

Asymptotic properties of the QR factorization of banded Hessenberg-Toeplitz matrices

Xiao-Wen Chang^{1,*}, Martin J. Gander² and Samir Karra³

¹*School of Computer Science, McGill University, Montreal, QC, Canada, H3A 2A7*

²*Section de Mathématiques, Université de Genève, CP 240, CH-1211 Genève*

³*Department of Mathematics and Statistics, Sultan Qaboos University, Muscat, Sultanate of Oman*

SUMMARY

We consider the Givens QR factorization of banded Hessenberg-Toeplitz matrices of large order and relatively small bandwidth. We investigate the asymptotic behavior of the R factor and the Givens rotation when the order of the matrix goes to infinity, and present some interesting convergence properties. These properties can lead to savings in the computation of the exact QR factorization and give insight for approximate QR factorizations of interest in preconditioning. The properties also reveal the relation between the limit of the main diagonal elements of R and the largest absolute root of a polynomial. Copyright © 2000 John Wiley & Sons, Ltd.

KEY WORDS: Banded Hessenberg-Toeplitz matrices; discrete and continuous QR factorization; convergence

1. INTRODUCTION

Structured matrices are encountered in various scientific and engineering applications, see for example [11] and [12]. One well-known class is the class of Toeplitz matrices, whose elements along each diagonal are the same constants. In many applications, such as image processing, the Toeplitz matrices are banded. In this paper, we consider the real $N \times N$ banded Hessenberg-

*Correspondence to: School of Computer Science, McGill University, Montreal, QC, Canada, H3A 2A7

Contract/grant sponsor: NSERC of Canada Grant for Xiao-Wen Chang; contract/grant number: RGPIN217191-03

differential operator in its approximate LU factors and leads to an efficient method for these types of problems in acoustics.

This motivated us to look if there is also a continuous limit of the QR factorization in such situations. Applying the QR factorization to the tridiagonal matrix representing the discrete Laplacian in one dimension, one observes that the Q factor converges to a matrix with 1 in the sub-diagonal and zero everywhere else (a shifting matrix), whereas the R factor converges to a matrix like the discretized Laplacian, but the stencil shifted to the upper triangle of the matrix, i.e., the QR factors converge to a factorization of the following form (only one row of each of the three matrices are displayed and the underlined entry is the diagonal entry):

$$\frac{1}{h^2}[\dots, 0, 1, \underline{-2}, 1, 0, \dots] = [\dots, 0, 1, \underline{0}, \dots] \frac{1}{h^2}[\dots, 0, \underline{1}, -2, 1, 0, \dots].$$

In this simple, one dimensional case, it seems that in the continuous limit, the QR factorization of $\frac{d^2}{dx^2}$ is simply $Id \cdot \frac{d^2}{dx^2}$, the orthogonal matrix Q is the identity and the upper triangular factor remains $\frac{d^2}{dx^2}$, which is a local operator and thus already diagonal in the continuous limit. Can such a result be proved for more general Toeplitz matrices? And if yes, is it possible to use the continuous limits of the QR factorization like in the case of the continuous limits of the LU factorization to construct a new class of preconditioners for Krylov methods of the type LSQR [14], which would have as the ideal preconditioner the QR factorization of the matrix A ?

We try to give an answer to the first question for the banded Hessenberg-Toeplitz matrix in (1) in this paper. We investigate the convergence of the diagonals of the R factor and the convergence of the Givens rotations which are used to compute the QR decomposition of the banded Hessenberg-Toeplitz matrix, when N becomes large. This convergence result can be useful by itself, even before it is used to find the continuous limit of the QR factorization of a differential operator and to construct a preconditioner. In fact, if one can show that the entries in the QR factorization converge rapidly to machine precision, one can avoid the computation of the entire factorization and use the limits directly instead in the factors, as soon as convergence is achieved, leading to significant computational savings in solving the linear system $Ax = c$, or the least squares problem $\min \|\bar{A}x - d\|_2$ (where $\bar{A} = \begin{bmatrix} A \\ be_N^T \end{bmatrix}$ with $e_N^T = [0, \dots, 0, 1]$ is still a Toeplitz matrix) by the QR factorization. If A is a symmetric, diagonally dominant, tridiagonal Toeplitz matrix, it is shown in [13] that the diagonals of the LU factors of A converge and computational savings are possible. Similar properties also hold for cyclic reduction, see [4]; but for the QR factorization we are not aware of any results in the literature.

Before proceeding, we introduce some basic notation. We use e_i to denote the unit vector whose i -th element is 1. For any matrix B , we denote by $B(i, :)$ the i -th row, by $B(:, j)$ the j -th column, and by $B(:, j_1 : j_2)$ the submatrix formed by column j_1 up to column j_2 . For a square matrix B , we denote its spectral radius by $\rho(B)$. For a complex number c , we use $\Re(c)$ to denote its real part.

The rest of this paper is organized as follows. In Section 2, we first introduce the Givens QR procedure and give iterative formulas which are the key to our later analysis. Then, after giving two lemmas, we present the main convergence results, followed by a convergence rate analysis. In Section 3, we use some numerical examples to illustrate our findings. Finally a brief summary is given in Section 4.

2. MAIN RESULTS

2.1. The QR factorization of A

For simplicity, we assume from now on that $b = 1$ in the banded Hessenberg-Toeplitz A unless we state otherwise, and say A is “normalized”. This is without loss of generality, since if $b \neq 1$, we can write $A = b \cdot (A/b)$ where the subdiagonal elements of A/b are 1s, and the Q factor of A can be taken as the same as that of A/b and the R factor of A is just b times that of A/b . The QR factorization $A = QR$ can be computed by a sequence of Givens rotations. In the general n -th step of the QR factorization process, a Givens rotation $Q_{n,n+1}$ is applied to rows n and $n+1$ of A to annihilate the n -th element of the subdiagonal. Let $\xi_i^{(n)}$ denote the n -th element of the $(i-1)$ -th upper diagonal of R . The n -th step can then be described as follows:

$$Q_{n,n+1} \begin{bmatrix} x_1^{(n-1)} & x_2^{(n-1)} & \cdots & x_m^{(n-1)} & 0 \\ 1 & a_1 & a_2 & \cdots & a_m \end{bmatrix} = \begin{bmatrix} \xi_1^{(n)} & \xi_2^{(n)} & \cdots & \xi_m^{(n)} & \xi_{m+1}^{(n)} \\ 0 & x_1^{(n)} & x_2^{(n)} & \cdots & x_m^{(n)} \end{bmatrix}, \quad (2)$$

where

$$x^{(0)} \equiv [x_1^{(0)}, x_2^{(0)}, \dots, x_m^{(0)}]^T = [a_1, a_2, \dots, a_m]^T, \quad (3)$$

$$Q_{n,n+1} = \begin{bmatrix} c_n & s_n \\ -s_n & c_n \end{bmatrix}, \quad c_n = \frac{x_1^{(n-1)}}{\sqrt{1 + (x_1^{(n-1)})^2}}, \quad s_n = \frac{1}{\sqrt{1 + (x_1^{(n-1)})^2}}. \quad (4)$$

From (2) and (4) it is straightforward to verify that $x^{(n)} \equiv [x_1^{(n)}, x_2^{(n)}, \dots, x_m^{(n)}]^T$, $n = 0, 1, \dots$, satisfy

$$x^{(n)} = \frac{1}{\sqrt{1 + (x_1^{(n-1)})^2}} G x^{(n-1)}, \quad G \equiv \begin{bmatrix} a_1 & -1 & & & \\ a_2 & & -1 & & \\ \vdots & & & \ddots & \\ \vdots & & & & -1 \\ a_m & & & & \end{bmatrix}, \quad (5)$$

and $\xi^{(n)} \equiv [\xi_1^{(n)}, \xi_2^{(n)}, \dots, \xi_{m+1}^{(n)}]^T$, $n = 1, 2, \dots$, satisfy

$$\begin{cases} \xi_1^{(n)} = \sqrt{1 + (x_1^{(n-1)})^2}, \\ \xi_i^{(n)} = (x_1^{(n-1)} x_i^{(n-1)} + a_{i-1}) / \sqrt{1 + (x_1^{(n-1)})^2}, & i = 2, \dots, m, \\ \xi_{m+1}^{(n)} = a_m / \sqrt{1 + (x_1^{(n-1)})^2}. \end{cases} \quad (6)$$

Notice here that the main diagonal elements $\xi_1^{(n)}$ ($n = 1, 2, \dots$) of R are taken to be positive. If A is not “normalized”, i.e., $b \neq 1$, the signs of the main diagonal elements of its R factor will be taken to be the same as that of b .

