Theorem 2.2 we show that the problem (2.3) has a uniquely defined minimizer when the value of this problem is positive (and not zero). In the previous lemma we have shown that $m + \ell + 1 < d(A)$ is necessary and sufficient so that the value of (2.3) is positive for all nonzero polynomials $g \in \mathcal{P}_\ell$. However, it is possible that for some nonzero polynomial $g \in \mathcal{P}_\ell$ the value of (2.3) is positive even when $m + 1 \leq d(A) \leq m + \ell + 1$. It is possible to further analyze this special case, but for the ease of the presentation we simply assume that the value of (2.3) is positive. The same assumption is made in Theorem 2.3 below, where we prove uniqueness of the minimizer of (2.4) (under the additional assumption that A is nonsingular).

THEOREM 2.2. *Let $A \in \mathbb{C}^{n \times n}$ be a given matrix, $\ell \geq 0$ and $m \geq 0$ be given integers, and $g \in \mathcal{P}_\ell$ be a given nonzero polynomial. If the value of (2.3) is positive, then this problem has a uniquely defined minimizer.*

Proof. The general strategy in the following is similar to the construction in [3, Section 5]. We suppose that $q_1 = z^{m+1}g + h_1 \in \mathcal{G}_{\ell,m}^{(g)}$ and $q_2 = z^{m+1}g + h_2 \in \mathcal{G}_{\ell,m}^{(g)}$ are two distinct solutions to (2.3) and derive a contradiction. Suppose that the minimal norm attained by the two polynomials is

$$C = \|q_1(A)\| = \|q_2(A)\|.$$

By assumption, $C > 0$. Define $q \equiv \frac{1}{2}(q_1 + q_2) \in \mathcal{G}_{\ell,m}^{(g)}$, then

$$\|q(A)\| \leq \frac{1}{2}(\|q_1(A)\| + \|q_2(A)\|) = C.$$

^bIf f and $g \neq 0$ are polynomials over a field \mathbb{F} , then there exist uniquely defined polynomials s and r over \mathbb{F} , such that (i) $f = g \cdot s + r$, and (ii) either $r = 0$ or $\deg r < \deg g$. If $\deg f \geq \deg g$, then $\deg f = \deg g + \deg s$. For a proof of this standard result, see, e.g., [5, Chapter 4].

Since C is assumed to be the minimal value of (2.3), we must have $\|q(A)\| = C$. Denote the singular value decomposition of $q(A)$ by

$$(2.5) \quad q(A) = V \operatorname{diag}(\sigma_1, \dots, \sigma_n) W^*.$$

Suppose that the maximal singular value $\sigma_1 = C$ of $q(A)$ is J -fold, with left and right singular vectors given by v_1, \dots, v_J and w_1, \dots, w_J , respectively.

It is well known that the 2-norm for vectors $v \in \mathbb{C}^n$, $\|v\| \equiv (v^*v)^{1/2}$, is strictly convex. For each w_j , $1 \leq j \leq J$, we have

$$C = \|q(A)w_j\| \leq \frac{1}{2} (\|q_1(A)w_j\| + \|q_2(A)w_j\|) \leq C,$$

which implies

$$\|q_1(A)w_j\| = \|q_2(A)w_j\| = C, \quad 1 \leq j \leq J.$$

By the strict convexity of the vector 2-norm,

$$q_1(A)w_j = q_2(A)w_j, \quad 1 \leq j \leq J.$$

Similarly, one can show that

$$q_1(A)^*v_j = q_2(A)^*v_j, \quad 1 \leq j \leq J.$$

Thus,

$$(2.6) \quad (q_2(A) - q_1(A))w_j = 0, \quad (q_2(A) - q_1(A))^*v_j = 0, \quad 1 \leq j \leq J.$$

By assumption, $q_2 - q_1 \in \mathcal{P}_m$ is a nonzero polynomial. By the division theorem for polynomials (see the footnote on p. 4), there exist uniquely defined polynomials s and r , with $\deg s \leq \ell + m + 1$ and $\deg r < \deg(q_2 - q_1) \leq m$ (or $r = 0$) so that

$$z^{m+1}g = (q_2 - q_1) \cdot s + r.$$

Hence we have shown that for the given polynomials $q_2 - q_1$ and g there exist polynomials s and r , such that

$$\tilde{q} \equiv (q_2 - q_1) \cdot s = z^{m+1}g - r \in \mathcal{G}_{\ell, m}^{(g)}.$$

