

## COMPUTABLE CONVERGENCE BOUNDS FOR GMRES\*

JÖRG LIESEN†

**Abstract.** The purpose of this paper is to derive new computable convergence bounds for GMRES. The new bounds depend on the initial guess and are thus conceptually different from standard “worst-case” bounds. Most importantly, approximations to the new bounds can be computed from information generated during the run of a certain GMRES implementation. The approximations allow predictions of how the algorithm will perform. Heuristics for such predictions are given. Numerical experiments illustrate the behavior of the new bounds as well as the use of the heuristics.

**Key words.** linear systems, convergence analysis, GMRES method, Krylov subspace methods, iterative methods

**AMS subject classification.** 65F10

**PII.** S0895479898341669

**1. Introduction.** The GMRES algorithm by Saad and Schultz [28] is a popular iterative method for solving systems of linear equations:

$$(1.1) \quad Ax = b, \quad A \in \mathbf{C}^{N \times N}, \quad b \in \mathbf{C}^N.$$

Given an initial guess  $x_0$  for the solution of (1.1), GMRES yields iterates  $x_n$  so that

$$(1.2) \quad \|r_n\| := \|b - Ax_n\| = \min_{p_n \in \pi_n} \|p_n(A)r_0\|,$$

where  $\|\cdot\|$  denotes the Euclidean vector and corresponding matrix norm,  $r_0 := b - Ax_0$  is the initial residual, and  $\pi_n$  denotes the set of polynomials of degree at most  $n$  with value 1 at the origin.

The convergence analysis of GMRES has been an active field of research in recent years. Among others, Saad and Schultz [28], Nachtigal, Reddy, and Trefethen [23], Eiermann [8], Campbell et al. [5], and Starke [30] derived upper bounds on the residual norms. An overview of GMRES convergence results is given in [21, section 5.2].

These bounds are typically derived from the “ideal GMRES” approximation problem [16]: Find  $p_n^* \in \pi_n$  so that

$$\|p_n^*(A)\| = \min_{p_n \in \pi_n} \|p_n(A)\|.$$

Being independent of the initial guess, “ideal GMRES” gives an upper bound for all possible GMRES convergence curves for the given matrix  $A$ . Hence, bounds on the GMRES residual norms derived from ideal GMRES intend to describe the algorithm’s worst-case behavior. This approach sometimes leads to sharp bounds, for example, when  $A$  is normal [13, 17]. But as demonstrated in an example by Toh [33], ideal GMRES can overestimate even the worst-case behavior of GMRES arbitrarily high. While the practical implication of Toh’s results and the similar results of others (see,

---

\*Received by the editors June 25, 1998; accepted for publication (in revised form) by A. Greenbaum June 19, 1999; published electronically February 18, 2000. This work was supported by the Studienstiftung des deutschen Volkes, Bonn, Germany.

<http://www.siam.org/journals/simax/21-3/34166.html>

†Fakultät für Mathematik, Universität Bielefeld, 33501 Bielefeld, Germany (liesen@mathematik.uni-bielefeld.de).

e.g., [10]) is not clear yet, none of the bounds derived from ideal GMRES can completely characterize the algorithm’s convergence behavior for any matrix  $A$ . Moreover, while some of the known bounds are only applicable to a restricted class of nonsingular linear systems, others include factors that are not computable. It is generally agreed that the convergence analysis of GMRES needs further work (see, e.g., [12, p. 55], [31, p. 187]).

In this paper we take an alternative approach to deriving convergence bounds for GMRES. Using characterizations of matrices that generate the same GMRES residuals—so-called GMRES-equivalent matrices—we derive bounds that depend on the initial guess. Our bounds are thus a posteriori and do not give new insight into how the convergence of GMRES depends on properties of  $A$ . However, all factors in our bounds can be computed from information generated during the run of WZ-GMRES, the GMRES implementation by Walker and Zhou [35]. Moreover, approximations to the factors in our bounds that are obtained early in the iteration often allow us to *predict how GMRES will perform* in later stages of the iteration. Our theory is applicable to any nonsingular linear system and can even be used in the singular case if a certain condition is satisfied.

The paper is organized as follows. In section 2 we introduce the GMRES algorithm and its mathematical properties. We generalize characterizations of GMRES-equivalent matrices by Greenbaum and Strakoš [15] in section 3. Section 4 includes the derivation of our new convergence bounds. In section 5 we explain how to obtain approximations to the factors in our bounds from information generated during the run of WZ-GMRES. We also devise heuristics for the behavior of these approximations and for the prediction of the performance of GMRES. Numerical examples are presented in section 6, and section 7 contains our concluding remarks.

By  $\Lambda(M)$  we denote the spectrum and by  $M^H$  the hermitian transpose of the matrix  $M \in \mathbf{C}^{N \times N}$ , respectively. If  $M$  is nonsingular, we define  $\kappa(M) := \|M\| \|M^{-1}\|$ , and  $\kappa(M) := \infty$ , otherwise. We define  $K_n(M, v) := \text{span}\{v, Mv, \dots, M^{n-1}v\}$ ,  $n = 1, 2, \dots$ , the  $n$ th Krylov space generated by the matrix  $M$  and the vector  $v \in \mathbf{C}^N$ . For technical reasons we also define  $K_0(A, v) := \{0\}$ .

**2. Mathematical properties of GMRES.** In this section we introduce the GMRES algorithm and state several of its properties. Note that the linear system (1.1) might be singular, if not assumed otherwise, and that we are not concerned with implementational details.

DEFINITION 2.1. *Suppose that  $x_0$  is an initial guess for the solution of a linear system (1.1) and define the initial residual  $r_0 := b - Ax_0$ . For  $n = 1, 2, \dots$ , let*

$$(2.1) \quad x_n \in x_0 + K_n(A, r_0)$$

be a vector, such that

$$(2.2) \quad r_n := b - Ax_n \in r_0 + AK_n(A, r_0)$$

satisfies

$$(2.3) \quad \|r_n\| = \min_{r \in r_0 + AK_n(A, r_0)} \|r\|;$$

then the vectors  $x_n$  and  $r_n$  are called  $n$ th GMRES iterate and residual, respectively.

As a comparison with [28, formulas (4) and (8)] shows, our definition of the GMRES iterates and residuals is mathematically equivalent to their original derivation

by Saad and Schultz. The following well-known result (see, e.g., [6, p. 8]) justifies Definition 2.1.

LEMMA 2.2. *Suppose that  $A \in \mathbf{C}^{N \times N}$  and  $r_0 \in \mathbf{C}^N$ . Then for every  $n = 1, 2, \dots$ , there exists a unique vector  $r_n \in r_0 + AK_n(A, r_0)$  that satisfies (2.3).*

Thus, for  $n = 1, 2, \dots$ , the GMRES residuals are uniquely defined. In case of a nonsingular linear system, uniqueness of the residuals does result in uniqueness of the iterates. Definition 2.1 leads to the following algorithm, which does not, however, require uniqueness of the iterates.

ALGORITHM 2.3 (GMRES).

*Input:* A linear system (1.1) and an initial guess  $x_0$ .

*Initialize:*  $r_0 := b - Ax_0$ ,  $n = 0$

*While*  $r_n \neq 0$

$n = n + 1$

*Construct*  $x_n \in x_0 + K_n(A, r_0)$  *so that*

$$\|r_n\| = \|b - Ax_n\| = \min_{r \in r_0 + AK_n(A, r_0)} \|r\|$$

*End While*

If  $r_n = 0$  for some  $n \geq 0$ , or, equivalently,  $x_n$  is a solution of (1.1), we say that GMRES terminates at step  $n$ . Before analyzing the termination properties of the algorithm, we point out an equivalent way to define the GMRES residuals.

LEMMA 2.4. *Let  $S \subset \mathbf{C}^N$  be a subspace and suppose that  $y_0 \in \mathbf{C}^N$  and  $y_n \in y_0 + S$ . Then  $\|y_n\| = \min_{y \in y_0 + S} \|y\|$  if and only if  $y_n \perp S$ .*

From this classical result (see, e.g., [6, p. 9]) it follows that our definition of the GMRES residuals via (2.2) and (2.3) is equivalent to defining them by (2.2) combined with

$$(2.4) \quad r_n \perp AK_n(A, r_0).$$

The following lemma is a direct consequence of (2.3).

LEMMA 2.5. *The GMRES algorithm applied to a linear system (1.1) and initial guess  $x_0$  terminates at step  $d \geq 0$  if and only if  $d$  is the smallest integer for which  $r_0 \in AK_d(A, r_0)$ .*

To further analyze this situation, we state an important property of the Krylov spaces, which is easy to prove.

LEMMA 2.6. *If  $A \in \mathbf{C}^{N \times N}$  and  $r_0 \in \mathbf{C}^N$ , then*

$$\dim AK_N(A, r_0) = d \text{ for some } d, \quad 0 \leq d \leq N,$$

*is equivalent to*

$$\dim AK_j(A, r_0) = j \text{ for } 0 \leq j \leq d \text{ and}$$

$$\dim AK_j(A, r_0) = d \text{ for all } j \geq d + 1.$$

Since the spaces  $AK_n(A, r_0)$  coincide for all  $n \geq \dim AK_N(A, r_0)$ , the following is an immediate consequence.

COROLLARY 2.7. *Suppose that a linear system (1.1) and an initial guess  $x_0$  are given and define  $d := \dim AK_N(A, r_0)$ . Then for the GMRES residuals  $r_n$ ,  $n \geq d$ ,  $r_n = r_d$  follows.*

Thus, either GMRES terminates at some step  $n \leq d := \dim AK_N(A, r_0) \leq N$ , which requires  $r_0 \in AK_N(A, r_0) = AK_d(A, r_0)$ , or the algorithm fails to terminate, but stagnates after step  $d$ . The next theorem can be shown easily using Lemma 2.6.

**THEOREM 2.8** (cf. [28, Proposition 2]). *Suppose that the matrix  $A$  in (1.1) is nonsingular. Then for any initial guess  $x_0$ , GMRES applied to (1.1) and  $x_0$  terminates at step  $d \geq 0$  if and only if  $\dim AK_N(A, r_0) = d$ .*

While GMRES must fail to terminate for inconsistent linear systems, for consistent singular systems the question of termination depends on the initial guess. To derive a reasonable bound on the residual norms, one has to assume that GMRES will terminate, which is equivalent to assuming  $r_0 \in AK_d(A, r_0) = AK_N(A, r_0)$ . In the derivation of our new bounds we will always make this assumption. For further analysis of the application of GMRES to singular systems, in particular the case  $r_0 \notin AK_d(A, r_0)$ , we refer to [4].