Since Givens rotations do not change the 2-norm of each column of A , we have

$$|\xi_1^{(n)}| \leq \|A(:, n)\|_2 \leq \sqrt{1 + \sum_{k=1}^m a_k^2},$$

$$\sqrt{(\xi_{i+1}^{(n)})^2 + (x_i^{(n)})^2} \leq \|A(:, n+i)\|_2 \leq \sqrt{1 + \sum_{k=1}^m a_k^2}, \quad i = 1, \dots, m, \quad n = 1, 2, \dots \quad (7)$$

Therefore all $|\xi_i^{(n)}|$ and $|x_i^{(n)}|$ are bounded.

Our goal is to study the asymptotic behavior of the diagonals of the R factor and the Givens rotations as n goes to infinity. From (4) and (6), we see that the sequences $\{Q_{n,n+1}\}$ and $\{\xi^{(n)}\}$ depend on the sequence $\{x^{(n)}\}$. Thus the key is to investigate the convergence of $\{x^{(n)}\}$ based on the iteration (5). Let the $m \times m$ matrix in (5) be denoted by G , and its largest eigenvalue in absolute value by λ_{\max} . Note that G is nonsingular and $-G$ is the companion matrix of the monic polynomial

$$p(x) = x^m + a_1 x^{m-1} + a_2 x^{m-2} + \cdots + a_{m-1} x + a_m. \quad (8)$$

Also notice that (5) is some kind of an iterative process of the power method for computing an eigenvector associated with λ_{\max} . But unlike two usual scaling factors $\|Gx^{(n-1)}\|_2$ and $\max(Gx^{(n-1)})$ (the entry of $Gx^{(n-1)}$ with largest absolute value) used in the power method (see for example [10, p.330] and [18, p.571]), the scaling factor in (5) is $\sqrt{1 + (x_1^{(n-1)})^2}$, which is equal to $\xi_1^{(n)}$. In the power method one is only interested in the limit of the direction of $\{x^{(n)}\}$, but here we want to know exactly what the limit of $\{x^{(n)}\}$ is. Due to the new scaling factor, our convergence analysis of $\{x^{(n)}\}$ is more complicated than the standard one.

2.2. Convergence results

We need the following two technical Lemmas to prove our main convergence result.

Lemma 1. *If λ is an eigenvalue of G , then the vector $u \equiv [u_1, u_2, \dots, u_m]^T$ defined by*

$$\begin{cases} u_1 &= \lambda, \\ u_i &= (-1)^{i-1} \lambda^i + \sum_{k=1}^{i-1} (-1)^{i-k-1} a_k \lambda^{i-k} = \lambda(a_{i-1} - u_{i-1}), \quad 2 \leq i \leq m-1, \\ u_m &= a_m = \lambda(a_{m-1} - u_{m-1}), \end{cases} \quad (9)$$

is an eigenvector of G associated with λ , and $[(-\lambda)^{m-1}, (-\lambda)^{m-2}, \dots, -\lambda, 1]^T$ is an eigenvector of G^T associated with λ .

Proof. This result is a direct consequence of the special structure of G . \square

Lemma 2. *Let the real sequence $\{x_n\}$ be defined by $x_n = \frac{\lambda x_{n-1}}{\sqrt{1+\gamma^2 x_{n-1}^2}} + y_{n-1}$, where $\lim_{n \rightarrow \infty} y_n \rightarrow 0$, $\lambda \geq 1$, $\gamma \neq 0$, $x_n \geq 0$ for all $n \geq 0$, and for any $n > 0$ there is always $n_0 > n$ such that $x_{n_0} \neq 0$. Then*

$$\lim_{n \rightarrow \infty} x_n = \sqrt{\lambda^2 - 1} / |\gamma|. \quad (10)$$

Proof. The proof is technical and given in Appendix II. \square

We now state and prove the main convergence result of this paper.

Theorem 1. *Let A be the banded Hessenberg-Toeplitz matrix given by (1) with $b = 1$. Let $x_i^{(n)}$ ($i = 1, \dots, m$) and $\xi_i^{(n)}$ ($i = 1, \dots, m+1$) be defined by (2), c_n and s_n by (4), and let G be the $m \times m$ matrix in (5).*

(i) If $\rho(G) < 1$, then

$$\lim_{n \rightarrow \infty} \xi_1^{(n)} = 1, \quad \lim_{n \rightarrow \infty} \xi_i^{(n)} = a_{i-1}, \quad i = 2, \dots, m+1, \quad (11)$$

and

$$\lim_{n \rightarrow \infty} c_n = 0, \quad \lim_{n \rightarrow \infty} s_n = 1. \quad (12)$$

(ii) Suppose G has a distinct largest eigenvalue λ_{\max} in modulus, i.e., λ_{\max} may be repeated, but there are no other different eigenvalues which have the same modulus as λ_{\max} . If $\rho(G) \equiv |\lambda_{\max}| \geq 1$, then

$$\begin{aligned} \lim_{n \rightarrow \infty} \xi_1^{(n)} &= |\lambda_{\max}|, \\ \lim_{n \rightarrow \infty} \xi_i^{(n)} &= \text{sign}(\lambda_{\max}) \left[\sum_{k=0}^{i-1} (-1)^{i-k-1} (a_k - a_{k-2}) \lambda_{\max}^{i-k} + a_{i-2} \right] \end{aligned} \quad (13)$$

$$= -\lambda_{\max} \lim_{n \rightarrow \infty} \xi_{i-1}^{(n)} + \text{sign}(\lambda_{\max}) (a_{i-1} \lambda_{\max} + a_{i-2}), \quad 2 \leq i \leq m, \quad (14)$$

$$\lim_{n \rightarrow \infty} \xi_{m+1}^{(n)} = a_m / |\lambda_{\max}|,$$

where $a_{-2} = a_{-1} \equiv 0$ and $a_0 \equiv 1$, and

$$\lim_{n \rightarrow \infty} |c_n| = \sqrt{\lambda_{\max}^2 - 1} / |\lambda_{\max}|, \quad \lim_{n \rightarrow \infty} s_n = 1 / |\lambda_{\max}|, \quad (15)$$

where when $\rho(G) > 1$ and n is large enough,

$$\text{sign}(c_n) = \text{sign}(\lambda_{\max}) \text{sign}(c_{n-1}). \quad (16)$$

Proof. (i) From (5), we have

$$x^{(n)} = \frac{1}{\prod_{i=0}^{n-1} \sqrt{1 + (x_1^{(i)})^2}} G^n x^{(0)},$$

which implies

$$\|x^{(n)}\|_2 \leq \|G^n\|_2 \|x^{(0)}\|_2.$$

If $\rho(G) < 1$, then $G^n \rightarrow 0$, as $n \rightarrow \infty$, see for example ([10], Lemma 7.3.2). Therefore $x^{(n)} \rightarrow 0$. Thus (11) follows from (6) and (12) follows from (4).

(ii) Since $G - \lambda I$ has rank at least $m - 1$ for any λ , G is nonderogatory, that is, in the Jordan form of G no eigenvalue appears in more than one Jordan block. Suppose that G has k distinct eigenvalues and they are arranged in decreasing order,

$$|\lambda_1| > |\lambda_2| \geq |\lambda_3| \geq \cdots \geq |\lambda_{k-1}| \geq |\lambda_k|,$$

where $\lambda_1 = \lambda_{\max}$, which is real according to the given assumption. Let G have the Jordan decomposition

$$G = SJS^{-1}, \quad J \equiv \text{diag}(J_1, J_2, \dots, J_k), \quad J_i \equiv \begin{bmatrix} \lambda_i & 1 & & \\ & \cdot & \cdot & \\ & & \lambda_i & 1 \\ & & & \lambda_i \end{bmatrix}_{l_i \times l_i}. \quad (17)$$

Defining the new sequence

$$\tilde{x}^{(n)} \equiv [\tilde{x}_1^{(n)}, \tilde{x}_2^{(n)}, \dots, \tilde{x}_m^{(n)}]^T = S^{-1} x^{(n)}, \quad n = 0, 1, \dots, \quad (18)$$

we observe that $\{\tilde{x}^{(n)}\}$ is bounded, since $\{x^{(n)}\}$ is bounded (see (7)). From (5) and (6), we obtain

$$\tilde{x}^{(n)} = \frac{1}{\xi_1^{(n)}} \begin{bmatrix} J_1 & & & \\ & J_2 & & \\ & & \ddots & \\ & & & J_k \end{bmatrix} \tilde{x}^{(n-1)}, \quad n = 1, 2, \dots, \quad (19)$$

where

$$\xi_1^{(n)} = \frac{1}{\sqrt{1 + (x_1^{(n-1)})^2}} = \frac{1}{\sqrt{1 + |\sum_{i=1}^m \gamma_i \tilde{x}_i^{(n-1)}|^2}}, \quad (20)$$

$$S(1, \cdot) \equiv [\gamma_1, \dots, \gamma_m]. \quad (21)$$

We would like to show $\tilde{x}_{l_1}^{(0)} \neq 0$. Since $\tilde{x}^{(0)T} = x^{(0)T} S^{-T}$, we have

$$\tilde{x}_{l_1}^{(0)} = (x^{(0)T} \cdot S^{-T}(:, l_1)). \quad (22)$$

Since $G^T = S^{-T} J^T S^T$, $S^{-T}(:, l_1)$ is an eigenvector of G^T associated with λ_1 . Since G^T is nonderogatory, we have, by Lemma 2.1,

$$(S^{-T}(:, l_1)) = \alpha [(-\lambda_1)^{m-1}, (-\lambda_1)^{m-2}, \dots, -\lambda_1, 1]^T$$

for a nonzero constant α . Combining this and $(x^{(0)T} = [a_1, a_2, \dots, a_m]$ (see (3)), we obtain from (22) that

$$\tilde{x}_{l_1}^{(0)} = \alpha \sum_{i=1}^m a_i (-\lambda_1)^{m-i}.$$

But since $-\lambda_1$ is an eigenvalue of $-G$, we obtain with (8)

$$\sum_{i=1}^m a_i (-\lambda_1)^{m-i} = -(-\lambda_1)^m.$$

Hence, $\tilde{x}_{l_1}^{(0)}$ must be nonzero.