Since $g \neq 0$, we must have $\tilde{q} \neq 0$. For a fixed $\epsilon \in (0, 1)$, consider the polynomial

$$q_\epsilon = (1 - \epsilon)q + \epsilon\tilde{q} \in \mathcal{G}_{\ell, m}^{(g)}.$$

By (2.6),

$$\tilde{q}(A)w_j = 0, \quad \tilde{q}(A)^*v_j = 0, \quad 1 \leq j \leq J,$$

and thus

$$\begin{aligned} q_\epsilon(A)^*q_\epsilon(A)w_j &= (1 - \epsilon)q_\epsilon(A)^*q(A)w_j = (1 - \epsilon)Cq_\epsilon(A)^*v_j \\ &= (1 - \epsilon)Cq(A)^*v_j = (1 - \epsilon)^2C^2w_j, \end{aligned}$$

which shows that w_1, \dots, w_J are right singular vectors of $q_\epsilon(A)$ corresponding to the singular value $(1 - \epsilon)C$. Note that $(1 - \epsilon)C < C$ since $C > 0$.

Now there are two cases: Either $\|q_\epsilon(A)\| = (1 - \epsilon)C$, or $(1 - \epsilon)C$ is not the largest singular value of $q_\epsilon(A)$. In the first case we have a contradiction to the fact that C is the minimal value of (2.3). Therefore the second case must hold. In that case, none of the vectors w_1, \dots, w_J corresponds to the largest singular value of $q_\epsilon(A)$. Using this fact and the singular value decomposition (2.5), we get

$$\begin{aligned}
 \|q_\epsilon(A)\| &= \|q_\epsilon(A)W\| \\
 &= \|q_\epsilon(A)[w_{J+1}, \dots, w_n]\| \\
 &= \|(1 - \epsilon)q(A)[w_{J+1}, \dots, w_n] + \epsilon\tilde{q}(A)[w_{J+1}, \dots, w_n]\| \\
 &\leq (1 - \epsilon)\|[v_{J+1}, \dots, v_n] \text{diag}(\sigma_{J+1}, \dots, \sigma_n)\| + \epsilon\|\tilde{q}(A)[w_{J+1}, \dots, w_n]\| \\
 (2.7) \quad &\leq (1 - \epsilon)\sigma_{J+1} + \epsilon\|\tilde{q}(A)[w_{J+1}, \dots, w_n]\|.
 \end{aligned}$$

Now note that the norm $\|\tilde{q}(A)[w_{J+1}, \dots, w_n]\|$ in (2.7) does not depend on the choice of ϵ , and that (2.7) goes to σ_{J+1} as ϵ goes to zero. Since $\sigma_J > \sigma_{J+1}$, one can find a positive $\epsilon_* \in (0, 1)$, such that (2.7) is less than σ_J for all $\epsilon \in (0, \epsilon_*)$. Any of the corresponding polynomials q_ϵ gives a matrix $q_\epsilon(A)$ whose norm is less than σ_J . This contradiction finishes the proof. \square

In the following theorem we prove that the problem (2.4), and hence in particular the problem (1.5), has a uniquely defined minimizer.

THEOREM 2.3. *Let $A \in \mathbb{C}^{n \times n}$ be a given nonsingular matrix, $\ell \geq 0$ and $m \geq 0$ be given integers, and $h \in \mathcal{P}_m$ be a given nonzero polynomial. If the value of (2.4) is positive, then this problem has a uniquely defined minimizer.*

Proof. Most parts of the following proof are analogous to the proof of Theorem 2.2, and are stated only briefly. However, the construction of the polynomial q_ϵ used to derive the contradiction is different.