**3. GMRES-equivalent matrices.** For a given matrix  $A \in \mathbf{C}^{N \times N}$  and vector  $r_0 \in \mathbf{C}^N$ , for which  $r_0 \in AK_N(A, r_0)$ , we now characterize the matrices  $B \in \mathbf{C}^{N \times N}$  that satisfy  $AK_n(A, r_0) = BK_n(B, r_0)$ , for  $1 \leq n \leq \dim AK_N(A, r_0)$ .

Since  $A$  might be singular and we allow  $\dim AK_N(A, r_0) < N$ , we obtain generalizations of results of Greenbaum and Strakoš [15, section 2] (also cf. [1, 14]). In the following, by  $r_n^{M,r}$  we denote the  $n$ th residual, when GMRES is applied to a linear system with the matrix  $M$  and the initial residual is  $r$ .

**General assumptions.** A linear system (1.1) and initial guess  $x_0$  are given such that  $2 \leq \dim AK_N(A, r_0) =: d \leq N$  and  $r_0 \in AK_d(A, r_0)$ . The latter assumption is always satisfied when  $A$  is nonsingular (cf. Theorem 2.8). Assuming  $d \geq 2$ , we exclude the trivial cases of GMRES termination at steps zero and one. Furthermore, let  $w_1, w_2, \dots, w_d$  be orthonormal vectors that satisfy

$$(3.1) \quad \text{span}\{w_1, w_2, \dots, w_n\} = AK_n(A, r_0) \text{ for } 1 \leq n \leq d.$$

Denoting  $W_d := [w_1, w_2, \dots, w_d] \in \mathbf{C}^{N \times d}$ , we have

$$(3.2) \quad r_0 = W_d h, \text{ where } h := [\eta_1, \dots, \eta_d]^T, \text{ and } \eta_n := w_n^H r_0 \text{ for } 1 \leq n \leq d.$$

It is easy to see that

$$(3.3) \quad AW_d = W_d H,$$

where  $H \in \mathbf{C}^{d \times d}$  is a nonsingular unreduced upper Hessenberg matrix. If  $H$  was not unreduced, then for some  $n$ ,  $1 \leq n \leq d - 1$ ,  $Aw_n$  would be a linear combination of  $w_1, \dots, w_n$ . This contradicts the dimension assumption on  $AK_N(A, r_0)$ . Also note that if we define

$$(3.4) \quad \hat{H} := \begin{bmatrix} 0 & \dots & 0 & 1/\eta_d \\ 1 & & & -\eta_1/\eta_d \\ & \ddots & & \vdots \\ & & 1 & -\eta_{d-1}/\eta_d \end{bmatrix},$$

then

$$(3.5) \quad H = \tilde{R}\hat{H},$$

where  $\tilde{R} \in \mathbf{C}^{d \times d}$  is some nonsingular and upper triangular matrix.

**DEFINITION 3.1.** *Suppose that  $A \in \mathbf{C}^{N \times N}$  and  $r_0 \in \mathbf{C}^N$ . Then  $B \in \mathbf{C}^{N \times N}$  is called GMRES-equivalent to  $A$  with respect to  $r_0$  if*

$$\dim BK_N(B, r_0) = d \text{ and } AK_n(A, r_0) = BK_n(B, r_0) \text{ for } 1 \leq n \leq d,$$

where  $d := \dim AK_N(A, r_0)$ .

Clearly, if  $B$  is GMRES-equivalent to  $A$  with respect to  $r_0$ , then

$$r_n^{A,r_0} = r_n^{B,r_0} \text{ for } 1 \leq n \leq d.$$

The following result is a generalization of [15, Theorems 2.1 and 2.2].

**THEOREM 3.2.** *Under our general assumptions, for a matrix  $B \in \mathbb{C}^{N \times N}$  the following three statements are equivalent:*

- (a)  $B$  is GMRES-equivalent to  $A$  with respect to  $r_0$ .
- (b)  $B$  satisfies  $BW_d = W_d \tilde{R} \hat{H}$ , where  $\hat{H}$  is defined in (3.4) and  $\tilde{R} \in \mathbb{C}^{d \times d}$  is a nonsingular upper triangular matrix.
- (c)  $B$  satisfies  $BW_d = W_d \bar{R} H$ , where  $H$  is defined in (3.3) and  $\bar{R} \in \mathbb{C}^{d \times d}$  is a nonsingular upper triangular matrix.

*Proof.* The equivalence of (b) and (c) is obvious from (3.5). We thus have to show only that (a) and (b) are equivalent. Suppose that (a) holds. Since  $r_0 = W_d h$ , we have

$$(3.6) \quad B[r_0, W_{d-1}] = BW_d \hat{H}^{-1}.$$

Furthermore, from (3.1) and the GMRES-equivalence of  $A$  and  $B$ , it follows that

$$(3.7) \quad B[r_0, W_{d-1}] = W_d \tilde{R}$$

for some upper triangular matrix  $\tilde{R}$ . It is easy to show that  $\tilde{R}$  must be nonsingular. Now (3.6) and (3.7) yield  $BW_d = W_d \tilde{R} \hat{H}$ .

Conversely, if (b) holds, then  $B^n r_0 = B^n W_d h = r_{11} W_d (\tilde{R} \hat{H})^{n-1} e_1$ . Since  $\tilde{R}$  is nonsingular and upper triangular,  $AK_n(A, r_0) = \text{span}\{w_1, w_2, \dots, w_n\} = BK_n(B, r_0)$  for  $1 \leq n \leq d$ . Finally we note that  $B^{d+1} r_0 \in BK_d(B, r_0)$ , so that  $\dim BK_N(B, r_0) = d$ .  $\square$

By varying  $\tilde{R}$  and  $\bar{R}$  in Theorem 3.2, we obtain a whole class of matrices which are GMRES-equivalent to the given matrix  $A$  with respect to the initial residual  $r_0 = W_d h$ . We obtain a useful corollary.

**COROLLARY 3.3.** *Suppose that, in the notation of Theorem 3.2,*

$$B = W_d \bar{R} H W_d^H \text{ and } \hat{B} = W_d \tilde{R} \hat{H} W_d^H,$$

where  $\bar{R}$  and  $\tilde{R}$  are arbitrary nonsingular upper triangular matrices. Then  $B$  and  $\hat{B}$  are both GMRES-equivalent to  $A$  with respect to  $r_0$  and, in particular,

$$\|r_n^{B,r_0}\| = \|r_n^{\hat{B},r_0}\| = \|r_n^{A,r_0}\| \text{ for } 1 \leq n \leq d.$$

We remark that if  $d < N$ , then the matrices  $B$  and  $\hat{B}$  in Corollary 3.3 must be singular, regardless of whether  $A$  is singular.

**4. New bounds on GMRES residual norms.** The results of this section are derived under the general assumptions stated in section 3. Using Theorem 3.2, we can relate the convergence of GMRES for the given linear system to the convergence for some other system, *which is easier to analyze*. Our starting point for the derivation of the new bounds therefore is the construction of suitable matrices that are GMRES-equivalent to  $A$  with respect to the given  $r_0$ .

As stated in Corollary 3.3, for arbitrary nonsingular upper triangular matrices  $\bar{R}, \tilde{R} \in \mathbb{C}^{d \times d}$ , the matrices

$$(4.1) \quad W_d \bar{R} H W_d^H \text{ and } W_d \tilde{R} \hat{H} W_d^H$$

are both GMRES-equivalent to  $A$  with respect to  $r_0$ . We consider the following decompositions of  $H$  and  $\hat{H}$ :

$$(4.2) \quad H = QR \quad \text{and} \quad \hat{H} = \hat{R}\hat{Q}.$$

The matrices  $Q$  and  $\hat{Q}$  are unitary and the matrices  $R$  and  $\hat{R}$  are nonsingular and upper triangular. All four are of size  $d \times d$ . These two decompositions always exist due to nonsingularity of  $H$  and  $\hat{H}$  (see, e.g., [11]).

Two other similar decompositions are possible:  $H = R_1Q_1$  and  $\hat{H} = Q_2R_2$ . We are particularly interested in the spectra of the  $Q$ -factors of these decompositions. Since  $Q_2$  is simply a unitary shift matrix and the results derived for  $\Lambda(\hat{Q})$  also hold for  $\Lambda(Q_1)$  (cf. (3.5)), we do not consider these two other decompositions.

For later use we note that  $W_d^HAW_d = QR$  yields

$$(4.3) \quad \kappa(R) \leq \kappa(A)$$

with equality if  $d = N$ . Often  $d = N$  is a reasonable assumption, in particular when the “generic” case of an initial residual without special properties is to be analyzed.

Choosing  $\tilde{R} = R^{-1}$  and  $\hat{\tilde{R}} = \hat{R}^{-1}$  in (4.1) and using (4.2), we define the matrices

$$(4.4) \quad B := W_dR^{-1}QRW_d^H \quad \text{and} \quad \hat{B} := W_d\hat{Q}W_d^H,$$

which are both GMRES-equivalent to  $A$  with respect to  $r_0$ .

*Remark 4.1.* Greenbaum and Strakoš [15, section 3.1] consider the  $RQ$ -decomposition of the matrix  $H$ . They point out that if  $d = N$ , then the matrix  $\hat{B} = W_N\hat{Q}W_N^H$  is unitary. Thus, whenever  $A$  and  $r_0$  are such that  $d = N$ , there exists a unitary matrix which is GMRES-equivalent to  $A$  with respect to  $r_0$ .

**THEOREM 4.2.** *Suppose that  $B$  and  $\hat{B}$  are defined as in (4.4). Then*

$$(4.5) \quad \frac{\|r_n\|}{\|r_0\|} \leq \kappa(R) \min_{p_n \in \pi_n} \max_{\lambda \in \Lambda(Q)} |p_n(\lambda)|$$

and

$$(4.6) \quad \frac{\|r_n\|}{\|r_0\|} \leq \min_{p_n \in \pi_n} \max_{\lambda \in \Lambda(\hat{Q})} |p_n(\lambda)|$$

for  $1 \leq n \leq d$ .

*Proof.* Corollary 3.3 yields for  $1 \leq n \leq d$ ,

$$\begin{aligned} \|r_n\| &\equiv \|r_n^{A,r_0}\| = \|r_n^{B,r_0}\| = \min_{p_n \in \pi_n} \|W_dR^{-1}p_n(Q)RW_d^H r_0\| \\ &\leq \|R^{-1}\| \|R\| \min_{p_n \in \pi_n} \|p_n(Q)\| \|r_0\|, \end{aligned}$$

from which (4.5) follows. The proof of (4.6) is similar.  $\square$

To study the approximation problem

$$(4.7) \quad \min_{p_n \in \pi_n} \max_{\lambda \in \Lambda(U)} |p_n(\lambda)|, \quad U \text{ unitary,}$$

we need the following definition.

