

## CHEBYSHEV APPROXIMATION VIA POLYNOMIAL MAPPINGS AND THE CONVERGENCE BEHAVIOUR OF KRYLOV SUBSPACE METHODS \*

BERND FISCHER <sup>†</sup> AND FRANZ PEHERSTORFER <sup>‡</sup>

**Abstract.** Let  $\varphi_m$  be a polynomial satisfying some mild conditions. Given a set  $R \subset \mathbb{C}$ , a continuous function  $f$  on  $R$  and its best approximation  $p_{n-1}^*$  from  $\Pi_{n-1}$  with respect to the maximum norm, we show that  $p_{n-1}^* \circ \varphi_m$  is a best approximation to  $f \circ \varphi_m$  on the inverse polynomial image  $S$  of  $R$ , i.e.  $\varphi_m(S) = R$ , where the extremal signature is given explicitly. A similar result is presented for constrained Chebyshev polynomial approximation. Finally, we apply the obtained results to the computation of the convergence rate of Krylov subspace methods when applied to a preconditioned linear system. We investigate pairs of preconditioners where the eigenvalues are contained in sets  $S$  and  $R$ , respectively, which are related by  $\varphi_m(S) = R$ .

**Key words.** Chebyshev polynomial, optimal polynomial, extremal signature, Krylov subspace method, convergence rate.

**AMS subject classifications.** 41A10, 30E10, 65F10.

**1. Notations and statement of the problem.** Let  $R \subset \mathbb{C}$  denote a compact subset of the complex plane and let  $C(R)$  be the set of continuous functions on  $R$ . For  $f \in C(R)$  we denote by  $\|f\|_R := \max_{z \in R} |f(z)|$  the uniform norm on  $R$ . Furthermore, let  $g_1, g_2, \dots, g_n \in C(R)$  be linearly independent functions with  $V_n := \text{span}\{g_1, g_2, \dots, g_n\}$ . Then the *best approximation*  $g^*$  of  $f$  with respect to  $V_n$  on  $R$  is the solution of the complex Chebyshev approximation problem

$$(1.1) \quad \|f - g^*\|_R = \min_{g \in V_n} \|f - g\|_R.$$

It is well-known, that  $g^*$  exists for any  $R$  and is unique provided that  $R$  contains at least  $n$  points.

Now, let  $\Pi_n := \{p(z) = \sum_{j=0}^n a_j z^j \mid a_j \in \mathbb{C}\}$  denote the set of all polynomials of degree up to  $n$  and let  $\varphi_m \in \Pi_m \setminus \Pi_{m-1}$  be a polynomial of exact degree  $m$ .

$R$  and  $\varphi_m$  may be used to define the set (cf. Figure 1.1)

$$S = S(R, \varphi_m) = \{s \in \mathbb{C} : \varphi_m(s) \in R\}.$$

In other words, we have  $\varphi_m(S) = R$  and  $S = \varphi_m^{-1}(R)$ , respectively.

By construction, it is clear that

$$\min_{g \in V_n} \|f \circ \varphi_m - g \circ \varphi_m\|_S = \min_{g \in V_n} \|f - g\|_R.$$

However, if  $V_n$  is a space of polynomials, e.g.,  $V_n = \Pi_{n-1}$  or  $V_n = (z - c)\Pi_{n-1}$ , one may ask the question whether even the following equations

$$(1.2) \quad \begin{aligned} \|f - p_{n-1}^*\|_R &= \min_{p_{n-1} \in \Pi_{n-1}} \|f - p_{n-1}\|_R \\ &= \min_{p_{mn-1} \in \Pi_{mn-1}} \|f \circ \varphi_m - p_{mn-1}\|_S = \|f \circ \varphi_m - p_{n-1}^* \circ \varphi_m\|_S \end{aligned}$$

\*Received May 18, 1999. Accepted for publication June 22, 2001. Recommended by R. Freund.

<sup>†</sup>Institute of Mathematics, Medical University of Lübeck, 23560 Lübeck, Wallstraße 40, Germany. E-mail: fischer@math.mu-luebeck.de

<sup>‡</sup>Institute for Analysis and Computational Mathematics, University of Linz, 4040 Linz, Altenbergerstraße, Austria. This research was supported by the Austrian Fonds zur Förderung der wissenschaftlichen Forschung, project-number P12985-TEC. E-mail: franz.peherstorfer@jk.uni-linz.ac.at

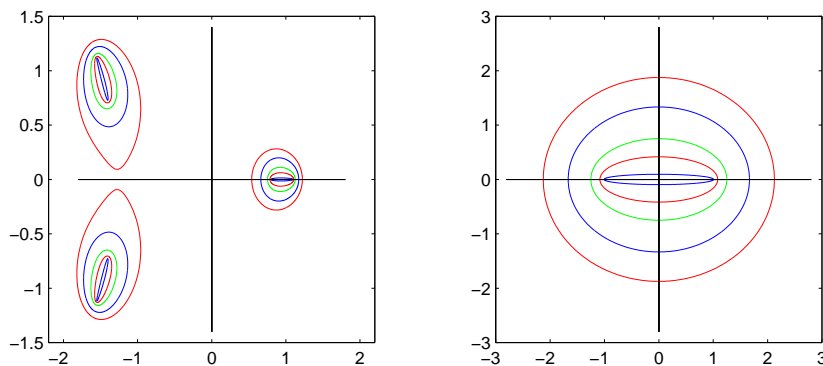


FIG. 1.1.  $(\varphi_3)^{-1}(R)$  (left) for various ellipses  $R$  (right) and  $\varphi_3(z) = z^3 + 2z^2 + 0.25z - 3$ .

hold true. In Section 2, we will give a complete answer to this question. We remark that the case where  $S$  is the inverse image of equipotential lines under a polynomial mapping has been considered in Peherstorfer [9].