Since for $n \geq l_i - 1$,

$$J_i^n = \begin{bmatrix} \lambda_i^n & \binom{n}{1} \lambda_i^{n-1} & \dots & \binom{n}{l_i-1} \lambda_i^{n-l_i+1} \\ & \lambda_i^n & \dots & \binom{n}{l_i-2} \lambda_i^{n-l_i+2} \\ & & \ddots & \vdots \\ & & & \lambda_i^n \end{bmatrix}$$

we obtain from (19) for $n \geq \max\{l_1 - 1, \dots, l_k - 1\}$,

$$\begin{aligned}
\tilde{x}^{(n)} &= \frac{1}{\xi_1^{(1)} \xi_1^{(2)} \dots \xi_1^{(n)}} \left\{ \left[\left(\lambda_1^n \tilde{x}_1^{(0)} + \binom{n}{1} \lambda_1^{n-1} \tilde{x}_2^{(0)} + \dots + \binom{n}{l_1-1} \lambda_1^{n-l_1+1} \tilde{x}_{l_1}^{(0)} \right) e_1 \right. \right. \\
&\quad + \left. \left(\lambda_1^n \tilde{x}_2^{(0)} + \dots + \binom{n}{l_1-2} \lambda_1^{n-l_1+2} \tilde{x}_{l_1}^{(0)} \right) e_2 + \dots + \lambda_1^n \tilde{x}_{l_1}^{(0)} e_{l_1} \right] \\
&\quad + \left[\left(\lambda_2^n \tilde{x}_{l_1+1}^{(0)} + \binom{n}{1} \lambda_2^{n-1} \tilde{x}_{l_1+2}^{(0)} + \dots + \binom{n}{l_2-1} \lambda_2^{n-l_2+1} \tilde{x}_{l_1+l_2}^{(0)} \right) e_{l_1+1} \right. \\
&\quad + \left. \left(\lambda_2^n \tilde{x}_{l_1+2}^{(0)} + \dots + \binom{n}{l_2-2} \lambda_2^{n-l_2+1} \tilde{x}_{l_1+l_2}^{(0)} \right) e_{l_1+2} + \dots + \tilde{x}_{l_1+l_2}^{(0)} \lambda_2^n e_{l_1+l_2} \right] \\
&\quad + \dots \\
&\quad + \left[\left(\lambda_k^n \tilde{x}_{m-l_k+1}^{(0)} + \binom{n}{1} \lambda_k^{n-1} \tilde{x}_{m-l_k+2}^{(0)} + \dots + \binom{n}{l_k-1} \lambda_k^{n-l_k+1} \tilde{x}_m^{(0)} \right) e_{m-l_k+1} \right. \\
&\quad + \left. \left(\lambda_k^n \tilde{x}_{m-l_k+2}^{(0)} + \dots + \binom{n}{l_k-2} \lambda_k^{n-l_k+1} \tilde{x}_m^{(0)} \right) e_{m-l_k+2} + \dots + \tilde{x}_m^{(0)} \lambda_k^n e_m \right] \left. \right\} \\
&= \binom{n}{l_1-1} \frac{\lambda_1^{n-l_1+1}}{\xi_1^{(1)} \xi_1^{(2)} \dots \xi_1^{(n)}} \left\{ \left[\left(\tilde{x}_{l_1}^{(0)} + \frac{l_1-1}{n-l_1+2} \lambda_1 \tilde{x}_{l_1-1}^{(0)} + \dots + \binom{n}{l_1-1}^{-1} \lambda_1^{l_1-1} \tilde{x}_1^{(0)} \right) e_1 \right. \right. \\
&\quad + \left. \left(\frac{l_1-1}{n-l_1+2} \lambda_1 \tilde{x}_{l_1}^{(0)} + \dots + \binom{n}{l_1-1}^{-1} \lambda_1^{l_1-1} \tilde{x}_2^{(0)} \right) e_2 + \dots + \binom{n}{l_1-1}^{-1} \lambda_1^{l_1-1} \tilde{x}_{l_1}^{(0)} e_{l_1} \right] \\
&\quad + \binom{n}{l_1-1}^{-1} \binom{n}{l_2-1} \left(\frac{\lambda_2}{\lambda_1} \right)^{n-l_1+1} \left[\left(\lambda_2^{l_1-l_2} \tilde{x}_{l_1+l_2}^{(0)} + \dots + \binom{n}{l_2-1}^{-1} \lambda_2^{l_1-1} \tilde{x}_{l_1+1}^{(0)} \right) e_{l_1+1} \right. \\
&\quad + \left. \left(\frac{l_2-1}{n-l_2+2} \lambda_2^{l_1-l_2+1} \tilde{x}_{l_1+l_2}^{(0)} + \dots \right) e_{l_1+2} + \dots + \binom{n}{l_2-1}^{-1} \lambda_2^{l_1-1} \tilde{x}_{l_1+l_2}^{(0)} e_{l_1+l_2} \right] \\
&\quad + \dots \\
&\quad + \binom{n}{l_1-1}^{-1} \binom{n}{l_k-1} \left(\frac{\lambda_k}{\lambda_1} \right)^{n-l_1+1} \left[\left(\lambda_1^{l_1-l_k} \tilde{x}_m^{(0)} + \dots + \binom{n}{l_k-1}^{-1} \lambda_k^{l_1-1} \tilde{x}_{m-l_k+1}^{(0)} \right) e_{m-l_k+1} \right. \\
&\quad + \left. \left(\frac{l_k-1}{n-l_k+2} \lambda_k^{l_1-l_k+1} \tilde{x}_m^{(0)} + \dots \right) e_{m-l_k+2} + \dots + \binom{n}{l_k-1}^{-1} \lambda_k^{l_1-1} \tilde{x}_m^{(0)} e_m \right] \left. \right\}. \quad (23)
\end{aligned}$$

Therefore, using the fact that $|\lambda_i|/|\lambda_1| < 1$ for $i = 2, \dots, m$ and the fact that $\{\tilde{x}_1^{(n)}\}$ is bounded, we can conclude that as $n \rightarrow \infty$,

$$\tilde{x}_1^{(n)} \sim \binom{n}{l_1-1} \lambda_1^{n-l_1+1} \tilde{x}_{l_1}^{(0)} / (\xi_1^{(1)} \xi_1^{(2)} \dots \xi_1^{(n)}), \quad (24)$$

$$\tilde{x}_i^{(n)} \rightarrow 0, \quad i = 2, 3, \dots, m. \quad (25)$$

Then we have with (20) that as $n \rightarrow \infty$,

$$\Re(\tilde{x}_1^{(n)}) \sim \tilde{x}_1^{(n)}, \quad |\tilde{x}_1^{(n)}| \sim \frac{n|\lambda_1| |\tilde{x}_1^{(n-1)}|}{(n-l_1+1)|\xi_1^{(n)}|} \sim \frac{|\lambda_1| |\tilde{x}_1^{(n-1)}|}{\sqrt{1+|\gamma_1|^2 |\tilde{x}_1^{(n-1)}|^2}}.$$

Since $\{|\tilde{x}_1^{(n)}|\}$ is bounded, it follows that we can write

$$|\tilde{x}_1^{(n)}| \equiv \frac{|\lambda_1||\tilde{x}_1^{(n-1)}|}{\sqrt{1 + |\gamma_1|^2|\tilde{x}_1^{(n-1)}|^2}} + y^{(n-1)} \quad (26)$$

for some sequence $\{y^{(n)}\}$ which converges to zero. In order to apply Lemma 2 to the above sequence, we need to show that γ_1 , the first entry of $S(1, \cdot)$ in (21), is nonzero. In fact, since $S(\cdot, 1)$ is the eigenvector of G associated with λ_1 and G is nonderogatory, by Lemma 1 we have

$$S(\cdot, 1) = \beta[u_1, u_1, \dots, u_m]^T, \quad (27)$$

where β is a nonzero real constant, and $u_i, i = 1, \dots, m$ are defined by (9) with $\lambda \equiv \lambda_1$. Thus

$$\gamma_1 = \beta u_1 = \beta \lambda_1 \neq 0. \quad (28)$$

Then applying Lemma 2 to (26), we obtain

$$\lim_{n \rightarrow \infty} |\tilde{x}_1^{(n)}| = (\sqrt{\lambda_1^2 - 1})/|\gamma_1|. \quad (29)$$

Recalling that $x^{(n)} = S\tilde{x}^{(n)}$ (see (18)), we have by (25), (27)–(29) that

$$\lim_{n \rightarrow \infty} |x_1^{(n)}| = \sqrt{\lambda_1^2 - 1}, \quad \lim_{n \rightarrow \infty} x_1^{(n)} x_i^{(n)} = u_i \frac{\lambda_1^2 - 1}{\lambda_1}, \quad i = 2, \dots, m. \quad (30)$$

With these relations, we can obtain from (6) and (9) with some algebraic manipulations (13) and (14). Combining the limit of $\{|x_1^{(n)}|\}$ with (4) leads to (15).