**DEFINITION 4.3.** *Suppose that the spectrum of the  $d \times d$  unitary matrix  $U$  is given by  $\Lambda(U) := \{e^{i\beta_j} : 0 \leq \beta_1 \leq \beta_2 \leq \dots \leq \beta_d < 2\pi\}$  and let  $\beta_{d+1} := 2\pi + \beta_1$ . We define the  $d$  gaps in  $\Lambda(U)$  by*

$$(4.8) \quad g_j := \beta_{j+1} - \beta_j \quad \text{for } 1 \leq j \leq d.$$

For  $1 \leq j \leq d$ , the complex number  $e^{i(\beta_{j+1}+\beta_j)/2}$  is called the center of the gap  $g_j$ .  $\phi := \max\{g_1, g_2, \dots, g_d\}$  denotes the largest gap in  $\Lambda(U)$ .

Note that for every unitary matrix  $U$  and every  $\alpha \in \mathbf{R}$ , we have

$$(4.9) \quad \min_{p_n \in \pi_n} \max_{\lambda \in \Lambda(U)} |p_n(\lambda)| = \min_{p_n \in \pi_n} \max_{\lambda \in \Lambda(e^{i\alpha}U)} |p_n(\lambda)|.$$

Therefore, in the study of (4.7) we can assume, without loss of generality, that 1 is the center of the largest gap  $\phi$  in  $\Lambda(U)$ . Let us define

$$(4.10) \quad \Omega_\phi := \left\{ e^{i\alpha} : \frac{\phi}{2} \leq \alpha \leq 2\pi - \frac{\phi}{2} \right\}.$$

Then (4.9) shows that

$$(4.11) \quad \min_{p_n \in \pi_n} \max_{\lambda \in \Lambda(U)} |p_n(\lambda)| \leq \min_{p_n \in \pi_n} \max_{z \in \Omega_\phi} |p_n(z)|.$$

Suppose that  $0 < \phi < 2\pi$ , or, equivalently, that  $\Lambda(U)$  contains at least two distinct points. Then for  $n \geq 1$ , we can introduce the polynomials

$$(4.12) \quad p_n(z) := \prod_{k=0}^{n-1} \left( 1 - z \frac{(1 - \gamma e^{i\psi_k})}{e^{i\psi_k} (\gamma - e^{i\psi_k})} \right), \text{ where}$$

$$(4.13) \quad \psi_k := \frac{2\pi k}{n} \text{ for } 0 \leq k \leq n-1,$$

$$(4.14) \quad \gamma := \left( \cos \frac{\phi}{4} \right)^{-1}.$$

By assumption,  $1 < \gamma < \infty$  and hence the polynomials (4.12) are well defined for every  $n \geq 1$ . Furthermore, we have  $p_n \in \pi_n$ .

*Remark 4.4.* The polynomials  $p_n$  in (4.12) are constructed using a conformal mapping technique from [19] (also cf. [21]). As shown there,  $\Psi(z) := \frac{z(\gamma-z)}{1-\gamma z}$  maps the exterior of the unit circle conformally onto the exterior of  $\Omega_\phi$ . Thus, the zeros of  $p_n$  are Fejér points (see, e.g., [29, Chapter 1]) associated with  $\Omega_\phi$ .

LEMMA 4.5. *The polynomials introduced in (4.12) satisfy*

$$(4.15) \quad \max_{z \in \Omega_\phi} |p_n(z)| \leq \frac{4}{\gamma^n - 1} \text{ for } n \geq 1.$$

*Proof.* Suppose that  $z = e^{i\alpha} \in \Omega_\phi$ , i.e.,  $\alpha \in \left[ \frac{\phi}{2}, 2\pi - \frac{\phi}{2} \right]$ . A simple manipulation yields

$$(4.16) \quad \begin{aligned} |p_n(e^{i\alpha})| &= \prod_{k=0}^{n-1} \left| 1 - e^{i\alpha} \frac{(1 - \gamma e^{i\psi_k})}{e^{i\psi_k} (\gamma - e^{i\psi_k})} \right| \\ &= \prod_{k=0}^{n-1} \left| \frac{e^{i\frac{\alpha}{2}} (\gamma (e^{i\frac{\alpha}{2}} + e^{-i\frac{\alpha}{2}}) - (e^{i(\psi_k - \frac{\alpha}{2})} + e^{-i(\psi_k - \frac{\alpha}{2})}))}{\gamma - e^{i\psi_k}} \right| \\ &= \frac{2^n}{\gamma^n - 1} \prod_{k=0}^{n-1} \left| \gamma \cos \frac{\alpha}{2} - \cos \left( \psi_k - \frac{\alpha}{2} \right) \right|. \end{aligned}$$

In the last equation, we utilized that  $\prod_{k=0}^{n-1} (z - e^{i\psi_k}) = z^n - 1$  for all  $z \in \mathbf{C}$ . Now note that  $-1 \leq \gamma \cos \frac{\alpha}{2} \leq 1$ . We define  $\beta := \arccos(\gamma \cos \frac{\alpha}{2})$ . Then (4.16) yields

$$\begin{aligned} |p_n(e^{i\alpha})| &= \frac{2^n}{\gamma^n - 1} \prod_{k=0}^{n-1} \left| \cos \beta - \cos \left( \psi_k - \frac{\alpha}{2} \right) \right| \\ &= \frac{4^n}{\gamma^n - 1} \prod_{k=0}^{n-1} \left| \sin \frac{\beta + \psi_k - \frac{\alpha}{2}}{2} \sin \frac{\beta - \psi_k + \frac{\alpha}{2}}{2} \right| \\ &= \frac{2^n}{\gamma^n - 1} \prod_{k=0}^{n-1} \left| \sqrt{1 - \cos \left( \beta + \psi_k - \frac{\alpha}{2} \right)} \sqrt{1 - \cos \left( \beta - \psi_k + \frac{\alpha}{2} \right)} \right| \\ &= \frac{1}{\gamma^n - 1} \prod_{k=0}^{n-1} \left| e^{i(\frac{\alpha}{2} - \beta)} - e^{i\psi_k} \right| \left| e^{i(\frac{\alpha}{2} + \beta)} - e^{i\psi_k} \right| \\ &= \frac{1}{\gamma^n - 1} \left| e^{i(\frac{\alpha}{2} - \beta)n} - 1 \right| \left| e^{i(\frac{\alpha}{2} + \beta)n} - 1 \right| \leq \frac{4}{\gamma^n - 1}. \quad \square \end{aligned}$$

As Lemma 4.5 shows, the polynomials (4.12) converge to zero on  $\Omega_\phi$ , with the speed of convergence being directly related to the gap  $\phi$ : the larger  $\phi$ , the faster the convergence, and vice versa. Since the polynomials (4.12) are constructed using Fejér points, they converge to zero on  $\Omega_\phi$  with the largest geometric degree of convergence (cf. [36, Chapter 4.7]). Because of the moderate constant 4, Lemma 4.5 is useful not only asymptotically, but also for small  $n$ .

We now combine (4.3), (4.11), Theorem 4.2, and Lemma 4.5 for our main convergence result. Note that from our general assumption  $d \geq 2$ , it follows that  $\Lambda(Q)$  and  $\Lambda(\hat{Q})$  both contain at least two distinct points, so that Lemma 4.5 is applicable.

**THEOREM 4.6.** *Suppose that a linear system (1.1) and an initial guess  $x_0$  are given, such that  $2 \leq d := \dim AK_N(A, r_0) \leq N$  and  $r_0 \in AK_d(A, r_0)$ . Define the matrices  $Q$ ,  $R$ , and  $\hat{Q}$  as in (4.2) and let  $\phi$  and  $\hat{\phi}$  denote the largest gaps in  $\Lambda(Q)$  and  $\Lambda(\hat{Q})$ , respectively. Define  $\gamma := (\cos \frac{\phi}{4})^{-1}$  and  $\hat{\gamma} := (\cos \frac{\hat{\phi}}{4})^{-1}$ ; then the GMRES residuals for (1.1) and  $x_0$  satisfy*

$$(4.17) \quad \frac{\|r_n\|}{\|r_0\|} \leq 4 \frac{\kappa(R)}{\gamma^n - 1}$$

and

$$(4.18) \quad \frac{\|r_n\|}{\|r_0\|} \leq \frac{4}{\hat{\gamma}^n - 1}$$

for  $1 \leq n \leq d$ . Furthermore,  $\kappa(R) \leq \kappa(A)$  with equality if  $d = N$ .

It is important to note that we cannot expect (4.17) and (4.18) to be sharp, except for special cases. There are several reasons for this: First, the proof of Lemma 4.5 uses the submultiplicative property of the Euclidean norm, which generally does not lead to sharp bounds. Second, while the polynomials used in Lemma 4.5 are asymptotically optimal, the GMRES algorithm actually constructs polynomials that are optimal in each step. Third, both our bounds decrease with constant rates, each being derived from only one largest spectral gap. But there might be several considerable gaps in  $\Lambda(Q)$  and  $\Lambda(\hat{Q})$  that have an impact on the convergence behavior.

As shown in the remainder of this paper, despite all these theoretical objections, (4.17) and (4.18) are valuable for describing and predicting the behavior of GMRES.



We first point out that by construction  $\gamma^{-1}$  and  $\hat{\gamma}^{-1}$ , the rates of reduction of our bounds, depend on the initial guess. Both rates lie strictly between zero and one. We have not found an algebraic relation between them.

At first sight, (4.18) seems to be more useful than (4.17), since it does not include the potentially large constant  $\kappa(R)$ . However, it is possible to implement GMRES in a way that approximations to (4.17) are easily computable in each step of the algorithm. On the other hand, computing approximations to (4.18) during the GMRES run is sometimes not possible and is always more expensive.

Furthermore, the appearance of the constant  $\kappa(R)$  leads to the following observation: Suppose that GMRES stagnates early in the iteration. Then  $\hat{\phi}$  must be small, because (4.18) involves no large constant. But because of  $\kappa(R)$  in (4.17), early stagnation of GMRES does not generally imply that  $\phi$  is small. Hence,  $\gamma^{-1}$  and  $\hat{\gamma}^{-1}$  might be of different sizes. In particular, when  $\kappa(R)$  is large, the rate  $\gamma^{-1}$  has the potential to describe the convergence of GMRES after an initial phase of stagnation. Numerical experiments confirming this observation are shown in section 6.

**5. Approximate bounds predict the behavior of GMRES.** Note that if close approximations to the right-hand sides of (4.17) and (4.18) are available in early stages of the computation, then these potentially *predict* the behavior of the algorithm in later stages. It is therefore useful to compute such approximations during the GMRES run.

