

# Recycling Givens rotations for the efficient approximation of pseudospectra of band-dominated operators

MARKO LINDNER\* and TORGE SCHMIDT†

June 14, 2016

**ABSTRACT.** We study spectra and pseudospectra of certain bounded linear operators on  $\ell^2(\mathbb{Z})$ . The operators are generally non-normal, and their matrix representation has a characteristic off-diagonal decay. Based on a result of Chandler-Wilde, Chonchaiya and Lindner for tridiagonal infinite matrices, we demonstrate an efficient algorithm for the computation of upper and lower bounds on the pseudospectrum of operators that are merely norm limits of band matrices – the so-called band-dominated operators. After approximation by a band matrix and fixing a parameter  $n \in \mathbb{N}$ , one looks at  $n$  consecutive columns  $\{k+1, \dots, k+n\}$ ,  $k \in \mathbb{Z}$ , of the corresponding matrix and computes the smallest singular value of that section via QR factorization. We here propose a QR factorization by a sequence of Givens rotations in such a way that a large part of the computation can be reused for the factorization of the next submatrix – when  $k$  is replaced by  $k+1$ . The computational cost for the next factorization(s) is  $\mathcal{O}(nd)$  as opposed to a naive implementation with  $\mathcal{O}(nd^2)$ , where  $d$  is the bandwidth. So our algorithm pays off for large bands, which is attractive when approximating band-dominated operators with a full (i.e. not banded) matrix.

*Mathematics subject classification (2000):* Primary 65J10; Secondary 47A10, 47B36, 65F15.

*Keywords:* pseudospectrum, QR factorization, band-dominated operators

## 1 Introduction, Notations, and Main Results

**Band-dominated operators.** We study bounded linear operators on the space  $\ell^2 := \ell^2(\mathbb{Z})$  of square-summable bi-infinite complex sequences  $x = (x_k)_{k \in \mathbb{Z}}$  with  $\|x\| = \sqrt{\sum_{k \in \mathbb{Z}} |x_k|^2} < \infty$ . Each linear operator  $A$  on  $\ell^2$  acts via matrix-vector multiplication with a bi-infinite matrix  $(a_{ij})_{i,j \in \mathbb{Z}}$  – and vice versa. We say that  $A$  is a *band operator* if its matrix  $(a_{ij})$  is banded (i.e. supported on only finitely many diagonals) and has uniformly bounded entries, so that  $A$  is a bounded linear operator. In that case,  $d := \max\{|i-j| : a_{ij} \neq 0\}$  is called the *bandwidth* of  $A$ . Moreover,  $A$  is called a *band-dominated operator* if it is the limit, in the induced operator norm on  $\ell^2$ , of a sequence of band operators; in particular it is a bounded operator, too, and its matrix entries decay with their distance from the main diagonal.

**Pseudospectra.** Because the spectrum of a non-normal operator  $A$  can be highly unstable under small perturbations of  $A$ , one is interested in the so-called  $\varepsilon$ -pseudospectrum of  $A$ , that is,

$$\text{spec}_\varepsilon A := \{\lambda \in \mathbb{C} : \|(A - \lambda I)^{-1}\| > 1/\varepsilon\} = \bigcup_{\|T\| < \varepsilon} \text{spec}(A + T), \quad \varepsilon > 0.$$

Here we agree upon writing  $\|B^{-1}\| = \infty$  if  $B$  is not invertible. The second equality sign (see e.g. [28]) shows that  $\text{spec}_\varepsilon A$  measures the sensitivity of  $\text{spec} A$  w.r.t. additive perturbations of  $A$

---

\*Email: lindner@tuhh.de

†Email: torge.schmidt@tuhh.de

of norm  $< \varepsilon$ . For normal operators  $A$ ,  $\text{spec}_\varepsilon A$  is the  $\varepsilon$ -neighbourhood of  $\text{spec} A$ ; otherwise it is generally larger (but never smaller). The interest in pseudospectra has been increasing over the last two decades. See [28] for many more reasons to study pseudospectra and for more references.

**The lower norm.** As a counterpart to the operator norm  $\|A\| = \sup_{\|x\|=1} \|Ax\|$ , we look at the quantity

$$\nu(A) := \inf_{\|x\|=1} \|Ax\|,$$

that is sometimes (by abuse of notation) called the *lower norm* of  $A$ . While  $\|A\|$  is the largest singular value of  $A$ ,  $\nu(A)$  is the smallest – provided maximum/minimum exist, such as in the case of finite matrices. It is well-known (see e.g. [19, p.69f]) that  $\nu(A) > 0$  holds iff  $A$  is injective and has a closed image; moreover, the equality

$$\|A^{-1}\| = 1/\min(\nu(A), \nu(A^*))$$

holds with  $1/0 := \infty$  indicating non-invertibility of  $A$ . In particular,  $A$  is invertible iff  $\nu(A)$  and  $\nu(A^*)$  are both nonzero, in which case they coincide. Together with the definition of  $\text{spec}_\varepsilon A$  it follows that

$$\text{spec}_\varepsilon A = \{\lambda \in \mathbb{C} : \min(\nu(A - \lambda I), \nu((A - \lambda I)^*)) < \varepsilon\}. \quad (1)$$

