# The Rate of Convergence of GMRES on a Tridiagonal Toeplitz Linear System

**Ren-Cang Li**
**Wei Zhang**

# The Rate of Convergence of GMRES
# on a Tridiagonal Toeplitz Linear System

Ren-Cang Li [*]         Wei Zhang [†]

February 2007

### Abstract

The Generalized Minimal Residual method (GMRES) is often used to solve a nonsymmetric linear system $Ax = b$. But its convergence analysis is a rather difficult task in general. A commonly used approach is to diagonalize $A = X\Lambda X^{-1}$ and then separate the study of GMRES convergence behavior into optimizing the condition number of $X$ and a polynomial minimization problem over $A$'s spectrum. This artificial separation could greatly overestimate GMRES residuals and likely yields error bounds that are too far from the actual ones. On the other hand, considering the effects of both $A$'s spectrum and the conditioning of $X$ at the same time poses a difficult challenge, perhaps impossible to deal with in general but only possible for certain particular linear systems. This paper will do so for a (nonsymmetric) tridiagonal Toeplitz system. Sharp error bounds on and sometimes exact expressions for residuals are obtained. These expressions and/or bounds are in terms of the three parameters that define $A$ and Chebyshev polynomials of the first kind.

## 1   Introduction

The Generalized Minimal Residual method (GMRES) is often used to solve a nonsymmetric linear system $Ax = b$. The basic idea is to seek approximate solutions, optimal in certain sense, from the so-called Krylov subspaces. Specifically, the $k$th approximation $x_k$ is sought so that the $k$th residual $r_k = b - Ax_k$ satisfies [20] (without loss of generality, we take initially $x_0 = 0$ and thus $r_0 = b$.)

$$\|r_k\|_2 = \min_{y \in \mathcal{K}_k} \|b - Ay\|_2,$$

where the *kth Krylov subspace* $\mathcal{K}_k \equiv \mathcal{K}_k(A, b)$ of $A$ on $b$ is defined as

$$\mathcal{K}_k \equiv \mathcal{K}_k(A, b) \stackrel{\text{def}}{=} \mathsf{span}\{b, Ab, \dots, A^{k-1}b\}, \tag{1.1}$$

and generic norm $\|\cdot\|_2$ is the usual $\ell_2$ norm of a vector or the spectral norm of a matrix.

This paper is concerned with the convergence analysis of GMRES on linear system $Ax = b$ whose coefficient matrix $A$ is a (nonsymmetric) tridiagonal Toeplitz coefficient matrix:

$$A = \begin{pmatrix} \lambda & \mu & & \\ \nu & \ddots & \ddots & \\ & \ddots & \ddots & \mu \\ & & \nu & \lambda \end{pmatrix},$$

where $\lambda$, $\mu$, $\nu$ are assumed nonzero and possibly complex. Linear systems as such have been studied quite extensively in the past. For the nonsymmetric case, i.e., $\mu \neq \nu$ as we are interested here, most update-to-date and detailed studies are due to Ernst [9] and Liesen and Strakoš [17]. Both papers, motivated to better understand the convergence behavior of GMRES on a convection-diffusion model problem [18], established various bounds on residual ratios. Ernst's bounds to which we shall return are comparable to ours, while most results in Liesen and Strakoš [17] are of qualitative nature, intended to explain GMRES convergence behaviors for such linear systems. In particular, Liesen and Strakoš showed that GMRES for tiny $|\mu|$ behaves much like GMRES after setting $\mu$ to 0.

Throughout this paper, exact arithmetic is assumed, $A$ is $N$-by-$N$, and $k$ is GMRES iteration index. Since in exact arithmetic GMRES computes the exact solution in at most $N$ steps, $r_N = 0$. For this reason, we restrict $k < N$ at all times. This restriction is needed to interpret our later results concerning (worst) asymptotic speed in terms of certain limits of $\|r_k\|^{1/k}$ as $k \to \infty$.

Our first main contribution in this paper is the following error bound (Theorem 2.1)

$$\frac{\|r_k\|_2}{\|r_0\|_2} \leq \sqrt{k+1} \left[ \sum_{j=0}^{k} \zeta^{2j} |T_j(\tau)|^2 \right]^{-1/2}, \tag{1.2}$$

where $T_j(t)$ is the $j$th Chebyshev polynomial of the first kind, and

$$\xi = -\frac{\sqrt{\mu\nu}}{\nu}, \quad \tau = \frac{\lambda}{2\sqrt{\mu\nu}}, \quad \zeta = \min\{|\xi|, |\xi|^{-1}\}.$$

We will also prove that this upper bound is nearly achieved by $b = e_1$ (the first column of the identity matrix) when $|\xi| \leq 1$ or by $b = e_N$ (the last column of the identity matrix) when $|\xi| \geq 1$. By "nearly achieved", we mean it is within a factor about at most $(k+1)^{3/2}$ of the exact residual ratios.

Our second main contribution is about the worst asymptotic speed of $\|r_k\|_2$ among all possible $r_0$. It is proven that (Theorem 2.2)

$$\lim_{k \to \infty} \left[ \sup_{r_0} \frac{\|r_k\|_2}{\|r_0\|_2} \right]^{1/k} = \min \left\{ (\zeta\rho)^{-1}, 1 \right\}, \tag{1.3}$$

where $\rho = \max \left\{ \left| \tau + \sqrt{\tau^2 - 1} \right|, \left| \tau - \sqrt{\tau^2 - 1} \right| \right\}$. Technically speaking, $\left[ \sup_{r_0} \|r_k\|_2 \|r_0\|_2 \right]^{1/k}$ is not a sequence because of the freedom in $N$ (except $N > k$), namely for each $N > k$ it

renders a number. Nonetheless, the limit in (1.3) can be meaningfully interpreted as follows: *for any given $\epsilon > 0$, there exists a positive integer $K$ such that*

$$\left| \left[ \sup_{r_0} \frac{\|r_k\|_2}{\|r_0\|_2} \right]^{1/k} - \min\left\{ (\zeta\rho)^{-1}, 1 \right\} \right| < \epsilon \quad \text{for all } N > k \geq K.$$

This interpretation will be adopted to understand similar limits later in this paper. In the case when $r_0$ is given, $\sup_{r_0}$ will be dropped. A related work that also studied asymptotic speed of convergence but for the conjugate gradient method (CG) and special right-hand sides and $\lambda = 2$ and $\mu = \nu = -1$ is [2], where that $N/k$ remains constant is required as $k \to \infty$.

A by-product of (1.3) is that the worst asymptotic speed can be separated into the factor $\zeta^{-1} \geq 1$ contributed by $A$'s departure from normality and the factor $\rho^{-1}$ contributed by $A$'s eigenvalue distribution. Take, for example, $\lambda = 0.5$, $\mu = -0.3$, and $\nu = 0.7$ which was used in [4, p.562] to explain the effect of nonnormality on GMRES convergence. We have $(\zeta\rho)^{-1} = 0.90672$, whereas in [4, p.562] it is implied $\|r_k\|_2/\|r_0\|_2 \leq (0.913)^k$ for $N = 50$, which is rather good, considering that $N = 50$ is rather small.

Ernst [9], in our notation, obtained the following inequality: *if $A$'s field of values does not contain the origin, then*

$$\frac{\|r_k\|_2}{\|r_0\|_2} \leq \left( |\xi|^k + |\xi|^{-k} \right) \frac{\widetilde{\rho}^k}{1 - \widetilde{\rho}^{2k}}, \tag{1.4}$$

*where $\widetilde{\rho} = \max\left\{ \left| \widetilde{\tau} + \sqrt{\widetilde{\tau}^2 - 1} \right|, \left| \widetilde{\tau} - \sqrt{\widetilde{\tau}^2 - 1} \right| \right\}$ and $\widetilde{\tau} = \left[ \cos \frac{\pi}{N+1} \right]^{-1} \tau$.* Our bound (1.2) is comparable to Ernst's bound for large $N$. This can be seen by noting that for $N$ large enough, $\widetilde{\tau} \approx \tau$ and $\widetilde{\rho} \approx \rho$, and that $T_j(\tau) \approx \frac{1}{2}\rho^j$ when $\rho > 1$ and $|\zeta|^{-k} \leq |\xi|^k + |\xi|^{-k} \leq 2|\zeta|^{-k}$. Ernst's bound also leads to

$$\limsup_{k \to \infty} \left[ \sup_{r_0} \frac{\|r_k\|_2}{\|r_0\|_2} \right]^{1/k} \leq \min\left\{ (\zeta\rho)^{-1}, 1 \right\}. \tag{1.5}$$

In differentiating our contributions here from Ernst's, we use a different technique to arrive at (1.2) and (1.3). While our approach is not as elegant as Ernst's which was based on $A$'s field of values (see also [5]), it allows us to establish both lower and upper bounds on relative residuals for special right-hand sides to conclude that our bound is nearly achieved. Also (1.3) is an equality while only an inequality (1.5) can be deduced from Ernst's bound and approach.

We also obtain residual bounds and exact expressions especially for right-hand sides $b = e_1$ and $b = e_N$ (Theorems 2.3 and 2.4). They suggest, besides the sharpness of (1.2), an interesting GMRES convergence behavior. For $b = e_1$, that $|\xi| > 1$ speeds up GMRES convergence, and in fact $\|r_k\|_2$ is roughly proportional to $|\xi|^{-k}$. So the bigger the $|\xi|$ is, the faster the convergence will be. Note as $|\xi|$ gets bigger, $A$ gets further away from a normal matrix. Thus, loosely speaking, the nonnormality contributes to the convergence rate in the positive way. Nonetheless this does not contradict our usual perception that high nonnormality is bad for GMRES if the worst behavior of GMRES among all $b$ is considered. This mystery can be best explained by looking at the extreme case: $|\xi| = \infty$, i.e., $\nu = 0$,

for which $b = e_1$ is an eigenvector (and convergence occurs in just one step). In general for $\nu \neq 0$, as $|\xi|$ gets bigger and bigger, roughly speaking $b = e_1$ comes closer and closer to $A$'s invariant subspaces of lower dimensions and consequently speedier convergence is witnessed. Similar comments apply to the case when $b = e_N$.

The rest of this paper is organized as follows. We state our main results in Section 2. Tedious proofs that rely on residual reformulation involving rectangular Vandermonde matrices and complicated analysis will be presented separately in Section 3. Exact residual norm formulas for two special right-hand sides $b = e_1$ and $e_N$ are established in Section 4. Finally in Section 5 we present our concluding remarks.

**Notation.** Throughout this paper, $\mathbb{K}^{n \times m}$ is the set of all $n \times m$ matrices with entries in $\mathbb{K}$, where $\mathbb{K}$ is $\mathbb{C}$ (the set of complex numbers) or $\mathbb{R}$ (the set of real numbers), $\mathbb{K}^n = \mathbb{K}^{n \times 1}$, and $\mathbb{K} = \mathbb{K}^1$. $I_n$ (or simply $I$ if its dimension is clear from the context) is the $n \times n$ identity matrix, and $e_j$ is its $j$th column. The superscript ".*" takes conjugate transpose while ".$^T$" takes transpose only. $\sigma_{\min}(X)$ denotes the smallest singular value of $X$.

We shall also adopt MATLAB-like convention to access the entries of vectors and matrices. The set of integers from $i$ to $j$ inclusive is $i : j$. For vector $u$ and matrix $X$, $u_{(j)}$ is $u$'s $j$th entry, $X_{(i,j)}$ is $X$'s $(i,j)$th entry, $\mathrm{diag}(u)$ is the diagonal matrix with $(\mathrm{diag}(u))_{(j,j)} = u_{(j)}$; $X$'s submatrices $X_{(k:\ell,i:j)}$, $X_{(k:\ell,:)}$, and $X_{(:,i:j)}$ consists of intersections of row $k$ to row $\ell$ and column $i$ to column $j$, row $k$ to row $\ell$ and all columns, and all rows and column $i$ to column $j$, respectively. Finally $\| \cdot \|_p$ $(1 \leq p \leq \infty)$ is the $\ell_p$ norm of a vector or the $\ell_p$ operator norm of a matrix, defined as

$$\|u\|_p = \left( \sum_j |u_{(j)}|^p \right)^{1/p}, \quad \|X\|_p = \max_{\|u\|_p=1} \|Xu\|_p.$$

$\lfloor \alpha \rfloor$ be the largest integer that is no bigger than $\alpha$, and $\lceil \alpha \rceil$ the smallest integer that is no less than $\alpha$.