**5.1. WZ-GMRES.** To obtain the approximations, we consider an implementation of GMRES first proposed by Walker and Zhou [35]. In the following we will refer to this implementation as WZ-GMRES. For a description of WZ-GMRES, we need to introduce the Arnoldi process [2]:

Suppose that  $M \in \mathbf{C}^{N \times N}$  and a unit vector  $w_1 \in \mathbf{C}^N$  satisfy  $\dim K_N(M, w_1) = d \leq N$ . Then the Arnoldi process applied to  $M$  and  $w_1$  is the construction of an orthonormal basis,  $w_1, w_2, \dots, w_d$ , of  $K_d(M, w_1)$  by means of any orthogonalization method. In many references (see, e.g., [11]), the orthogonalization performed by the Arnoldi process is implemented using the classical or modified Gram–Schmidt procedure. As pointed out by Walker [34], it is also possible to use Householder transformations for this purpose. For a collection of different implementations of the Arnoldi process, we refer to [7] (see also [24]). Our development is independent of the particular implementation, so that any (numerically stable) one will serve our purpose.

In practice,  $d$  is not known, and the basis vectors are constructed recursively. After  $n < d$  steps of the Arnoldi process, we have

$$(5.1) \quad MW_n = W_{n+1}H_n^e, \quad W_n := [w_1, \dots, w_n],$$

where

$$(5.2) \quad H_n^e = \begin{bmatrix} h_{11} & \cdots & h_{1n} \\ h_{21} & \ddots & \vdots \\ & \ddots & h_{nn} \\ & & & h_{n+1,n} \end{bmatrix}$$

is an  $(n+1) \times n$  unreduced upper Hessenberg matrix. In exact arithmetic, the process has to terminate at step  $d$  with the final matrix relation

$$(5.3) \quad MW_d = W_d H_d, \quad W_d \in \mathbf{C}^{N \times d}, \quad H_d := W_d^H A W_d \in \mathbf{C}^{d \times d}.$$

The original implementation of GMRES [28] is based on the Arnoldi process applied to the matrix  $A$  and the vector  $r_0/\|r_0\|$ . WZ-GMRES uses

$$(5.4) \quad w_1 := Ar_0/\|Ar_0\|$$

as the starting vector for the Arnoldi process. Recall that for our given linear system and initial guess, we have assumed  $\dim AK_N(A, r_0) = d$ . Then, since

$$K_n(A, w_1) = AK_n(A, r_0) \text{ for } 1 \leq n \leq d,$$

the application of the Arnoldi process to  $A$  and  $w_1$  yields orthonormal vectors that satisfy (3.1). Using these vectors, (2.1) can be rewritten as

$$(5.5) \quad x_n = x_0 + [r_0/\|r_0\|, W_{n-1}]y_n$$

for some  $y_n \in \mathbf{C}^n$ ,  $1 \leq n \leq d$ . To implement GMRES, we use the orthogonality relation (2.4) in the computation of  $y_n$ . Together with (5.1) and (5.5), this yields

$$0 = W_n^H r_n = W_n^H r_0 - W_n^H A[r_0/\|r_0\|, W_{n-1}]y_n,$$

which is equivalent to

$$(5.6) \quad [(\|Ar_0\|/\|r_0\|)e_1, H_{n-1}^e]y_n = W_n^H r_0.$$

Thus,  $y_n$  is the solution of an upper triangular system. Since  $H_{n-1}^e$  is unreduced,  $y_n$  is well defined. It is important to note that  $r_n$  can be obtained without explicitly computing  $x_n$  [25],

$$(5.7) \quad r_n = r_{n-1} - (w_n^H r_0)w_n.$$

The upper triangular solve for  $y_n$  therefore has to be performed only if the residual norm computed from (5.7) is small enough. Clearly, in exact arithmetic, the termination of the Arnoldi process at step  $d$  is equivalent to the termination of WZ-GMRES at step  $d$ . We remark that from (5.5) and (5.6) it is easy to see that WZ-GMRES is (mathematically) independent of the orthogonalization method used for the Arnoldi process. For more details about WZ-GMRES, we refer to [21, 25, 35].

**5.2. Computing approximations to (4.17).** If (5.4) holds, then the vectors  $w_1, w_2, \dots, w_d$  produced by the Arnoldi process satisfy (3.1). Consequently, the matrix  $H_d$  in (5.3) can be identified with the matrix  $H$  in (4.2). We can therefore compute the numbers  $\phi$  and  $\kappa(R)$  used in (4.17) from the  $QR$ -decomposition of  $H_d$ . Thus, the right-hand side of (4.17) is computable from information generated during the run of WZ-GMRES. Moreover, because of the recursive nature of the Arnoldi process, the matrix  $H_n^e$  is the leading  $(n+1) \times n$  submatrix of  $H_d$ . For  $n = 1, 2, \dots, d-1$ , we can therefore decompose

$$(5.8) \quad H_n^e = Q_n \begin{bmatrix} R_n \\ 0 \end{bmatrix}, \quad Q_n \in \mathbf{C}^{(n+1) \times (n+1)}, \quad R_n \in \mathbf{C}^{n \times n},$$

and use  $\phi_n$ , the largest gap in  $\Lambda(Q_n)$ , as an approximation to  $\phi$  and  $\kappa(R_n)$  as an approximation to  $\kappa(R)$ . Note that because  $H_n^e$  is unreduced,  $R_n$  is nonsingular and (5.8) is well defined.

The decomposition (5.8) is particularly easy to compute. It is well known that because of the special structure of  $H_n^e$ , the matrix  $Q_n$  can be derived as a product of  $n$  Givens transformations,  $Q_n = G_1 \cdots G_n$ , where

$$(5.9) \quad G_j = \begin{bmatrix} I_{j-1} & & & & \\ & c_j & -\bar{s}_j & & \\ & s_j & c_j & & \\ & & & I_{n-j} & \\ & & & & \end{bmatrix} \in \mathbf{C}^{(n+1) \times (n+1)}, \quad 1 \leq j \leq n.$$

Suppose that for some  $n$ ,  $2 \leq n \leq d - 1$ ,

$$(5.10) \quad Q_{n-1}^H H_{n-1}^e = \begin{bmatrix} R_{n-1} \\ 0 \end{bmatrix}, \quad Q_{n-1}^H = G_{n-1}^H \cdots G_1^H.$$

We define

$$h^n := [h_{1n}, \dots, h_{nn}]^T, [\rho_1, \dots, \rho_n]^T := Q_{n-1}^H h^n, \text{ and } \nu_n := [\rho_1, \dots, \rho_{n-1}]^T.$$

Then

$$(5.11) \quad \begin{aligned} Q_n^H H_n^e &= G_n^H \begin{bmatrix} Q_{n-1}^H & \\ & 1 \end{bmatrix} \begin{bmatrix} H_{n-1}^e & h^n \\ & h_{n+1,n} \end{bmatrix} \\ &= G_n^H \begin{bmatrix} R_{n-1} & \nu_n \\ & \rho_n \\ & & h_{n+1,n} \end{bmatrix}. \end{aligned}$$

If  $\rho_n = |\rho_n|e^{i\delta_n}$ , we define  $\omega_n := (|\rho_n|^2 + |h_{n+1,n}|^2)^{1/2}$  and choose

$$(5.12) \quad c_n = \frac{|\rho_n|}{\omega_n},$$

$$(5.13) \quad s_n = \frac{e^{-i\delta_n} h_{n+1,n}}{\omega_n}$$

to give (5.11) the required upper triangular form. Hence, to compute  $Q_n$  and  $R_n$  from  $Q_{n-1}$  and  $R_{n-1}$ , respectively, only one matrix-vector product and a few flops to compute  $c_n$  and  $s_n$  are required.

To compute  $\phi_n$  we have to solve an  $n \times n$  unitary eigenproblem. For the efficient computation of  $\kappa(R_n)$  we can use the method of incremental condition estimation [3]. With this method, the computation of  $\kappa(R_n)$  costs only  $\mathcal{O}(n)$  operations, provided  $\kappa(R_{n-1})$  is known.

**5.3. Computing approximations to (4.18).** Suppose that we apply WZ-GMRES to (1.1) and  $x_0$  and use the Arnoldi vectors  $w_1, w_2, \dots, w_d$  to set up the matrix  $\hat{H}$  as defined in (3.4). Note that  $\hat{H}$  can be obtained at no additional cost if the updated residuals are computed using (5.7). The bound (4.18) can then be computed easily because it depends only on the largest gap  $\hat{\phi}$  in  $\Lambda(\hat{Q})$ , where  $\hat{Q}$  is defined by the  $RQ$ -decomposition  $\hat{H} = \hat{R}\hat{Q}$  (cf. (4.2)).

For approximations to (4.18), we consider (5.7), which yields

$$|\eta_n| = |w_n^H r_0| = \|r_{n-1} - r_n\| \text{ for } 1 \leq n \leq d.$$

Hence, if GMRES does not stagnate in step  $n$ ,  $2 \leq n \leq d - 1$ , then  $\eta_n \neq 0$  and the decomposition

$$(5.14) \quad \hat{H}_n := \begin{bmatrix} 0 & \dots & 0 & 1/\eta_n \\ 1 & & & -\eta_1/\eta_n \\ & \ddots & & \vdots \\ & & 1 & -\eta_{n-1}/\eta_n \end{bmatrix} = \hat{R}_n \hat{Q}_n \in \mathbf{C}^{n \times n}$$

is defined. In this case,  $\hat{\phi}_n$ , the largest gap in  $\Lambda(\hat{Q}_n)$ , can be used as an approximation to  $\hat{\phi}$ . However, we have no easy way to update the decomposition (5.14).  $\hat{R}_n$  and  $\hat{Q}_n$  have to be recomputed from scratch in each step. As shown above, the situation is drastically different for the decomposition (5.8).

An approximation  $\hat{\phi}_n$  can also be computed from the  $RQ$ -decomposition of the matrix  $H_n := W_n^H A W_n$  (cf. (5.1)). In exact arithmetic, stagnation of WZ-GMRES in step  $n$  leads to singularity of  $H_n$ . Hence  $\hat{\phi}_n$  is not well defined in this case. But on the contrary to (5.14),  $\hat{\phi}_n$  might still be computable. This is particularly true in finite precision arithmetic.