If $|\lambda_1| > 1$, then $\lim_{n \rightarrow \infty} |\Re(\tilde{x}_1^{(n)})| = \lim_{n \rightarrow \infty} |\tilde{x}_1^{(n)}| \neq 0$. It follows from (24) that when n is large enough,

$$\text{sign}(\Re(\tilde{x}_1^{(n)})) = \text{sign}(\lambda_1) \text{sign}(\Re(\tilde{x}_1^{(n-1)})),$$

i.e., the sign of $\Re(\tilde{x}_1^{(n)})$ will change alternatively when n is large enough. But since $x_1^{(n)} \sim \gamma_1 \tilde{x}_1^{(n)} \sim \gamma_1 \Re(\tilde{x}_1^{(n)})$ when n is large enough,

$$\text{sign}(x_1^{(n)}) = \text{sign}(\gamma_1) \text{sign}(\Re(\tilde{x}_1^{(n)})) = \text{sign}(\lambda_1) \text{sign}(x_1^{(n-1)}).$$

Hence (16) follows from this and (4). \square

Remark 1. *Theorem 1 can be extended to a broader class of matrices. Notice that the initial vector $x^{(0)}$ in the iteration (5) is defined in (3), but this is not necessary for having the convergence results. In case (i), we see from the proof that $x^{(0)}$ can be arbitrary. In case (ii), the sequence $\{\tilde{x}^{(n)}\}$ starts with the initial vector $\tilde{x}^{(0)} = S^{-1}x^{(0)}$. We observe that the arguments of the proof can still be carried out if $\tilde{x}^{(0)}$ has a nonzero component $\tilde{x}_i^{(0)}$ for some i with $1 \leq i \leq l_1$. Since $x^{(0)} = S\tilde{x}^{(0)}$, this means that the results in Theorem 1 (ii) still hold, if the initial vector $x^{(0)}$ has at least one nonzero component in the invariant subspace of G associated with λ_{\max} . This is exactly a condition assumed for the convergence of the power method, see [18, p.583]. Thus, even if the upper-left elements of the banded Hessenberg Toeplitz matrix A are changed, we still have the convergence results of Theorem 1 as long as in case (ii) after some finite steps of the QR factorization, we can get a vector $x^{(k)}$ which satisfies the above condition. Such a class of matrices may arise in certain applications, see the second example given in Section 3.*

Remark 2. From case (i), we observe that the limits of the rows of R can be obtained by just simply shifting the rows of A one position to the right. Case (ii) establishes an interesting relation between the limit of the main diagonal elements of the R factor and the largest root in magnitude of the polynomial equation $p(x) = 0$ (see (8)). If we know one of them, the other will be known. When the bandwidth of A is small, $\lambda_{\max}(G)$ can easily be computed, and then the limits of $\{\xi_i^{(n)}\}$ ($i = 1, \dots, m + 1$) can easily be obtained from the recursion (14).

Remark 3. If the discretization A of a differential operator A leads to case (i), then the limits of the QR factorization can be interpreted, when the mesh is refined, as the continuous QR factorization of the differential operator $A = Id \cdot A$, as it was shown for the simple example in the introduction. Furthermore, the limits can be used to construct an approximate QR factorization of A , $\hat{A} = \hat{Q}\hat{R}$ where $\hat{Q} = [e_2, \dots, e_N, e_1]$ and $\hat{R} = [e_1, A(:, 1:N - 1)]$, i.e., \hat{R} is a right-shift of A to upper triangular form (note that the subdiagonal of A has been normalized to one). It is easy to verify that the difference $A - \hat{A}$ is a rank one matrix, so the approximate QR factorization is only a rank one change distant from the exact one. A similar low rank distance property was also shown for approximate LU factorizations of banded Toeplitz matrices, see [9], where the low rank difference is treated by the Sherman-Morrison-Woodbury formula in the solution process.

There are however discretizations of differential operators which lead to case (ii); an example of which is given in Section 3. In that case, the rows of the R factor do not converge to a simple shift of the rows of the discretized operator A . Nevertheless, even in that case, $\hat{A} = \hat{Q}\hat{R}$ defined above constitutes an approximate QR factorization of A , since $A - \hat{A}$ is still a rank one matrix. It might well be this approximate factorization which has good preconditioning properties, even in case (ii).

Remark 4. Theorem 1 is not only interesting in theory but also useful in computations. If the dimension of the banded Hessenberg-Toeplitz matrix is large, we may not need to carry out the entire QR factorization process to obtain the QR factors, because the limits of the R factor and the Givens rotations may have been reached to machine precision before the QR factorization process is completed, which can lead to significant savings in the computation.

A special case of Theorem 1 is when the matrix is tridiagonal,

$$A = \begin{bmatrix} a & c & & & \\ b & a & c & & \\ & \ddots & \ddots & \ddots & \\ & & b & a & c \\ & & & b & a \end{bmatrix}. \quad (31)$$

Here, we prefer not to “normalize” the subdiagonal elements of A . The n -th step of the QR factorization of A can be described by

$$Q_{n,n+1} \begin{bmatrix} x_1^{(n-1)} & x_2^{(n-1)} & 0 \\ b & a & c \end{bmatrix} = \begin{bmatrix} \xi_1^{(n)} & \xi_2^{(n)} & \xi_3^{(n)} \\ 0 & x_1^{(n)} & x_2^{(n)} \end{bmatrix},$$

where (cf. (3)–(6))

$$\begin{aligned} x^{(0)} &\equiv [x_1^{(0)}, x_2^{(0)}]^T = [a, c]^T, \\ Q_{n,n+1} &= \begin{bmatrix} c_n & s_n \\ -s_n & c_n \end{bmatrix}, \quad c_n = \frac{\text{sign}(b) x_1^{(n-1)}}{\sqrt{b^2 + (x_1^{(n-1)})^2}}, \quad s_n = \frac{\text{sign}(b)b}{\sqrt{b^2 + (x_1^{(n-1)})^2}}, \\ \begin{bmatrix} x_1^{(n)} \\ x_2^{(n)} \end{bmatrix} &= \frac{\text{sign}(b)}{\sqrt{b^2 + (x_1^{(n-1)})^2}} \begin{bmatrix} a & -b \\ c & 0 \end{bmatrix} \begin{bmatrix} x_1^{(n-1)} \\ x_2^{(n-1)} \end{bmatrix}, \end{aligned} \quad (32)$$

$$\xi_1^{(n)} = \text{sign}(b)\sqrt{b^2 + (x_1^{(n-1)})^2}, \quad \xi_2^{(n)} = \frac{\text{sign}(b)(x_1^{(n-1)} x_2^{(n-1)} + ab)}{\sqrt{b^2 + (x_1^{(n-1)})^2}}, \quad \xi_3^{(n)} = \frac{\text{sign}(b)bc}{\sqrt{b^2 + (x_1^{(n-1)})^2}}. \quad (33)$$

The eigenvalues of the matrix $G = \begin{bmatrix} a & -b \\ c & 0 \end{bmatrix}$ corresponding to the “unnormalized” A are

$$\lambda_{\max, \min} = \frac{1}{2}(a \pm \sqrt{a^2 - 4bc}), \quad |\lambda_{\max}| \geq |\lambda_{\min}|. \quad (34)$$

So the eigenvalues of G/b corresponding to the “normalized” A/b are $\lambda_{\max, \min}/b$. By Theorem 1, we have

Corollary 1. *Let A be the tridiagonal Toeplitz matrix in (31), and let $\lambda_{\max, \min}$ be defined by (34).*