**Approximating the lower norm of band-dominated operators.** For  $x \in \ell^2$ , we denote its support by  $\text{supp } x := \{j \in \mathbb{Z} : x_j \neq 0\}$ , and we say that a bounded set  $J \subset \mathbb{Z}$  has diameter  $\text{diam } J := \max\{|i - j| : i, j \in J\}$ . One of the main observations of [11] (also see [13, 21]) is that the lower norm<sup>1</sup> of a band-dominated operator  $A$  can be realized, up to a given  $\delta > 0$ , by a unit element  $x \in \ell^2$  with bounded support, say of diameter less than  $n \in \mathbb{N}$  (dependent on  $\delta$ , of course). So one has

$$\nu(A) \leq \|Ax\| \leq \nu(A) + \delta \quad (2)$$

for a particular  $x \in \ell^2$  with  $\|x\| = 1$  and  $\text{diam}(\text{supp } x) < n$ . If  $\text{supp } x$  were known to be contained in the discrete interval  $J_k^n := \{k + 1, \dots, k + n\}$  with a given  $k \in \mathbb{Z}$ , then the optimal term  $\|Ax\|$  in (2) could be practically computed as the lower norm / smallest singular value of the restriction of  $A$  to  $\ell^2(J_k^n)$ . Since  $\text{diam}(\text{supp } x) < n$ , the support must be contained in some interval  $J_k^n$  with  $k \in \mathbb{Z}$ . Unfortunately, this  $k$  is in general not known. It “remains” to look at – and minimize over – all  $k \in \mathbb{Z}$ :

$$\nu(A) \leq \inf_{k \in \mathbb{Z}} \nu(A|_{\ell^2(J_k^n)}) \leq \nu(A) + \delta \quad (3)$$

If  $A$  is a band operator then  $A|_{\ell^2(J_k^n)}$  corresponds to a finite rectangular matrix (containing columns  $k + 1, \dots, k + n$  of the infinite matrix, truncated to their joint support that is finite – due to the band structure), so that the smallest singular value,  $\nu(A|_{\ell^2(J_k^n)})$ , can be computed effectively. However, consideration of all  $k \in \mathbb{Z}$  is, in general, of course practically impossible – unless the set of all restrictions  $\{A|_{\ell^2(J_k^n)} : k \in \mathbb{Z}\}$  is finite, e.g. when  $A$  is eventually periodic or otherwise structured.

It is clear that the size  $n$  has to be increased in order to decrease the error  $\delta$  in (2) and (3). The analysis in [11] (also see §3 and 4 in [13]) shows, for the particular case of tridiagonal (bandwidth  $d = 1$ ) bi-infinite matrices  $(a_{ij})_{i,j \in \mathbb{Z}}$ , that  $\delta$  is of the order  $1/n$ ; more precisely,

$$\delta \leq 2 \left( \sup_{j \in \mathbb{Z}} |a_{j+1,j}| + \sup_{j \in \mathbb{Z}} |a_{j-1,j}| \right) \sin \frac{\pi}{2n+2} \in \mathcal{O} \left( \frac{1}{n} \right), \quad (4)$$

The constant turns out to be optimal. We make use of that result by two simple steps of reduction:

- (i) Given an accuracy  $\eta > 0$ , approximate our band-dominated operator  $A$  (with a generally full matrix) by a band operator  $B$  with  $\|A - B\| \leq \eta$  and use the contractivity of  $\nu(\cdot)$ , so that

$$|\nu(A) - \nu(B)| \leq \|A - B\| \leq \eta, \quad \text{as well as} \quad |\nu(A^*) - \nu(B^*)| \leq \|A^* - B^*\| \leq \eta. \quad (5)$$

---

<sup>1</sup>A symmetric result holds for the norm,  $\|A\|$ , see Proposition 3.4 and inequality (ONL) in [16].

- (ii) Use that the matrix of the band operator  $B$  is block-tridiagonal (with block size equal to the band width of  $B$ , see Figure 1) and that the results of [11, 13] even apply to tridiagonal matrices with operator entries<sup>2</sup> – hence to block-tridiagonal matrices.

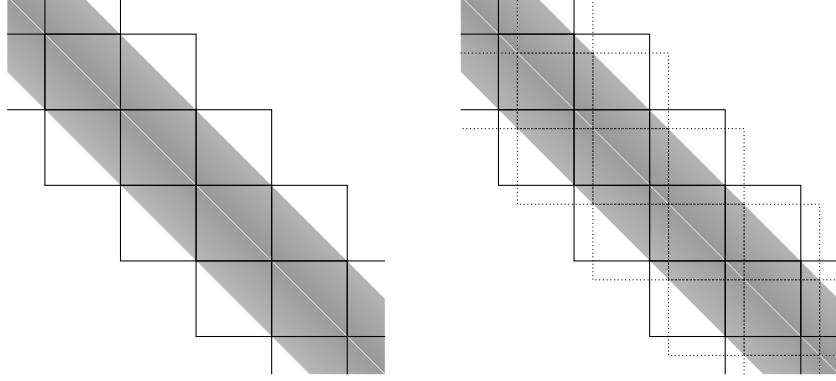


Figure 1.1: Left: A banded matrix (support shown in gray) is turned into block-tridiagonal form with blocks of according size. Right: The dotted blocks equally do the job of turning the banded matrix into block-tridiagonal form. There are  $b$  different ways of positioning a  $b \times b$  grid along the main diagonal. Two of them are depicted here (solid and dotted lines).

We discuss further details of steps (i) and (ii) in Section 2.

**Approximating pseudospectra of band-dominated operators.** From (1) and the above approximations and bounds on the lower norm we conclude approximations and bounds on the pseudospectrum:

Inequality (3) and its counterpart for the adjoint,  $A^*$ , lead to

$$\min(\nu(A), \nu(A^*)) \leq \inf_{k \in \mathbb{Z}} \min(\nu(A|_{\ell^2(J_k^n)}), \nu(A^*|_{\ell^2(J_k^n)})) \leq \min(\nu(A), \nu(A^*)) + \delta,$$

from which we conclude the implications

$$\begin{aligned} \inf_{k \in \mathbb{Z}} \min(\nu(A|_{\ell^2(J_k^n)}), \nu(A^*|_{\ell^2(J_k^n)})) < \varepsilon &\Rightarrow \min(\nu(A), \nu(A^*)) < \varepsilon \\ &\Rightarrow \inf_{k \in \mathbb{Z}} \min(\nu(A|_{\ell^2(J_k^n)}), \nu(A^*|_{\ell^2(J_k^n)})) < \varepsilon + \delta \end{aligned}$$

for all  $\varepsilon > 0$ , and consequently

$$\Gamma_\varepsilon^n(A) \subset \text{spec}_\varepsilon A \subset \Gamma_{\varepsilon+\delta}^n(A), \quad (6)$$

where

$$\Gamma_\varepsilon^n(A) := \bigcup_{k \in \mathbb{Z}} \left\{ \lambda \in \mathbb{C} : \min(\nu((A - \lambda I)|_{\ell^2(J_k^n)}), \nu((A - \lambda I)^*|_{\ell^2(J_k^n)})) < \varepsilon \right\}. \quad (7)$$