Furthermore, we will investigate the convergence behavior of Krylov subspace methods when applied to the linear system

$$\mathcal{A}\mathbf{x} = \mathbf{b}, \quad \mathcal{A} \in \mathbb{C}^{N \times N}.$$

Actually, this application, to a certain extent, provided the motivation for the discussion of Chebyshev approximation problems connected via the polynomial mapping  $\varphi_m$ . As we will describe in Section 3, the convergence rate of Krylov subspace methods may be given in terms of so-called *optimal polynomials*. For a given parameter  $c \notin R$  (typically  $c = 0$ ), such a polynomial  $P_n^R$  is the solution of the constrained Chebyshev approximation problem

$$(1.3) \quad \|P_n^R\|_R = \|1 - (z - c)q_{n-1}^*(z)\|_R = \min_{q_{n-1} \in \Pi_{n-1}} \|1 - (z - c)q_{n-1}(z)\|_R,$$

where typically  $R$  is a set which contains all eigenvalues of the given matrix  $\mathcal{A}$ .

When using iterative methods, preconditioning is an important issue. Here, one is looking for a preconditioning matrix  $\mathcal{M}$  such that the new system  $\mathcal{M}^{-1}\mathcal{A}\mathbf{x} = \mathcal{M}^{-1}\mathbf{b}$  is easier to solve than the original system  $\mathcal{A}\mathbf{x} = \mathbf{b}$ . Let  $S$  denote a set which contains all eigenvalues of  $\mathcal{M}^{-1}\mathcal{A}$ . Naturally, one is interested in studying the convergence properties of the preconditioned system as compared to the original system. It turns out that for certain classes of preconditioners and certain linear systems the eigenvalue inclusion sets and sometimes even the eigenvalues themselves are related via  $R = \varphi_m(S)$ . Thus, our analysis enables one to relate the corresponding convergence rates to each other. See Section 3 for details.

Finally, we note that it is often easier to compute the *Chebyshev Polynomial*  $T_n^R$  than the optimal polynomial. It is the solution of the Chebyshev approximation problem

$$(1.4) \quad \|T_n^R\|_R = \|z^n - p_{n-1}^*\|_R = \min_{p_{n-1} \in \Pi_{n-1}} \|z^n - p_{n-1}(z)\|_R.$$

Observe that the scaled Chebyshev polynomial  $T_n^R(z)/T_n^R(c)$  always provides an upper bound for the norm of the optimal polynomial

$$\|P_n^R\|_R \leq \left\| \frac{T_n^R}{T_n^R(c)} \right\|_R,$$

where equality holds for some sets  $R$ . Examples of such sets include single intervals (cf. Markoff [8]) and certain ellipses (cf. Fischer and Freund [4]).

**2. Best approximations and extremal signatures on inverse images of polynomial mappings.** Throughout this section we assume that  $\varphi_m \in \Pi_m \setminus \Pi_{m-1}$  is a polynomial of degree  $m$  with leading coefficient  $a_m$ . Moreover, let  $z_j \in \mathbb{C}$  be a given point, then we denote by  $z_{j,1}, z_{j,2}, \dots, z_{j,m}$  the zeros of  $\varphi_m(z) - z_j$ . We always assume that these zeros are simple, i.e.,  $\varphi_m$  has a full set of inverse branches  $z_{j,l} = \varphi_{m,l}^{-1}(z_j)$ . This assumption, for example, implies that we have the partial fraction expansion

$$(2.1) \quad \frac{1}{\varphi_m(z) - z_j} = \sum_{l=1}^m \frac{1}{\varphi'_m(z_{j,l})} \frac{1}{(z - z_{j,l})}.$$

For the proof of our first theorem, the following well-known characterization of best approximation will be useful (cf. Rivlin and Shapiro [11]). The function  $g^*$  is a best approximation of  $f$  with respect to  $V_n$  (cf. (1.1)) if, and only if, there exist  $r$  extremal points  $z_1, z_2, \dots, z_r \in \{z \in R : |(f - g^*)(z)| = \|f - g^*\|_R\}$  and positive numbers  $\mu_1, \mu_2, \dots, \mu_r \in \mathbb{R}_+$  ( $r \leq 2n + 1$  in the complex case and  $r \leq n + 1$  in the real case) such that

$$(2.2) \quad \sum_{j=1}^r \mu_j \operatorname{sgn} \overline{(f - g^*)(z_j)} g_k(z_j) = 0 \quad \text{for } k = 1, 2, \dots, n.$$

The set  $\{(z_j, \mu_j \operatorname{sgn} \overline{(f - g^*)(z_j)}) \mid j = 1, 2, \dots, r\}$  is called an *extremal signature* for  $f - g^*$  on  $R$  with respect to  $V_n$ .

**THEOREM 2.1.** *Let  $R \subset \mathbb{C}$ ,  $f \in C(R)$ , and  $S = S(R, \varphi_m)$  be given. If  $p_{n-1}^*$  is the best approximation of  $f$  with respect to  $\Pi_{n-1}$  on  $R$*

$$\|f - p_{n-1}^*\|_R = \min_{p_{n-1} \in \Pi_{n-1}} \|f - p_{n-1}\|_R,$$

then  $p_{n-1}^* \circ \varphi_m$  is the best approximation of  $f \circ \varphi_m$  with respect to  $\Pi_{mn-1}$  on  $S$

$$\|f \circ \varphi_m - p_{n-1}^* \circ \varphi_m\|_S = \min_{p_{mn-1} \in \Pi_{mn-1}} \|f \circ \varphi_m - p_{mn-1}\|_S.$$