## 2 Main Results

An $N \times N$ tridiagonal Toeplitz $A$ takes this form

$$A = \begin{pmatrix} \lambda & \mu & & \\ \nu & \ddots & \ddots & \\ & \ddots & \ddots & \mu \\ & & \nu & \lambda \end{pmatrix} \in \mathbb{C}^{N \times N}. \tag{2.1}$$

Throughout the rest of this paper, $\nu$, $\lambda$, and $\mu$ are reserved as the defining parameters of $A$ in (2.1), and set

$$\xi = -\frac{\sqrt{\mu\nu}}{\nu}, \quad \tau = \frac{\lambda}{2\sqrt{\mu\nu}}, \quad \zeta = \min\left\{ |\xi|, \frac{1}{|\xi|} \right\}, \tag{2.2}$$

$$\rho = \max\left\{ \left| \tau + \sqrt{\tau^2 - 1} \right|, \left| \tau - \sqrt{\tau^2 - 1} \right| \right\}. \tag{2.3}$$

Any branch of $\sqrt{\mu\nu}$, once picked and fixed, is a valid choice in this paper. Note $\rho \geq 1$ always because $(\tau + \sqrt{\tau^2 - 1})(\tau - \sqrt{\tau^2 - 1}) = 1$. In particular if $\lambda \in \mathbb{R}$, $\mu < 0$ and $\nu > 0$, then $\rho = |\tau| + \sqrt{|\tau|^2 + 1}$.

Recall Chebyshev polynomials of the first kind:

$$
\begin{aligned}
T_m(t) &= \cos(m \arccos t) & \text{for } |t| \leq 1, & \quad (2.4) \\
&= \frac{1}{2}\left(t + \sqrt{t^2 - 1}\right)^m + \frac{1}{2}\left(t - \sqrt{t^2 - 1}\right)^m & \text{for } |t| \geq 1, & \quad (2.5)
\end{aligned}
$$

and define

$$
\Phi_{k+1}^{(+)}(\tau, \xi) \overset{\text{def}}{=} \sum_{j=0}^{k}{}' |\xi|^{2j}\, |T_j(\tau)|^2 , \tag{2.6}
$$

$$
\Phi_{k+1}^{(-)}(\tau, \xi) \overset{\text{def}}{=} \sum_{j=0}^{k}{}' |\xi|^{-2j}\, |T_j(\tau)|^2 , \tag{2.7}
$$

$$
\Phi_{k+1}(\tau, \xi) \overset{\text{def}}{=} \sum_{j=0}^{k}{}' \zeta^{2j}\, |T_j(\tau)|^2 \equiv \min\left\{\Phi_{k+1}^{(+)}(\tau, \xi), \Phi_{k+1}^{(-)}(\tau, \xi)\right\}, \tag{2.8}
$$

where $\sum_j'$ means the first term is halved.

## 2.1 General Right-hand Sides

Our first main result is given in Theorem 2.1 whose proof, along with the proofs of other results in the section, involve complicated computations and will be postponed to Section 3.

**Theorem 2.1** *For $Ax = b$, where $A$ is tridiagonal Toeplitz as in* (2.1) *with nonzero (real or complex) parameters $\nu$, $\lambda$, and $\mu$. Then the $k$th GMRES residual $r_k$ satisfies for $1 \leq k < N$*

$$
\frac{\|r_k\|_2}{\|r_0\|_2} \leq \sqrt{k+1} \left[\frac{1}{2} + \Phi_{k+1}(\tau, \xi)\right]^{-1/2}. \tag{2.9}
$$

Figure 2.1 plots residual histories for several examples of GMRES with each of $b$'s entries being uniformly random in $[-1, 1]$. In each of the plots in Figure 2.1, as well as in Figures 2.2 and 2.3 below, we fix $|\tau|$ and $|\xi|$, take $\lambda = 1$ always, and then take

$$
|\mu| = \frac{|\xi|}{2|\tau|}, \quad \mu = \pm|\mu|, \quad \text{and} \quad \nu = |\nu| = \frac{1}{2|\tau\xi|}.
$$

Thus $\mu, \nu \in \mathbb{R}$, and in fact $\nu > 0$ always. When $\mu > 0$, $\xi = -\sqrt{\mu/\nu} < 0$ and $\tau = 1/(2\sqrt{\mu\nu}) > 0$, but when $\mu < 0$, both $\xi = -\iota\sqrt{|\mu/\nu|}$ and $\tau = -\iota/(2\sqrt{|\mu\nu|})$ are imaginary, where $\iota = \sqrt{-1}$ is the imaginary unit. Figure 2.1 indicates that GMRES converges much faster for $\mu < 0$ than for $\mu > 0$ in each of the plots. There is a simple explanation for this: the eigenvalues of $A$ (see (3.11) below) are further away from the origin for a pure imaginary $\tau$ than for a real $\tau$ for any fixed $|\tau|$.

Our next main result given in Theorem 2.2 tells the worst asymptotic speed for $\|r_k\|_2$.

**Theorem 2.2** *Under the conditions of* Theorem 2.1,

$$
\lim_{k\to\infty}\left[\sup_{r_0} \frac{\|r_k\|_2}{\|r_0\|_2}\right]^{1/k} = \lim_{k\to\infty}\left[\max_{r_0 \in \{e_1, e_N\}} \frac{\|r_k\|_2}{\|r_0\|_2}\right]^{1/k} = \min\{(\zeta\rho)^{-1}, 1\}. \tag{2.10}
$$

$||r_k||_2/||r_0||_2$, upper bounds ($|\tau|=0.8$, $|\xi|=0.7$)

N =50

○   $||r_k||_2/||r_0||_2$ for $\mu=0.4375$, $\nu=0.89286$
△   $||r_k||_2/||r_0||_2$ for $\mu=-0.4375$, $\nu=0.89286$

$k+1$

$||r_k||_2/||r_0||_2$, upper bounds ($|\tau|=0.8$, $|\xi|=1.2$)

N =50

○   $||r_k||_2/||r_0||_2$ for $\mu=0.75$, $\nu=0.52083$
△   $||r_k||_2/||r_0||_2$ for $\mu=-0.75$, $\nu=0.52083$

$k+1$

$||r_k||_2/||r_0||_2$, upper bounds ($|\tau|=1$, $|\xi|=0.7$)

N =50

○   $||r_k||_2/||r_0||_2$ for $\mu=0.35$, $\nu=0.71429$
△   $||r_k||_2/||r_0||_2$ for $\mu=-0.35$, $\nu=0.71429$

$k+1$

$||r_k||_2/||r_0||_2$, upper bounds ($|\tau|=1$, $|\xi|=1.2$)

N =50

○   $||r_k||_2/||r_0||_2$ for $\mu=0.6$, $\nu=0.41667$
△   $||r_k||_2/||r_0||_2$ for $\mu=-0.6$, $\nu=0.41667$

$k+1$

$||r_k||_2/||r_0||_2$, upper bounds ($|\tau|=1.2$, $|\xi|=0.7$)

N =50

○   $||r_k||_2/||r_0||_2$ for $\mu=0.29167$, $\nu=0.59524$
△   $||r_k||_2/||r_0||_2$ for $\mu=-0.29167$, $\nu=0.59524$

$k+1$

$||r_k||_2/||r_0||_2$, upper bounds ($|\tau|=1.2$, $|\xi|=1.2$)

N =50

○   $||r_k||_2/||r_0||_2$ for $\mu=0.5$, $\nu=0.34722$
△   $||r_k||_2/||r_0||_2$ for $\mu=-0.5$, $\nu=0.34722$
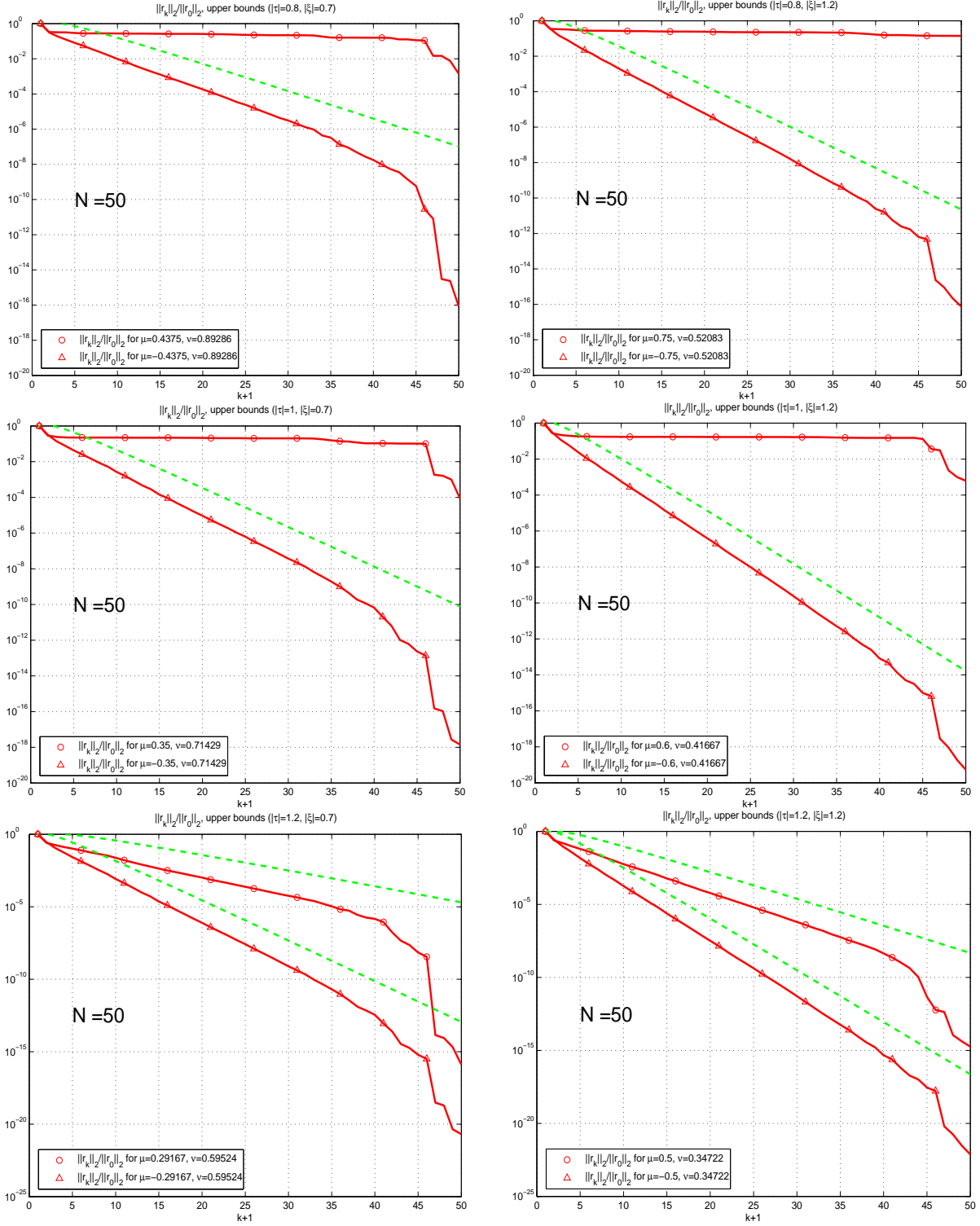
$k+1$

Figure 2.1: GMRES residuals for random $b$ uniformly in $[-1, 1]$, and their upper bounds (dashed lines) by (2.9). All indicate that our upper bounds are tight, except for the last few steps. Upper bounds for the case $\mu > 0$ in the top and bottom two plots are visually indistinguishable from the horizonal line $10^0$, suggesting slow convergence.

## 2.2 Special Right-hand Sides

We now consider three special right-hand sides: $b = e_1$ or $e_N$ or $b_{(1)}e_1 + b_{(N)}e_N$. In particular they show that the upper bound in Theorem 2.1 is within a factor about at most $(k+1)^{3/2}$ of the true residual for $b = e_1$ or $e_N$, depending on whether $|\xi| \leq 1$ or $|\xi| \geq 1$.