(i) *If $|\lambda_{\max}| < |b|$, then*

$$\lim_{n \rightarrow \infty} \xi_1^{(n)} = b, \quad \lim_{n \rightarrow \infty} \xi_2^{(n)} = a, \quad \lim_{n \rightarrow \infty} \xi_3^{(n)} = c, \quad \lim_{n \rightarrow \infty} c_n = 0, \quad \lim_{n \rightarrow \infty} s_n = 1.$$

(ii) *If $|\lambda_{\max}| \geq |b|$, $a^2 - 4bc > 0$, and $a \neq 0$, then*

$$\begin{aligned} \lim_{n \rightarrow \infty} \xi_1^{(n)} &= \text{sign}(b\lambda_{\max})\lambda_{\max}, \quad \lim_{n \rightarrow \infty} \xi_2^{(n)} = \text{sign}(b\lambda_{\max})(b+c), \quad \lim_{n \rightarrow \infty} \xi_3^{(n)} = \text{sign}(b\lambda_{\max})\lambda_{\min}, \\ \lim_{n \rightarrow \infty} |c_n| &= \sqrt{\lambda_{\max}^2 - b^2}/|\lambda_{\max}|, \quad \lim_{n \rightarrow \infty} s_n = |b|/|\lambda_{\max}|, \end{aligned}$$

where, if $|\lambda_{\max}| > b$ and n is large enough,

$$\text{sign}(c_n) = \text{sign}(b\lambda_{\max})\text{sign}(c_{n-1}).$$

Remark 5. *In Theorem 1 (ii) and Corollary 1 (ii), we (implicitly) assume that G has only one distinct dominant eigenvalue. This is also an assumption for the power method. The following example shows that if this is not true, then convergence may not be guaranteed. Suppose in the tridiagonal A , $a = 0$. In this case, $G = \begin{bmatrix} 0 & -b \\ c & 0 \end{bmatrix}$ has exactly two distinct eigenvalues with the same modulus. Since $x_1^{(0)} = 0$, it is easy to verify from (32) that $x_1^{(2n)} = 0$ for all n , and*

$$x_1^{(1)} = -\text{sign}(b)c, \quad x_1^{(2n+1)} = -\text{sign}(b) \frac{c x_1^{(2n-1)}}{\sqrt{b^2 + (x_1^{(2n-1)})^2}}, \quad n = 1, 2, \dots$$

If $|c| > |b|$, then by Lemma 2 the sequence $\{|x_1^{(2n+1)}|\}$ converges to $\sqrt{c^2 - b^2}$. Hence, from (33),

$$\lim_{n \rightarrow \infty} \xi_1^{(2n)} = b, \quad \lim_{n \rightarrow \infty} \xi_1^{(2n+1)} = \text{sign}(b)|c|,$$

which shows that $\{\xi_1^{(n)}\}$ does not converge.

2.3. Evaluation of the convergence rates

In Section 2.2, we presented the convergence results for the diagonals of the R factor and the Givens rotations. How fast do these sequences converge? This is the question we would like to address in this section. Since the power method with standard scaling factors usually converges linearly, we can expect that these sequences usually have the same order of convergence. But since the situation here is more complicated than the power method with standard scaling factors, we perform a convergence rate analysis here, giving not only the convergence order, but also convergence factors.

For a sequence $\{y^{(n)}\}$ which converges to the limit y^* , we define

$$r\{y^{(n)}\} \equiv \limsup_{n \rightarrow \infty} \|y^{(n)} - y^*\|^{1/n},$$

where we assume that the limit exists. If $0 < r\{y^{(n)}\} < 1$ or $r\{y^{(n)}\} = 1$, we say the sequence $\{y^{(n)}\}$ has r-linear or r-sublinear convergence, respectively (here “r” stands for “root”); and $r\{y^{(n)}\}$ is called the convergence factor. For various measures of efficiency of iterative processes, see [15, Chap.3].

In our analysis, we will use the following lemmas.

Lemma 3. Let $\{y^{(n)}\}$ be a sequence defined by the iterative process

$$y^{(n)} = F(y^{(n-1)}), \quad n = 1, 2, \dots, \quad (35)$$

where $F : \mathbb{R}^p \rightarrow \mathbb{R}^p$ is a mapping. Let y^* be the limit of $\{y^{(n)}\}$. If F is Fréchet differentiable at y^* , then

$$r\{y^{(n)}\} = \rho(DF(y^*)),$$

where $DF(y^*)$ is the Jacobian matrix of F at y^* .

Proof. It is straightforward to see that the conclusion is true from the proof of [15, Theorem 3.5]. But we would like to make a remark here to avoid possible confusion. In [15, Theorem 3.5], it is assumed that $\rho(DF(y^*)) < 1$. But this condition was only used to show the local convergence of $\{y^{(n)}\}$ (in our notation), while here we have assumed that $\{y^{(n)}\}$ converges. So we do not need the condition $\rho(DF(y^*)) < 1$. \square

Lemma 4. Let two sequences $\{y^{(n)}\} \in \mathbb{R}^p$ and $\{z^{(n)}\} \in \mathbb{R}^q$ have limits y^* and z^* , respectively, and be related by

$$z^{(n)} = F(y^{(n)}),$$

where the mapping $F : \mathbb{R}^p \rightarrow \mathbb{R}^q$ is Fréchet-differentiable at y^* . Then

$$r\{z^{(n)}\} \leq r\{y^{(n)}\}. \quad (36)$$

Proof. Since F is Fréchet-differentiable at y^* , it is locally Lipschitz-continuous at y^* . So there exists a constant $c > 0$ and an integer n_0 such that

$$\|z^{(n)} - z^*\| \leq c \|y^{(n)} - y^*\| \quad \text{for } n > n_0,$$

which leads to (36). \square

We now analyze the convergence rates of $\{\xi^{(n)}\}$, $\{c_n\}$ and $\{s_n\}$ based on different cases in Theorem 1.

Case (i): $\rho(G) \equiv |\lambda_{\max}| < 1$. Write (5) as

$$x^{(n)} = F(x^{(n-1)}), \quad (37)$$

where $F: \mathbb{R}^m \rightarrow \mathbb{R}^m$ is Fréchet-differentiable. A direct calculation shows that the Jacobian of F at the limit zero is given by $DF(0) = G$. Thus by Lemma 3, we have

$$r\{x^{(n)}\} = \rho(G). \quad (38)$$

This result can also be obtained from (23), where $|\lambda_2|$ can be equal to $|\lambda_1|$ and $r\{x^{(n)}\} = r\{\tilde{x}^{(n)}\}$.

Since all c_n , s_n and $\xi^{(n)}$ involve $x_1^{(n)}$, we are interested in obtaining $r\{x_1^{(n)}\}$. Using the equations in (5) from the bottom to the top, one can show by Lemma 4 that

$$r\{x_i^{(n)}\} \leq r\{x_1^{(n)}\}, \quad i = m, m-1, \dots, 2. \quad (39)$$

On the other hand, since $x_1^{(n)}$ is an element of the vector $x^{(n)}$, we have

$$r\{x_1^{(n)}\} \leq r\{x^{(n)}\}.$$

Therefore we must have

$$r\{x_1^{(n)}\} = r\{x^{(n)}\}. \quad (40)$$

In the derivation of (40), we did not use the condition $\rho(G) < 1$, so it is always true. Then from (40) and (38), we obtain

$$r\{x_1^{(n)}\} = \rho(G). \quad (41)$$

Using the expressions of c_n and s_n in (4) one can show that

$$r\{c_n\} = r\{x_1^{(n)}\} = \rho(G), \quad r\{s_n\} = (r\{x_1^{(n)}\})^2 = \rho^2(G).$$

Therefore both $\{c_n\}$ and $\{s_n\}$ have r-linear convergence, but the latter converges faster than the former.

Now let us consider the convergence rate of $\{\xi^{(n)}\}$. From (6), it is easy to show with (39) and (41) that

$$\begin{aligned} r\{\xi_1^{(n)}\} &= r\{\xi_m^{(n)}\} = (r\{x_1^{(n)}\})^2 = \rho^2(G), \\ r\{\xi_i^{(n)}\} &\leq r\{\max\{|x_1^{(n)}|, |x_i^{(n)}|\}\} \cdot r\{x_1^{(n)}\} \leq (r\{x_1^{(n)}\})^2 = \rho^2(G). \end{aligned}$$

Then it follows that

$$r\{\xi^{(n)}\} = \rho^2(G).$$

Thus $\{\xi^n\}$ also has r-linear convergence.

Case (ii) (a): $\rho(G) \equiv |\lambda_{\max}| = 1$. This is the special case in Theorem 1 (ii). By similar arguments as those in Case (i), we can show that $\{c_n\}$, $\{s_n\}$ and $\{\xi^{(n)}\}$ have r-sublinear convergence. Therefore all these sequences converge very slowly.