Concerning the approximation step (i) above, by (5), we have the implications

$$\nu(B) < \varepsilon - \eta \quad \Rightarrow \quad \nu(A) < \varepsilon \quad \Rightarrow \quad \nu(B) < \varepsilon + \eta,$$

and the same holds for the adjoints. Subtracting  $\lambda I$  from  $A$  and  $B$  and using (1), this shows that

$$\text{spec}_{\varepsilon-\eta} B \subset \text{spec}_\varepsilon A \subset \text{spec}_{\varepsilon+\eta} B, \quad 0 < \eta < \varepsilon, \quad (8)$$

so that upper and lower bounds on certain pseudospectra of  $B$  yield bounds on  $\text{spec}_\varepsilon A$ . Moreover, the inclusions (8) are as tight as desired (in the Hausdorff distance) by sending  $\eta \rightarrow 0$ .

<sup>2</sup>In that case,  $|a_{j+1,j}|$  and  $|a_{j-1,j}|$  in (4) are interpreted as operator norms.

**Existing results.** The probably most natural idea to approximate  $\text{spec}_\varepsilon A$  is to look at the pseudospectra  $\text{spec}_\varepsilon A_n$  of the finite sections  $A_n = (a_{ij})_{i,j=-n}^n$  of  $A = (a_{ij})_{i,j \in \mathbb{Z}}$  as  $n \rightarrow \infty$ . In some rare cases (Toeplitz operators [24, 8], random Jacobi operators [12]), the sets  $\text{spec}_\varepsilon A_n$  indeed converge to  $\text{spec}_\varepsilon A$  w.r.t. the Hausdorff distance – but in general, the sequence  $\text{spec}_\varepsilon A_n$  does not converge at all; its cluster points usually contain  $\text{spec}_\varepsilon A$  but also further points (see e.g. [27], one speaks of spectral pollution). Even in a simple selfadjoint example such as  $A = \text{diag}(\dots, B, B, B, \dots)$  with  $B = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}$ , one has<sup>3</sup>  $\text{spec} A = \{-1, 1\}$ , while  $\text{spec} A_n$  repeatedly switches between  $\{-1, 1\}$  and  $\{-1, 0, 1\}$  as  $n$  grows. As an alternative that is somewhere between spectra and pseudospectra, [17, 26] study so-called  $(N, \varepsilon)$ -pseudospectra, where  $2^N$ -th powers of the resolvent and of  $1/\varepsilon$  are compared to each other. In [5] the lower norms of rectangular submatrices are suggested for the approximation of the spectrum and the  $(N, \varepsilon)$ -pseudospectrum. Needless to say, there is a large amount of literature on the selfadjoint case (see e.g. [1, 14] and the references therein).

One major problem in approximating  $\mathbb{Z}$  by the intervals  $\{-n, \dots, n\}$  is (besides the potential of spectral pollution) that generally, huge values of  $n$  are required to capture spectral properties of  $A$  properly. (Think of an infinite diagonal matrix with distinguished entries in very remote locations.) From a computational perspective, such huge sections  $\{-n, \dots, n\}$  are too expensive. The approach of [11] (also see §3 and 4 in [13]) – that is very much in the spirit of Gershgorin and that we adopt here – replaces  $\{-n, \dots, n\}$  with  $n \rightarrow \infty$  by the family of intervals  $J_k^n = \{k+1, \dots, k+n\}$  for all  $k \in \mathbb{Z}$  but with  $n$  of moderate size. The price that is obviously paid is the infinite amount of positions  $k$  that one has to look at, so that a certain structural simplicity of the infinite matrix is required to make the approach practically feasible. The other major plus of the [11] approach is that it comes with sharp and explicit bounds (4) on the accuracy of the approximation (6), while working for the general non-normal case.

**What is new here?** The tridiagonal results and the ideas of transferring them to band-dominated operators via (i) and (ii) are from [11], therefore not new. But there are two degrees of freedom in the choice of the blocks in step (ii): Firstly, the size of the blocks, say  $b \in \mathbb{N}$ , could be any number greater than or equal to the bandwidth  $d \in \mathbb{N}$ . Secondly, once this size  $b$  is fixed, there are  $b$  different choices for the position of the blocks inside the infinite matrix (see Figure 1.1).

We play with that second degree of freedom, arguing that there is usually no best choice (in terms of sharpness of (6)) of block positioning, and instead we consider all  $b$  possibilities, thereby improving sharpness of the bounds on  $\text{spec}_\varepsilon A$ . (We take the union of the  $b$  different lower bounds and the intersection of the  $b$  upper bounds.) Naively implemented, this increases the computational cost by the factor  $b$ . However, we present an algorithm that compensates for this increase by reusing much of the effort that was put into the computation of  $\nu(A|_{\ell^2(J_k^n)})$  for the computation of  $\nu(A|_{\ell^2(J_{k+1}^n)})$ . This is possible due to the large overlap between the two matrices  $A|_{\ell^2(J_k^n)}$  and  $A|_{\ell^2(J_{k+1}^n)}$ . We cannot see a similar idea to work for the  $b$ -sized step from  $\nu(A|_{\ell^2(J_k^n)})$  to  $\nu(A|_{\ell^2(J_{k+b}^n)})$  in the block matrix, though.