Furthermore, if  $\{(z_j, \mu_j \operatorname{sgn} \overline{(f - p_{n-1}^*)(z_j)}) \mid j = 1, 2, \dots, r\}$  is an extremal signature for  $f - p_{n-1}^*$  on  $R$ , i.e.,

$$(2.3) \quad \sum_{j=1}^r \mu_j \operatorname{sgn} \overline{(f - p_{n-1}^*)(z_j)} z_j^k = 0, \quad \text{for } k = 0, 1, \dots, n-1,$$

then

$$\sum_{j=1}^r \sum_{l=1}^m \mu_j \operatorname{sgn} \overline{(f - p_{n-1}^*)(\varphi_m(z_{j,l}))} z_{j,l}^k = 0, \quad \text{for } k = 0, 1, \dots, nm-1,$$

where  $z_{j,1}, z_{j,2}, \dots, z_{j,m}$  denote the zeros of  $\varphi_m(z) - z_j$ , for  $j = 1, 2, \dots, r$ . That is,  $\{(z_{j,l}, \mu_{j,l} \operatorname{sgn} \overline{(f - p_{n-1}^*)(\varphi_m(z_{j,l}))}) \mid j = 1, 2, \dots, r, l = 1, 2, \dots, m\}$ , with  $\mu_{j,l} := \mu_j$ , is an extremal signature for  $(f - p_{n-1}^*) \circ \varphi_m$  on  $S$  with respect to  $\Pi_{nm-1}$ .

*Proof.* For convenience we set

$$\hat{\mu}_j := \mu_j \operatorname{sgn} \overline{(f - p_{n-1}^*)(\varphi_m(z_{j,l}))} = \mu_j \operatorname{sgn} \overline{(f - p_{n-1}^*)(z_j)}.$$

We then have to show that

$$\sum_{j=1}^r \hat{\mu}_j \sum_{l=1}^m z_{j,l}^k = 0, \quad \text{for } k = 0, 1, \dots, mn - 1.$$

To start with, we note that by construction the points  $z_{j,l}$  are extremal points for  $(f - p_{n-1}^*) \circ \varphi_m$ . Furthermore, in view of (2.1) and (2.3) we obtain

$$\begin{aligned} \sum_{j=1}^r \hat{\mu}_j \sum_{l=1}^m \frac{1}{\varphi_m'(z_{j,l})} \frac{1}{(z - z_{j,l})} &= \sum_{j=1}^r \frac{\hat{\mu}_j}{\varphi_m(z) - z_j} \\ &= \sum_{j=1}^r \hat{\mu}_j \left( \sum_{k=0}^{\infty} \frac{z_j^k}{\varphi_m(z)^{k+1}} \right) \\ &= \sum_{k=0}^{\infty} \left( \sum_{j=1}^r \hat{\mu}_j z_j^k \right) \frac{1}{\varphi_m(z)^{k+1}} \\ &= \mathcal{O} \left( \frac{1}{z^{m(n+1)}} \right), \end{aligned}$$

as  $z \rightarrow \infty$ . On the other hand, we have

$$\sum_{j=1}^r \hat{\mu}_j \sum_{l=1}^m \frac{1}{\varphi_m'(z_{j,l})} \frac{1}{(z - z_{j,l})} = \sum_{k=0}^{\infty} \left( \sum_{j=1}^r \hat{\mu}_j \sum_{l=1}^m \frac{z_{j,l}^k}{\varphi_m'(z_{j,l})} \right) \frac{1}{z^{k+1}},$$

and consequently

$$\sum_{j=1}^r \hat{\mu}_j \sum_{l=1}^m \frac{q(z_{j,l})}{\varphi_m'(z_{j,l})} = 0, \quad \text{for all } q \in \Pi_{m(n+1)-2}.$$

For the special choice  $q(z) = z^k \varphi_m'(z)$  we obtain

$$\sum_{j=1}^r \hat{\mu}_j \sum_{l=1}^m z_{j,l}^k = 0, \quad \text{for } k = 0, 1, \dots, mn - 1,$$

which concludes the proof.  $\square$

With  $f(z) = z^n$  we immediately arrive at the following corollary.

**COROLLARY 2.2.** *Let  $R \subset \mathbb{C}$ , and  $S = S(R, \varphi_m)$  be given. If  $T_n^R$  denotes the Chebyshev polynomial with respect to  $R$  (cf. (1.4)), then the Chebyshev polynomial of degree  $mn$  with respect to  $S$  is given by*

$$T_{mn}^S(z) = \frac{1}{a_m^n} T_n^R(\varphi_m(z)).$$

We remark that one may also find Corollary 2.2 in Kamo and Borodin [6]. Their proof, however, is based on Kolmogorov's criterion and therefore does not provide an extremal

signature for  $T_n^R \circ \varphi_m$ . Furthermore, we note that the above theorem includes the result of Lebedev [7] for the case where  $S$  is the union of finitely many intervals on the real line.

Next, we turn our attention to the computation of optimal polynomials or, more general, to the case of constrained Chebyshev approximation. For constrained Chebyshev problems on two and several intervals see, for example, Fischer [2], Peherstorfer and Schiefermayr [10] and references therein.

The proof of the next theorem is along the lines of the proof of Theorem 2.1 and is therefore omitted.