**5.4. Behavior of the approximations and prediction of the performance of GMRES.**  $Q_n$  in (5.8) is obtained from  $Q_{n-1}$  by a low-rank modification. If  $q_1, \dots, q_n$  denote the columns of  $Q_{n-1}$ , then

$$Q_n = \begin{bmatrix} q_1 & \dots & q_{n-1} & c_n q_n & -\bar{s}_n q_n \\ 0 & \dots & 0 & s_n & c_n \end{bmatrix}.$$

By (5.12),  $c_n \neq -1$ , so that we can readily apply results of Elsner and He [9, section 5]. We partition

$$Q_n = \begin{bmatrix} U_{11} & U_{12} \\ U_{21} & U_{22} \end{bmatrix}, \quad U_{11} = [q_1, \dots, q_{n-1}, c_n q_n] \in \mathbf{C}^{n \times n}, \quad U_{22} = c_n$$

to get, in the notation of [9, Theorem 5.3]:  $U_n := U_{11} - U_{12}(I + U_{22})^{-1}U_{21} = Q_{n-1}$ . From [9, Theorem 5.3] it follows that if

$$\begin{aligned} \Lambda(Q_{n-1}) &= \Lambda(U_n) = \{e^{i\beta_j} : 0 \leq \beta_1 \leq \dots \leq \beta_n < 2\pi\}, \\ \Lambda(Q_n) &= \{e^{i\tau_j} : 0 \leq \tau_1 \leq \dots \leq \tau_{n+1} < 2\pi\}, \end{aligned}$$

then  $\tau_j \leq \beta_j \leq \tau_{j+1}$  for  $1 \leq j \leq n$ .

Because of this interlacing property, we can expect a “smooth” behavior of the eigenvalues of the matrices  $Q_n$  and, in particular, of their largest gaps. This is not necessarily the case for the spectra of the matrices  $\hat{Q}_n$  in (5.14), since we know of no easy algebraic relation between them.

The matrix  $R_n$  in (5.8) is the leading principal submatrix  $R$  in (4.2). Thus, we have  $\kappa(R_{n-1}) \leq \kappa(R_n) \leq \kappa(R)$  for  $n = 2, \dots, d - 1$ . Monitoring the numbers  $\kappa(R_n)$  is a useful check on the numerical accuracy of the WZ-GMRES iterates. As (5.6) shows, the conditioning of  $R_{n-1}$  determines the accuracy of the vector  $y_n$ . In the derivation of our bounds we have assumed  $r_0 \in AK_d(A, r_0)$  for some  $d \leq N$ . From this assumption follows that  $\kappa(R_n) \leq \kappa(R) < \infty$ . But for ill-conditioned (or singular)  $A$ ,  $\kappa(R)$  and thus  $\kappa(R_n)$  might still be very large. In such a case it might be useful to terminate WZ-GMRES early.

In numerical experiments we often observed that for small  $n$ ,  $\phi_n$ ,  $\kappa(R_n)$ , and  $\hat{\phi}_n$  were close to  $\phi$ ,  $\kappa(R)$ , and  $\hat{\phi}$ , respectively. Such observations become plausible following this heuristic reasoning: In many cases even small leading submatrices of  $H$  convey much of the information obtained by the Arnoldi process. For example, close approximations to  $\kappa(A)$  and  $\Lambda(A)$  can typically be obtained from  $H_n$  for small  $n$  (see, e.g., [26]). In the same fashion, information about properties of the factors in the  $QR$ - and  $RQ$ -factorizations of  $H$  should be obtainable from its leading submatrices. Hence, close approximations to quantities associated with  $Q$ ,  $R$ , and  $\hat{Q}$  should for small  $n$  typically be obtainable from  $Q_n$ ,  $R_n$ , and  $\hat{Q}_n$ , respectively.

With close approximations to  $\phi$ ,  $\kappa(R)$ , and  $\hat{\phi}$ , predictions can be made on how GMRES will perform in later stages of the iteration. Theorem 4.6 implies the following heuristic for such predictions.

**HEURISTIC 5.1.** *When the approximations  $\phi_n$ ,  $\kappa(R_n)$ , and  $\hat{\phi}_n$  are close to  $\phi$ ,  $\kappa(R)$ , and  $\hat{\phi}$ , respectively, their sizes predict the behavior of GMRES as follows:*

1. *A large  $\phi_n$  predicts fast GMRES convergence, possibly after an initial phase of stagnation in case  $\kappa(R_n)$  is large. A small  $\phi_n$  predicts slow convergence over the whole course of the iteration.*
2. *A large  $\hat{\phi}_n$  predicts fast GMRES convergence over the whole course of the iteration. A small  $\hat{\phi}_n$  predicts (at least early) stagnation of GMRES.*

Our experiments indicate that stagnation of the approximations for several steps typically indicates that they are close to their final values. This is not, however, a reliable method of checking the accuracy of approximations. For additional information it is possible to check the convergence of the Arnoldi process itself, for example, by estimating how closely  $\Lambda(H_n)$  approximates  $\Lambda(A)$ . Finding direct bounds on the errors  $\|\phi - \phi_n\|$ , etc., is an interesting open problem which we plan to study in the future.

**6. Numerical experiments.** We present numerical experiments run in MATLAB 5.0 [32] using five different nonsymmetric (A199: nonhermitian) test matrices:

$A$	$N$	$\kappa(A)$	Data type
A199	199	153.5	complex
PDE225	225	39.1	real
FS1836	183	1.7e+11	real
STEAM1	240	2.8e+07	real
ISING100	100	1.0	real

The matrices PDE225, FS1836, and STEAM1 can be obtained from the MatrixMarket website at <http://math.nist.gov/MatrixMarket>. PDE225 comes from the finite difference discretization of a two-dimensional elliptic partial differential operator, FS1836 represents a problem in chemical kinetics, and STEAM1 comes from an oil recovery problem. More information on these matrices can be obtained from the MatrixMarket website. More information on A199 and ISING100 is given below.

For our experiments we use a Householder version [34] of WZ-GMRES. The approximations  $\phi_n$  are computed as described in section 5.2 and the  $\hat{\phi}_n$  are derived from the  $RQ$ -decomposition of  $H_n = W_n^H A W_n$  (cf. section 5.3). To check the numerical accuracy of our computations, we compare the norms of the directly computed residuals  $(b - Ax_n)$  and the updated residuals (cf. (5.7)).

Before presenting our results we point out that the bounds (4.17) and (4.18), in particular the approximations  $\phi_n$  and  $\hat{\phi}_n$ , are sensitive to the scaling of the diagonal entries of the  $R$ -factors in the  $QR$ - and  $RQ$ -decompositions. The results of this section

are derived with the following choices:

- The  $QR$ -decompositions (5.8) are computed using Givens transformations with nonnegative cosines, i.e.,  $c_n \geq 0$  for all  $n$  (cf. (5.12)).
- The  $RQ$ -decompositions  $H_n = \hat{R}_n \hat{Q}_n$  are computed so that the  $\hat{R}_n$  have positive diagonal entries.

Of all the scalings we tested, these choices typically led to the largest gaps  $\phi$  and  $\hat{\phi}$  and thus yielded the smallest upper bounds (4.17) and (4.18).

*Experiment 6.1.* We use the matrices A199, PDE225, and FS1836 to show typical types of behavior of the quantities in our bounds and of the computed approximations. A199 is a complex matrix of the form

$$A199 = \begin{bmatrix} 1 & & & \\ & 1.5 & & \\ & & \ddots & \\ & & & 100 \end{bmatrix} + iB \in \mathbf{C}^{199 \times 199},$$

where  $B$  is a random matrix generated with the MATLAB command `B=rand(199)`.

For  $A = A199$  and  $A = PDE225$ , we use  $b = A[1, \dots, 1]^T$  and  $x_0 = 0$ . We terminate WZ-GMRES at step  $d$  when the updated scaled residual norm satisfies  $\|r_d\|/\|r_0\| < 10^{-14}$ . Since the residual norm has decreased almost to machine precision,  $\mathbf{u} = 10^{-16}$ , we identify  $d$  with  $\dim AK_N(A, r_0)$ . Hence, we consider last-computed approximations as equal to the quantities used in Theorem 4.6:

	$d$	$\phi$	$\hat{\phi}$	$\kappa(R)$
A199	94	2.1325	2.1601	150.7
PDE225	93	1.9675	1.4250	34.9

Figure 6.1 shows the convergence curve of WZ-GMRES for A199 as well as our convergence bounds. We stress that the actual right-hand sides of (4.17) and (4.18) are plotted and not just the “rates of convergence”  $\gamma^{-1}$  and  $\hat{\gamma}^{-1}$ . Both bounds describe the actual rate of convergence well. Due to the constant  $\kappa(R)$ , (4.17) is slightly weaker than (4.18). Figures 6.2 and 6.4 show that the approximations  $\kappa(R_n)$  and  $\phi_n$  converge “smoothly” and reach their final levels after approximately 40 steps.

Note that, as shown in Figure 6.3, there is only one significant gap in  $\Lambda(Q)$  (the same holds for  $\Lambda(\hat{Q})$ ). The number of significant gaps in  $\Lambda(Q)$  and  $\Lambda(\hat{Q})$  has an important influence on our bounds: The rates of reduction of both bounds depend only on the largest gap in the two respective spectra. Thus, the occurrence of several large gaps might cause our bounds to overestimate the actual convergence behavior. As mentioned in section 4, this is one reason we do not generally expect our bounds to be sharp. However, our bounds, as well as Heuristic 5.1, will not lose their predictive value in cases of several large gaps, as the occurrence of such gaps will still correspond to and predict fast GMRES convergence. Experiment 6.3 gives an example of our bounds in the context of two large gaps.

Results for PDE225 are shown in Figures 6.5–6.8, which have the same notation as Figures 6.1–6.4, respectively. Note the considerable difference between  $\phi$  and  $\hat{\phi}$ . As Theorem 4.6 implies,  $\hat{\phi}$  must be small due to slow initial GMRES convergence. Because of the constant  $\kappa(R)$  in (4.17),  $\phi$  might be—and actually is—larger than  $\hat{\phi}$ .

For  $A = FS1836$  (see Figures 6.9–6.12) we use two different right-hand sides,

$$b^{(1)} = [1, 1, \dots, 1]^T, \quad b^{(2)} = A[1, 1, \dots, 1]^T,$$

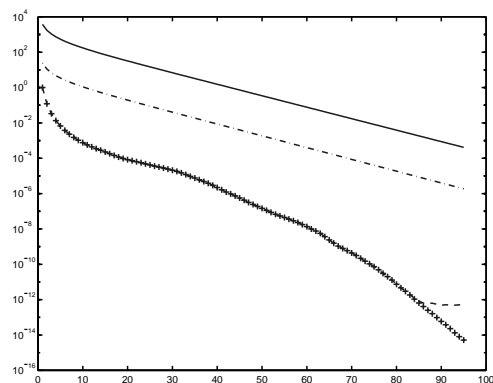


FIG. 6.1. A199: Computed (dashed) and updated (pluses) residual norms, bounds (4.17) (solid) and (4.18) (dashdot).