In a nutshell, the smallest singular value of  $A|_{\ell^2(J_k^n)}$  coincides with that of the upper triangular matrix<sup>4</sup>  $R_k$  from the factorization  $A|_{\ell^2(J_k^n)} = Q_k R_k$  with a unitary  $Q_k$  that results from a sequence of Givens rotations. The key idea is now to rearrange and reuse most of these Givens rotations for the next step when  $k$  is replaced by  $k+1$ . With this algorithm, the complexity of the computation of  $\nu(A|_{\ell^2(J_{k+1}^n)}) = \nu(R_{k+1})$  decreases from  $\mathcal{O}(nd^2)$  to just  $\mathcal{O}(nd)$ , thereby compensating for the increase by a factor of  $b \approx d$  that was mentioned above. The same recycling idea and the same complexity then also apply to the computation of  $\nu(A|_{\ell^2(J_{k+2}^n)})$ ,  $\nu(A|_{\ell^2(J_{k+3}^n)})$ , etc.

**Contents of the paper.** In Section 2 we show the details of both reduction steps (i) and (ii). The heart of the paper is Section 3, where we present the algorithm for the computation of  $\nu(A|_{\ell^2(J_{k+1}^n)})$  from  $\nu(A|_{\ell^2(J_k^n)})$  by appropriately reordering Givens rotations. In Section 4 we

<sup>3</sup>The  $\varepsilon$ -pseudospectra are the  $\varepsilon$ -neighbourhoods of the spectra in this selfadjoint example.

<sup>4</sup>For the computation of the smallest singular value of  $R_k$ , one can use an inverse Lanczos method.

illustrate our results in two examples with non-trivial pseudospectra. Moreover, we compare the efficiency of our algorithm with the standard QR decomposition in each step.

## 2 From band-dominated to tridiagonal operators

Recall that we call an operator  $A$  on  $\ell^2$  band-dominated if it is the limit, in the operator norm, of a sequence of band operators (which are bounded operators with a banded matrix representation). Let us denote the sets of all band and all band-dominated operators on  $\ell^2$  by BO and BDO, respectively. We make use of the results from [11, 13] for tridiagonal operators by two steps of reduction:

### 2.1 Step (i): From band-dominated to banded

Let  $A \in \text{BDO}$  and  $\eta > 0$  be given. There are different approaches of constructing a band operator  $B$  with  $\|A - B\| \leq \eta$ , leading to (5) and (8):

**Case 1.** If  $A$  is in the so-called Wiener algebra, the problem is simple. To explain this, let  $(b_{ij})_{i,j \in \mathbb{Z}}$  be the matrix representation of some  $B \in \text{BO}$  and let  $d_k := (b_{j+k,j})_{j \in \mathbb{Z}}$  be its  $k$ -th diagonal, where  $k \in \mathbb{Z}$ . Then

$$B = \sum_{k \in \mathbb{Z}} M_{d_k} V_k,$$

where  $M_f$  refers to the operator on  $\ell^2$  of entrywise multiplication with a sequence  $f \in \ell^\infty$  and  $V_k$  is the forward shift on  $\ell^2$  by  $k$  positions. (Note that the sum is actually finite, by  $B \in \text{BO}$ .) It now follows that

$$\|B\| = \left\| \sum_{k \in \mathbb{Z}} M_{d_k} V_k \right\| \leq \sum_{k \in \mathbb{Z}} \|M_{d_k}\| \|V_k\| = \sum_{k \in \mathbb{Z}} \|d_k\|_\infty =: \llbracket B \rrbracket. \quad (9)$$

The new expression  $\llbracket \cdot \rrbracket$  indeed defines a norm on BO. The completion of BO with respect to  $\llbracket \cdot \rrbracket$  is the so-called *Wiener algebra*  $\mathcal{W}$ . By (9),  $\mathcal{W}$  is contained in the completion of BO w.r.t.  $\|\cdot\|$ , that is BDO. Moreover,  $(\mathcal{W}, \llbracket \cdot \rrbracket)$  is a Banach algebra (see §1.6.8 in [18] or §3.7.3 in [20]).

So if  $A \in \mathcal{W} \subset \text{BDO}$  and  $d_k$  refers to its  $k$ -th diagonal for all  $k \in \mathbb{Z}$ , then

$$B_n := \sum_{k=-n}^n M_{d_k} V_k \in \text{BO} \quad (10)$$

is the desired approximation of  $A$  if  $n \in \mathbb{N}$  is chosen large enough for

$$\|A - B_n\| \leq \llbracket A - B_n \rrbracket = \sum_{|k| > n} \|d_k\|_\infty \leq \eta. \quad (11)$$

Such an  $n$  exists since  $\sum_{n \in \mathbb{Z}} \|d_k\|_\infty < \infty$ , by  $A \in \mathcal{W}$ .

**Case 2.** If  $A \in \text{BDO} \setminus \mathcal{W}$ , the simple approach (10) of restriction to a finite subset of diagonals need not lead to a sequence  $B_n$  that converges to  $A$  in the operator norm. A simple example is shown in Remark 1.40 of [19]. The example relies on the fact that, for a continuous  $2\pi$ -periodic function  $f$  on  $\mathbb{R}$ , the partial sums of the Fourier series need not converge uniformly to  $f$ . This is repaired by looking at Fejer-Cesaro means instead, and the same trick works for the approximation of band-dominated operators:

$$C_n := \frac{B_0 + \dots + B_n}{n+1} = \sum_{k=-n}^n \left(1 - \frac{|k|}{n+1}\right) M_{d_k} V_k \in \text{BO} \quad (12)$$

with  $B_n$  from (10) can be shown to converge to  $A$  in the operator norm as  $n \rightarrow \infty$ , see e.g. the proof of the implication (e)  $\Rightarrow$  (a) in Theorem 2.1.6 of [23].

Another way to explicitly approximate  $A \in \text{BDO}$  by band operators is shown in (1.18) of [25].