**Theorem 2.3** *In* Theorem 2.1, *if* $b = e_1$, *then the $k$th GMRES residual $r_k$ satisfies for* $1 \leq k < N$

$$\frac{1}{2}\left[\sum_{j=0}^{\lceil\frac{k+1}{2}\rceil-1}|\xi|^{2j}\right]^{-1}\left[\Phi_{k+1}^{(+)}(\tau,\xi)-\frac{1}{4}\right]^{-1/2} \leq \frac{\|r_k\|_2}{\|r_0\|_2} \leq \frac{1}{2}(1+|\xi|^2)\left[\Phi_{k+1}^{(+)}(\tau,\xi)-\frac{1}{4}\right]^{-1/2}. \tag{2.11}$$

*In particular,*

$$\frac{1}{2\lceil\frac{k+1}{2}\rceil}\left[\Phi_{k+1}^{(+)}(\tau,\xi)-\frac{1}{4}\right]^{-1/2} \leq \frac{\|r_k\|_2}{\|r_0\|_2} \leq \left[\Phi_{k+1}^{(+)}(\tau,\xi)-\frac{1}{4}\right]^{-1/2} \qquad \text{for } |\xi| \leq 1. \tag{2.12}$$

**Theorem 2.4** *In* Theorem 2.1, *if* $b = e_N$, *then the $k$th GMRES residual $r_k$ satisfies for* $1 \leq k < N$

$$\frac{1}{2}\left[\sum_{j=0}^{\lceil\frac{k+1}{2}\rceil-1}|\xi|^{-2j}\right]^{-1}\left[\Phi_{k+1}^{(-)}(\tau,\xi)-\frac{1}{4}\right]^{-1/2} \leq \frac{\|r_k\|_2}{\|r_0\|_2} \leq \frac{1}{2}(1+|\xi|^{-2})\left[\Phi_{k+1}^{(-)}(\tau,\xi)-\frac{1}{4}\right]^{-1/2}. \tag{2.13}$$

*In particular,*

$$\frac{1}{2\lceil\frac{k+1}{2}\rceil}\left[\Phi_{k+1}^{(-)}(\tau,\xi)-\frac{1}{4}\right]^{-1/2} \leq \frac{\|r_k\|_2}{\|r_0\|_2} \leq \left[\Phi_{k+1}^{(-)}(\tau,\xi)-\frac{1}{4}\right]^{-1/2} \qquad \text{for } |\xi| \geq 1. \tag{2.14}$$

The upper bound and the lower bound in (2.12) and these in (2.14) differ by a factor roughly $(k+1)$, and thus they are rather sharp; so are the bounds in (2.11) for $|\xi| \leq 1$ and these in (2.13) for $|\xi| \geq 1$. Comparing them to (2.9), we conclude that the upper bound by (2.9) is fairly sharp for worst possible $b$.

But the bounds in (2.11) differ by a factor $\mathcal{O}(|\xi|^{k+1})$ for $|\xi| > 1$, and thus at least one of them (upper or lower bound) is bad. Similar comments apply to the bounds in (2.13) for $|\xi| < 1$. Our numerical examples indicate that the upper bounds are rather good regardless of the magnitude of $|\xi|$ for both cases $b = e_1$ and $b = e_N$. See Figure 2.2, where only the case $b = e_1$ is presented, since the case $b = e_N$ is similar.

Given Theorems 2.3 and 2.4, it would not be unreasonable to expect that the upper bound in Theorem 2.1 would be sharp for very large or tiny $|\xi|$ within a factor possibly about at most $(k+1)^{3/2}$ for right-hand side $b$ with $b_{(i)} = 0$ for $2 \leq i \leq N-1$ and $|b_{(1)}| = |b_{(N)}| > 0$. The following theorem indeed confirms this but only for $k \leq N/2$. Our numerical examples even support that the lower bounds by (2.15) would be good for $k > N/2$ (see Figure 2.3), too, but we do not have a way to mathematically justify it yet.
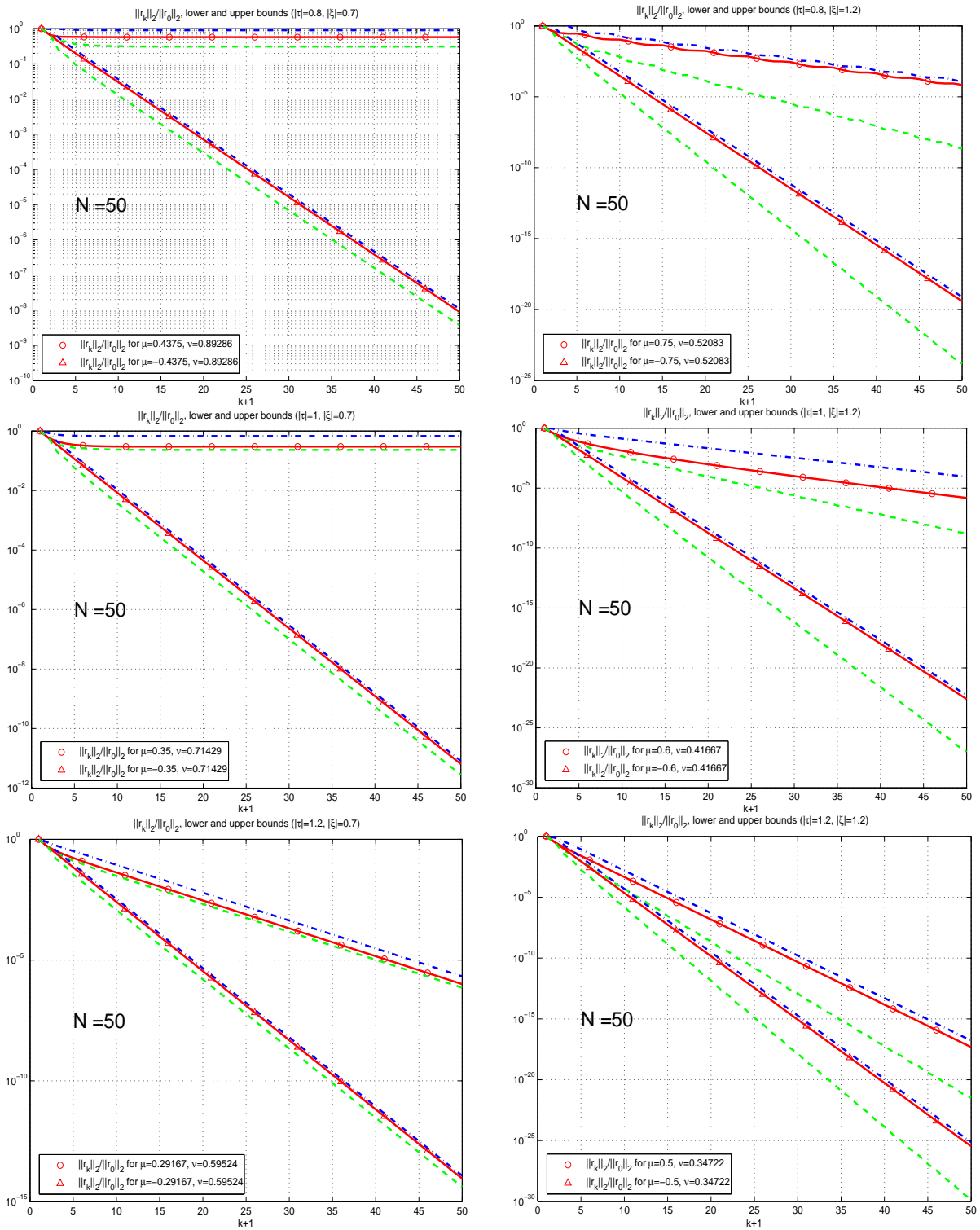
Figure 2.2: GMRES residuals for $b = e_1$, sandwiched by their lower and upper bounds by (2.11). All lower and upper bounds are very good for $|\xi| \leq 1$ as expected, but only upper bounds are good when $|\xi| > 1$. We also ran GMRES for $b = e_N$ and obtained residual history that is very much the same as for $b = e_1$ with $|\xi|$ replaced by $|\xi|^{-1}$.
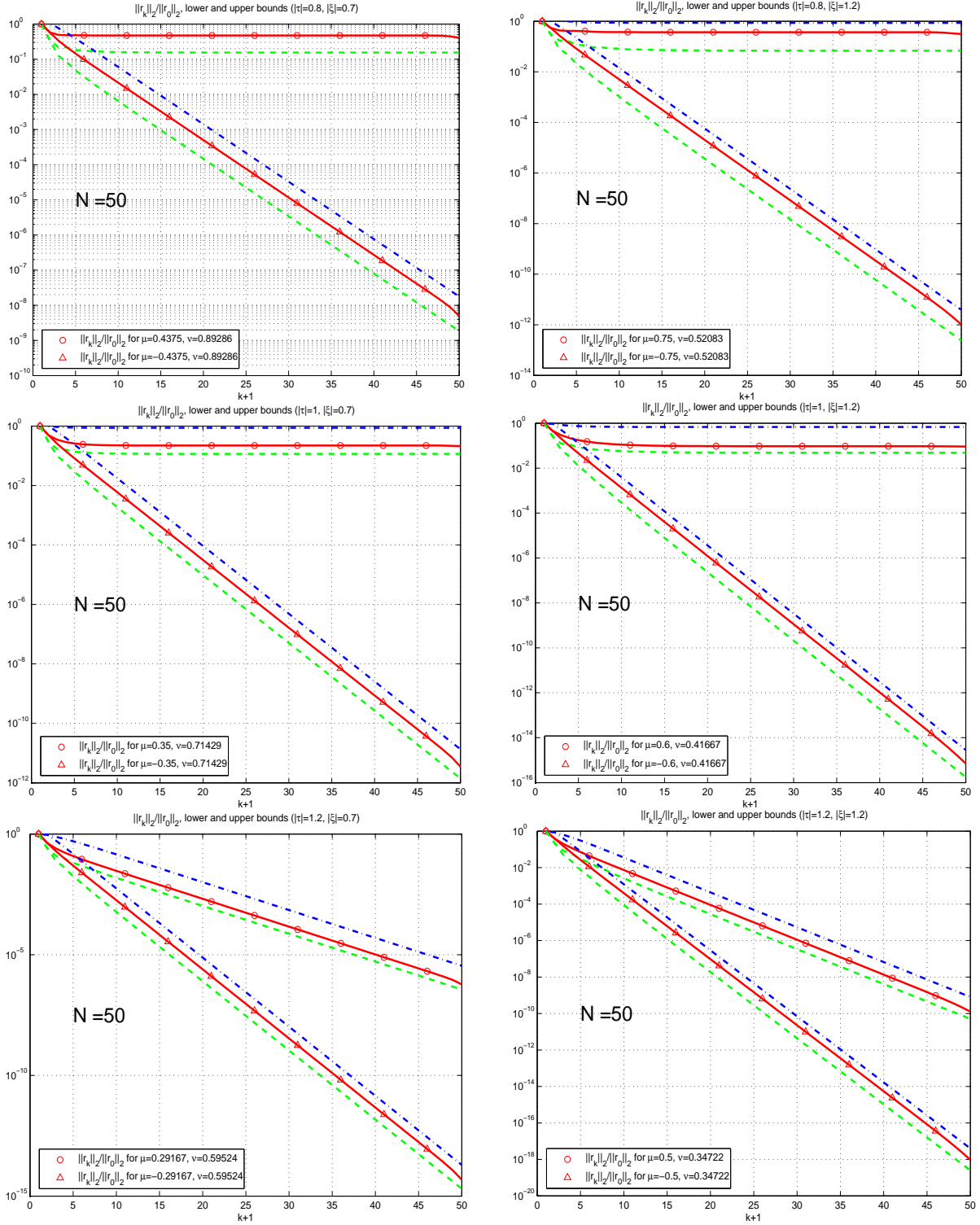
Figure 2.3: GMRES residuals for $b = e_1 + e_N$, sandwiched by their lower and upper bounds by (2.15) and (2.16). Strictly speaking, (2.15) is only proved for $k \leq N/2$, but it seems to be very good even for $k > N/2$ as well. We also ran GMRES for $b = e_1 - e_N$ and obtained residual history that is very much the same.