We suppose that $q_1 = z^{m+1}g_1 + h \in \mathcal{H}_{\ell, m}^{(h)}$ and $q_2 = z^{m+1}g_2 + h \in \mathcal{H}_{\ell, m}^{(h)}$ are two distinct solutions to (2.4), and that the minimal norm attained by them is $C = \|q_1(A)\| = \|q_2(A)\|$. By assumption, $C > 0$. Define $q \equiv \frac{1}{2}(q_1 + q_2) \in \mathcal{H}_{\ell, m}^{(h)}$, then $\|q(A)\| = C$. Denote the singular value decomposition of $q(A)$ by $q(A) = V \text{diag}(\sigma_1, \dots, \sigma_n)W^*$, and suppose that the maximal singular value $\sigma_1 = C$ of $q(A)$ is J -fold, with left and right singular vectors given by v_1, \dots, v_J and w_1, \dots, w_J , respectively. As previously, we can show that

$$(q_2(A) - q_1(A))w_j = 0, \quad (q_2(A) - q_1(A))^*v_j = 0, \quad 1 \leq j \leq J.$$

Since A is nonsingular, and $q_2 - q_1 = z^{m+1}(g_2 - g_1)$, these relations imply that

$$(2.8) \quad (g_2(A) - g_1(A))w_j = 0, \quad (g_2(A) - g_1(A))^*v_j = 0, \quad 1 \leq j \leq J.$$

By assumption, $0 \neq g_2 - g_1 \in \mathcal{P}_\ell$. Hence there exists an integer d , $0 \leq d \leq \ell$, so that

$$g_2 - g_1 = \sum_{i=d}^{\ell} \gamma_i z^i, \quad \text{with } \gamma_d \neq 0.$$

Now define

$$\tilde{g} \equiv z^{-d}(g_2 - g_1) \in \mathcal{P}_{\ell-d}.$$

By construction, \tilde{g} is a polynomial with a nonzero linear term. Furthermore, define

$$\hat{h} \equiv z^{-m-1-\ell+d}h \quad \text{and} \quad \hat{g} \equiv z^{-\ell+d}\tilde{g}.$$

After a formal change of variables $z^{-1} \mapsto y$, we obtain

$$\widehat{h}(y) \in \mathcal{P}_{m+1+\ell-d} \quad \text{and} \quad \widehat{g}(y) \in \mathcal{P}_{\ell-d} \setminus \mathcal{P}_{\ell-d-1}.$$

By the division theorem for polynomials (see the footnote on p. 4), there exist uniquely defined polynomials $s(y)$ and $r(y)$ with $\deg s \leq m+1$ (since $\widehat{g} \neq 0$ is of exact degree $\ell-d$) and $\deg r < \ell-d$ (or $r = 0$) such that

$$\widehat{h}(y) = \widehat{g}(y) \cdot s(y) - r(y).$$

We now multiply the preceding equation by $y^{-m-1-\ell+d}$, which gives

$$y^{-m-1-\ell+d} \widehat{h}(y) = (y^{-\ell+d} \widehat{g}(y)) \cdot (y^{-m-1} s(y)) - y^{-m-1} (y^{-\ell+d} r(y)).$$

Since $y^{-1} = z$, this equation is equivalent to

$$h = \widetilde{g} \cdot \widetilde{s} - z^{m+1} \widetilde{r},$$

where $\widetilde{s} \in \mathcal{P}_{m+1}$ and $\widetilde{r} \in \mathcal{P}_{\ell-d-1}$. Hence we have shown that for the given polynomials h and \widetilde{g} there exist polynomials $\widetilde{s} \in \mathcal{P}_{m+1}$ and $\widetilde{r} \in \mathcal{P}_{\ell-d-1}$, such that

$$\widetilde{q} \equiv \widetilde{g} \cdot \widetilde{s} = z^{m+1} \widetilde{r} + h \in \mathcal{H}_{\ell,m}^{(h)}.$$

For a fixed $\epsilon \in (0, 1)$, consider

$$q_\epsilon = (1 - \epsilon)q + \epsilon \widetilde{q} \in \mathcal{H}_{\ell,m}^{(h)}.$$

Since $\widetilde{q} = \widetilde{r}z^{-d}(g_2 - g_1)$, (2.8) implies that

$$\widetilde{q}(A)w_j = 0, \quad \widetilde{q}(A)^*v_j = 0, \quad 1 \leq j \leq J,$$

which can be used to show that

$$q_\epsilon(A)^*q_\epsilon(A)w_j = (1 - \epsilon)^2 C^2 w_j, \quad 1 \leq j \leq J.$$

Now the same argument as in the proof of Theorem 2.2 gives a contradiction to the original assumption that $q_2 \neq q_1$. \square