## 2.2 Step (ii): From banded to tridiagonal

Now we can assume  $A \in \text{BO}$ . Let  $d$  denote its bandwidth. The idea is captured by Figure 1.1 above:  $A$  can be expressed as a block-tridiagonal matrix with block size  $b \geq d$ . Besides the choice of  $b$ , there is another degree of freedom in this identification. If the blocks are centered on the main diagonal, there are still  $b$  different positions at which to start, see Figure 1.1.

Precisely, each block is of the form

$$\begin{pmatrix} a_{i+1,j+1} & \cdots & a_{i+1,j+b} \\ \vdots & & \vdots \\ a_{i+b,j+1} & \cdots & a_{i+b,j+b} \end{pmatrix} \in \mathbb{C}^{b \times b} \quad \text{with} \quad i, j \in c + b\mathbb{Z} := \{c + bz : z \in \mathbb{Z}\}, \quad (13)$$

where  $c \in \{0, \dots, b-1\}$  is this second degree of freedom. This leads to  $b$  different ways (one for each choice of the offset  $c$ ) of turning  $A$  into a tridiagonal matrix.

For the moment, fix one choice of  $c \in \{0, \dots, b-1\}$ . To apply the results on the block-tridiagonal matrix behind  $A$ , we have to adjust the intervals  $J_k^n := \{k+1, \dots, k+n\}$  (of matrix columns under current investigation in (3)) with the blocks. Therefore, we restrict ourselves to positions  $k \in c + b\mathbb{Z}$  and to interval lengths  $n = Nb$ , where  $N \in \mathbb{N}$  is the number of blocks to be considered in  $J_k^n$ .

Now we slightly modify (7) to

$$\Gamma_\varepsilon^{n,M}(A) := \bigcup_{k \in M} \left\{ \lambda \in \mathbb{C} : \min(\nu((A - \lambda I)|_{\ell^2(J_k^n)}), \nu((A - \lambda I)^*|_{\ell^2(J_k^n)})) < \varepsilon \right\} \quad (14)$$

for any set  $M \subset \mathbb{Z}$ , where our particular interest is in sets of the form  $M = c + b\mathbb{Z}$ . Assuming  $b$  as given and fixed, we abbreviate  $\Gamma_\varepsilon^{n,c+b\mathbb{Z}}(A) =: \Gamma_\varepsilon^{n,c}(A)$ .

Because each offset  $c \in \{0, \dots, b-1\}$  yields a tridiagonal representation of  $A$ , we get from (6) that all inclusions

$$\left. \begin{array}{l} \Gamma_\varepsilon^{n,0}(A) \subset \text{spec}_\varepsilon A \subset \Gamma_{\varepsilon+\delta}^{n,0}(A) \\ \Gamma_\varepsilon^{n,1}(A) \subset \text{spec}_\varepsilon A \subset \Gamma_{\varepsilon+\delta}^{n,1}(A) \\ \vdots \\ \Gamma_\varepsilon^{n,b-1}(A) \subset \text{spec}_\varepsilon A \subset \Gamma_{\varepsilon+\delta}^{n,b-1}(A) \end{array} \right\} \quad (15)$$

hold. Here, by evaluating (4) for the block tridiagonal matrix,

$$\delta \leq 2 \left( \sup_{l \in \mathbb{Z}} \|A_{l+1,l}\| + \sup_{l \in \mathbb{Z}} \|A_{l-1,l}\| \right) \sin \frac{\pi}{2N+2} \in \mathcal{O}\left(\frac{1}{N}\right) = \mathcal{O}\left(\frac{1}{n}\right), \quad (16)$$

where we denote the block (13) by  $A_{kl}$  if  $i = c + bk$  and  $j = c + bl$  with  $k, l \in \mathbb{Z}$ .

Taking unions on the left and intersections on the right of (15), we conclude

$$\Gamma_\varepsilon^{n,0}(A) \cup \dots \cup \Gamma_\varepsilon^{n,b-1}(A) \subset \text{spec}_\varepsilon A \subset \Gamma_{\varepsilon+\delta}^{n,0}(A) \cap \dots \cap \Gamma_{\varepsilon+\delta}^{n,b-1}(A). \quad (17)$$

In examples one observes that the bound (17) on  $\text{spec}_\varepsilon A$  is sharper than any of (15).

**Example 2.1** We look at the following 2-periodic bi-infinite matrix with bandwidth  $d = 2$ :

$$A = \left( \begin{array}{c|c|c|c|c} \dots & \dots & \dots & & \\ \dots & 0 & 9 & 4 & \\ \dots & 9 & 0 & 2 & 0 \\ \hline & 0 & 2 & 0 & 9 & 4 \\ & & 0 & 9 & 0 & 2 & \dots \\ \hline & & & 0 & 2 & 0 & \dots & \dots \\ & & & & \dots & \dots & \dots & \dots \end{array} \right) = \left( \begin{array}{c|c|c|c|c|c} \dots & \dots & \dots & & & \\ \dots & 0 & 9 & 4 & & \\ \dots & 9 & 0 & 2 & 0 & \\ \hline & 0 & 2 & 0 & 9 & 4 \\ & & 0 & 9 & 0 & 2 & \dots \\ \hline & & & 0 & 2 & 0 & \dots & \dots \\ & & & & \dots & \dots & \dots & \dots \end{array} \right)$$

The block size was chosen to be  $b = d = 2$ , leading to the two different possibilities of block positioning ( $c = 0$  and  $c = 1$ ) shown above. Figure 2.1 below shows a plot of  $\Gamma_\varepsilon^{n,0}(A)$  and, for comparison, of  $\Gamma_\varepsilon^{n,1}(A)$ , as well as  $\Gamma_\varepsilon^{n,0}(A) \cap \Gamma_\varepsilon^{n,1}(A)$  for  $n = 6$ .  $\square$