**THEOREM 2.3.** *Let  $R \subset \mathbb{C}$ ,  $c \notin R$ ,  $f \in C(R)$ , and  $S = S(R, \varphi_m)$  be given. If  $Q_{n-1}^* = (z - c)q_{n-1}^*$  is the best approximation of  $f$  with respect to  $(z - c)\Pi_{n-1}$  on  $R$*

$$\|f - Q_{n-1}^*\|_R = \min_{q_{n-1} \in \Pi_{n-1}} \|f - (z - c)q_{n-1}\|_R,$$

*then  $Q_{n-1}^* \circ \varphi_m = (\varphi_m(z) - c)q_{n-1}^* \circ \varphi_m$  is the best approximation of  $f \circ \varphi_m$  with respect to  $(\varphi_m(z) - c)\Pi_{mn-1}$  on  $S$*

$$\|f \circ \varphi_m - Q_{n-1}^* \circ \varphi_m\|_S = \min_{q_{mn-1} \in \Pi_{mn-1}} \|f \circ \varphi_m - (\varphi_m(z) - c)q_{mn-1}\|_S.$$

*Furthermore, if  $\{(z_j, \mu_j \operatorname{sgn} \overline{(f(z_j) - Q_{n-1}^*(z_j))}) \mid j = 1, 2, \dots, r\}$  is an extremal signature for  $f - Q_{n-1}^*$  on  $R$ , i.e.,*

$$\sum_{j=1}^r \mu_j \operatorname{sgn} \overline{(f(z_j) - Q_{n-1}^*(z_j))} (z_j - c) z_j^k = 0, \quad \text{for } k = 0, 1, \dots, n-1,$$

*then, for  $k = 0, 1, \dots, nm - 1$ ,*

$$\sum_{j=1}^r \sum_{l=1}^m \mu_j \operatorname{sgn} \overline{(f(\varphi_m(z_{j,l})) - Q_{n-1}^*(\varphi_m(z_{j,l})))} (\varphi_m(z_{j,l}) - c) z_{j,l}^k = 0,$$

*where  $z_{j,1}, z_{j,2}, \dots, z_{j,m}$  denote the zeros of  $\varphi_m(z) - z_j$ , for  $j = 1, 2, \dots, r$ . That is,  $\{(z_{j,l}, \mu_{j,l} \operatorname{sgn} \overline{(f(\varphi_m(z_{j,l})) - Q_{n-1}^*(\varphi_m(z_{j,l})))}) \mid j = 1, 2, \dots, r, l = 1, 2, \dots, m\}$ , with  $\mu_{j,l} := \mu_j$ , is an extremal signature for  $f \circ \varphi_m - Q_{n-1}^* \circ \varphi_m$  on  $S$  with respect to  $\Pi_{nm-1}$ .*

The special case  $f(z) = 1$  gives rise to the following corollary.

**COROLLARY 2.4.** *Let  $R \subset \mathbb{C}$ ,  $c \notin R$ , and  $S = S(R, \varphi_m)$  be given. If  $P_n^R = 1 - (z - c)q_{n-1}^*$  denotes the optimal polynomial with respect to  $R$  and to  $c$  (cf. (1.3)), then we have*

$$\|P_n^R \circ \varphi_m\|_S = \min_{q_{mn-1} \in \Pi_{mn-1}} \|1 - (\varphi_m(z) - c)q_{mn-1}\|_S.$$

It is worth noticing that the corollary above in general does not imply that  $P_n^R \circ \varphi_m$  is the optimal polynomial for the set  $S$ . This would be the case, if

$$\|P_n^R \circ \varphi_m\|_S = \|1 - (\varphi_m(z) - c)q_{n-1}^* \circ \varphi_m\|_S = \min_{q_{mn-1} \in \Pi_{mn-1}} \|1 - (z - e)q_{mn-1}\|_S,$$

where  $e$  is a zero of  $\varphi_m(z) - c$ .

In the remaining part of this section we will further investigate polynomial mappings  $\varphi_2$  of degree two. We start with the case where  $S$  is the union of two disjoint intervals on the real line.

**THEOREM 2.5.** *Let  $R = [a, b]$ ,  $a < b \in \mathbb{R}$ , and  $c \in \mathbb{R} \setminus [a, b]$  be given. Furthermore, let  $\varphi_2$  be such that  $S = S(R, \varphi_2)$  is the union of two intervals and let  $e_1, e_2$  denote the*

zeros of  $\varphi_2(z) - c$ . If  $P_n^R$  denotes the optimal polynomial with respect to  $R$  and to  $c$ , then  $P_{2n}^S = P_n^R \circ \varphi_2$  is the optimal polynomial with respect to  $S$  and to  $e_1$  and  $e_2$ , respectively,

$$\begin{aligned} \|P_n^R \circ \varphi_2\|_S &= \min_{q_{2n-1} \in \Pi_{2n-1}} \|1 - (z - e_1)q_{2n-1}\|_S \\ &= \min_{q_{2n-1} \in \Pi_{2n-1}} \|1 - (z - e_2)q_{2n-1}\|_S. \end{aligned}$$

Moreover,  $P_{2n}^S$  is even optimal for  $\Pi_{2n+1}$ , i.e., we have  $P_{2n+1}^S = P_{2n}^S$ .

*Proof.* Let  $T_n^R$  denote the Chebyshev polynomial with respect to  $R$ . By Corollary 2.2 we have that the Chebyshev polynomial with respect to  $S$  is given by  $T_{2n}^S(z) = T_n^R(\varphi_2(z))/a_2^2$ . Since  $T_n^R$  is a shifted version of the classical Chebyshev polynomials on the unit interval, it has precisely  $n + 1$  extremal points on  $R$ . Now, Theorem 2.1 implies that  $T_{2n}^S$  has  $2(n + 1)$  extremal points on  $S$ . Finally, the assertion follows from Theorem 4.3 in Fischer [2], which states that the optimal polynomial with respect to  $S$  is a scaled Chebyshev polynomial

$$P_{2n}^S(z) = \frac{T_{2n}^S(z)}{T_{2n}^S(e_1)} = \frac{T_{2n}^S(z)}{T_{2n}^S(e_2)},$$

if  $T_{2n}^S$  has  $2n + 2$  extremal points on  $S$ . Actually, the fact that  $T_{2n}^S$  has one additional extremal point also implies that  $P_{2n+1}^S = P_{2n}^S$ , cf. Corollary 3.3.6(b) in Fischer [3].  $\square$

Next, we analyze the special mappings  $\varphi_2^\pm(z) = \pm z(z - 1)$ .