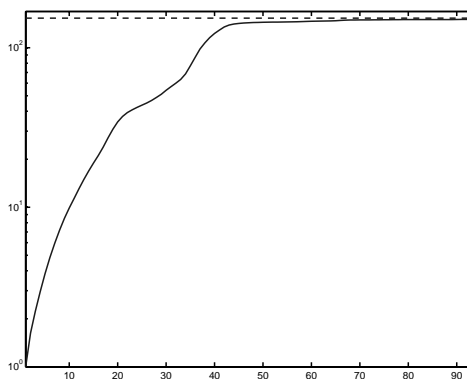


FIG. 6.2. A199:  $\kappa(R_n)$ .

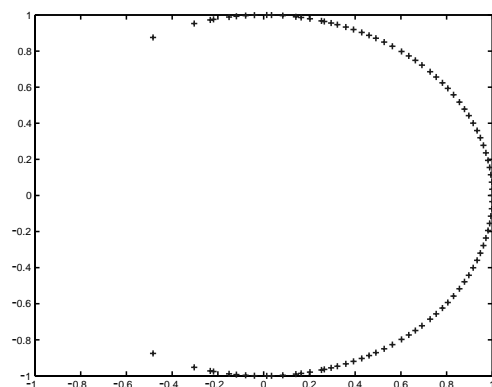


FIG. 6.3. A199:  $\Lambda(Q)$ .

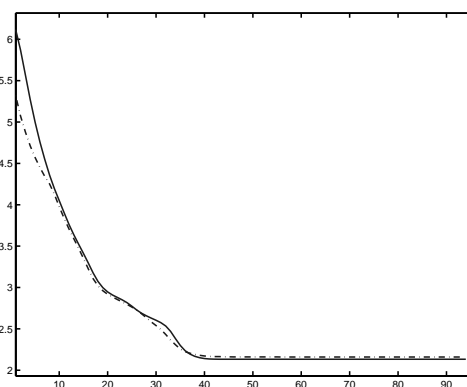


FIG. 6.4. A199:  $\phi_n$  (solid) and  $\hat{\phi}_n$  (dashdot).

and  $x_0 = 0$  in both cases. FS1836 is ill conditioned, and the performance of GMRES for this matrix is highly sensitive to the initial residual. Furthermore, due to the ill-conditioning, for both right-hand sides a finite precision accuracy of little more than  $\mathbf{u} * \kappa(A)$  is achievable by WZ-GMRES. We terminate WZ-GMRES at step  $n$  shortly after we notice a considerable difference between the computed and updated residuals. After termination of WZ-GMRES we continue the Arnoldi process to compute  $\phi$  and  $\hat{\phi}$ . For both right-hand sides these are close to the computed approximations at step  $n$ :

	$n$	$\phi_n$	$\phi$	$\hat{\phi}_n$	$\hat{\phi}$	$\kappa(R_n)$
FS1836 ( $b^{(1)}$ )	50	2.2045	2.2045	0.2497	0.2497	1.2e+11
FS1836 ( $b^{(2)}$ )	23	2.0876	2.0421	1.4332	1.2793	3.4e+09

Our bounds, in particular (4.18), largely overestimate the actual convergence curves. We therefore have not included them in Figure 6.9, which shows the GMRES convergence curves for both right-hand sides. Still, our theory provides useful information because of the different largest gaps in  $\Lambda(Q)$  and  $\Lambda(\hat{Q})$ . GMRES stagnates for  $b^{(1)}$ , so  $\hat{\phi}$  must be small. On the other hand,  $\phi$  for  $b^{(1)}$  is large and the behavior of the approximations (cf. Figure 6.11) predicts for small  $n$  that GMRES will eventually

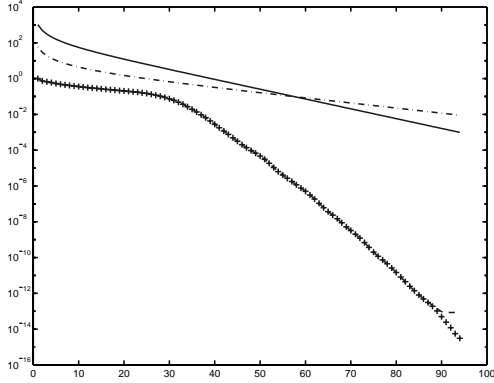


FIG. 6.5. PDE225: Computed (dashed) and updated (pluses) residual norms, bounds (4.17) (solid) and (4.18) (dashdot).

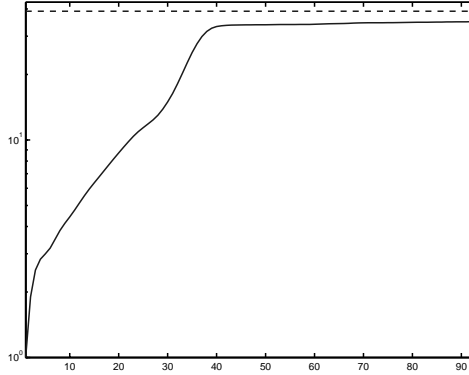


FIG. 6.6. PDE225:  $\kappa(R_n)$ .

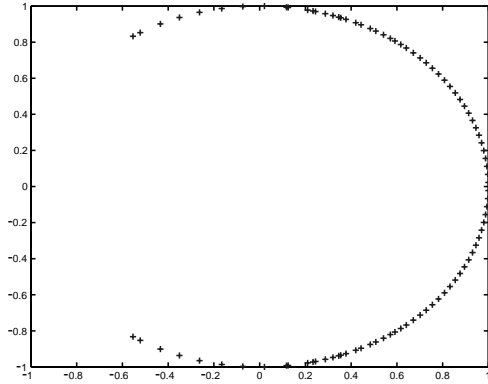


FIG. 6.7. PDE225:  $\Lambda(Q)$ .

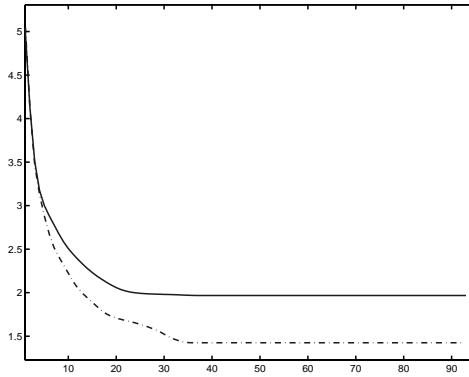


FIG. 6.8. PDE225:  $\phi_n$  (solid) and  $\hat{\phi}_n$  (dashdot).

stop to stagnate. If interpreted as “rates of convergence,”  $\hat{\gamma}^{-1} \approx 0.99$  describes the first phase of the iteration and  $\gamma^{-1} \approx 0.85$  describes the second phase. For  $b^{(2)}$ , there is no stagnation and  $\phi$  and  $\hat{\phi}$  are closer to each other (cf. Figure 6.12).

*Experiment 6.2.* We show the application of our heuristics in the context of preconditioning. We also identify situations in which our heuristics work well and situations in which they yield unsatisfactory results.

We use  $A = \text{STEAM1}$  and  $A = \text{PDE225}$  and apply WZ-GMRES with  $x_0 = 0$  to preconditioned linear systems

$$M^{-1}A\hat{x} = M^{-1}A[1, \dots, 1]^T.$$

For the preconditioner  $M$  we use

- the diagonal of  $A$  (Jacobi preconditioner),
- the lower triangular part of  $A$  (Gauss–Seidel preconditioner), and
- an incomplete LU factorization of  $A$  (ILU(0) preconditioner).

See, e.g., [27, Chapter 10] for more information on these preconditioners. We terminate WZ-GMRES at step  $d$  when the updated scaled residual norm is less than  $10^{-13}$ . We identify  $d$  with  $\dim AK_N(A, r_0)$  and thus consider the last computed approximations as equal to the quantities used in Theorem 4.6.



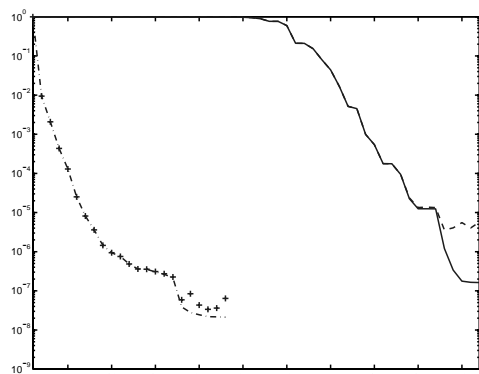


FIG. 6.9. FS1836: Computed (dashed) and updated (solid) residual norms for  $b^{(1)}$ , computed (pluses) and updated (dashdot) residual norms for  $b^{(2)}$ .

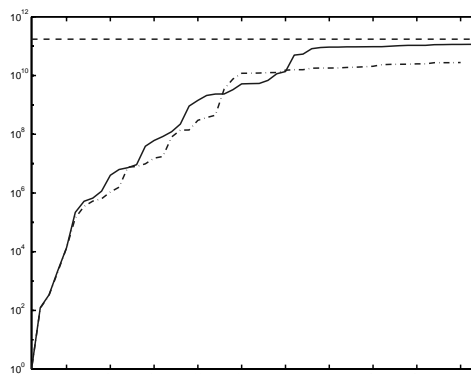


FIG. 6.10. FS1836:  $\kappa(R_n)$  for  $b^{(1)}$  (solid) and  $b^{(2)}$  (dashdot).

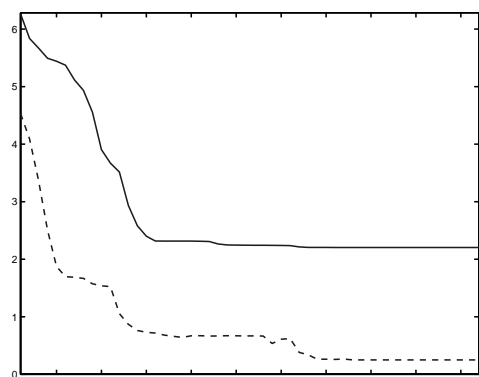


FIG. 6.11. FS1836:  $\phi_n$  (solid) and  $\hat{\phi}_n$  (dashed) for  $b^{(1)}$ .

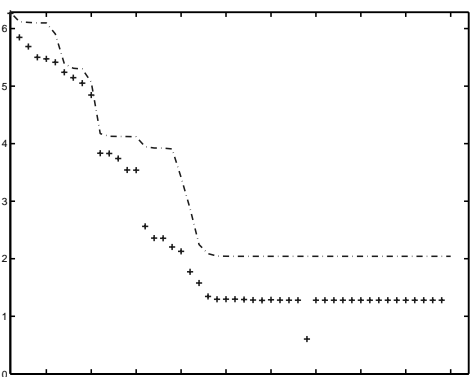


FIG. 6.12. FS1836:  $\phi_n$  (dashdot) and  $\hat{\phi}_n$  (pluses) for  $b^{(2)}$ .