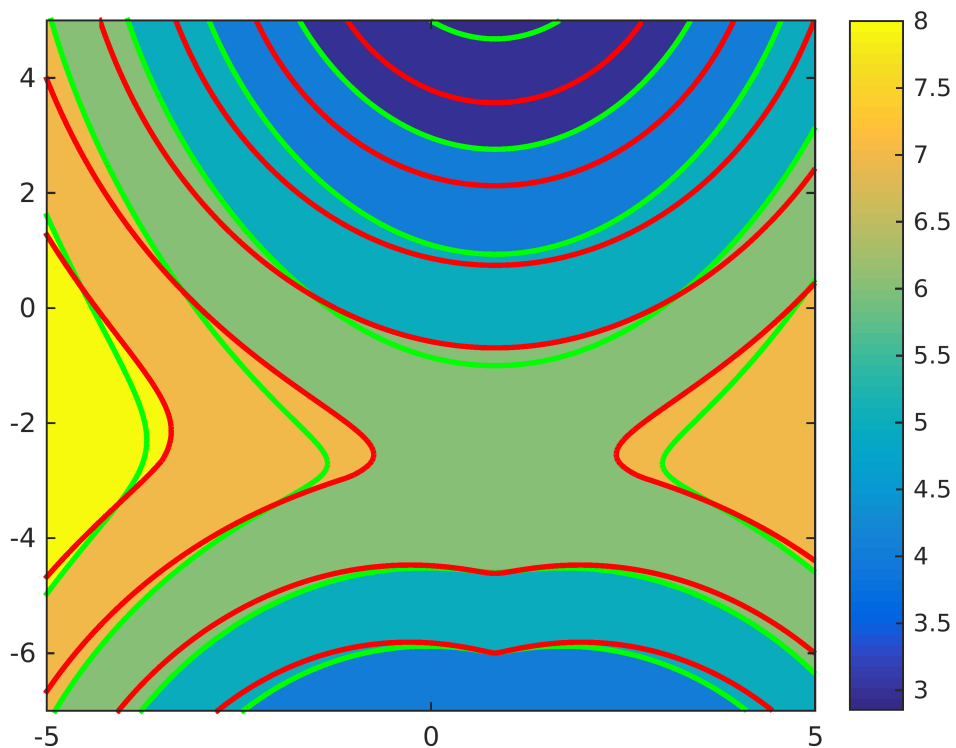


Figure 2.1: Regarding Example 2.1, we see the boundaries of the sets  $\Gamma_\varepsilon^{n,0}(A)$  (dark/red line) and, for comparison, of  $\Gamma_\varepsilon^{n,1}(A)$  (light/green line), both for  $n = 6$  and  $\varepsilon = 2, 3, \dots, 8$ . The colored areas denote  $\Gamma_\varepsilon^{n,0}(A) \cap \Gamma_\varepsilon^{n,1}(A)$ .

This is why we suggest to look at all (instead of just one) of the inclusions (15). Of course this improvement in quality of the bound on  $\text{spec}_\varepsilon A$  increases the numerical costs by a factor of  $b$ . The next section shows how to compensate for that.

### 3 The Algorithm

To simplify notation abbreviate, for  $k \in \mathbb{Z}, n \in \mathbb{N}$  and  $\lambda \in \mathbb{C}$ ,

$$A_\lambda^k := (A - \lambda I)|_{\ell^2(J_k^n)} : \begin{array}{c} \ell^2(J_k^n) \\ \cong \\ \mathbb{C}^n \end{array} \rightarrow \begin{array}{c} \ell^2(J_{k-d}^{n+2d}) \\ \cong \\ \mathbb{C}^{n+2d} \end{array}$$

and treat  $A_\lambda^k$  as a finite rectangular matrix. We define  $\bar{A}_\lambda^k := (A - \lambda I)^*|_{\ell^2(J_k^n)}$  analogously.

As has been described in the previous section, we need to approximate  $\nu(A_\lambda^k)$  and  $\nu(\bar{A}_\lambda^k)$  for different values  $\lambda \in \mathbb{C}$  and multiple consecutive values of  $k$ . This can be done by computing the smallest singular values  $\sigma_n(A_\lambda^k)$  and  $\sigma_n(\bar{A}_\lambda^k)$  which is strongly related to pseudospectra of rectangular matrices ([31]) and similar computational problems arise.

If the considered matrices  $A_\lambda^k$  were square, we could compute the Schur decomposition of  $A_0^k$  – thus transforming  $A_0^k$  into upper rectangular form – while preserving the shift by  $\lambda$ . Afterwards we could compute  $\sigma_n(A_\lambda^k)$  for multiple values of  $\lambda \in \mathbb{C}$  using a bidiagonalization method [9] on  $(A_\lambda^k)^{-1}$ . In the rectangular case though, no shift-preserving method to reduce  $A_0^k$  to a simple form appears to be known and an inverse iteration is more difficult to implement for rectangular matrices.

We can however use the fact, that for each  $\lambda \in \mathbb{C}$  we have two sequences  $(A_\lambda^k)_k$  and  $(\bar{A}_\lambda^k)_k$  each of which contains large overlaps between consecutive matrices. We will introduce an algorithm that takes advantage of this property.

We fix  $\lambda \in \mathbb{C}$  and  $n \in \mathbb{N}$  and abbreviate  $A^k := A_\lambda^k \in \mathbb{C}^{(n+2d) \times n}$  for  $k \in \mathbb{Z}$ .