**THEOREM 2.6.** *Let  $R \subset \mathbb{C}$ ,  $0 \notin R$ , and  $S^\pm = S(R, \varphi_2^\pm)$  with  $\varphi_2^\pm(z) = \pm z(z - 1)$  be given. If  $P_n^R$  denotes the optimal polynomial with respect to  $R$  and to 0, then the optimal polynomials of degree  $2n$  and  $2n + 1$  with respect to  $S^\pm$  and to 0 and 1, respectively, are given by*

$$P_{2n+1}^{S^\pm} = P_{2n}^{S^\pm} = P_n^R \circ \varphi_2^\pm$$

with

$$\|P_{2n+1}^{S^+}\|_{S^+} = \|P_{2n}^{S^+}\|_{S^+} = \|P_n^R\|_R = \|P_{2n}^{S^-}\|_{S^-} = \|P_{2n+1}^{S^-}\|_{S^-}.$$

*Proof.* We only consider the case  $\varphi_2 = \varphi_2^+$  with  $S = S^+$ . We start by noting that the polynomial  $\varphi_2$  satisfies the symmetry relation  $\varphi_2(l(z)) = \varphi_2(z)$  with  $l(z) = 1 - z$ . Hence, by construction, we have  $S = l(S)$ . From the fact that the optimal polynomial is uniquely defined and

$$\max_{z \in S} |P_k^S(z)| = \max_{z \in S} |P_k^S(l(z))|,$$

we obtain the identity  $P_k^S(z) = P_k^S(l(z))$ . For odd degree polynomials this symmetry relation implies that the leading coefficient vanishes, because  $P_{2n+1}^S(z) = \alpha_{2n+1}z^{2n+1} + \dots$  and  $P_{2n+1}^S(l(z)) = \alpha_{2n+1}(1-z)^{2n+1} + \dots = -\alpha_{2n+1}z^{2n+1} + \dots$ . Now, let us consider the even degree case  $k = 2n$ . Here, we deduce from the symmetry relation  $P_{2n}^S(z) = P_{2n}^S(l(z))$  that all zeros of  $P_{2n}^S$  come in pairs  $P_{2n}^S(z_j) = P_{2n}^S(l(z_j)) = 0$ ,  $j = 1, 2, \dots, n$ . Let  $q_n$  denote the polynomial with the zeros  $y_j = \varphi_2(z_j)$ ,  $j = 1, 2, \dots, n$ , and constant term 1. Since  $P_{2n}^S - q_n \circ \varphi_2 \in \Pi_{2n}$  has  $2n + 1$  zeros, namely  $z_j, l(z_j), j = 1, 2, \dots, n$  and  $z_{n+1} = 0$ , we conclude that  $P_{2n}^S = q_n \circ \varphi_2$ , which proves the assertion.  $\square$

**3. Application.** In this section we will analyze the convergence of iterative methods when applied to the solution of a non-symmetric, nonsingular linear system

$$\mathcal{A}\mathbf{x} = \mathbf{b}, \quad \mathcal{A} \in \mathbb{C}^{N \times N}.$$

Throughout this section we assume for convenience that  $\mathcal{A}$  is diagonalizable. Some of the most effective iterative methods available are those of Krylov subspace type which have in-built minimization properties. Here we consider minimal residual methods, whose iterates minimize the Euclidian norm of the residual  $\mathbf{r}_n$  at each step. More precisely, for the methods under consideration the  $n$ th residual  $\mathbf{r}_n$  may be written as

$$\mathbf{r}_n = \mathbf{b} - \mathcal{A}\mathbf{x}_n = p_n(\mathcal{A})\mathbf{r}_0,$$

where  $p_n \in \Pi_n$  is a polynomial of exact degree  $n$  satisfying the interpolatory constraint  $p_n(0) = 1$ . The polynomial  $p_n$  and consequently the iterate  $\mathbf{x}_n$  are uniquely determined by the minimization property

$$\|\mathbf{r}_n\|_2 = \min\{\|p(\mathcal{A})\mathbf{r}_0\|_2; p \in \Pi_n, p(0) = 1\}.$$

The actual implementation of such a method depends on the properties of the given coefficient matrix  $\mathcal{A}$ . Among the most well-known methods belonging to this class are CR (for  $\mathcal{A}$  symmetric and positive definite), MINRES (for  $\mathcal{A}$  symmetric but indefinite) and GMRES (for  $\mathcal{A}$  non-symmetric) (cf., e.g., Saad [12]).

Let  $V$  denote the eigenvector matrix of  $\mathcal{A}$ . Then the 2-norm of the residual may be bounded by the following standard estimate

$$(3.1) \quad \frac{\|\mathbf{r}_n\|_2}{\|\mathbf{r}_0\|_2} \leq \kappa_2(V) \min_{p \in \Pi_n, p(0)=1} \max_{\lambda \in \sigma(\mathcal{A})} |p(\lambda)| = \kappa_2(V) \|P_n^{\sigma(\mathcal{A})}\|_{\sigma(\mathcal{A})},$$

where  $\sigma(\mathcal{A})$  denotes the spectrum of  $\mathcal{A}$  and  $\kappa_2(V)$  denotes the condition number of the eigenvector matrix  $V$ . In conclusion, to bound the convergence rate of minimal residual methods one is tempted to compute the norm of the optimal polynomial  $P_n^{\sigma(\mathcal{A})}$  (cf. (1.3)) with respect to  $\sigma(\mathcal{A})$  and to 0.