First, consider the results for STEAM1 (see Figures 6.13–6.15):

$M$	$d$	$\phi$	$\hat{\phi}$	$\kappa(R)$	$\kappa(M^{-1}A)$
Jacobi	22	4.1638	2.5861	2.4e+06	3.4e+06
Gauss–Seidel	15	5.4312	5.4452	2.16	2.25
ILU(0)	7	6.2691	6.2692	1.02	1.02

As a comparison of Figures 6.13 and 6.14 shows, the difference in the three convergence curves is reflected very well by the computed approximations  $\phi_n$ . Heuristic 5.1 correctly predicts the differences in the GMRES convergence behavior for the three preconditioned systems. Of particular interest is the large difference between the  $\phi_n$  for Gauss–Seidel and ILU(0): It implies after only 6 steps that GMRES will converge faster for ILU(0). In contrast, after 10 steps the actual residual norms for these two preconditioners are still approximately equal. Figure 6.15 shows that for Gauss–Seidel and ILU(0) the approximations  $\hat{\phi}_n$  behave similar to the  $\phi_n$ . For Jacobi preconditioning we note that  $\phi$  is much larger than  $\hat{\phi}$ , which is reflected early by the approximations. This difference occurs simultaneously with the appearance of the

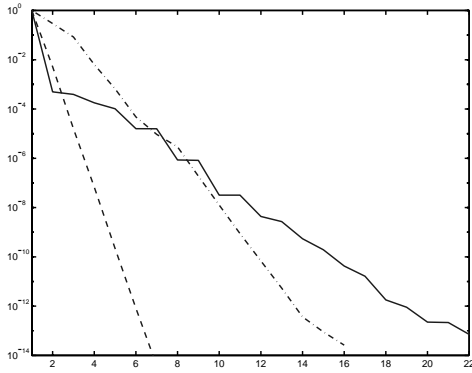


FIG. 6.13. STEAM1: GMRES convergence with Jacobi (solid), Gauss–Seidel (dashdot), and ILU(0) (dashed) preconditioning.

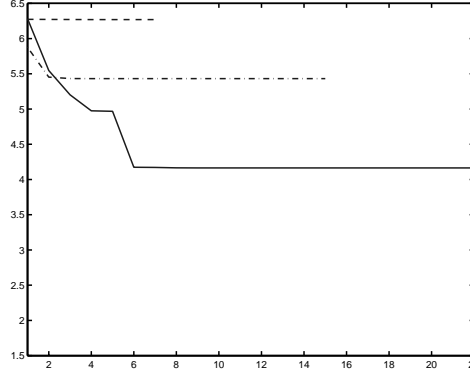


FIG. 6.14. STEAM1:  $\phi_n$  for Jacobi (solid), Gauss–Seidel (dashdot), and ILU(0) (dashed) preconditioning.

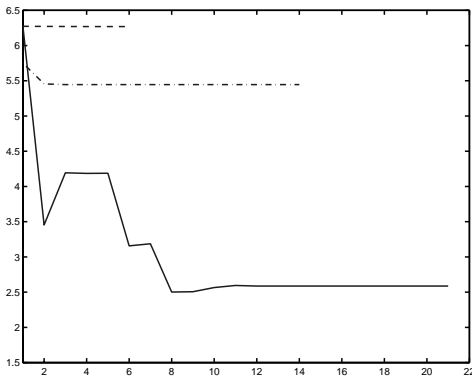


FIG. 6.15. STEAM1:  $\hat{\phi}_n$  for Jacobi (solid), Gauss–Seidel (dashdot), and ILU(0) (dashed) preconditioning.

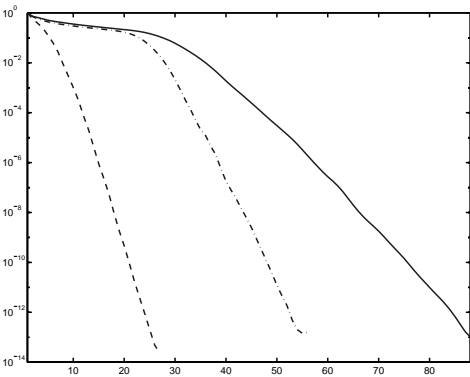


FIG. 6.16. PDE225: GMRES convergence with Jacobi (solid), Gauss–Seidel (dashdot), and ILU(0) (dashed) preconditioning.

large constant  $\kappa(R)$  in (4.17) (cf. Experiment 6.1). Figure 6.15 also shows that, while due to interlacing the  $\phi_n$  usually decrease monotonically, the  $\hat{\phi}_n$  might behave more erratically.

Experiments with PDE225 (see Figures 6.16–6.18) show limitations of the predictive value of our bounds. As we see in the following table and in Figures 6.17 and 6.18, the difference in the GMRES behavior for Jacobi and Gauss–Seidel preconditioning (cf. Figure 6.16) is neither reflected in nor predicted by our bounds.

$M$	$d$	$\phi$	$\hat{\phi}$	$\kappa(R)$	$\kappa(M^{-1}A)$
Jacobi	88	1.9289	1.3852	29.8	33.0
Gauss–Seidel	55	2.0605	1.4552	21.7	23.7
ILU(0)	26	4.0565	3.9351	3.7	4.3

In agreement with our earlier observations, the outcome of these experiments can be explained as follows:

In case of ILU(0), our bounds work very well. We have not plotted them, but the reader might check that they both closely approximate the actual convergence curve.

For Jacobi and Gauss–Seidel, GMRES has an initial phase of stagnation. Thus,  $\hat{\phi}$

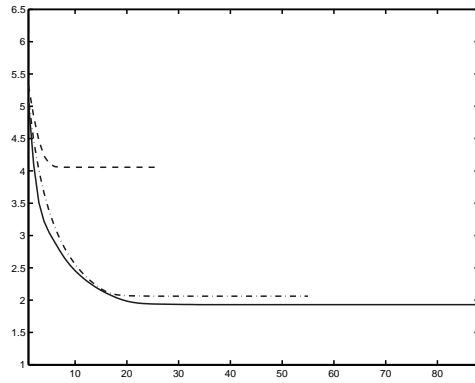


FIG. 6.17. *PDE225*:  $\phi_n$  for Jacobi (solid), Gauss–Seidel (dashdot), and ILU(0) (dashed) preconditioning.

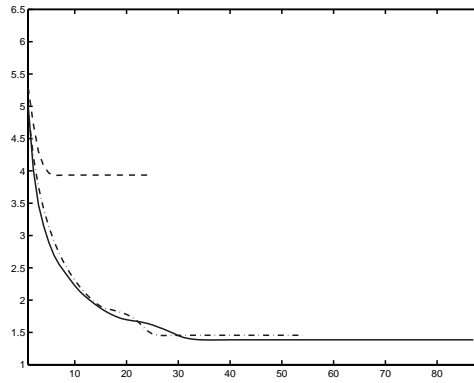


FIG. 6.18. *PDE225*:  $\hat{\phi}_n$  for Jacobi (solid), Gauss–Seidel (dashdot), and ILU(0) (dashed) preconditioning.

must in both cases be relatively small and cannot reflect later phase of faster convergence. Because of the constant  $\kappa(R)$ ,  $\phi$  in both cases might be—and actually is—larger than  $\hat{\phi}$ . But, because of the initial stagnation, the values of  $\phi$  cannot possibly differ much from each other, because the two respective constants  $\kappa(R)$  are only of a similar and moderate size.

*Experiment 6.3.* We study our bounds in the case of two significantly large gaps in the spectra of the matrices  $Q$  and  $\hat{Q}$ . It was difficult to find a suitable matrix  $A$ : With the imposed diagonal scalings of the  $R$ -factors, no nonunitary (or nonorthogonal) matrix  $A$  we tested led to matrices  $Q$  or  $\hat{Q}$  having more than one significant gap in their spectra. We therefore had to construct an orthogonal matrix  $A$  to start with.

One source of orthogonal matrices is the well-known Ising model (see, e.g., [22, p. 26]). The Ising matrices are products of two matrices  $K$  and  $L$  of the form

$$K = \begin{bmatrix} E(\alpha) & & \\ & \ddots & \\ & & E(\alpha) \end{bmatrix} \in \mathbf{R}^{2s \times 2s}, \quad E(\alpha) = \begin{bmatrix} \cos \alpha & \sin \alpha \\ -\sin \alpha & \cos \alpha \end{bmatrix}, \quad \text{and}$$

$$L = \begin{bmatrix} \cos \beta & & & -\sin \beta \\ & E(\beta) & & \\ \vdots & & \ddots & \vdots \\ \sin \beta & & & E(\beta) & \cos \beta \end{bmatrix} \in \mathbf{R}^{2s \times 2s}.$$

We choose  $A = KL$  with  $N = 2s = 100$ ,  $\alpha = \pi/4$ ,  $\beta = \pi/6$  and call this matrix ISING100. We use a random right-hand side generated with the MATLAB call `b=rand(100,1)`, and  $x_0 = 0$ . We terminate WZ-GMRES at step  $d$  when the updated scaled residual norm is less than  $10^{-13}$ . As above, we consider the last computed approximations as equal to the quantities used in Theorem 4.6:

	$d$	$\phi$	$\hat{\phi}$	$\kappa(R)$
ISING100	52	3.6652	3.6652	1.0

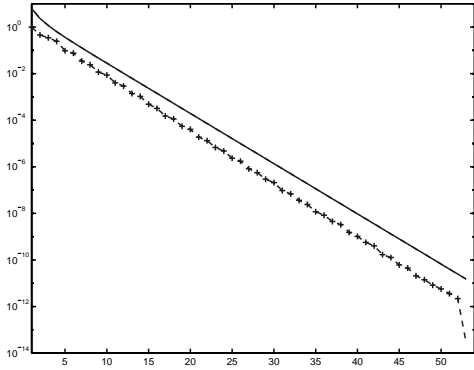


FIG. 6.19. *ISING100*: Computed (dashed) and updated (pluses) residual norms, bounds (4.17) (solid) and (4.18) (dashdot).

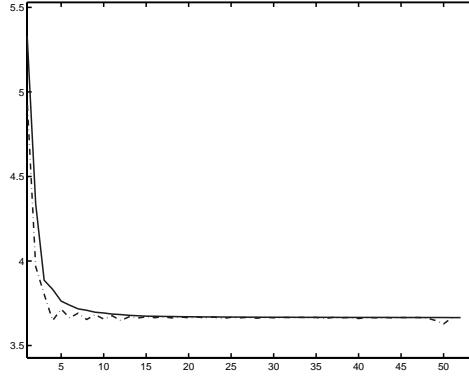


FIG. 6.20. *ISING100*:  $\phi_n$  (solid) and  $\hat{\phi}_n$  (dashdot).