REMARK 2.4. *Similarly as in Lemma 2.1, the assumption of nonsingularity in the previous theorem is in general necessary. In other words, when A is singular the approximation problem (2.4) might have more than one solution even when the value of (2.4) is positive. The following example demonstrating this fact was pointed out to us by Krystyna Ziętak:*

Consider a normal matrix $A = U\Lambda U^$, where $U^*U = I$ and $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_n)$. Suppose that A is singular with n distinct eigenvalues and $\lambda_1 = 0$. Furthermore, suppose that $h \in \mathcal{P}_m$ is any given polynomial that satisfies $h(0) \neq 0$ and $|h(0)| > |h(\lambda_j)|$ for $j = 2, \dots, n$. Then for any integer $\ell \geq 0$,*

$$\min_{p \in \mathcal{H}_{\ell,m}^{(h)}} \|p(A)\| = \min_{g \in \mathcal{P}_\ell} \max_j |\lambda_j^{m+1} g(\lambda_j) + h(\lambda_j)| = |h(0)| > 0.$$

One solution of this problem is given by the polynomial $g = 0$. Moreover, the minimum value is attained for any polynomial $g \in \mathcal{P}_\ell$ that satisfies

$$\min_{g \in \mathcal{P}_\ell} \max_{2 \leq j \leq n} |\lambda_j^{m+1} g(\lambda_j) + h(\lambda_j)| < |h(0)|,$$

i.e., for any polynomial $g \in \mathcal{P}_\ell$ that is close enough to the zero polynomial.

3. Characterization of best approximation with respect to the matrix 2-norm. In this section we discuss general characterizations of best approximation in linear spaces of matrices with respect to the matrix 2-norm obtained by Ziętak [14, 15], and give an example from our specific problem. To state Ziętak's results, we need some notation. Suppose that we are given m matrices $A_1, \dots, A_m \in \mathbb{C}^{n \times n}$ that are linearly independent in $\mathbb{C}^{n \times n}$. We assume that $m < n^2$ to avoid trivialities. Denote $\mathbb{A} \equiv \text{span} \{A_1, \dots, A_m\}$, which is an m -dimensional subspace of $\mathbb{C}^{n \times n}$. As above, let $\|\cdot\|$ denote the matrix 2-norm. For a given matrix $B \in \mathbb{C}^{n \times n} \setminus \mathbb{A}$, we consider the best approximation (or matrix nearness) problem

$$(3.1) \quad \min_{M \in \mathbb{A}} \|B - M\|.$$

A matrix $A_* \in \mathbb{A}$ for which this minimum is achieved (such a matrix exists since \mathbb{A} is finite dimensional) is called a *spectral approximation of B from the subspace \mathbb{A}* . The corresponding matrix $R(A_*) = B - A_*$ is called a *residual matrix*.

The approximation problems (2.3) and (2.4) studied in the previous section are both special cases of (3.1). In case of (2.3),

$$B = A^{m+1}g(A), \text{ where } g \in \mathcal{P}_\ell \text{ is given, and } \mathbb{A} = \{I, A, \dots, A^m\},$$

while in case of (2.4),

$$B = h(A), \text{ where } h \in \mathcal{P}_m \text{ is given, and } \mathbb{A} = \{A^{m+1}, \dots, A^{\ell+m+1}\}.$$

We have shown that when the values of these approximation problems are positive (which is true if $\ell + m + 1 < d(A)$), for both these problems there exists a uniquely defined spectral approximation A_* of B from the subspace \mathbb{A} (in case of (2.4), we have assumed that A is nonsingular).

In general, however, the spectral approximation of a matrix $B \in \mathbb{C}^{n \times n}$ from a subspace $\mathbb{A} \subset \mathbb{C}^{n \times n}$ is not unique. Ziętak [14] gives a general characterization of spectral approximations based on the singular value decomposition of the residual matrices. In particular [14, Theorem 4.3] contains the following sufficient condition for uniqueness of the spectral approximation (that theorem is formulated for real matrices only, but with results from [15] it is not hard to generalize it to the complex case).