**Theorem 2.5** *In* Theorem 2.1, *if* $b_{(i)} = 0$ *for* $2 \leq i \leq N-1$, *then the kth GMRES residual* $r_k$ *satisfies*

$$\frac{\|r_k\|_2}{\|r_0\|_2} \geq \frac{\min_{i \in \{1,N\}} |b_{(i)}|}{2\chi \|r_0\|_2} \left[ \Phi_{k+1}(\tau,\xi) - \frac{1}{4} \right]^{-1/2} \qquad \text{for } 1 \leq k \leq N/2, \tag{2.15}$$

$$\frac{\|r_k\|_2}{\|r_0\|_2} \leq \sqrt{3} \left[ \frac{1}{2} + \Phi_{k+1}(\tau,\xi) \right]^{-1/2}, \tag{2.16}$$

*where*

$$1 < \chi = \sum_{j=0}^{\lceil \frac{k+1}{2} \rceil - 1} \zeta^{2j} \leq \left\lceil \frac{k+1}{2} \right\rceil.$$

Figures 2.2 and 2.3 plot residual histories for several examples of GMRES with $b = e_1$ and $b = e_1 + e_N$, respectively. Finally we have the following theorem about the asymptotic speeds of $\|r_k\|_2$ for $b = e_1$ and $b = e_N$.

**Theorem 2.6** *Assume the conditions of* Theorem 2.1 *hold.*

1. *Let* $b = e_1$. *If* $\rho > 1$, *then*

$$\min\{(|\xi|^2\rho)^{-1}, (|\xi|\rho)^{-1}, 1\} \leq \liminf_{k \to \infty} \left[ \frac{\|r_k\|_2}{\|r_0\|_2} \right]^{1/k} \leq \limsup_{k \to \infty} \left[ \frac{\|r_k\|_2}{\|r_0\|_2} \right]^{1/k} \leq \min\{(|\xi|\rho)^{-1}, 1\}. \tag{2.17}$$

*If* $\rho = 1$ *(which happens when and only when* $\tau \in [-1,1]$*), then*

$$\min\{|\xi|^{-1}, 1\} \times \eta \leq \liminf_{k \to \infty} \left[ \frac{\|r_k\|_2}{\|r_0\|_2} \right]^{1/k} \leq \limsup_{k \to \infty} \left[ \frac{\|r_k\|_2}{\|r_0\|_2} \right]^{1/k} \leq \eta, \tag{2.18}$$

*where* $\eta = \limsup_{k \to \infty} \left[ 1/4 + \sum_{j=1}^{k} |\xi|^{2j} (\cos j\theta)^2 \right]^{-1/(2k)}$ *and* $\theta = \arccos \tau$. *Regardless of* $\rho > 1$ *or* $\rho = 1$,

$$\lim_{k \to \infty} \left[ \frac{\|r_k\|_2}{\|r_0\|_2} \right]^{1/k} = \min\{(|\xi|\rho)^{-1}, 1\} \quad \text{for } |\xi| \leq 1. \tag{2.19}$$

2. *Let* $b = e_N$. *If* $\rho > 1$, *then*

$$\min\{(|\xi|^{-2}\rho)^{-1}, (|\xi|^{-1}\rho)^{-1}, 1\} \leq \liminf_{k \to \infty} \left[ \frac{\|r_k\|_2}{\|r_0\|_2} \right]^{1/k} \leq \limsup_{k \to \infty} \left[ \frac{\|r_k\|_2}{\|r_0\|_2} \right]^{1/k} \leq \min\{(|\xi|^{-1}\rho)^{-1}, 1\}. \tag{2.20}$$

*If* $\rho = 1$ *(which happens when and only when* $\tau \in [-1,1]$*), then*

$$\min\{|\xi|, 1\} \times \eta \leq \liminf_{k \to \infty} \left[ \frac{\|r_k\|_2}{\|r_0\|_2} \right]^{1/k} \leq \limsup_{k \to \infty} \left[ \frac{\|r_k\|_2}{\|r_0\|_2} \right]^{1/k} \leq \eta, \tag{2.21}$$

*where* $\eta = \limsup_{k \to \infty} \left[ 1/4 + \sum_{j=1}^{k} |\xi|^{-2j} (\cos j\theta)^2 \right]^{-1/(2k)}$ *and* $\theta = \arccos \tau$. *Regardless of* $\rho > 1$ *or* $\rho = 1$,

$$\lim_{k \to \infty} \left[ \frac{\|r_k\|_2}{\|r_0\|_2} \right]^{1/k} = \min\{(|\xi|^{-1}\rho)^{-1}, 1\} \quad \text{for } |\xi| \geq 1. \tag{2.22}$$

REMARK **2.1** As we commented before, our numerical examples indicate that the upper bounds in Theorems 2.3 and 2.4 are rather accurate regardless of the magnitude of $|\xi|$ for both cases $b = e_1$ and $b = e_N$ (see Figure 2.2) and the lower bound in Theorem 2.5 is also accurate regardless of whether $k \leq N/2$ or not (see Figure 2.3). This leads us to conjecture that the following equations would hold.

$$\lim_{k\to\infty} \|r_k\|_2^{1/k} = \min\{(|\xi|\rho)^{-1}, 1\} \qquad \text{for } b = e_1, \tag{2.23}$$

$$\lim_{k\to\infty} \|r_k\|_2^{1/k} = \min\{(|\xi|^{-1}\rho)^{-1}, 1\} \quad \text{for } b = e_N, \tag{2.24}$$

where no constraint is assumed between $k$ and $N$, except $k < N$ as usual.

## 3 Proofs

We starting by reformulating the computation of GMRES residuals into an optimization problem involving rectangular Vandermonde matrices when $A$ is diagonalizable but otherwise general, i.e., not necessarily tridiagonal Toeplitz.

Recall we assumed, without loss of generality, the initial approximation $x_0 = 0$ and thus the initial residual $r_0 = b - Ax_0 = b$.

Let $N$-by-$N$ matrix $A$ have eigendecomposition

$$A = X\Lambda X^{-1}, \quad \Lambda = \mathrm{diag}(\lambda_1, \lambda_2, \ldots, \lambda_N), \tag{3.1}$$

and let $V_{k+1,N}$ be the $(k+1) \times N$ rectangular Vandermonde matrix

$$V_{k+1,N} \overset{\text{def}}{=} \begin{pmatrix} 1 & 1 & \cdots & 1 \\ \lambda_1 & \lambda_2 & \cdots & \lambda_N \\ \vdots & \vdots & \ddots & \vdots \\ \lambda_1^k & \lambda_2^k & \cdots & \lambda_N^k \end{pmatrix} \tag{3.2}$$

having nodes $\{\lambda_j\}_{j=1}^N$, and

$$Y = X\mathrm{diag}(X^{-1}b). \tag{3.3}$$

Using $(b, Ab, \ldots, A^{k-1}b) = Y V_{k+1,N}^T$ [24, Lemma 2.1], we have for GMRES

$$\|r_k\|_2 = \min_{u_{(1)}=1} \|(b, Ab, \ldots, A^{k-1}b)u\|_2 = \min_{u_{(1)}=1} \|Y V_{k+1,N}^T u\|_2. \tag{3.4}$$

It can be seen that for $Y$ as in (3.3)

$$\begin{aligned}
\|Y\|_2 &\leq \|X\|_2 \max_i |(X^{-1}b)_{(i)}| \\
&\leq \|X\|_2\|X^{-1}b\|_2 \\
&\leq \|X\|_2\|X^{-1}\|_2\|b\|_2, \\
\min_{u_{(1)}=1} \|Y V_{k+1,N}^T u\|_2 &\leq \|Y\|_2 \min_{u_{(1)}=1} \|V_{k+1,N}^T u\|_2 \\
&\leq \kappa(X) \min_{\phi_k(0)=1} \max_i |\phi_k(\lambda_i)| \|b\|_2,
\end{aligned} \tag{3.5}$$

11

where $\phi_k$ is a polynomial of degree no higher than $k$. Thus, together with (3.4), they imply

$$\|r_k\|_2/\|r_0\|_2 \le \kappa(X) \min_{\phi_k(0)=1} \max_i |\phi_k(\lambda_i)|. \tag{3.6}$$

Inequality (3.6) is often the starting point in existing quantitative analysis on GMRES convergence [11, Page 54], as it seems that there is no easy way to do otherwise. It simplifies the analysis by separating the study of GMRES convergence behavior into optimizing the condition number of $X$ and a polynomial minimization problem over $A$'s spectrum, but it could potentially overestimate GMRES residuals. This is partly because, as observed by Liesen and Strakoš [17], possible cancelations of huge components in $X$ and/or $X^{-1}$ were artificially ignored for the sake of the convergence analysis. For tridiagonal Toeplitz matrix $A$ we are interested in here, however, rich structure allows us to do differently, namely starting with (3.4) directly.

Switch to tridiagonal Toeplitz matrix $A$ as in (2.1) which is diagonalizable when $\mu \ne 0$ and $\nu \ne 0$. In fact [21, pp.113-115] (see also [9, 17]),

$$A = X\Lambda X^{-1}, \quad X = \Xi Z, \quad \Lambda = \mathrm{diag}(\lambda_1,\dots,\lambda_N), \tag{3.7}$$

$$\lambda_j = \lambda - 2\sqrt{\mu\nu}\, t_j, \quad t_j = \cos\theta_j, \quad \theta_j = \frac{j\pi}{N+1}, \tag{3.8}$$

$$Z_{(:,j)} = \sqrt{\frac{2}{N+1}}\, (\sin j\theta_1,\dots,\sin j\theta_N)^T, \tag{3.9}$$

$$\Xi = \mathrm{diag}(1, \xi^{-1},\dots,\xi^{-N+1}). \tag{3.10}$$

It can be verified that $Z^T Z = I_N$; So $A$ is normal if $|\xi| = 1$, i.e., $|\mu| = |\nu| > 0$. Set $\omega = -2\sqrt{\mu\nu}$. By (3.8), we have

$$\lambda_j = \omega(t_j - \tau), \quad 1 \le j \le N. \tag{3.11}$$

We define the $m$th *Translated Chebyshev Polynomial* in $z$ of degree $m$ as

$$\begin{aligned} T_m(z;\omega,\tau) &\stackrel{\text{def}}{=} T_m(z/\omega + \tau) && \text{(3.12)}\\ &= a_{mm}z^m + a_{m-1\,m}z^{m-1} + \cdots + a_{1m}z + a_{0m}, && \text{(3.13)} \end{aligned}$$

where $a_{jm} \equiv a_{jm}(\omega,\tau)$ are functions of $\omega$ and $\tau$, and upper triangular $R_m \in \mathbb{C}^{m\times m}$, a matrix-valued function in $\omega$ and $\tau$, too, as

$$R_m \equiv R_m(\omega,\tau) \stackrel{\text{def}}{=} \begin{pmatrix} a_{00} & a_{01} & a_{02} & \cdots & a_{0\,m-1} \\ & a_{11} & a_{12} & \cdots & a_{1\,m-1} \\ & & a_{22} & \cdots & a_{2\,m-1} \\ & & & \ddots & \vdots \\ & & & & a_{m-1\,m-1} \end{pmatrix}, \tag{3.14}$$

i.e., the $j$th column consists of the coefficients of $T_{j-1}(z;\omega,\tau)$. Set

$$\boldsymbol{T}_N \stackrel{\text{def}}{=} \begin{pmatrix} T_0(t_1) & T_0(t_2) & \cdots & T_0(t_N) \\ T_1(t_1) & T_1(t_2) & \cdots & T_1(t_N) \\ \vdots & \vdots & & \vdots \\ T_{N-1}(t_1) & T_{N-1}(t_2) & \cdots & T_{N-1}(t_N) \end{pmatrix} \tag{3.15}$$

12

and $V_N = V_{N,N}$ for short. Then

$$V_N^T R_N = \boldsymbol{T}_N^T. \tag{3.16}$$

Equation (3.16) yields $V_N^T = \boldsymbol{T}_N^T R_N^{-1}$. Extracting the first $k+1$ columns from both sides of $V_N^T = \boldsymbol{T}_N^T R_N^{-1}$ yields

$$V_{k+1,N}^T = \boldsymbol{T}_{k+1,N}^T R_{k+1}^{-1}, \tag{3.17}$$

where $\boldsymbol{T}_{k+1,N} = (\boldsymbol{T}_N)_{(1:k+1,:)}$.

In what follows, we will prove theorems in the previous section in the order of their appearance, except Theorem 2.2 whose proof requires Theorem 2.6 will be proved last.