In the remaining part of this section we will identify matrices  $\mathcal{A}$  and preconditioner  $\mathcal{M}$  such that the spectrum of  $\mathcal{A}$  and its preconditioned version  $\mathcal{M}^{-1}\mathcal{A}$  are related by a polynomial mapping

$$\varphi_k(\sigma(\mathcal{A})) = \sigma(\mathcal{M}^{-1}\mathcal{A}),$$

with

$$P_{kn}^{\sigma(\mathcal{A})} = P_n^{\sigma(\mathcal{M}^{-1}\mathcal{A})} \circ \varphi_k,$$

and

$$\|P_{kn}^{\sigma(\mathcal{A})}\|_{\sigma(\mathcal{A})} = \|P_n^{\sigma(\mathcal{M}^{-1}\mathcal{A})} \circ \varphi_k\|_{\sigma(\mathcal{A})} = \|P_n^{\sigma(\mathcal{M}^{-1}\mathcal{A})}\|_{\sigma(\mathcal{M}^{-1}\mathcal{A})}.$$

These equations together with (3.1) have an important consequence. They indicate, at least in the absence of roundoff errors, that a minimal residual method converges about  $k$ -times faster for the preconditioned system than for the original system.

Let us now consider linear systems of the form

$$(3.2) \quad \mathcal{A} = \begin{bmatrix} A & B^T \\ B & 0 \end{bmatrix} \in \mathbb{C}^{(n+m) \times (n+m)},$$

where the  $n \times n$  matrix  $A$  is regular and the  $m \times n$ ,  $m < n$ , matrix  $B$  has full rank. The efficient solution of such systems plays an important role in many different applications. For example, the mixed finite element approximation of the Stokes problem (cf. [5]) leads to a system of the form (3.2) with symmetric positive definite  $A$ , whereas the discretization of the Oseen problem (cf. [1]) results in a matrix problem with non-symmetric  $A$ . Other examples include linear least squares problems and linear KKT systems.

Given the structure of  $\mathcal{A}$ , an obvious choice is to use block preconditioning. Furthermore, often a fast solver for the top-left block  $A$  is available. Therefore, preconditioners of the form

$$(3.3) \quad \mathcal{M}(F, M) = \begin{bmatrix} A & F \\ 0 & M \end{bmatrix},$$

where  $M$  is nonsingular are quite popular (cf. [1]). To provide an efficient implementation of the preconditioner, common choices for the top-right block  $F$  are  $F = 0$  and  $F = B^T$ , respectively.

In general, there is some flexibility in applying the preconditioner. One may consider a left  $\mathcal{M}^{-1}\mathcal{A}$ , a right  $\mathcal{A}\mathcal{M}^{-1}$ , or a centrally preconditioned coefficient matrix  $\mathcal{M}^{-1/2}\mathcal{A}\mathcal{M}^{-1/2}$ , respectively. For the latter case the matrix  $\mathcal{M}$  is assumed to be positive definite. However, it is easy to see, that the three preconditioned matrices share the same eigenvalues. They are given by the solution of the generalized eigenvalue problem

$$(3.4) \quad \mathcal{A}\mathbf{v} = \lambda\mathcal{M}\mathbf{v}.$$

The next lemma relates the wanted solutions  $\lambda$  of (3.4) for the specific choices  $\mathcal{M}(0, \pm M)$  and  $\mathcal{M}(B^T, -M)$  to the generalized eigenvalues  $\mu$  of the so-called Schur complement

$$(3.5) \quad BA^{-1}B^T\mathbf{p} = \mu M\mathbf{p}.$$

In particular, it shows that the two sets of eigenvalues are related by a polynomial mapping. The proof is straightforward and might be found in Elman and Silvester [1].

LEMMA 3.1. *Let  $\mu_k$ ,  $k = 1, 2, \dots, m$ , denote the solutions of the generalized eigenvalue problem (3.5).*

(i) *Let  $\varphi_2^+(z) = z(z-1)$ . The solutions of  $\mathcal{A}\mathbf{v} = \lambda\mathcal{M}(0, M)\mathbf{v}$  are  $\lambda_0 = 1$ , of multiplicity  $n-m$ , and*

$$\lambda_k^\pm = \frac{1 \pm \sqrt{1 + 4\mu_k}}{2}, \quad \text{i.e., } \varphi_2^+(\lambda_k^\pm) = \mu_k, \quad k = 1, 2, \dots, m.$$

(ii) *Let  $\varphi_2^-(z) = -z(z-1)$ . The solutions of  $\mathcal{A}\mathbf{v} = \lambda\mathcal{M}(0, -M)\mathbf{v}$  are  $\lambda_0 = 1$ , of multiplicity  $n-m$ , and*

$$\lambda_k^\pm = \frac{1 \pm \sqrt{1 - 4\mu_k}}{2}, \quad \text{i.e., } \varphi_2^-(\lambda_k^\pm) = \mu_k, \quad k = 1, 2, \dots, m.$$

(iii) *The solutions of  $\mathcal{A}\mathbf{v} = \lambda\mathcal{M}(B^T, -M)\mathbf{v}$  are  $\lambda_0 = 1$ , of multiplicity  $n$ , and*

$$\lambda_k = \mu_k, \quad k = 1, 2, \dots, m.$$

Let us now investigate the case that the matrices  $A$  and  $M$  are symmetric and positive definite. Here, it is easy to see that the corresponding matrix  $\mathcal{A}$  is symmetric but indefinite and that the preconditioner  $\mathcal{M}(0, M)$  is symmetric and positive definite. It follows from Lemma 3.1(i) that the preconditioned matrix  $\mathcal{A}^+ = \mathcal{M}(0, M)^{-1/2}\mathcal{A}\mathcal{M}(0, M)^{-1/2}$



is indefinite and that its eigenvalues are, apart from  $\lambda = 1$ , symmetric about the point  $1/2$ . On the other hand, the alternative preconditioner  $\mathcal{M}(0, -M)$  is itself indefinite and therefore may not be applied centrally. As a consequence, the preconditioned system  $\mathcal{A}^- = \mathcal{M}(0, -M)^{-1}\mathcal{A}$  (or  $\mathcal{A}^- = \mathcal{A}\mathcal{M}(0, -M)^{-1}$ ) is no longer symmetric. However, all its eigenvalues are located in the right-half plane. Lemma 3.1(ii) shows that apart from  $\lambda = 1$ , all of them are either located on a cross with center  $1/2$  or on a real interval.