LEMMA 3.1. *In the notation established above, let A_* be a spectral approximation of B from the subspace \mathbb{A} . If the residual matrix $R(A_*) = B - A_*$ has an n -fold maximal singular value, then the spectral approximation A_* of B from the subspace \mathbb{A} is unique.*

It is quite obvious that this sufficient condition is, in general, not necessary. To construct a nontrivial counterexample, we recall that the dual norm to the matrix 2-norm is the *Frobenius norm*,

$$(3.2) \quad \|A\|_F \equiv \left(\sum_{j=1}^n \sigma_j(A)^2 \right)^{1/2} = \langle A, A \rangle^{1/2},$$

where $\sigma_1(A), \dots, \sigma_n(A)$ denote the singular values of A , the trace of $A = [a_{ij}]$ is defined by $\text{tr}(A) \equiv a_{11} + \dots + a_{nn}$, and $\langle A, X \rangle \equiv \text{tr}(A^*X)$. Using this notation, we can state the following result, which is given in [15, p. 173].

LEMMA 3.2. *In the notation established above, $A_* \in \mathbb{A}$ is a spectral approximation of B from the subspace \mathbb{A} if and only if there exists a matrix $Z \in \mathbb{C}^{n \times n}$ with $\|Z\|_F = 1$, such that*

$$(3.3) \quad \langle Z, X \rangle = 0, \quad \forall X \in \mathbb{A}, \quad \text{and} \quad \operatorname{Re} \langle Z, B - A_* \rangle = \|B - A_*\|.$$

THEOREM 3.3. *For $\lambda \in \mathbb{C}$, consider the $n \times n$ Jordan block*

$$J_\lambda \equiv \begin{pmatrix} \lambda & 1 & & \\ & \ddots & \ddots & \\ & & \ddots & 1 \\ & & & \lambda \end{pmatrix}.$$

Then for any nonnegative integer m with $m+1 \leq n$, the solution to the approximation problem (1.4) with $A = J_\lambda$, i.e. the $(m+1)$ st ideal Arnoldi (or Chebyshev) polynomial of J_λ , is uniquely defined and given by $(z - \lambda)^{m+1}$.

Proof. With $A = J_\lambda$, the approximation problem (1.4) reads

$$(3.4) \quad \min_{p \in \mathcal{P}_m} \|J_\lambda^{m+1} - p(J_\lambda)\|.$$

In the notation established in this section, we seek a spectral approximation A_* of $B = J_\lambda^{m+1}$ from the subspace $\mathbb{A} = \operatorname{span}\{I, J_\lambda, \dots, J_\lambda^m\}$. We claim that the uniquely defined solution is given by the matrix $A_* = J_\lambda^{m+1} - (J_\lambda - \lambda I)^{m+1}$. For this matrix A_* we get

$$B - A_* = J_\lambda^{m+1} - A_* = (J_\lambda - \lambda I)^{m+1} = J_0^{m+1}.$$

For $m+1 = n$, $A_* = J_\lambda^n - (J_\lambda - \lambda I)^n = J_\lambda^n$ yields $B - A_* = J_0^n = 0$. The corresponding ideal Arnoldi polynomial of J_λ is uniquely defined and equal to $(z - \lambda)^n$, the minimal polynomial of J_λ .

For $m+1 < n$, the value of (3.4) is positive, and hence Theorem 2.2 ensures that the spectral approximation of J_λ^{m+1} from the subspace \mathbb{A} is uniquely defined. We prove our claim using Lemma 3.2. Define $Z \equiv e_1 e_{m+2}^T$, then

$$\langle Z, J_\lambda^k \rangle = 0, \quad \text{for } k = 0, \dots, m,$$

and $\|B - A_*\| = \|J_0^{m+1}\| = 1$, so that

$$\langle Z, B - A_* \rangle = \langle Z, J_0^{m+1} \rangle = 1 = \|B - A_*\|,$$

which shows (3.3) and completes the proof. \square

The proof of this theorem shows that the residual matrix of the spectral approximation A_* of $B = J_\lambda^{m+1}$ from the subspace $\mathbb{A} = \operatorname{span}\{I, J_\lambda, \dots, J_\lambda^m\}$ is given by $R(A_*) = J_0^{m+1}$. This matrix $R(A_*)$ has $m+1$ singular values equal to zero, and $n - m - 1$ singular values equal to one. Hence, for $m+1 < n$, the maximal singular value of the residual matrix is not n -fold, and the sufficient condition of Lemma 3.1 does not hold. Nevertheless, the spectral approximation of B from the subspace \mathbb{A} is unique whenever $m+1 < n$.