We need to estimate GMRES residual

$$\|r_k\|_2 = \min_{u_{(1)}=1} \|Y V_{k+1,N}^T u\|_2$$

for $Ax = b$ here, where $Y = X\mathrm{diag}(X^{-1}b)$. Now notice $Y = X\,\mathrm{diag}(X^{-1}b)$ and $X = \Xi Z$ with $Z$ in (3.9) being real and orthogonal to get

$$
\begin{aligned}
Y V_{k+1,N}^T &= \Xi Z\,\mathrm{diag}(Z^T\Xi^{-1}b)\,(\boldsymbol{T}_N^T)_{(:,1:k+1)}R_{k+1}^{-1} \\
&= \Xi M_{(:,1:k+1)}R_{k+1}^{-1} \tag{3.18} \\
&= \Xi M_{(:,1:k+1)}\Xi_{k+1}^{-1}\Xi_{k+1}R_{k+1}^{-1}. \tag{3.19}
\end{aligned}
$$

where $\Xi_{k+1} = \Xi_{(1:k+1,1:k+1)}$, the $(k+1)$th leading submatrix of $\Xi$,

$$M = Z\,\mathrm{diag}(Z^T\Xi^{-1}b)\,\boldsymbol{T}_N^T. \tag{3.20}$$

It follows from (3.4) and (3.19) that

$$\sigma_{\min}(\Xi M_{(:,1:k+1)}\Xi_{k+1}^{-1}) \le \frac{\|r_k\|_2}{\min_{u_{(1)}=1}\|\Xi_{k+1}R_{k+1}^{-1}u\|_2} \le \|\Xi M_{(:,1:k+1)}\Xi_{k+1}^{-1}\|_2. \tag{3.21}$$

The second inequality in (3.21) is our foundation to prove Theorem 2.1. There are two quantities to deal with

$$\min_{u_{(1)}=1} \|\Xi_{k+1}R_{k+1}^{-1}u\|_2 \quad \text{and} \quad \|\Xi M_{(:,1:k+1)}\Xi_{k+1}^{-1}\|_2. \tag{3.22}$$

We shall now do so. In its present general form, the next lemma was proven in [14, 15]. It was also implied by the proof of [13, Theorem 2.1]. See also [16].

**Lemma 3.1** *If $W$ has full column rank, then*

$$\min_{u_{(1)}=1} \|Wu\|_2 = \left[e_1^T(W^*W)^{-1}e_1\right]^{-1/2}. \tag{3.23}$$

*In particular if $W$ is nonsingular, $\min_{u_{(1)}=1}\|Wu\|_2 = \|W^{-*}e_1\|_2^{-1}$.*

By this lemma, we have (note $a_{00} = 1$)

$$\min_{u_{(1)}=1} \|\Xi_{k+1}R_{k+1}^{-1}u\|_2 = \|\Xi_{k+1}^{-*}R_{k+1}^*e_1\|_2^{-1} = \left[\frac{1}{2} + \Phi_{k+1}^{(+)}(\tau,\xi)\right]^{-1/2}. \tag{3.24}$$

This gives the first quantity in (3.22). We now turn to the second one there. It can be seen that $\Xi M_{(:,1:k+1)}\Xi_{k+1}^{-1} = (\Xi M\Xi^{-1})_{(:,1:k+1)}$ since $\Xi$ is diagonal. To compute $\Xi M\Xi^{-1}$, we shall investigate $M$ in (3.20) first.

$$
\begin{aligned}
M &= \sum_{\ell=1}^{N} Z \operatorname{diag}(Z\Xi^{-1}b_{(\ell)}e_\ell)\,\boldsymbol{T}_N^T \\
&= \sum_{\ell=1}^{N} b_{(\ell)}\xi^{\ell-1} Z \operatorname{diag}(Ze_\ell)\,\boldsymbol{T}_N^T \\
&= \sum_{\ell=1}^{N} b_{(\ell)}\xi^{\ell-1} Z \operatorname{diag}(Z_{(:,\ell)})\,\boldsymbol{T}_N^T.
\end{aligned}
\tag{3.25}
$$

In Lemma 3.2 and in the proof of Lemma 3.3 below, without causing notational conflict, we will temporarily use $k$ as a running index, as opposed to the rest of the paper where $k$ is reserved for GMRES step index.

**Lemma 3.2** *For $\theta_j = \frac{j}{N+1}\pi$ and integer $\ell$,*

$$
\sum_{k=1}^{N} \cos \ell\theta_k = \begin{cases} N, & \text{if } \ell = 2m(N+1) \text{ for some integer } m, \\ 0, & \text{if } \ell \text{ is odd,} \\ -1, & \text{if } \ell \text{ is even, but } \ell \neq 2m(N+1) \text{ for any integer } m. \end{cases}
\tag{3.26}
$$

*Proof:* If $\ell = 2m(N+1)$ for some integer $m$, then $\ell\theta_k = 2mk\pi$ and thus $\cos \ell\theta_k = 1$. Assume that $\ell \neq 2m(N+1)$ for any integer $m$. Set $\phi = \ell\pi/(N+1)$. We have [10, p.30]

$$
\sum_{k=1}^{N} \cos \ell\theta_k = \sum_{k=1}^{N} \cos k\phi = \cos \frac{N+1}{2}\phi \times \frac{\sin \frac{N\phi}{2}}{\sin \frac{\phi}{2}}.
$$

Now notice $\cos \frac{N+1}{2}\phi = \cos \frac{\ell}{2}\pi = 0$ for odd $\ell$ and $(-1)^{\ell/2}$ for even $\ell$, and $\sin \frac{N\phi}{2} = \sin(\frac{\ell}{2}\pi - \phi) = -(-1)^{\ell/2} \sin \phi$ for even $\ell$ to conclude the proof. ∎

**Lemma 3.3** *Let $M_\ell \overset{\text{def}}{=} Z \operatorname{diag}(Z_{(:,\ell)})\boldsymbol{T}_N^T$ for $1 \leq \ell \leq N$. Then the entries of $M_\ell$ are zeros, except at those positions $(i,j)$, graphically forming four straight lines:*

$$
\begin{aligned}
&\text{(a)} \quad i+j = \ell+1, \\
&\text{(b)} \quad i-j = \ell-1, \\
&\text{(c)} \quad j-i = \ell+1, \\
&\text{(d)} \quad i+j = 2(N+1) - \ell+1.
\end{aligned}
\tag{3.27}
$$

*$(M_\ell)_{(i,j)} = 1/2$ for (a) and (b), except at their intersection $(\ell,1)$ for which $(M_\ell)_{(\ell,1)} = 1$. $(M_\ell)_{(i,j)} = -1/2$ for (c) and (d). Notice no valid entries for (c) if $\ell \geq N-1$ and no valid entries for (d) if $\ell \leq 2$.*

*Proof:* For $1 \leq i, j \leq N$,

$$
2(N+1)\cdot(M_\ell)_{(i,j)} = 4\sum_{k=1}^{N} \sin k\theta_i \, \sin \ell\theta_k \, \cos(j-1)\theta_k
$$

$$= 4\sum_{k=1}^{N} \sin i\theta_k \, \sin \ell\theta_k \, \cos(j-1)\theta_k$$

$$= 2\sum_{k=1}^{N} [\cos(i-\ell)\theta_k - \cos(i+\ell)\theta_k] \cos(j-1)\theta_k$$

$$= \sum_{k=1}^{N} [\cos(i+j-\ell-1)\theta_k + \cos(i-j-\ell+1)\theta_k$$
$$- \cos(i+j+\ell-1)\theta_k - \cos(i-j+\ell+1)\theta_k].$$

Since all

$$
\begin{aligned}
i_1 &= i+j-\ell-1, \\
i_2 &= i-j-\ell+1, \\
i_3 &= i+j+\ell-1, \\
i_4 &= i-j+\ell+1
\end{aligned}
$$

are either even or odd at the same time, Lemma 3.2 implies $(M_\ell)_{(i,j)} = 0$ unless one of them takes the form $2m(N+1)$ for some integer $m$. We now investigate all possible situations as such, keeping in mind that $1 \le i, j, \ell \le N$.

1. $i_1 = i+j-\ell-1 = 2m(N+1)$. This happens if and only if $m = 0$, and thus $i+j = \ell+1$. Then

$$i_2 = -2j+2, \; i_3 = 2\ell, \; i_4 = -2j+2\ell+2.$$

They are all even. $i_3$ and $i_4$ do not take the form $2m(N+1)$ for some integers $m$. This is obvious for $i_3$, while $i_4 = 2m(N+1)$ implies $m = 0$ and $j = \ell+1$, and thus $i = 0$ which cannot happen. However if $i_2 = 2m(N+1)$, then $m = 0$ and $j = 1$, and thus $i = \ell$.

So Lemma 3.2 implies $(M_\ell)_{(i,j)} = 1/2$ for $i+j = \ell+1$ and $i \ne \ell$, while $(M_\ell)_{(\ell,1)} = 1$.

2. $i_2 = i-j-\ell+1 = 2m(N+1)$. This happens if and only if $m = 0$, and thus $i-j = \ell-1$. Then

$$i_1 = 2j-2, \; i_3 = 2j+2\ell-2, \; i_4 = 2\ell.$$

They are all even. $i_3$ and $i_4$ do not take the form $2m(N+1)$ for some integers $m$. This is obvious for $i_4$, while $i_3 = 2m(N+1)$ implies $m = 1$ and $j = N+2-\ell$, and thus $i = N+1$ which cannot happen. However if $i_1 = 2m(N+1)$, then $m = 0$ and thus $j = 1$ and $i = \ell$ which has already been considered in Item 1.

So Lemma 3.2 implies $(M_\ell)_{(i,j)} = 1/2$ for $i-j = \ell-1$ and $i \ne \ell$, while $(M_\ell)_{(\ell,1)} = 1$.

3. $i_3 = i+j+\ell-1 = 2m(N+1)$. This happens if and only if $m = 1$, and thus $i+j = 2(N+1) - \ell + 1$. Then

$$i_1 = 2(N+1) - 2\ell, \; i_2 = 2(N+1) - 2j - 2\ell + 2, \; i_4 = 2(N+1) - 2j + 2.$$

They are all even. $i_1$ and $i_2$ do not take the form $2m(N+1)$ for some integers $m$. This is obvious for $i_1$, while $i_2 = 2m(N+1)$ implies $m = 0$ and $j = N+2-\ell$, and thus

15

$i = N + 1$ which cannot happen. However if $i_4 = 2m(N + 1)$, then $m = 1$ and thus $j = 1$ and $i = 2(N + 1) - \ell$ which is bigger than $N + 2$ and not possible.

So Lemma 3.2 implies $(M_\ell)_{(i,j)} = -1/2$ for $i + j = 2(N + 1) - \ell + 1$.

4. $i_4 = i - j + \ell + 1 = 2m(N+1)$. This happens if and only if $m = 0$, and thus $j - i = \ell + 1$. Then
$$i_1 = 2j - 2\ell - 2, \quad i_2 = -2\ell, \quad i_3 = 2j - 2.$$

They are all even, and do not take the form $2m(N + 1)$ for some integers $m$. This is obvious for $i_2$. $i_1 = 2m(N + 1)$ implies $m = 0$ and $j = \ell + 1$, and thus $i = 0$ which cannot happen. $i_3 = 2m(N+1)$ implies $m = 0$ and thus $j = 1$ and $i = -\ell$ which cannot happen either.

So Lemma 3.2 implies $(M_\ell)_{(i,j)} = -1/2$ for $j - i = \ell + 1$.