Case (ii) (b): $|\lambda_{\max}| > 1$. From the proof of Theorem 1 (ii), we see that the sign of λ_1 does not affect the convergence rates of the sequences in question, so we assume here without loss of generality that λ_{\max} is positive. Then the sequence $\{x^{(n)}\}$ itself converges. Again writing (5) like (37), denoting the limit of $x^{(n)}$ by x^* , some calculations using (30), (9), (27) and (28) show that

$$DF(x^*) = \frac{1}{\lambda_1}G + \frac{1 - \lambda_1^2}{\lambda^2 \gamma_1} [S(:, 1), 0].$$

Then using the Jordan decomposition (17) and (21) one can verify that

$$DF(x^*) = S \left(\frac{1}{\lambda_1}J + \frac{1 - \lambda_1^2}{\lambda^2 \gamma_1} \begin{bmatrix} S(1, \cdot) \\ 0 \end{bmatrix} \right) S^{-1},$$

where the middle matrix on the righthand side is upper triangular and its diagonal part is

$$\text{diag}(DF(x^*)) = \text{diag}(\underbrace{1/\lambda_1^2, 1, \dots, 1}_{l_1}, \underbrace{\lambda_2/\lambda_1, \dots, \lambda_2/\lambda_1}_{l_2}, \dots, \underbrace{\lambda_k/\lambda_1, \dots, \lambda_k/\lambda_1}_{l_k}).$$

Then it follows that

$$r\{x^{(n)}\} = \rho(DF(x^*)) = \begin{cases} 1, & \text{if } l_1 > 1, \\ \max\{1/\lambda_1^2, |\lambda_2/\lambda_1|\} < 1, & \text{if } l_1 = 1. \end{cases} \quad (42)$$

From (23) we see that for $2 \leq i \leq m$ the absolute value of the (i, i) entry of $\text{diag}(DF(x^*))$ is just the convergence factor of the sequence $\{\tilde{x}_i^{(n)}\}$, which converges to zero. The absolute value of the $(1, 1)$ entry of $\text{diag}(DF(x^*))$ can be understood as the convergence factor of $\{\tilde{x}_1^{(n)}\}$ (see (23) and (26)). If the Jordan block associated with the largest eigenvalue of G has dimension larger than 1, (42) indicates that $\{x^{(n)}\}$ has only r-sublinear convergence. This can easily be observed from (23), where the coefficient of e_2 converges to zero sublinearly—very slowly.

It is not difficult to show from the expressions of c_n and s_n in (4) with (40) and (42) that

$$r\{c_n\} = r\{s_n\} = r\{x_1^{(n)}\} = \begin{cases} 1, & \text{if } l_1 > 1, \\ \max\{1/\lambda_1^2, |\lambda_2/\lambda_1|\}, & \text{if } l_1 = 1. \end{cases}$$

Therefore $\{c_n\}$ and $\{s_n\}$ have r-sublinear or r-linear convergence.

Applying Lemma 4 to (6), we obtain

$$r\{\xi^{(n)}\} \leq r\{x^{(n)}\}.$$

On the other hand, from (6) we have

$$\begin{cases} (x_1^{(n-1)})^2 = (\xi_1^{(n)})^2 - 1, \\ (x_i^{(n-1)})^2 = (\xi_i^{(n)} \xi_1^{(n)} - a_{i-1}) / ((\xi_1^{(n)})^2 - 1), & i = 2, \dots, m, \end{cases} \quad (43)$$

where $\xi_1^* \equiv \lim_{n \rightarrow \infty} \xi_1^{(n)} = \sqrt{1 + x_1^*} > 1$, and applying Lemma 4 to (43) gives

$$r\{(x^{(n)})^2\} \leq r\{\xi^{(n)}\}, \quad \text{where } (x^{(n)})^2 \equiv [(x_1^{(n)})^2, \dots, (x_m^{(n)})^2]^T.$$

Table I. Results for the first example

n	$\xi_1^{(n)}$	$\xi_2^{(n)}$	$\xi_3^{(n)}$	c_n	s_n
1	5.830951894845299e+00	3.086974532565159e+00	8.574929257125441e-01	5.144957554275265e-01	8.574929257125441e-01
2	5.046839430306270e+00	3.042090280161516e+00	9.907190567575820e-01	1.359255332061169e-01	9.907190567575820e-01
3	5.001039152986085e+00	2.996605851029752e+00	9.997922125873653e-01	-2.038459343868916e-02	9.997922125873653e-01
4	5.003881405883999e+00	2.998475700943611e+00	9.992243209682320e-01	-3.937964430484279e-02	9.992243209682320e-01
5	5.000955808892974e+00	3.000196507856399e+00	9.998088747572463e-01	-1.955029300672493e-02	9.998088747572463e-01
6	5.000037166925046e+00	3.000053080536299e+00	9.999925666702456e-01	-3.855723570804071e-03	9.999925666702456e-01
7	5.000006372788071e+00	2.999990020309384e+00	9.999987254440106e-01	1.596593359157377e-03	9.999987254440106e-01
8	5.000007474459249e+00	2.999998275996201e+00	9.999985051103849e-01	1.729097161974871e-03	9.999985051103849e-01
9	5.000001289312853e+00	3.000000468146182e+00	9.999997421374963e-01	7.181399175182089e-04	9.999997421374963e-01
10	5.000000018089969e+00	3.000000050234309e+00	9.99999963820063e-01	8.506460698205245e-05	9.99999963820063e-01
11	5.000000021431909e+00	2.999999979264790e+00	9.999999957136183e-01	-9.258921839789244e-05	9.999999957136183e-01
12	5.000000013164724e+00	2.99999998820036e+00	9.999999973670549e-01	-7.256645217115860e-05	9.999999973670549e-01
13	5.000000001565255e+00	3.00000000876607e+00	9.99999996869491e-01	-2.502202766403997e-05	9.99999996869491e-01
14	5.00000000000624e+00	3.00000000012134e+00	9.9999999998750e-01	-4.999261687355971e-07	9.9999999998750e-01
15	5.000000000055329e+00	2.99999999964450e+00	9.99999999889340e-01	4.704449831513951e-06	9.99999999889340e-01
16	5.000000000021355e+00	3.00000000000937e+00	9.99999999957291e-01	2.922655132641901e-06	9.99999999957291e-01
17	5.000000000001651e+00	3.00000000001385e+00	9.9999999996698e-01	8.127031132861006e-07	9.9999999996698e-01
18	5.000000000000024e+00	2.9999999999907e+00	9.9999999999953e-01	-9.690915855652647e-08	9.9999999999953e-01
19	5.000000000000122e+00	2.9999999999948e+00	9.9999999999757e-01	-2.206861177911299e-07	9.9999999999757e-01
20	5.000000000000032e+00	3.00000000000005e+00	9.9999999999936e-01	-1.130298389633724e-07	9.9999999999936e-01
21	5.000000000000001e+00	3.00000000000002e+00	9.999999999997e-01	-2.368067981979774e-08	9.999999999997e-01
22	5.00000000000000e+00	3.00000000000000e+00	1.00000000000000e+00	8.397559900795826e-09	1.00000000000000e+00
23	5.00000000000000e+00	3.00000000000000e+00	1.00000000000000e+00	9.774671904437044e-09	1.00000000000000e+00
24	5.00000000000000e+00	3.00000000000000e+00	1.00000000000000e+00	4.185291162503062e-09	1.00000000000000e+00
25	5.00000000000000e+00	3.00000000000000e+00	1.00000000000000e+00	5.562403166144281e-10	1.00000000000000e+00

numerical boundary conditions, see [2] and [3] (here for simplicity we have already normalized the subdiagonal elements by multiplying the original matrix by -3). If we delete the first and last rows, this matrix becomes a Toeplitz matrix, so it belongs to a class of matrices called quasi-Toeplitz matrices, see [3]. Note that the first row of A contains four nonzero elements instead of three and after the first step of the QR factorization the second row has 3 nonzero elements, so we consider $x^{(1)}$ as the initial vector of the regular QR procedure (2). The matrix

$$G = \begin{bmatrix} 3/2 & -1 & 0 \\ -3 & 0 & -1 \\ 1/2 & 0 & 0 \end{bmatrix}$$

corresponding to A has the eigenvalues $\lambda_1 = \frac{5+\sqrt{33}}{4} > 1$, $\lambda_2 = -1$, $\lambda_3 = \frac{5-\sqrt{33}}{4}$. Direct computations show that this initial vector is not orthogonal to the eigenvector of G associated with λ_1 . Thus by Theorem 1 (ii) and Remark 1, we have

$$\begin{aligned} \lim_{n \rightarrow \infty} \xi_1^{(n)} &= |\lambda_1| = 2.686140661634507e+00, \\ \lim_{n \rightarrow \infty} \xi_2^{(n)} &= \text{sign}(\lambda_1)(-\lambda_1^2 + a_1\lambda_1 + 1) = -2.186140661634507e+00, \\ \lim_{n \rightarrow \infty} \xi_3^{(n)} &= \text{sign}(\lambda_1)(\lambda_1^3 - a_1\lambda_1^2 + (a_2 - 1)\lambda_1 + a_1) = -6.861406616345072e-01, \\ \lim_{n \rightarrow \infty} \xi_4^{(n)} &= a_3/\lambda_1 = 1.861406616345072e-01 \\ \lim_{n \rightarrow \infty} |c_n| &= \sqrt{\lambda_1^2 - 1}/|\lambda_1| = 9.281199364010406e-01, \\ \lim_{n \rightarrow \infty} s_n &= 1/\lambda_1 = 3.722813232690143e-01. \end{aligned}$$