As shown above, for $m = 0, 1, \dots, n-1$ the polynomial $(z - \lambda)^{m+1}$ solves the ideal Arnoldi approximation problem (1.4) for $A = J_\lambda$. For $\lambda \neq 0$, we can write

$$(z - \lambda)^{m+1} = (-\lambda)^{m+1} \cdot (1 - \lambda^{-1}z)^{m+1}.$$

Note that the rightmost factor is a polynomial that has value one at the origin. Hence it is a candidate for the solution of the ideal GMRES approximation problem (1.5) for $A = J_\lambda$. More generally, it is tempting to assume that for a given matrix A its $(m+1)$ st ideal GMRES polynomial is equal to a scaled version of its $(m+1)$ st ideal Arnoldi (or Chebyshev) polynomial. However, this assumption is false, as we can already see in case $A = J_\lambda$. As shown in [9], the determination of the ideal GMRES polynomials for a Jordan block is an intriguing problem, since these polynomials can become quite complicated. They are of the simple form $(1 - \lambda^{-1}z)^{m+1}$ if and only if $0 \leq m+1 < n/2$ and $|\lambda| \geq \varrho_{m+1, n-m-1}^{-1}$, cf. [9, Theorem 3.2]. Here $\varrho_{k,n}$ denotes the radius of the polynomial numerical hull of degree k of an $n \times n$ Jordan block (this radius is independent of the eigenvalue λ).

Now let n be even and consider $m+1 = n/2$. If $|\lambda| \leq 2^{-2/n}$, the ideal GMRES polynomial of degree $n/2$ of J_λ is equal to the constant polynomial 1. If $|\lambda| \geq 2^{-2/n}$, the ideal GMRES polynomial of degree $n/2$ of J_λ is equal to

$$(3.5) \quad \frac{2}{4\lambda^n + 1} + \frac{4\lambda^n - 1}{4\lambda^n + 1} (1 - \lambda^{-1}z)^{n/2},$$

cf. [9, p. 465]. Obviously, neither the polynomial 1 nor the polynomial (3.5) are scalar multiples of $(z - \lambda)^{n/2}$, the ideal Arnoldi polynomial of degree $n/2$ of J_λ .

4. Matrix norms and strict convexity. We conclude the paper with an informal discussion of the theoretical background of our exposition, namely matrix norms and strict convexity.

A norm $|\cdot|$ on a vector space \mathcal{V} is called *strictly convex*, when for all vectors $v_1, v_2 \in \mathcal{V}$ the equation $|v_1| = |v_2| = \frac{1}{2}|v_1 + v_2|$ implies that $v_1 = v_2$. A geometric interpretation of strict convexity is that the unit sphere in \mathcal{V} with respect to the norm $|\cdot|$ does not contain any line segments. Strictly convex norms are of interest, since they imply uniqueness of best approximation problems from finite dimensional subspaces of \mathcal{V} . More precisely, if $\mathcal{S} \subset \mathcal{V}$ is a finite dimensional subspace, then for any given $v \in \mathcal{V}$ there exists a *unique* $s_* \in \mathcal{S}$ so that

$$|v - s_*| = \min_{s \in \mathcal{S}} |v - s|.$$

A proof of this classical result can be found in most books on approximation theory; see, e.g., [2, Chapter 1].

In this paper we have studied best approximation problems in the space $\mathcal{V} = \mathbb{C}^{n \times n}$ and with respect to the matrix 2-norm. This norm is *not* strictly convex, as can be seen from the following simple argument: Suppose that we have two matrices $A_1, A_2 \in \mathbb{C}^{n \times n}$ of the form

$$A_1 = \begin{bmatrix} B & 0 \\ 0 & C \end{bmatrix}, \quad A_2 = \begin{bmatrix} B & 0 \\ 0 & D \end{bmatrix},$$

where $\|A_1\| = \|A_2\| = \sigma_1(B) \geq \frac{1}{2}\|C+D\|$. Then $\frac{1}{2}\|A_1 + A_2\| = \sigma_1(B)$, but whenever $C \neq D$, we have $A_1 \neq A_2$.