There has been some discussion in the literature (cf. [5]) on whether MINRES applied to the symmetric indefinite system  $\mathcal{A}^+$  or GMRES applied to one of the non-symmetric systems  $\mathcal{A}^-$  will converge faster. The next theorem shows that at least the convergence rates are identical. The proof is a direct consequence of (3.1), Lemma 3.1, and Theorem 2.6.

**THEOREM 3.2.** *Let  $V^-$  denote the eigenvector matrix of  $\mathcal{A}^-$  and let  $\mathbf{r}_0 = (0, v)^T$  with  $v \in \mathbb{R}^m$ . Furthermore, let  $R = \{\mu_j : j = 1, 2, \dots, m\}$  denote the generalized eigenvalues of the Schur complement (3.5) and let  $S^\pm = \sigma(\mathcal{A}^\pm) \setminus \{1\}$ . Then the 2-norm of the residuals  $\mathbf{r}_k^+ = p_k^+(\mathcal{A}^+)\mathbf{r}_0$  and  $\mathbf{r}_k^- = p_k^-(\mathcal{A}^-)\mathbf{r}_0$  may be estimated as follows*

$$\frac{\|\mathbf{r}_k^+\|_2}{\|\mathbf{r}_0\|_2} \leq \|P_{2n}^{S^+}\|_{S^+} = \|P_n^R\|_R,$$

$$\frac{\|\mathbf{r}_k^-\|_2}{\|\mathbf{r}_0\|_2} \leq \kappa_2(V^-) \|P_{2n}^{S^-}\|_{S^-} = \kappa_2(V^-) \|P_n^R\|_R,$$

for  $k \in \{2n, 2n + 1\}$ .

Some remarks are in place. The above theorem is essentially a result on polynomial approximation. It says nothing about the performance of MINRES and GMRES in finite precision arithmetic. However, all our numerical test runs (cf. also [5]) have shown precisely the predicted convergence behavior. Apart from the constant factors  $\kappa_2(V^-)$  the bounds of the preceding theorem are the same. Note that  $\kappa_2(V^+) = 1$ . Actually, an analysis based on orthogonal polynomials shows that even the iterates  $\mathbf{x}_{2n+1}^+ = \mathbf{x}_{2n}^+ = \mathbf{x}_{2n}^- = \mathbf{x}_{2n+1}^-$  are the same. Furthermore, the special choice  $\mathbf{r}_0 = (0, v)^T$  ensures that no eigendirection associated with the eigenvalue 1 exists in the initial residual (for details see [5]).

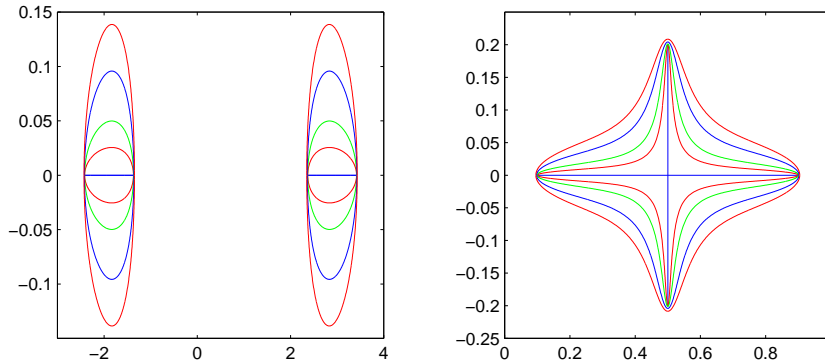


FIG. 3.1.  $(\varphi_2^+)^{-1}(E)$  (left) and  $(\varphi_2^-)^{-1}(E)$  (right) for various ellipses  $E$ .

The above analysis is based on the fact that the top-left block  $A$  in (3.2) is inverted in exact arithmetic. In practice, however, only an approximation to the inverse is available. Here the spectra of the preconditioned matrices are contained in regions which do include  $S^\pm$ . By Lemma 3.1 we have  $S^\pm = (\varphi_2^\pm)^{-1}(R)$ . Figure 3.1 shows, apart from  $S^\pm$ , the preimage of

$\varphi_2^\pm$  of ellipses  $E$  containing the set  $R$ . With the help of Theorem 2.6 it is again possible to explicitly compute the optimal polynomials with respect to the depicted enlarged versions of  $S^\pm$ .

Finally, we turn our attention to the non-symmetric case. That is, we assume that the matrix  $A$  in (3.2) is non-symmetric. Following Elman and Silvester [1] we investigate the block diagonal preconditioned system  $\mathcal{A}^{BD} = \mathcal{A}\mathcal{M}(0, M)^{-1}$  and the block triangular preconditioned system  $\mathcal{A}^{BT} = \mathcal{A}\mathcal{M}(B^T, -M)^{-1}$ , respectively. We learn from Lemma 3.1 that, apart from the eigenvalue  $\lambda = 1$ , the eigenvalues  $\lambda$  of  $\mathcal{A}^{BD}$  and the eigenvalues  $\mu$  of  $\mathcal{A}^{BT}$  are related by the quadratic mapping  $\mu = \varphi_2(\lambda)$  with  $\varphi_2(z) = z(z - 1)$ .