Let  $A^{k_0+1}, A^{k_0+2}, \dots, A^{k_0+k_{\max}}$  be a finite sequence of matrices given by (14). W.l.o.g. we consider  $k_0 = 0$ . We can describe the overlapping property of these matrices by

$$A_{i,j}^k = A_{i-1,j-1}^{k+1}, \text{ for all } \begin{cases} 1 \leq k \leq k_{\max} - 1 \\ 2 \leq i \leq n + 2d =: m \\ 2 \leq j \leq n. \end{cases} \quad (18)$$

Since  $\nu(A^k) = \sigma_n(A^k)$ , we are interested in computing the set

$$\{\sigma_n(A^k)\}_{1 \leq k \leq k_{\max}},$$

where  $\sigma_n$  denotes the smallest singular value, which can be approximated using a QR decomposition

$$Q^k A^k = R^k = \begin{pmatrix} \tilde{R}^k \\ \mathbf{0} \end{pmatrix}, \text{ with } \tilde{R}^k \in \mathbb{C}^{n \times n} \text{ upper triangular, } Q^k \in \mathbb{C}^{m \times m} \text{ unitary} \quad (19)$$

and applying an inverse Golub-Kahan-Lanczos-Bidiagonalization method ([2, 15]), from now on abbreviated as GKLB method, to  $\tilde{R}^k$  (i.e. applying the GKLB method to  $(\tilde{R}^k)^{-1}$ ). Since this is a unitary transformation, the singular values of  $A^k$  and  $\tilde{R}^k$  are the same. The inverse GKLB method requires solving two linear systems of equations in each iteration which can be achieved using backward-substitution, since  $\tilde{R}^k$  is upper triangular.

Note that unlike convention we write  $Q^k A^k = R^k$  instead of  $(Q^k)^H A^k = R^k$  to simplify notation. It is possible to compute a QR decomposition such that the banded structure of  $A^k$  is preserved in  $\tilde{R}^k$ , i.e.  $\tilde{R}^k$  has at most  $2d + 1$  consecutive non-zero diagonals. Therefore solving a linear system of equations involving  $\tilde{R}^k$  requires only  $\mathcal{O}(nd)$  flops. The QR decomposition (19) itself however requires  $\mathcal{O}(nd^2)$  operations and is therefore the bottleneck of the algorithm for large  $d$ .

This bottleneck is addressed in the QH-shift-algorithm which we will develop in this section. The idea of the algorithm is to use Givens rotations to compute the factorization  $Q^1 A^1 = H^1$ , where  $H^1 \in \mathbb{C}^{m \times n}$  is an upper Hessenberg-matrix<sup>5</sup> with  $2d + 1$  consecutive non-zero diagonals, and then

<sup>5</sup>We say a matrix  $H \in \mathbb{C}^{m \times n}$  is an upper Hessenberg-matrix if  $H_{i,j} = 0$  for all  $i > j + 1$



reuse these rotations to factorize  $A^2, A^3, \dots$  the same way.

Having factorized  $A^1, A^2, \dots$  into Hessenberg form using unitary transformations, we only need to apply  $n$  additional Givens rotations to each matrix to arrive at the QR decomposition (19). The total effort for each QR decomposition of  $A^2, A^3, \dots$  is only  $\mathcal{O}(nd)$  instead of  $\mathcal{O}(nd^2)$ .

**Preliminaries.** We will only use Givens rotations acting on consecutive rows and define a rotation on the  $i$ th and  $(i+1)$ st row by the mapping

$$G_i : \mathbb{D} \times \mathbb{D} \rightarrow \mathbb{C}^{m \times m} \\ (c, s) \mapsto G_i(c, s).$$

and

$$G_i(c, s) = \begin{matrix} & & & i & i+1 & & \\ & & & & & & \\ & & & & & & \\ & & & & & & \\ i & & & & & & \\ i+1 & & & & & & \end{matrix} \begin{pmatrix} 1 & \cdots & 0 & 0 & \cdots & 0 \\ \vdots & \ddots & \vdots & \vdots & & \vdots \\ 0 & \cdots & c & s & \cdots & 0 \\ 0 & \cdots & -\bar{s} & \bar{c} & \cdots & 0 \\ \vdots & & \vdots & \vdots & \ddots & \vdots \\ 0 & \cdots & 0 & 0 & \cdots & 1 \end{pmatrix}$$

where  $\mathbb{D} := \{z \in \mathbb{C} : |z| \leq 1\}$  is the closed complex unit disc. Details on the choice of  $c, s$  can be found in standard literature [15, 29]. To simplify notation we will, in most cases, write  $G_i \equiv G_i(c, s)$ , if the choice of  $c, s$  is clear from the context. This naturally leads to the problem of possibly having multiple rotations on the same row, each having different entries  $c, s$  and we hope that it will be clear from the context that these Givens rotations are not the same.

In the interest of readability we will mainly use the arrow-notation introduced by Raf Vandebril et al. in [29, 30]:

$$\begin{array}{c|c} 1 & \\ 2 & \curvearrowright \\ 3 & \curvearrowright \curvearrowright \\ 4 & \curvearrowright \curvearrowright \curvearrowright \\ 5 & \curvearrowright \curvearrowright \curvearrowright \curvearrowright \\ \hline & 4\ 3\ 2\ 1 \end{array} \quad (20)$$

The arrows in (20) each depict a Givens rotation operation, acting on the two rows in which the arrow is drawn (see axis of ordinates). The order of application of these rotations is described in the abscissa, i.e. from right to left, so that (20) represents the product  $G_4 G_3 G_2 G_1$ . It is important to note the order of application of the Givens rotations, since they do not commute in general unless they act on disjoint couples of rows:

$$\begin{matrix} \curvearrowright \\ \curvearrowright \end{matrix} \begin{bmatrix} \times & \times \\ \times & \times \\ \times & \times \\ \times & \times \end{bmatrix} = \begin{matrix} \curvearrowright \\ \curvearrowright \end{matrix} \begin{bmatrix} \times & \times \\ \times & \times \\ \times & \times \\ \times & \times \end{bmatrix}, \text{ but } \begin{matrix} \curvearrowright \\ \curvearrowright \end{matrix} \begin{bmatrix} \times & \times \\ \times & \times \\ \times & \times \\ \times & \times \end{bmatrix} \neq \begin{matrix} \curvearrowright \\ \curvearrowright \end{matrix} \begin{bmatrix} \times & \times \\ \times & \times \\ \times & \times \\ \times & \times \end{bmatrix} \quad (21)$$

We say that a product  $G_{i_1}, G_{i_2}, \dots, G_{i_l}$  is a *descending*, respectively *ascending*, *sequence of Givens rotations* of length  $l$ , if  $i_{p+1} = i_p - 1$ , respectively  $i_{p+1} = i_p + 1$ , for  $p = 1, \dots, l-1$ . (20) is an example of a descending sequence of length 4.