According to Case (ii) (b) in Section 2.3, all of the sequences $\{\xi^{(n)}\}$, $\{c_n\}$ and $\{s_n\}$ have r-linear convergence with the same convergence factor $|\lambda_2/\lambda_1| \approx 0.372$. In our computation, we took $\alpha = 6/5$ and $\beta = -4/5$ as in [3]. The computed results for the R factor and the Givens rotations of the first 22 steps of the QR factorization procedure are displayed in Tables 2 and 3, respectively. The results show that all sequences converge to their corresponding limits to machine precision in 21 steps and the convergence is linear with almost the same speed. So if the order of A is larger than 22, we do not need to compute the general rows of R any more after 21 steps, except the last two rows of R (because the last row of A is special), which can be computed separately.

Of course, one can find examples like the well-known second order difference matrix having the form of (31) with $a = -2$, $b = 1$ and $c = 1$, for which the convergence is extremely slow, since $\rho(G) = 1$.

4. SUMMARY AND FUTURE WORK

We analyzed the convergence of the Givens QR factorization of a banded Hessenberg-Toeplitz matrix when its dimension goes to infinity, and presented some interesting properties of these infinite factorizations. An immediate practical use of the results given here is that one will not need to carry out the entire QR factorization procedure for a given banded Hessenberg-Toeplitz matrix with large dimensions, allowing some saving in the computation. More importantly

Table II. Results for the second example— R factor

n	$\xi_1^{(n)}$	$\xi_2^{(n)}$	$\xi_3^{(n)}$	$\xi_4^{(n)}$
1	8.161494961096285e+00	-1.648962606011617e+01	1.183974265261573e+01	-3.511611553595848e+00
2	3.685408036776648e+00	-3.909220055578865e+00	8.814183112950758e-02	1.356701876727096e-01
3	2.842268895951231e+00	-2.413162633998660e+00	-6.050220596141659e-01	1.759157976615944e-01
4	2.701381390350119e+00	-2.210711477879456e+00	-6.757604005395086e-01	1.850904880688455e-01
5	2.688549175508812e+00	-2.189834721708226e+00	-6.846883629091571e-01	1.859739091085712e-01
6	2.686452196633474e+00	-2.186630209037751e+00	-6.859410633889107e-01	1.861190757931872e-01
7	2.686185340933632e+00	-2.186210009400613e+00	-6.861128970914022e-01	1.861375655583825e-01
8	2.686146749059763e+00	-2.186150167889759e+00	-6.861368209670350e-01	1.861402397970311e-01
9	2.686141512561523e+00	-2.186141986387669e+00	-6.861401288419524e-01	1.861406026680986e-01
10	2.686140779064683e+00	-2.186140844733832e+00	-6.861405878278367e-01	1.861406534969848e-01
11	2.686140677944381e+00	-2.186140687045691e+00	-6.861406514029780e-01	1.861406605042869e-01
12	2.686140663892535e+00	-2.186140665153915e+00	-6.861406602166520e-01	1.861406614780333e-01
13	2.686140661947622e+00	-2.186140662122441e+00	-6.861406614379904e-01	1.861406616128094e-01
14	2.686140661677891e+00	-2.186140661702120e+00	-6.861406616072721e-01	1.861406616315008e-01
15	2.686140661640521e+00	-2.186140661643879e+00	-6.861406616307326e-01	1.861406616340904e-01
16	2.686140661635341e+00	-2.186140661635806e+00	-6.861406616339838e-01	1.861406616344494e-01
17	2.686140661634623e+00	-2.186140661634687e+00	-6.861406616344347e-01	1.861406616344992e-01
18	2.686140661634523e+00	-2.186140661634532e+00	-6.861406616344973e-01	1.861406616345061e-01
19	2.686140661634509e+00	-2.186140661634510e+00	-6.861406616345058e-01	1.861406616345070e-01
20	2.686140661634507e+00	-2.186140661634508e+00	-6.861406616345070e-01	1.861406616345072e-01
21	2.686140661634507e+00	-2.186140661634507e+00	-6.861406616345072e-01	1.861406616345072e-01
22	2.686140661634507e+00	-2.186140661634507e+00	-6.861406616345072e-01	1.861406616345072e-01

Table III. Results for the second example—Givens rotation

n	c_n	s_n
1	9.924652332214365e-01	1.225265720026465e-01
2	9.624834547707337e-01	2.713403753454192e-01
3	9.360633143822804e-01	3.518315953231888e-01
4	9.289596573079728e-01	3.701809761376910e-01
5	9.282536401886650e-01	3.719478182171425e-01
6	9.281372519749216e-01	3.722381515863744e-01
7	9.281224201321695e-01	3.722751311167650e-01
8	9.281202748099058e-01	3.722804795940623e-01
9	9.281199837054626e-01	3.722812053361973e-01
10	9.281199429291787e-01	3.722813069939696e-01
11	9.281199373077336e-01	3.722813210085738e-01
12	9.281199365265680e-01	3.722813229560666e-01
13	9.281199364184471e-01	3.722813232256187e-01
14	9.281199364034524e-01	3.722813232630016e-01
15	9.281199364013749e-01	3.722813232681809e-01
16	9.281199364010869e-01	3.722813232688988e-01
17	9.281199364010471e-01	3.722813232689983e-01
18	9.281199364010414e-01	3.722813232690121e-01
19	9.281199364010407e-01	3.722813232690140e-01
20	9.281199364010406e-01	3.722813232690143e-01
21	9.281199364010406e-01	3.722813232690143e-01
22	9.281199364010406e-01	3.722813232690143e-01

however, the understanding of the limits of the QR factorization as the matrix dimension grows large gives us insight into the limits of the QR factorization of the discretization of differential operators. Theorem 1 formally shows that the limit of the QR factorization of the one dimensional Laplacian operator shown in the introduction is indeed correct. To generalize this idea to higher dimensional differential operators however, we need to take two more steps: we need to investigate the convergence of the QR factorization and its block variants of general banded Toeplitz matrices; and we need to estimate how well the QR factors based on the limits only are approximating the true factors, in order to understand the preconditioning qualities of the approximate factors for Krylov methods like LSQR. Preliminary numerical experiments suggest that one can still get convergence for general banded Toeplitz matrices under certain conditions and that the approximate QR factors are very well suited to precondition LSQR.

APPENDIX

II. APPENDIX: Proof of Lemma 2

We deal with the case $\lambda > 1$ and the case $\lambda = 1$ separately.

Case 1: $\lambda > 1$. If there exists $n_0 > 0$ such that when $n \geq n_0$, x_n is a constant, say c , then $c = \lambda c / \sqrt{1 + \gamma^2 c^2} + y_{n-1}$. Since $\lim_{n \rightarrow \infty} y_n = 0$, we have $c = \lambda c / \sqrt{1 + \gamma^2 c^2}$. Thus

$$c = 0, \quad \text{or} \quad c = \sqrt{\lambda^2 - 1} / |\gamma|.$$

But according to the assumption, $c = 0$ is impossible, thus in this situation the conclusion

(10) is true.