When we consider all singular values instead of just the largest one, we receive the Frobenius norm (3.2). One can easily show that $\|A\|_F^2 = \sum_{i=1}^n \sum_{j=1}^n |a_{ij}|^2$. Hence the Frobenius norm is nothing but the vector 2-norm in the space \mathbb{C}^{n^2} , and this norm is strictly convex. For any $\epsilon > 0$, the norm $\|A\|_\epsilon \equiv (1 - \epsilon)\|A\| + \epsilon\|A\|_F$ is strictly convex (adding a convex and a strictly convex function yields a strictly convex function). Hence, arbitrarily close to the matrix 2-norm there exists a matrix norm that is strictly convex. However, this norm is of rather theoretical interest only.

Now consider a Hilbert space \mathbb{H} with inner product $\langle \cdot, \cdot \rangle$. Then it is easily shown that the associated norm $\|\cdot\| = \langle \cdot, \cdot \rangle^{1/2}$ is strictly convex. In case of $\mathbb{H} = \mathbb{C}^{n \times n}$ we may for example consider the trace inner product $\langle A, B \rangle \equiv \text{tr}(A^* B)$. Then the associated norm is just the Frobenius norm, i.e. $\|A\|_F = \langle A, A \rangle^{1/2}$, which gives another proof of strict convexity of the Frobenius norm. Moreover, the lack of strict convexity of the matrix 2-norm $\|\cdot\|$ shows that this norm is not induced from any inner product on $\mathbb{C}^{n \times n}$.

Acknowledgements. We thank Krystyna Ziętak for many discussions and suggestions that helped to improve the content of this paper. We also thank Shmuel Friedland, Anne Greenbaum, and Nick Trefethen for their helpful comments.

REFERENCES

- [1] W. E. ARNOLDI, *The principle of minimized iteration in the solution of the matrix eigenvalue problem*, Quart. Appl. Math., 9 (1951), pp. 17–29.
- [2] E. W. CHENEY, *Introduction to Approximation Theory*, McGraw-Hill Book Co., New York, 1966.
- [3] A. GREENBAUM AND L. N. TREFETHEN, *GMRES/CR and Arnoldi/Lanczos as matrix approximation problems*, SIAM J. Sci. Comput., 15 (1994), pp. 359–368.
- [4] N. J. HIGHAM, *Functions of Matrices: Theory and Computation*, Society for Industrial and Applied Mathematics (SIAM), Philadelphia, PA, 2008.
- [5] K. HOFFMAN AND R. KUNZE, *Linear Algebra*, 2nd edition, Prentice-Hall Inc., Englewood Cliffs, N.J., 1971.
- [6] K. K. LAU AND W. O. J. RIHA, *Characterization of best approximations in normed linear spaces of matrices by elements of finite-dimensional linear subspaces*, Linear Algebra Appl., 35 (1981), pp. 109–120.
- [7] Y. SAAD AND M. H. SCHULTZ, *GMRES: a generalized minimal residual algorithm for solving nonsymmetric linear systems*, SIAM J. Sci. Statist. Comput., 7 (1986), pp. 856–869.
- [8] I. SINGER, *Best approximation in normed linear spaces by elements of linear subspaces*, Translated from the Romanian by Radu Georgescu. Die Grundlehren der mathematischen Wissenschaften, Band 171, Publishing House of the Academy of the Socialist Republic of Romania, Bucharest, 1970.
- [9] P. TICHÝ, J. LIESEN, AND V. FABER, *On worst-case GMRES, ideal GMRES, and the polynomial numerical hull of a Jordan block*, Electron. Trans. Numer. Anal., 26 (2007), pp. 453–473 (electronic).
- [10] K.-C. TOH, *Matrix Approximation Problems and Nonsymmetric Iterative Methods*, PhD thesis, Cornell University, Ithaca, N.Y., 1996.
- [11] K.-C. TOH, *GMRES vs. ideal GMRES*, SIAM J. Matrix Anal. Appl., 18 (1997), pp. 30–36.
- [12] K.-C. TOH AND L. N. TREFETHEN, *The Chebyshev polynomials of a matrix*, SIAM J. Matrix Anal. Appl., 20 (1998), pp. 400–419.
- [13] L. N. TREFETHEN AND M. EMBREE, *Spectra and Pseudospectra*, Princeton University Press, Princeton, N.J., 2005.
- [14] K. ZIĘTAK, *Properties of linear approximations of matrices in the spectral norm*, Linear Algebra Appl., 183 (1993), pp. 41–60.
- [15] K. ZIĘTAK, *On approximation problems with zero-trace matrices*, Linear Algebra Appl., 247 (1996), pp. 169–183.