Actually, for the Oseen operator Elman and Silvester [1] proved that all eigenvalues of  $\mathcal{A}^{BT}$  are contained in a rectangular box  $R$  in the right half plane. Consequently, the eigenvalues of  $\mathcal{A}^{BD}$  are contained in sets  $S$  which are the preimages of  $R$ . The next figure displays some typical inclusion sets.

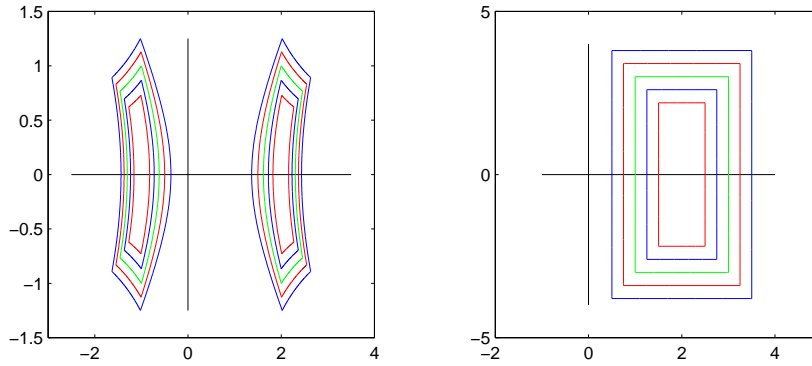


FIG. 3.2.  $S$  (left) and  $R = \varphi_2(S)$  (right) for different sets  $R$ .

Also, Elman and Silvester reported on some numerical test runs with these preconditioners. They observed that GMRES applied to  $\mathcal{A}^{BD}$  took twice as many iterations as GMRES applied to  $\mathcal{A}^{BT}$  to reach the same reduction of the initial norm. Moreover, the norm of the odd iterates stagnates for the diagonal preconditioning.

The next theorem provides an explanation for their observations in terms of the respective convergence rates. Again, the proof follows directly from (3.1), Lemma 3.1, and Theorem 2.6.

**THEOREM 3.3.** *Let  $V^{BD}$  and  $V^{BT}$  denote the eigenvector matrix of  $\mathcal{A}^{BD}$  and  $\mathcal{A}^{BT}$ , respectively and let  $\mathbf{r}_0 = (0, v)^T$  with  $v \in \mathbb{R}^m$ . Furthermore, let  $R = \{\mu_j : j = 1, 2, \dots, m\}$  denote the generalized eigenvalues of the Schur complement (3.5) and let  $S^{BD} = \sigma(\mathcal{A}^{BD}) \setminus \{1\}$ . Then the 2-norm of the residuals  $\mathbf{r}_k^{BD} = p_k^{BD}(\mathcal{A}^{BD})\mathbf{r}_0$  and  $\mathbf{r}_k^{BT} = p_k^{BT}(\mathcal{A}^{BT})\mathbf{r}_0$  may be estimated as follows*

$$\frac{\|\mathbf{r}_k^{BD}\|_2}{\|\mathbf{r}_0\|_2} \leq \kappa_2(V^{BD}) \|P_{2n}^{S^{BD}}\|_{S^{BD}} = \kappa_2(V^{BD}) \|P_n^R\|_R, \quad k \in \{2n, 2n + 1\},$$

$$\frac{\|\mathbf{r}_n^{BT}\|_2}{\|\mathbf{r}_0\|_2} \leq \kappa_2(V^{BT}) \|P_n^R\|_R.$$

**Acknowledgments.** Part of this work was done while the first author was visiting the Laboratoire d'Analyse Numérique et d'Optimisation at the Université des Sciences et Tech-

nologies de Lille. He would like to thank all members of the Institute for their warm hospitality. Finally, we thank the reviewers for helpful suggestions and comments.

## REFERENCES

- [1] H.C. ELMAN AND D.J. SILVESTER, *Fast Nonsymmetric iterations and preconditioning for Navier-Stokes equations*, SIAM J. Sci. Comput., 17 (1996), pp. 33-46.
- [2] B. FISCHER, *Chebyshev polynomials for disjoint compact sets*, Constr. Approx., 8 (1992), pp. 309-329.
- [3] B. FISCHER, *Polynomial Based Iteration Methods for Symmetric Linear Systems*, Wiley-Teubner, Chicester, Stuttgart, 1996.
- [4] B. FISCHER AND R. FREUND, *On the constraint Chebyshev approximation problem on ellipses*, J. Approx. Theory, 62 (1990), pp. 297-315.
- [5] B. FISCHER, A. RAMAGE, D.J. SILVESTER, AND A.J. WATHEN, *Minimum residual methods for augmented systems*, BIT, 38 (1998), pp. 527-543.
- [6] S.O. KAMO AND P.A. BORODIN, *Chebyshev polynomials for Julia sets*, Moscow Univ. Math. Bull., 49 (1994), pp. 44-45.
- [7] V.I. LEBEDEV, *Iterative methods for solving operator equations with a spectrum contained in several intervals*, Comput. Math. Math. Phys., 9 (1969), pp. 17-24.
- [8] W. MARKOFF, *Über Polynome, die in einem gegebenen Intervalle möglichst wenig von Null abweichen*, Math. Ann., 77 (1916), pp. 213-258.
- [9] F. PEHERSTORFER, *Minimal polynomials for compact sets of the complex plane*, Constr. Approx., 12 (1996), pp. 481-488.
- [10] F. PEHERSTORFER AND K. SCHIEFERMAYR, *Theoretical and numerical description of extremal polynomials on several intervals II*, Acta Math. Hungar., 83 (1999), pp. 103-128.
- [11] T.J. RIVLIN AND H.S. SHAPIRO, *A unified approach to certain problems of approximation and minimization*, J. Soc. Indust. Appl. Math., 9 (1961), pp. 670-699.
- [12] Y. SAAD, *Iterative Methods for Sparse Linear Systems*, PWS, Boston, 1996.