In the following we assume that x_n is not a constant for large enough n . Let $g(x) \equiv \lambda x / \sqrt{1 + \gamma^2 x^2}$ for $x \geq 0$. Then $g'(x) = \lambda / (1 + \gamma^2 x^2)^{3/2}$. Notice $g(x)$ is strictly monotonically increasing and $g'(x)$ is strictly monotonically decreasing. These two properties will be used a few times in our proof. One can verify that $g(x)$ has the two fixed points 0 and $x^* = \sqrt{\lambda^2 - 1} / |\gamma|$, i.e., $g(0) = 0$ and $g(x^*) = x^*$, and

$$g'(0) = \lambda > 1, \quad g'(x^*) = 1/\lambda^2 < 1, \quad g(x) > x \text{ for } x \in (0, x^*).$$

Let $g'(z) = 1$ for some z ; then we must have $0 < z < x^*$. Since $g(z) > z$, there exists a small positive number $\delta < z$, such that $\eta \equiv g(z - \delta) - z > 0$. Thus

$$z + \eta = g(z - \delta) < g(x^*) = x^*. \quad (\text{II.1})$$

We now show that for any $n > 0$, there exists some $n_0 > n$ such that $x_{n_0} \geq z - \delta$. In fact, if this is not true, then there exists some $n_0 > 0$ such that $x_n < z - \delta$ for any $n \geq n_0$. Since $g'(x)$ is monotonically decreasing, $g'(x) \geq g'(z - \delta) > 1$ for $0 \leq x \leq z - \delta$. Since $\lim_{n \rightarrow \infty} y_n = 0$, without loss of generality we assume $|y_{n+1} - y_n| \leq \omega \equiv \frac{1}{2}[g'(z - \delta) - 1]|x_{n_0+1} - x_{n_0}| \neq 0$ when $n \geq n_0$. Therefore we have for $n > n_0$,

$$\begin{aligned} |x_n - x_{n-1}| &\geq |g(x_{n-1}) - g(x_{n-2})| - |y_{n-1} - y_{n-2}| \\ &\geq g'(z - \delta)|x_{n-1} - x_{n-2}| - \omega \\ &\geq \dots\dots \\ &\geq g'(z - \delta)^{n-n_0-1}|x_{n_0+1} - x_{n_0}| - [\omega + g'(z - \delta)\omega + \dots + g'(z - \delta)^{n-n_0-2}\omega] \\ &= \omega[g'(z - \delta)^{n-n_0-1} - 1]/[g'(z - \delta) - 1]. \end{aligned}$$

Since $g'(z - \delta) > 1$, this shows that x_n is unbounded, contradicting the assumption that $x_n < z - \delta$ for any $n \geq n_0$.

Now we assume for some n_0 , $x_{n_0} \geq z - \delta$ and

$$|y_n| \leq \epsilon \equiv \min\{\eta/2, [1 - g'(z + \eta/2)](x^* - z - \eta/2)\} \text{ for any } n \geq n_0, \quad (\text{II.2})$$

where $g'(z + \eta/2) < g'(z) = 1$ and $x^* - z - \eta/2 > 0$ from (II.1). We show $x_n \geq z + \eta/2$ for any $n \geq n_0 + 1$ by induction. From (II.1) and (II.2), we have

$$x_{n_0+1} = g(x_{n_0}) + y_{n_0} \geq g(z - \delta) - \epsilon \geq (z + \eta) - \eta/2 = z + \eta/2.$$

Assuming that $x_{n-1} \geq z + \eta/2$ for some $n \geq n_0 + 2$, we need to show that $x_n \geq z + \eta/2$. If $x_{n-1} \geq x^*$, then from (II.2) and (II.1), we obtain

$$x_n = g(x_{n-1}) + y_{n-1} \geq g(x^*) - \epsilon \geq x^* - \eta/2 \geq z + \eta/2.$$

If $x_{n-1} < x^*$, then, since $x_{n-1} \geq z + \eta/2$ and $x^* > z + \eta/2$ (see (II.1)), we have with (II.2) that

$$\begin{aligned} |x^* - x_n| &= |g(x^*) - g(x_{n-1}) - y_{n-1}| \leq g'(z + \eta/2)(x^* - x_{n-1}) + \epsilon \\ &\leq g'(z + \eta/2)(x^* - z - \eta/2) + (1 - g'(z + \eta/2))(x^* - z - \eta/2) = x^* - z - \eta/2, \end{aligned}$$

which leads to $x_n \geq z + \eta/2$ again. Therefore for any $n \geq n_0 + 2$, we have

$$|x^* - x_n| = |g(x^*) - g(x_{n-1}) - y_{n-1}| \leq g'(z + \eta/2)|x^* - x_{n-1}| + |y_{n-1}|.$$

Since $g'(z + \eta/2) < 1$ and $y_{n-1} \rightarrow \infty$, it is easy to show that $\{x_n\}$ converges to x^* .

Case 2: $\lambda = 1$. Let h denote the inverse function of g , i.e., $h(x) = x/\sqrt{1 - \gamma^2 x^2}$, $x \geq 0$. One can verify that h' is monotonically increasing and $h'(x) > 1$ for $x \neq 0$. Suppose $\{x_n\}$ does not converge to 0, then, since $\{x_n\}$ is bounded, there exists a subsequence $\{x_{n_k}\}$ which converges to some $x^* \neq 0$. From $x_{n_k} = g(x_{n_k-1}) + y_{n_k-1}$ we have $x_{n_k-1} = h(x_{n_k} - y_{n_k-1})$. Thus

$$\lim_{k \rightarrow \infty} x_{n_k-1} = h(x^*) = h'(\zeta)x^*, \quad \text{where } 0 < \zeta < x^*.$$

Then we have

$$\lim_{k \rightarrow \infty} x_{n_k-2} = \lim_{k \rightarrow \infty} h(x_{n_k-1} - y_{n_k-2}) = h(h'(\zeta)x^*) \geq h'(\zeta)h(x^*) = h'(\zeta)^2 x^*.$$

Continuing the above steps, we arrive after j steps at

$$\lim_{k \rightarrow \infty} x_{n_k-j} \geq h'(\zeta)^j x^*.$$

Since $h'(\zeta) > 1$, this leads to a contradiction of the fact that $\{x_n\}$ is bounded. Thus $\{x_n\}$ converges to 0.

Acknowledgment

We wish to thank Chris Paige, Sanzheng Qiao and Gilbert Strang for their helpful suggestions and comments, and Gene Golub for pointing out the low rank distance of the approximate QR factorization to the underlying matrix A .

REFERENCES

1. Bauer FL. Ein direktes Iterationsverfahren zur Hurwitz-Zerlegung eines Polynoms. *Archiv der Elektrischen Übertragung, Stuttgart* 1995; 9:285-290.
2. Beam RM, Warming RF. An eigenvalue analysis of finite-difference approximations for hyperbolic IBVPs II: the auxiliary Dirichlet problem. In *Third International Conference on Hyperbolic Problems, Vol. II*, Engquist B, Gustafsson B (eds). Studentlitteratur, Lund, 1991; 923-937.
3. Beam RM, Warming RF. The asymptotic spectra of banded Toeplitz and quasi-Toeplitz matrices. *SIAM J. Sci. Comput.* 1993; 14:971-1006.
4. Bondeli S, Gander W. Cyclic reduction for special tridiagonal systems. *SIAM J. Matrix Anal. Appl.* 1994; 15:321-330.
5. Buzdin A. Tangential decomposition. *Computing* 1998; 61:257-276.
6. Gander MJ, Nataf F. AILU: A preconditioner based on the analytic factorization of the elliptic operator. *Numerical Linear Algebra with Applications* 2000; 7(7-8):543-567.
7. Gander MJ, Nataf F. AILU for Helmholtz problems: a new preconditioner based on the analytic parabolic factorization. *Journal of Computational Acoustics* 2001; 9(4):1499-1506.
8. Gander MJ. Preconditioners for acoustic problems: AILU. *Special issue of the journal Revista de Acústica*, Vol. XXXIII, 2002, [CDROM].
9. Fischer D, Golub GH, Hald O, Leiva C, Widlund O. *On Fourier-Toeplitz Methods for Separable Elliptic Problems*, Math. Comp. 1974; 28(126):349-368.
10. Golub GH, Van Loan CF. *Matrix Computations* (3rd ed). Johns Hopkins University Press, Baltimore, 1996.
11. Grenander A, Szegő G. *Toeplitz Forms and Their Applications*, University of California Press, Berkeley, 1958.
12. Kailath T, Sayed AH. *Fast Reliable Algorithms for Matrices with Structure*, SIAM, Philadelphia, PA, 1999.
13. Malcolm MA, Palmer J. A fast method for solving a class of tridiagonal linear system. *Comm. ACM* 1974; 17:14-17.

14. Paige CC, Saunders MA. LSQR: An algorithm for sparse linear equations and sparse least squares, *ACM Trans. Math. Soft.* 1982; **8**:43–71.
15. Rheinboldt WC. *Methods for Solving Systems of Nonlinear Equations* (2nd edn). SIAM, Philadelphia, PA, 1998.
16. Wagner C. Tangential frequency filtering decompositions for symmetric matrices. *Numer. Math.* 1997; **78**(1):119–142.
17. Wagner C. Tangential frequency filtering decompositions for unsymmetric matrices. *Numer. Math.* 1997; **78**(1):143–163.
18. Wilkinson JH. *The Algebraic Eigenvalue Problem*. Clarendon Press, Oxford, England, 1965.
19. Wittum G. An ILU-based smoothing correction scheme. *Parallel algorithms for partial differential equations*. Proc. 6th GAMM-Semin., Kiel/Ger., Notes Numer. Fluid Mech. 1991; **31**:228–240.
20. Wittum G. Filternde Zerlegungen. Schnelle Löser für grosse Gleichungssysteme, *Teubner Skripten zur Numerik*, Stuttgart, 1992.