This completes the proof. ∎

Now we know $M_\ell$. We still need to find out $\Xi M \Xi^{-1}$. Let us examine it for $N = 5$ in order to get some idea about what it may look like. $\Xi M \Xi^{-1}$ for $N = 5$ is

$$
\begin{pmatrix}
b_{(1)} & \frac{1}{2}\xi^2 b_{(2)} & -\frac{1}{2}\xi^2 b_{(1)} + \frac{1}{2}\xi^4 b_{(3)} & -\frac{1}{2}\xi^4 b_{(2)} + \frac{1}{2}\xi^6 b_{(4)} & -\frac{1}{2}\xi^6 b_{(3)} + \frac{1}{2}\xi^8 b_{(5)} \\
b_{(2)} & \frac{1}{2}b_{(1)} + \frac{1}{2}b_{(3)}\xi^2 & \frac{1}{2}b_{(4)}\xi^4 & -\frac{1}{2}\xi^2 b_{(1)} + \frac{1}{2}\xi^6 b_{(5)} & -\frac{1}{2}\xi^4 b_{(2)} \\
b_{(3)} & \frac{1}{2}b_{(2)} + \frac{1}{2}\xi^2 b_{(4)} & \frac{1}{2}b_{(1)} + \frac{1}{2}b_{(5)}\xi^4 & 0 & -\frac{1}{2}\xi^2 b_{(1)} - \frac{1}{2}\xi^6 b_{(5)} \\
b_{(4)} & \frac{1}{2}b_{(3)} + \frac{1}{2}b_{(5)}\xi^2 & \frac{1}{2}b_{(2)} & \frac{1}{2}b_{(1)} - \frac{1}{2}b_{(5)}\xi^4 & -\frac{1}{2}b_{(4)}\xi^4 \\
b_{(5)} & \frac{1}{2}b_{(4)} & \frac{1}{2}b_{(3)} - \frac{1}{2}b_{(5)}\xi^2 & \frac{1}{2}b_{(2)} - \frac{1}{2}\xi^2 b_{(4)} & \frac{1}{2}b_{(1)} - \frac{1}{2}b_{(3)}\xi^2
\end{pmatrix}.
$$

We observe that for $N = 5$, the entries of $\Xi M \Xi^{-1}$ are polynomials in $\xi$ with at most two terms. This turns out to be true for all $N$.

**Lemma 3.4** *The following statements hold.*

1. *The first column of $\Xi M \Xi^{-1}$ is $b$. Entries in every other columns taking one of the three forms: $(b_{(n_1)}\xi^{m_1} + b_{(n_2)}\xi^{m_2})/2$ with $n_1 \neq n_2$, $b_{(n_1)}\xi^{m_1}/2$, and $0$, where $1 \leq n_1, n_2 \leq N$ and $m_i \geq 0$ are nonnegative integer.*

2. *In each given column of $\Xi M \Xi^{-1}$, any particular entry of $b$ appears at most twice.*

*As the consequence, we have $\|\Xi M_{(:,1:k+1)} \Xi_{k+1}^{-1}\|_2 \leq \sqrt{k+1}\|b\|_2$ if $|\xi| \leq 1$.*

*Proof:* Notice $M = \sum_{\ell=1}^{N} b_{(\ell)} \xi^{\ell-1} M_\ell$ and consider $M$'s $(i, j)$th entry which comes from the contributions from all $M_\ell$. But not all of $M_\ell$ contribute as most of them are zero at the position. Precisely, with the help of Lemma 3.3, those $M_\ell$ that contribute nontrivially to the $(i, j)$th position are the following ones subject to satisfying the given inequalities.

(a) if $1 \leq i + j - 1 \leq N$ or equivalently $i + j \leq N + 1$, $M_{i+j-1}$ gives a $1/2$.

(b) if $1 \leq i - j + 1 \leq N$ or equivalently $i \geq j$, $M_{i-j+1}$ gives a $1/2$.

(c) if $1 \leq j - i - 1 \leq N$ or equivalently $j \geq i + 2$, $M_{j-i-1}$ gives a $-1/2$.
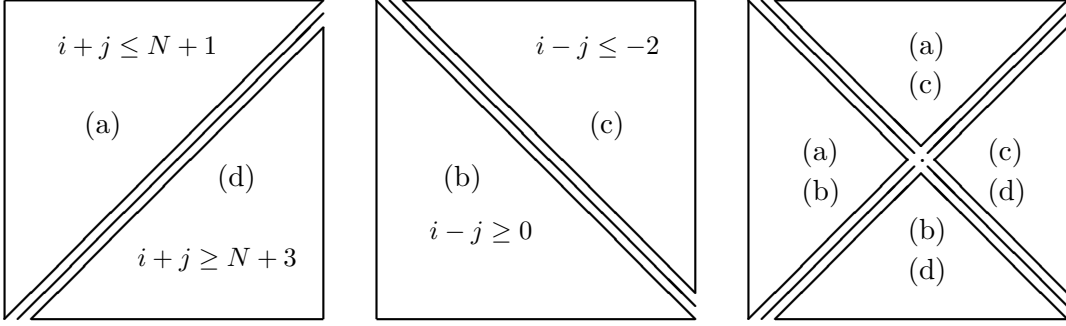
**Figure 3.1:** Computation of $M_{(i,j)}$. *Left:* Regions of entries as divided by inequalities in (a) and (d); *Middle:* Regions of entries as divided by inequalities in (b) and (c); *Right:* Regions of entries as divided by all inequalities in (a), (b), (c), and (d).

(d) if $1 \leq 2(N+1) - (i+j) + 1 \leq N$ or equivalently $i + j \geq N + 3$, $M_{2(N+1)-(i+j)+1}$ gives a $-1/2$.

These inequalities, effectively 4 of them, divided entries of $M$ into nine possible regions as detailed in Figure 3.1. We shall examine each region one by one. Recall

$$(\Xi M \Xi^{-1})_{(i,j)} = \xi^{-i+1} M_{(i,j)} \xi^{j-1} = \xi^{j-i} M_{(i,j)},$$

and let

$$
\begin{aligned}
\gamma_a &= \frac{1}{2} b_{(i+j-1)} \xi^{2j-2}, \\
\gamma_b &= \frac{1}{2} b_{(i-j+1)}, \\
\gamma_c &= -\frac{1}{2} b_{(j-i-1)} \xi^{2(j-i-1)}, \\
\gamma_d &= -\frac{1}{2} b_{(2(N+1)-(i+j)+1)} \xi^{2(N+1)-(i+j)}.
\end{aligned}
$$

Each entry in the 9 possible regions in the rightmost plot of Figure 3.1 is as follows.

1. (a) and (b): $(\Xi M \Xi^{-1})_{(i,j)} = \gamma_a + \gamma_b$.

2. (a) and (c): $(\Xi M \Xi^{-1})_{(i,j)} = \gamma_a + \gamma_c$.

3. (b) and (d): $(\Xi M \Xi^{-1})_{(i,j)} = \gamma_b + \gamma_d$.

4. (c) and (d): $(\Xi M \Xi^{-1})_{(i,j)} = \gamma_c + \gamma_d$.

5. (a) and $i - j = -1$: $(\Xi M \Xi^{-1})_{(i,j)} = \gamma_a$.

6. (b) and $i + j = N + 2$: $(\Xi M \Xi^{-1})_{(i,j)} = \gamma_b$.

7. (c) and $i + j = N + 2$: $(\Xi M \Xi^{-1})_{(i,j)} = \gamma_c$.

8. (d) and $i - j = -1$: $(\Xi M \Xi^{-1})_{(i,j)} = \gamma_d$.

9. $i - j = -1$ and $i + j = N + 2$: $(\Xi M \Xi^{-1})_{(i,j)} = 0$. In this case, $i = (N + 1)/2$ and $j = (N + 3)/2$. So there is only one such entry when $N$ is odd, and none when $N$ is even.

With this profile on the entries of $\Xi M \Xi^{-1}$, we have Item 1 of the lemma immediately. Item 2 is the consequence of $M = \sum_{\ell=1}^{N} b_{(\ell)} \xi^{\ell-1} M_\ell$ and Lemma 3.3 which implies that there are at most two nonzero entries in each column of $M_\ell$.

As the consequence of Item 1 and Item 2, each column of $\Xi M \Xi^{-1}$ can be expressed as the sum of two vectors $w$ and $v$ such that $\|w\|_2, \|v\|_2 \leq \|b\|_2/2$ when $|\xi| \leq 1$, and thus $\|(\Xi M \Xi^{-1})_{(:,j)}\|_2 \leq \|b\|_2$ for all $1 \leq j \leq N$. Therefore

$$\|\Xi M_{(:,1:k+1)} \Xi_{k+1}^{-1}\|_2 \leq \sqrt{\sum_{j=1}^{k+1} \|(\Xi M \Xi^{-1})_{(:,j)}\|_2^2} \leq \sqrt{k+1}\|b\|_2,$$

as expected. ∎

*Proof* of **Theorem 2.1**. We shall only prove

$$\|r_k\|_2 \leq \|b\|_2 \sqrt{k+1} \left[ \frac{1}{2} + \Phi_{k+1}^{(+)}(\tau, \xi) \right]^{-1/2} \qquad \text{for } |\xi| \leq 1 \tag{3.28}$$

since the other case when $|\xi| > 1$ can be turned into this case as follows. Let $\Pi = (e_N, \ldots, e_2, e_1) \in \mathbb{R}^{N \times N}$ be the permutation matrix. Notice $\Pi^T A \Pi = A^T$ and thus $Ax = b$ is equivalent to

$$A^T \Pi^T x = (\Pi^T A \Pi)(\Pi^T x) = \Pi^T b. \tag{3.29}$$

Note $\mathcal{K}_k(A^T, \Pi^T b) = \mathcal{K}_k(\Pi^T A \Pi, \Pi^T b) = \Pi^T \mathcal{K}_k(A, b)$, and

$$\|r_k\|_2 = \min_{y \in \mathcal{K}_k(A,b)} \|b - Ay\|_2 = \min_{\Pi^T y \in \Pi^T \mathcal{K}_k(A,b)} \|\Pi^T(b - A\Pi\,\Pi^T y)\|_2$$

$$= \min_{w \in \mathcal{K}_k(A^T, \Pi^T b)} \|\Pi^T b - A^T w\|_2.$$

If (3.28) is proven true, then for $|\xi| > 1$ we have

$$\|r_k\|_2 \leq \|\Pi^T b\|_2 \sqrt{k+1} \left[ \frac{1}{2} + \Phi_{k+1}^{(-)}(\tau, \xi) \right]^{-1/2}$$

because the $\xi$ for $A^T$ is the reciprocal of the one for $A$.

Assume $|\xi| \leq 1$. Inequality (3.28) is the consequence of (3.21), (3.24), and Lemma 3.4. ∎

REMARK **3.1** The leftmost inequality in (3.21) gives a lower bound on $\|r_k\|_2$ in terms of $\sigma_{\min}(\Xi M_{(:,1:k+1)} \Xi_{k+1}^{-1})$ which, however, is hard to bound from below because it can be as small as zero, unless we know more about $b$ such as $b = e_1$ or $e_N$ as in Theorems 2.3 and 2.4.

*Proof* of **Theorem 2.3**: If $b = e_1$, then $M = M_1$ is upper triangular. More specifically

$$M = M_1 = \begin{pmatrix} 1 & 0 & -1/2 & & \\ & 1/2 & 0 & \ddots & \\ & & 1/2 & & -1/2 \\ & & & \ddots & 0 \\ & & & & 1/2 \end{pmatrix} \tag{3.30}$$

18

and, by (3.18),

$$YV_{k+1,N} = \begin{pmatrix} \Xi_{k+1}\widetilde{M}R_{k+1}^{-1} \\ 0 \end{pmatrix} = \begin{pmatrix} \Xi_{k+1}\widetilde{M}\Xi_{k+1}^{-1}D^{-1} \times D\Xi_{k+1}R_{k+1}^{-1} \\ 0 \end{pmatrix},$$

where $\widetilde{M} = M_{(1:k+1,1:k+1)}$ and $D = \mathrm{diag}(2,1,1,\ldots,1)$. Therefore

$$\sigma_{\min}(\Xi_{k+1}\widetilde{M}\Xi_{k+1}^{-1}D^{-1}) \le \frac{\min_{u_{(1)}=1}\|Y V_{k+1,N}^T u\|_2}{\min_{u_{(1)}=1}\|D\Xi_{k+1}R_{k+1}^{-1}u\|_2} \le \|\Xi_{k+1}\widetilde{M}\Xi_{k+1}^{-1}D^{-1}\|_2. \qquad (3.31)$$