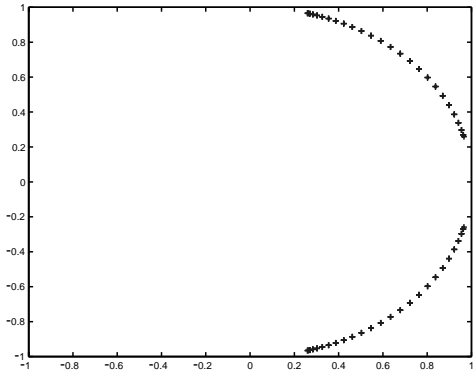


FIG. 6.21. *ISING100*:  $\Lambda(A)$ .

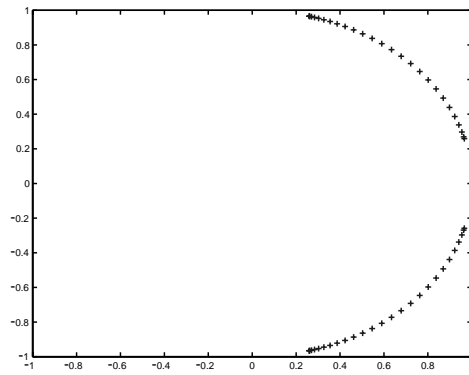


FIG. 6.22. *ISING100*:  $\Lambda(Q)$ .

Figure 6.19 shows the computed (pluses) and updated (dashed) scaled residual norms produced by WZ-GMRES. These two curves are identical, which is typical for well-conditioned matrices. We also plot our bounds (4.17) and (4.18). The bounds coincide as well, because  $\phi = \hat{\phi}$  and  $\kappa(R) = 1$ . The experiment shows that, despite the occurrence of two significant spectral gaps, our bounds describe the actual convergence curve very well. Figure 6.20 shows the approximations  $\phi_n$  and  $\hat{\phi}_n$ . Note that, due to interlacing, the curve of the  $\phi_n$  is smoother than the curve of the  $\hat{\phi}_n$ . Finally, Figures 6.21 and 6.22 show that the spectra of  $A$  and  $Q$  almost coincide.

**7. Conclusions.** We have derived two new bounds on the residual norms of GMRES, which are conceptually different from the standard worst-case analysis. Our bounds depend on the initial guess and can be approximated from information generated during the run of a certain GMRES implementation. The approximations allow predictions of how the algorithm will perform in the iteration. We have presented heuristics for such predictions, which are often confirmed in numerical experiments. We have not studied the use of predictions rigorously in the context of large-scale applications. However, we believe that the tools presented here might prove useful in some areas.

Numerical experiments have clearly demonstrated the dependence of our bounds on  $A$  and  $r_0$ . Hence, it will be generally difficult to derive a priori estimates on  $\phi$  and  $\hat{\phi}$  based solely on properties of  $A$ . Because information about the convergence behavior of GMRES seems to be “distilled” into  $\phi$  and  $\hat{\phi}$ , a priori estimates for these quantities are of great interest and should be the subject of further investigations.

Interesting related work was recently done by Knizhnerman [18]. He gives bounds for the largest gap in  $\Lambda(\hat{Q})$ , which depend on the GMRES convergence curve for the given  $A$  and  $r_0$ . In particular, the bounds show that fast GMRES convergence without steps close to stagnation implies large gaps. Our results, which show that large gaps in  $\Lambda(\hat{Q})$  imply fast GMRES convergence, are therefore in some sense dual to Knizhnerman’s. It is also interesting to note that, in accordance with our observations, the matrices  $\hat{Q}$  in Knizhnerman’s examples have only one significant gap in their spectra (cf. [18, Figures 2.1 and 2.2]).

**Acknowledgments.** Parts of this paper first appeared in [20] and my Ph.D. thesis [21]. Thanks to my advisor Ludwig Elsner for his careful reading of the manuscript and several suggestions that improved the paper. I am happy to acknowledge many fruitful discussions with Zdeněk Strakoš, who made very helpful comments and pointed out several references. Part of the work was performed while I was visiting the Institute of Computer Science of the Academy of Sciences of the Czech Republic. Thanks to Leonid Knizhnerman and Paul Saylor for careful reading of the manuscript and many comments, and to Tino Koch, Reinhard Nabben, and Miro Rozložník for helpful discussions.

#### REFERENCES

- [1] M. ARIOLI, V. PTÁK, AND Z. STRAKOŠ, *Krylov sequences of maximal length and convergence of GMRES*, BIT, 38 (1998), pp. 636–643.
- [2] W. ARNOLDI, *The principle of minimized iterations in the solution of the matrix eigenvalue problem*, Quart. Appl. Math., 9 (1951), pp. 17–29.
- [3] C. H. BISCHOF, *Incremental condition estimation*, SIAM J. Matrix Anal. Appl., 11 (1990), pp. 312–322.
- [4] P. N. BROWN AND H. F. WALKER, *GMRES on (nearly) singular systems*, SIAM J. Matrix Anal. Appl., 18 (1997), pp. 37–51.
- [5] S. L. CAMPBELL, I. C. F. IPSEN, C. T. KELLEY, AND C. D. MEYER, *GMRES and the minimal polynomial*, BIT, 36 (1996), pp. 664–675.
- [6] J. B. CONWAY, *A Course in Functional Analysis*, 2nd ed., Springer-Verlag, New York, 1990.
- [7] J. DRKOŠOVÁ, A. GREENBAUM, M. ROZLOŽNÍK, AND Z. STRAKOŠ, *Numerical stability of GMRES*, BIT, 35 (1995), pp. 309–330.
- [8] M. EIERMANN, *Field of Values and Iterative Methods*, manuscript, 1996.
- [9] L. ELSNER AND C. HE, *Perturbation and interlace theorems for the unitary eigenvalue problem*, Linear Algebra Appl., 188 (1993), pp. 207–229.
- [10] V. FABER, W. JOUBERT, E. KNILL, AND T. MANTEUFFEL, *Minimal residual method stronger than polynomial preconditioning*, SIAM J. Matrix Anal. Appl., 17 (1996), pp. 707–729.
- [11] G. H. GOLUB AND C. F. VAN LOAN, *Matrix Computations*, 3rd ed., The Johns Hopkins University Press, Baltimore, MD, 1996.
- [12] A. GREENBAUM, *Iterative Methods for Solving Linear Systems*, SIAM, Philadelphia, 1997.
- [13] A. GREENBAUM AND L. GURVITS, *MAX-MIN properties of matrix factor norms*, SIAM J. Sci. Comput., 15 (1994), pp. 348–358.
- [14] A. GREENBAUM, V. PTÁK, AND Z. STRAKOŠ, *Any nonincreasing convergence curve is possible for GMRES*, SIAM J. Matrix Anal. Appl., 17 (1996), pp. 465–469.
- [15] A. GREENBAUM AND Z. STRAKOŠ, *Matrices that generate the same Krylov residual spaces*, in Recent Advances in Iterative Methods, G. H. Golub, A. Greenbaum, and M. Luskin, eds., IMA Vol. Math. Appl. 60, Springer-Verlag, New York, 1994, pp. 95–118.
- [16] A. GREENBAUM AND L. N. TREFETHEN, *GMRES/CR and Arnoldi/Lanczos as matrix approximation problems*, SIAM J. Sci. Comput., 15 (1994), pp. 359–368.

- [17] W. JOUBERT, *A robust GMRES-based adaptive polynomial preconditioning algorithm for non-symmetric linear systems*, SIAM J. Sci. Comput., 15 (1994), pp. 427–439.
- [18] L. KNIZHNERMAN, *On GMRES-Equivalent Bounded Operators*, preprint, 1998.
- [19] T. KOCH AND J. LIESEN, *The conformal bratwurst maps and associated Faber polynomials*, Numer. Math., to appear.
- [20] J. LIESEN, *Computable Convergence Bounds for GMRES*, Preprint 98-043, Sonderforschungsbereich 343, Universität Bielefeld, Bielefeld, Germany, 1998.
- [21] J. LIESEN, *Construction and Analysis of Polynomial Iterative Methods for Non-Hermitian Systems of Linear Equations*, Ph.D. thesis, Fakultät für Mathematik, Universität Bielefeld, Bielefeld, Germany, 1998; also available from <http://archiv.ub.uni-bielefeld.de/disshabi/mathe.htm>.
- [22] M. MARCUS AND H. MINC, *A Survey of Matrix Theory and Matrix Inequalities*, Dover, New York, 1992.
- [23] N. M. NACHTIGAL, S. REDDY, AND L. N. TREFETHEN, *How fast are nonsymmetric matrix iterations?*, SIAM J. Matrix Anal. Appl., 13 (1992), pp. 778–795.
- [24] M. ROZLOŽNÍK, *Numerical Stability of the GMRES Method*, Ph.D. thesis, Institute of Computer Science, Academy of Sciences of the Czech Republic, 1997.
- [25] M. ROZLOŽNÍK AND Z. STRAKOŠ, *Variants of the residual minimizing Krylov space methods*, in Proceedings of the Eleventh Summer School on Software and Algorithms of Numerical Mathematics, I. Marek, ed., 1995, pp. 208–225.
- [26] Y. SAAD, *Variations on Arnoldi's method for computing eigenelements of large unsymmetric matrices*, Linear Algebra Appl., 34 (1980), pp. 269–295.
- [27] Y. SAAD, *Iterative Methods for Solving Sparse Linear Systems*, PWS Publishing Company, Boston, MA, 1996.
- [28] 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.
- [29] V. SMIRNOV AND N. LEBEDEV, *Functions of a Complex Variable—Constructive Theory*, MIT Press, Cambridge, MA, 1968.
- [30] G. STARKE, *Field-of-values analysis of preconditioned iterative methods for nonsymmetric elliptic problems*, Numer. Math., 78 (1997), pp. 103–117.
- [31] Z. STRAKOŠ, *Convergence and numerical behavior of the Krylov space methods*, in Proceedings of the NATO ASI Institute “Algorithms for Large Sparse Linear Algebraic Systems: The State of the Art and Applications in Science and Engineering” (invited lectures), Kluwer Academic Publishers, 1998, pp. 175–196.
- [32] THE MATHWORKS, *Using MATLAB (December 1996)*, The MathWorks Inc., Natick, MA, 1996.
- [33] K.-C. TOH, *GMRES vs. ideal GMRES*, SIAM J. Matrix Anal. Appl., 18 (1997), pp. 30–36.
- [34] H. F. WALKER, *Implementation of the GMRES method using Householder transformations*, SIAM J. Sci. Statist. Comput., 9 (1988), pp. 152–163.
- [35] H. F. WALKER AND L. ZHOU, *A simpler GMRES*, Numer. Linear Algebra Appl., 1 (1994), pp. 571–581.
- [36] J. WALSH, *Interpolation and Approximation by Rational Functions in the Complex Domain*, 5th ed., AMS, Providence, RI, 1969.