Recall $\|r_k\|_2 = \min_{u_{(1)}=1}\|Y V_{k+1,N}^T u\|_2$. Let $P_{k+1} = (e_1, e_3, \ldots, e_2, e_4, \ldots) \in \mathbb{R}^{(k+1)\times(k+1)}$. It can be seen that

$$P_{k+1}^T(\Xi_{k+1}\widetilde{M}\Xi_{k+1}^{-1}D^{-1})P_{k+1} = \frac{1}{2}\begin{pmatrix} E_1 & \\ & E_2 \end{pmatrix},$$

where $E_1 \in \mathbb{R}^{\lceil\frac{k+1}{2}\rceil\times\lceil\frac{k+1}{2}\rceil}$, $E_2 \in \mathbb{R}^{\lfloor\frac{k+1}{2}\rfloor\times\lfloor\frac{k+1}{2}\rfloor}$, and

$$E_i = \begin{pmatrix} 1 & -\xi^2 & & \\ & 1 & \ddots & \\ & & \ddots & -\xi^2 \\ & & & 1 \end{pmatrix}, \quad E_i^{-1} = \begin{pmatrix} 1 & \xi^2 & \cdots & \xi^{2(m-1)} \\ & 1 & \ddots & \vdots \\ & & \ddots & \xi^2 \\ & & & 1 \end{pmatrix}.$$

Hence $\|E_i\|_2 \le \sqrt{\|E_i\|_1\|E_i\|_\infty} = 1 + |\xi|^2$. Therefore

$$\|\Xi_{k+1}\widetilde{M}\Xi_{k+1}^{-1}D^{-1}\|_2 = \frac{1}{2}\max\{\|E_1\|_2, \|E_2\|_2\} \le \frac{1}{2}(1 + |\xi|^2).$$

Similarly use $\|E_i^{-1}\|_2 \le \sqrt{\|E_i^{-1}\|_1\|E_i^{-1}\|_\infty}$ to get

$$\|E_1^{-1}\|_2 \le \sum_{j=0}^{\lceil\frac{k+1}{2}\rceil-1}|\xi|^{2j}, \quad \|E_2^{-1}\|_2 \le \sum_{j=0}^{\lfloor\frac{k+1}{2}\rfloor-1}|\xi|^{2j}.$$

Therefore

$$\begin{aligned}
\sigma_{\min}(\Xi_{k+1}\widetilde{M}\Xi_{k+1}^{-1}D^{-1}) &= \frac{1}{2}\min\{\sigma_{\min}(E_1), \sigma_{\min}(E_2)\} \\
&= \frac{1}{2}\min\{\|E_1^{-1}\|_2^{-1}, \|E_2^{-1}\|_2^{-1}\} \\
&\ge \frac{1}{2}\left[\sum_{j=0}^{\lceil\frac{k+1}{2}\rceil-1}|\xi|^{2j}\right]^{-1}.
\end{aligned}$$

Finally, by Lemma 3.1, we have

$$\min_{u_{(1)}=1}\|D\Xi_{k+1}R_{k+1}^{-1}u\|_2 = \|D^{-*}\Xi_{k+1}^{-*}R_{k+1}^*e_1\|_2^{-1} = \left[\Phi_{k+1}^{(+)}(\tau,\xi) - \frac{1}{4}\right]^{-1/2}.$$

19

This, together with (3.31), lead to (2.11). ∎

As in the proof of Theorem 2.1, by applying Theorem 2.3 to the permuted system (3.29), we get Theorem 2.4 for $b = e_N$.

*Proof* of **Theorem 2.5**: Now $b = b_{(1)}e_1 + b_{(N)}e_N$. Notice the form of $M_1$ in (3.30), and that $M_N$ is $M_1$ after its rows reordered from the last to the first. For the case $M = b_{(1)}M_1 + \xi^{N-1}b_{(N)}M_N$, and also Lemma 3.4 implies that only positive powers of $\xi$ appear in the entries of $\Xi M \Xi^{-1}$. Therefore when $|\xi| \le 1$,

$$
\begin{aligned}
\|\Xi M_{(:,1:k+1)}\Xi_{k+1}^{-1}\|_2 &\le \|\Xi M \Xi^{-1}\|_2 \\
&\le |b_{(1)}| \, \| \, |M_1| \, \|_2 + |b_{(N)}| \, \| \, |M_N| \, \|_2 \\
&\le |b_{(1)}| \sqrt{3/2} + |b_{(N)}| \sqrt{3/2} \\
&\le \sqrt{3}\|b\|_2,
\end{aligned}
\tag{3.32}
$$

where $|M_\ell|$ takes entrywise absolute value, and we have used

$$
\| \, |M_N| \, \|_2 = \| \, |M_1| \, \|_2 \le \sqrt{\|M_1\|_1 \|M_1\|_\infty} = \sqrt{3/2}.
$$

Inequality (2.16) for $|\xi| \le 1$ is the consequence of (3.21), (3.24), and (3.32). Inequality (2.16) for $|\xi| \ge 1$ follows from itself for $|\xi| \le 1$ applied to the permuted system (3.29).

To prove (2.15), we use the lines of arguments in the proof for Theorem 2.3 and notice that for $1 \le k \le N/2$

$$
Y V_{k+1,N} = \begin{matrix} k \\ N-2k \\ k \end{matrix} \begin{pmatrix} W_1 \\ 0 \\ W_2 \end{pmatrix}^{\!\!k}
$$

It can be seen from the proof for Theorem 2.3 that

$$
\min_{u_{(1)}=1} \|W_1 u\|_2 \ge \frac{|b_{(1)}|}{2} \left[ \sum_{j=0}^{\lceil \frac{k+1}{2}\rceil - 1} |\xi|^{2j} \right]^{-1} \left[ \Phi_{k+1}^{(+)}(\tau, \xi) - \frac{1}{4} \right]^{-1/2},
$$

$$
\min_{u_{(1)}=1} \|W_2 u\|_2 \ge \frac{|b_{(N)}|}{2} \left[ \sum_{j=0}^{\lceil \frac{k+1}{2}\rceil - 1} |\xi|^{-2j} \right]^{-1} \left[ \Phi_{k+1}^{(-)}(\tau, \xi) - \frac{1}{4} \right]^{-1/2}.
$$

Finally use

$$
\min_{u_{(1)}=1} \|Y V_{k+1,N} u\|_2 \ge \max \left\{ \min_{u_{(1)}=1} \|W_1 u\|_2, \min_{u_{(1)}=1} \|W_2 u\|_2 \right\}
$$

to complete the proof. ∎

*Proof* of **Theorem 2.6**: We note that

$$
\limsup_{k\to\infty} \left[ \frac{\|r_k\|_2}{\|r_0\|_2} \right]^{1/k} \le 1, \quad \limsup_{k\to\infty} \left[ \sup_{r_0} \frac{\|r_k\|_2}{\|r_0\|_2} \right]^{1/k} \le 1
$$

for any $b$ because $\|r_k\|_2$ is nonincreasing.

20

Suppose $b = e_1$. Consider first $\rho > 1$. Then $|T_j(\tau)| \sim \frac{1}{2}\rho^j$, and thus

$$\Phi_{k+1}^{(+)}(\tau, \xi) - \frac{1}{4} \sim \frac{1}{4} \sum_{j=0}^{k} (|\xi|\rho)^{2j} = \frac{1}{4} \cdot \frac{(|\xi|\rho)^{2(k+1)} - 1}{(|\xi|\rho)^2 - 1}. \tag{3.33}$$

If $|\xi|\rho > 1$, then (3.33) and Theorem 2.3 imply

$$\limsup_{k\to\infty} \left[\frac{\|r_k\|_2}{\|r_0\|_2}\right]^{1/k} \leq \lim_{k\to\infty} \left[\frac{1}{2}(1 + |\xi|^2)\right]^{1/k} \left[\Phi_{k+1}^{(+)}(\tau, \xi) - \frac{1}{4}\right]^{-1/(2k)} \tag{3.34}$$

$$= (|\xi|\rho)^{-1},$$

$$\liminf_{k\to\infty} \left[\frac{\|r_k\|_2}{\|r_0\|_2}\right]^{1/k} \geq \lim_{k\to\infty} \frac{1}{2^{1/k}} \left[\sum_{j=0}^{\lceil \frac{k+1}{2}\rceil - 1} |\xi|^{2j}\right]^{-1/k} \left[\Phi_{k+1}^{(+)}(\tau, \xi) - \frac{1}{4}\right]^{-1/(2k)} \tag{3.35}$$

$$= \begin{cases} (|\xi|\rho)^{-1}, & \text{for } |\xi| \leq 1, \\ (|\xi|^2\rho)^{-1}, & \text{for } |\xi| > 1. \end{cases}$$

They together give (2.17) for the case $|\xi|\rho > 1$. If $|\xi|\rho \leq 1$, then must $|\xi| < 1$ and $\min\{(|\xi|\rho)^{-1}, 1\} = 1$, $\min\{(|\xi|^2\rho)^{-1}, (|\xi|\rho)^{-1}, 1\} = 1$, and

$$\liminf_{k\to\infty} \left[\frac{\|r_k\|_2}{\|r_0\|_2}\right]^{1/k} \geq 1$$

by (3.35) because $\Phi_{k+1}^{(+)}(\tau, \xi) - \frac{1}{4}$ is approximately bounded by $(k+1)/4$ by (3.33). So (2.17) holds for the case $|\xi|\rho \leq 1$, too. Now consider $\rho = 1$. Then $\tau + \sqrt{\tau^2 - 1} = e^{\iota\theta}$ for some $0 \leq \theta \leq \pi$, where $\iota = \sqrt{-1}$ is the imaginary unit. Thus $\tau \in [-1, 1]$ and in fact

$$2\tau = (\tau + \sqrt{\tau^2 - 1}) + (\tau - \sqrt{\tau^2 - 1}) = 2\cos\theta, \quad T_j(\tau) = \cos j\theta.$$

Therefore $\Phi_{k+1}^{(+)}(\tau, \xi) - \frac{1}{4} \sim \frac{1}{4} + \sum_{j=1}^{k} |\xi|^{2j}(\cos j\theta)^2$ which implies

$$\lim_{k\to\infty} \left[\Phi_{k+1}^{(+)}(\tau, \xi) - \frac{1}{4}\right]^{-1/(2k)} = \eta.$$

Inequalities (3.34) and (3.35) remain valid and yield (2.18). Finally regardless of $\rho > 1$ or $\rho = 1$, if $|\xi| \leq 1$, then all leftmost sides and rightmost sides in (2.17) and (2.18) are equal to $\min\{(|\xi|\rho)^{-1}, 1\}$. This proves (2.19). The proof for the case $b = e_1$ is done.

The case for $b = e_N$ can be dealt with by applied the results for $b = e_1$ to the permuted system (3.29). ∎

*Proof* of **Theorem 2.2**: Note again that $\limsup_{k\to\infty} \left(\sup_{r_0} \|r_k\|_2/\|r_0\|_2\right)^{1/k} \leq 1$.
First we prove

$$\limsup_{k\to\infty} \left[\max_{r_0 \in \{e_1, e_N\}} \frac{\|r_k\|_2}{\|r_0\|_2}\right]^{1/k} \leq \limsup_{k\to\infty} \left[\sup_{r_0} \frac{\|r_k\|_2}{\|r_0\|_2}\right]^{1/k} \leq \min\{(\zeta\rho)^{-1}, 1\}. \tag{3.36}$$

21

The first inequality is obvious because $\{e_1, e_N\} \in \{r_0\}$. We now prove the second one. If $\rho = 1$, then $\min\{(\zeta\rho)^{-1}, 1\} = 1$ because $\zeta^{-1} \geq 1$; no proof is needed. If $\rho > 1$, then $|T_j(\tau)| \sim \frac{1}{2}\rho^j$, and thus (3.33). Now if $\zeta\rho > 1$, then (3.33) and Theorem 2.1 imply $\limsup_{k\to\infty} \left(\sup_{r_0} \|r_k\|_2/\|r_0\|_2\right)^{1/k} \leq (\zeta\rho)^{-1}$ which also holds if $\zeta\rho \leq 1$ because then $(\zeta\rho)^{-1} \geq 1$.

Next we prove

$$\liminf_{k\to\infty} \left[\max_{r_0\in\{e_1, e_N\}} \frac{\|r_k\|_2}{\|r_0\|_2}\right]^{1/k} \geq \min\{(\zeta\rho)^{-1}, 1\}. \tag{3.37}$$

If $|\xi| \leq 1$, then $\zeta = |\xi|$ and thus

$$\liminf_{k\to\infty} \left[\max_{r_0\in\{e_1, e_N\}} \frac{\|r_k\|_2}{\|r_0\|_2}\right]^{1/k} \geq \liminf_{k\to\infty} \left[\max_{r_0=e_1} \frac{\|r_k\|_2}{\|r_0\|_2}\right]^{1/k} = \min\{(\xi\rho)^{-1}, 1\}$$

by (2.19) in Theorem 2.6. This is (3.37) for $|\xi| \leq 1$. For the case $|\xi| \geq 1$, we also have (3.37) similarly by (2.22). The proof is completed by combining (3.36) and (3.37). ∎

# 4 Exact Residual Norms for $b = e_1$ and $e_N$

In this section we present two theorems in which exact formulas for $\|r_k\|_2$ for $b = e_1$ and $b = e_N$ are established. Let $\Xi$ and $\Xi_{k+1}$ have their assignments as in Section 3.

**Theorem 4.1** *In Theorem 2.1, if $b = e_1$, then the kth GMRES residual $r_k$ satisfies for $1 \leq k < N$*

$$\|r_k\|_2 = \|2\Xi_{k+1}^{-*}y_{(1:k+1)}\|_2^{-1}, \tag{4.1}$$

*where $y \in \mathbb{C}^{2\lceil N/2 \rceil}$ is defined as*

$$y_{(2j-1)} = \sum_{i=1}^{j}{}' \bar{T}_{2i-2}(\tau), \quad y_{(2j)} = \sum_{i=1}^{j} \bar{T}_{2i-1}(\tau) \quad for \; j = 1, 2, \ldots, \lceil N/2 \rceil,$$

*and $\bar{T}_j(\tau)$ is the complex conjugate of $T_j(\tau)$.*

*Proof:* We still have (3.30), and $YV_{k+1,N} = \begin{pmatrix} \Xi_{k+1}\widetilde{M}R_{k+1}^{-1} \\ 0 \end{pmatrix}$, where $\widetilde{M} = M_{(1:k+1,1:k+1)}$ as in the proof of Theorem 2.3. Let $D = \text{diag}(2, 1, 1, \ldots, 1)$. Noticing $\Xi_{k+1}\widetilde{M}R_{k+1}^{-1} = \Xi_{k+1}\widetilde{M}D^{-1} \times DR_{k+1}^{-1}$ is nonsingular, we have by Lemma 3.1

$$\min_{u_{(1)}=1} \|YV_{k+1,N}^T u\|_2 = \|w\|_2^{-1},$$

where $w = \Xi_{k+1}^{-*}(\widetilde{M}D^{-1})^{-T}D^{-T}R_{k+1}^* e_1$, or equivalently $(\widetilde{M}D^{-1})^T\Xi_{k+1}^* w = D^{-T}R_{k+1}^* e_1$. We shall now solve it for $w$. Let $P_{k+1} = (e_1, e_3, \ldots, e_2, e_4, \ldots) \in \mathbb{R}^{(k+1)\times(k+1)}$. It can be verified that

$$P_{k+1}^T(\widetilde{M}D^{-1})P_{k+1} = \frac{1}{2}\begin{pmatrix} G_1 \\ & G_2 \end{pmatrix},$$

22

where $G_1 \in \mathbb{R}^{\lceil \frac{k+1}{2} \rceil \times \lceil \frac{k+1}{2} \rceil}$, $G_2 \in \mathbb{R}^{\lfloor \frac{k+1}{2} \rfloor \times \lfloor \frac{k+1}{2} \rfloor}$, and

$$
G_i = \begin{pmatrix} 1 & -1 & & \\ & 1 & \ddots & \\ & & \ddots & -1 \\ & & & 1 \end{pmatrix}, \quad G_i^{-1} = \begin{pmatrix} 1 & 1 & \cdots & 1 \\ & 1 & \ddots & \vdots \\ & & \ddots & 1 \\ & & & 1 \end{pmatrix}.
$$

Solve $(P_{k+1}^T \widetilde{M} D^{-1} P_{k+1})^T P_{k+1}^T \Xi_{k+1}^* w = P_{k+1}^T D^{-T} R_{k+1}^* e_1 \equiv P_{k+1}^T z$ for $w$ to get

$$
w = 2 \Xi_{k+1}^{-*} P_{k+1} \begin{pmatrix} G_1^{-T} & \\ & G_2^{-T} \end{pmatrix} P_{k+1}^T z,
$$

where $z = (\frac{1}{2} T_0(\tau), T_1(\tau), T_2(\tau), \ldots, T_k(\tau))^*$. Finally notice $w = 2 \Xi_{k+1}^{-*} y_{(1:k+1)}$ to complete the proof. ∎

Apply Theorem 4.1 to the permuted system (3.29) to get

**Theorem 4.2** *In Theorem 2.1, if $b = e_N$, then the $k$th GMRES residual $r_k$ satisfies for $1 \le k < N$*

$$
\|r_k\|_2 = \|2\Xi_{k+1}^* y_{(1:k+1)}\|_2^{-1}, \tag{4.2}
$$

*where $y \in \mathbb{C}^{2\lceil N/2 \rceil}$ is the same as the one in* Theorem 4.1.

# 5  Concluding Remarks

There are a few GMRES error bounds with simplicity comparable to the well-known bound for the conjugate gradient method [3, 11, 19, 23]. In [6, Section 6], Eiermann and Ernst proved

$$
\frac{\|r_k\|_2}{\|r_0\|_2} \le \left[1 - \gamma(A)\, \gamma(A^{-1})\right]^{k/2} \tag{5.1}
$$

where $\gamma(A) = \inf\{|z^* A z| : \|z\|_2 = 1\}$ is the distance from the origin to $A$'s field of values. When $A$'s Hermitian part, $H = (A + A^*)/2$, is positive definite, it yields a bound by Elman [8] (see also [7])

$$
\frac{\|r_k\|_2}{\|r_0\|_2} \le \left[1 - \left(\frac{1}{\|H^{-1}\|_2 \|A\|_2}\right)^2\right]. \tag{5.2}
$$

As observed in [1], this bound of Elman can be easily extended to cover the case when only $\gamma(A) > 0$

$$
\frac{\|r_k\|_2}{\|r_0\|_2} \le (\sin\theta)^k, \quad \theta = \arccos \frac{\gamma(A)}{\|A\|_2}. \tag{5.3}
$$

Recently Beckermann, Goreinov, and Tyrtyshnikov [1] improved (5.3) to

$$
\frac{\|r_k\|_2}{\|r_0\|_2} \le (2 + 2/\sqrt{3})(2 + \delta)\delta^k, \quad \delta = 2\sin \frac{\theta}{4 - 2\theta/\pi}. \tag{5.4}
$$

All three bounds (5.1), (5.3), and (5.4) yield meaningful estimates only when $\gamma(A) > 0$, i.e., $A$'s field of values does not contain the origin.

However in general, there is not much concrete quantitative results for the convergence rate of GMRES, based on limited information on $A$ and/or $b$. In part, it is a very difficult problem, and such a result most likely does not exist, thanks to the negative result of Greenbaum, Pták, and Strakoš [12] which says that *"Any Nonincreasing Convergence Curve is Possible for GMRES"*. A commonly used approach, as a step towards getting a feel of how fast GMRES may be, is through assuming that $A$ is diagonalizable to arrive at (3.6):

$$\|r_k\|_2/\|r_0\|_2 \leq \kappa(X) \min_{\phi_k(0)=1} \max_i |\phi_k(\lambda_i)|, \tag{5.5}$$

and then putting aside the effect of $\kappa(X)$ to study only the effect in the factor of the associated minimization problem. This approach does not always yield satisfactory results, especially when $\kappa(X) \gg 1$ which occurs when $A$ is highly nonnormal. Getting a fairly accurate quantitative estimate for the convergence rate of GMRES for a highly nonnormal case is likely to be very difficult. Trefethen [22] established residual bounds based on pseudospectra, which sometimes is more realistic than (5.5) but is very expensive to compute. In [4], Driscoll, Toh, and Trefethen provided an nice explanation on this matter.

Our analysis here on tridiagonal Toeplitz $A$ represents one of few diagonalizable cases where one can analyze $r_k$ directly to arrive at simple quantitative results such as (5.1) – (5.4). Previous other results except those in Ernst [9], while helpful in explaining and understanding various convergence behaviors, are more qualitative than quantitative.

Two conjectures are made in Remark 2.1.

# References

[1] B. BECKERMANN, S. A. GOREINOV, AND E. E. TYRTYSHNIKOV, *Some remarks on the Elman estimate for GMRES*, SIAM J. Matrix Anal. Appl., 27 (2006), pp. 772–778.

[2] B. BECKERMANN AND A. B. J. KUIJLAARS, *Superlinear CG convergence for special right-hand sides*, Electron. Trans. Numer. Anal., 14 (2002), pp. 1–19.

[3] J. DEMMEL, *Applied Numerical Linear Algebra*, SIAM, Philadelphia, 1997.

[4] T. A. DRISCOLL, K.-C. TOH, AND L. N. TREFETHEN, *From potential theory to matrix iterations in six steps*, SIAM Rev., 40 (1998), pp. 547–578.

[5] M. EIERMANN, *Fields of values and iterative methods*, Linear Algebra Appl., 180 (1993), pp. 167–197.

[6] M. EIERMANN AND O. G. ERNST, *Geometric aspects in the theory of Krylov subspace methods*, Acta Numer., 10 (2001), pp. 251–312.

[7] S. C. EISENSTAT, H. C. ELMAN, AND M. H. SCHULTZ, *Variational iterative methods for nonsymmetric systems of linear equations*, SIAM J. Numer. Anal., 20 (1983), pp. 345–357.

[8] H. C. ELMAN, *Iterative Methods for Large, Sparse Nonsymmetric Systems of Linear Equations*, PhD thesis, Department of Computer Science, Yale University, 1982.

[9] O. G. ERNST, *Residual-minimizing Krylov subspace methods for stabilized discretizations of convection-diffusion equations*, SIAM J. Matrix Anal. Appl., 21 (2000), pp. 1079–1101.

[10] I. S. GRADSHTEYN AND I. M. RYZHIK, *Table Of Integrals, Series, and Products*, Academic Press, New York, 1980. Corrected and Enlarged Edition prepared by A. Jeffrey, incorporated the fourth edition prepared by Yu. V. Geronimus and M. Yu. Tseytlin, translated from the Russian by Scripta Technica, Inc.

[11] A. GREENBAUM, *Iterative Methods for Solving Linear Systems*, SIAM, Philadelphia, 1997.

[12] 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.

[13] I. C. F. IPSEN, *Expressions and bounds for the GMRES residual*, BIT, 40 (2000), pp. 524–535.

[14] R.-C. LI, *Sharpness in rates of convergence for CG and symmetric Lanczos methods*, Technical Report 2005-01, Department of Mathematics, University of Kentucky, 2005. Avaliable at http://www.ms.uky.edu/∼math/MAreport/.

[15] ——, *On Meinardus' examples for the conjugate gradient method*, Math. Comp., (2006). to appear.

[16] J. LIESEN, M. ROZLOZNÍK, AND Z. STRAKOŠ, *Least squares residuals and minimal residual methods*, SIAM J. Sci. Comput., 23 (2002), pp. 1503–1525.

[17] J. LIESEN AND Z. STRAKOŠ, *Convergence of GMRES for tridiagonal Toeplitz matrices*, SIAM J. Matrix Anal. Appl., 26 (2004), pp. 233–251.

[18] ——, *GMRES convergence analysis for a convection-diffusion model problem*, SIAM J. Sci. Comput., 26 (2005), pp. 1989–2009.

[19] Y. SAAD, *Iterative Methods for Sparse Linear Systems*, SIAM, Philadelphia, 2nd ed., 2003.

[20] Y. SAAD AND M. SCHULTZ, *GMRES: A generalized minimal residual algorithm for solving nonsymmetric linear systems*, SIAM J. Sci. Statist. Comput., 7 (1986), pp. 856–869.

[21] G. D. SMITH, *Numerical Solution of Partial Differential Equations*, Clarendon Press, Oxford, UK, 2nd ed., 1978.

[22] L. N. TREFETHEN, *Pseudospectra of matrices*, in Numerical Analysis 1991:Proceedings of the 14th Dundee Conference, June, 1991, D. F. Griffiths and G. A. Watson, eds., Research Notes in Mathematics Series, Longman Press, 1992.

[23] L. N. TREFETHEN AND D. BAU, III, *Numerical Linear Algebra*, SIAM, Philadelphia, 1997.

[24] I. ZAVORIN, D. P. O'LEARY, AND H. ELMAN, *Complete stagnation of GMRES*, Linear Algebra Appl., 367 (2003), pp. 165–183.