Notice that, since the right hand side is of Hessenberg form, no rotations acting on the first row are required and the first row of  $Q^1$  is the unit vector  $e_1^T$ . Since  $Q^1$  is unitary it has the form

$$Q^1 = \begin{pmatrix} 1 & \mathbf{0} \\ \mathbf{0} & \tilde{Q}^1 \end{pmatrix}, \quad \tilde{Q}^1 \in \mathbb{C}^{m-1 \times m-1} \quad (23)$$

The computational effort of this step consists of the computation and application of  $n(2d-1)$  Givens rotations. Because of the band structure the number of flops required is  $\mathcal{O}(nd^2)$ .

We can now easily compute a QR decomposition by applying  $n$  additional Givens rotations to the matrix  $H^1$ . This is done in  $\mathcal{O}(nd)$  flops.

**Step 2:** Since  $A^1$  and  $A^2$  overlap in all but one row and column each, we can derive  $A^2$  from  $A^1$  by cutting off the first row and column, shifting all values by one entry to the top left (as in (18)) and add a new row and column at the end. More precisely, let

$$C_p := \begin{pmatrix} \mathbf{0} & I_{p-1} \\ 1 & \mathbf{0} \end{pmatrix} : \begin{pmatrix} x_1 \\ \vdots \\ x_{p-1} \\ x_p \end{pmatrix} \mapsto \begin{pmatrix} x_2 \\ \vdots \\ x_p \\ x_1 \end{pmatrix}$$

denote the circulant backward shift of size  $p$ . Then  $\hat{A}^2 := C_m A^1 C_n^{-1}$  differs from  $A^2$  only in the last column (the first  $n-1$  entries in the last row are zero in both matrices) and satisfies (18). We illustrate this step  $A^1 \rightarrow \hat{A}^2 \rightarrow A^2$  as follows, where  $-$  and  $+$  denote the entries lost and gained respectively:

$$\underbrace{\begin{bmatrix} - & 0 & 0 & 0 & 0 & 0 & 0 \\ - & \times & 0 & 0 & 0 & 0 & 0 \\ - & \times & \times & 0 & 0 & 0 & 0 \\ - & \times & \times & \times & 0 & 0 & 0 \\ - & \times & \times & \times & \times & 0 & 0 \\ 0 & \times & \times & \times & \times & \times & 0 \\ 0 & 0 & \times & \times & \times & \times & \times \\ 0 & 0 & 0 & \times & \times & \times & \times \\ 0 & 0 & 0 & 0 & \times & \times & \times \\ 0 & 0 & 0 & 0 & 0 & \times & \times \\ 0 & 0 & 0 & 0 & 0 & 0 & \times \\ 0 & 0 & 0 & 0 & 0 & 0 & \times \end{bmatrix}}_{=A^1} \rightarrow \underbrace{\begin{bmatrix} \times & 0 & 0 & 0 & 0 & 0 & - \\ \times & \times & 0 & 0 & 0 & 0 & - \\ \times & \times & \times & 0 & 0 & 0 & - \\ \times & \times & \times & \times & 0 & 0 & - \\ \times & \times & \times & \times & \times & 0 & 0 \\ 0 & \times & \times & \times & \times & \times & 0 \\ 0 & 0 & \times & \times & \times & \times & 0 \\ 0 & 0 & 0 & \times & \times & \times & 0 \\ 0 & 0 & 0 & 0 & \times & \times & 0 \\ 0 & 0 & 0 & 0 & 0 & \times & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & - \end{bmatrix}}_{=\hat{A}^2} \rightarrow \underbrace{\begin{bmatrix} \times & 0 & 0 & 0 & 0 & 0 & 0 \\ \times & \times & 0 & 0 & 0 & 0 & 0 \\ \times & \times & \times & 0 & 0 & 0 & 0 \\ \times & \times & \times & \times & 0 & 0 & 0 \\ \times & \times & \times & \times & \times & 0 & 0 \\ 0 & \times & \times & \times & \times & \times & 0 \\ 0 & 0 & \times & \times & \times & \times & + \\ 0 & 0 & 0 & \times & \times & \times & + \\ 0 & 0 & 0 & 0 & \times & \times & + \\ 0 & 0 & 0 & 0 & 0 & \times & + \\ 0 & 0 & 0 & 0 & 0 & 0 & \times & + \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & + \end{bmatrix}}_{=A^2}$$

We apply the same transformation to the factorization  $Q^1 A^1 = H^1$ :

$$Q^1 A^1 = H^1 \Rightarrow \underbrace{C_m Q^1 C_m^{-1}}_{=: \hat{Q}^2} \underbrace{C_m A^1 C_n^{-1}}_{=A^2} = \underbrace{C_m H^1 C_n^{-1}}_{=: \hat{H}^2}. \quad (24)$$

Notice that  $\hat{Q}^2$  is again unitary and can be written as

$$\hat{Q}^2 = C_m Q^1 C_m^{-1} = C_m \begin{pmatrix} 1 & \mathbf{0} \\ \mathbf{0} & \tilde{Q}^1 \end{pmatrix} C_m^{-1} = \begin{pmatrix} \tilde{Q}^1 & \mathbf{0} \\ \mathbf{0} & 1 \end{pmatrix}.$$

The matrix  $\hat{Q}^2$  consists of the same sequences of Givens rotations as before, where all Givens rotations have been shifted up by one row. We write the factorization (24) as

$$\hat{Q}^2 \hat{A}^2 = \left( \begin{array}{c|c} \tilde{Q}^1 & \mathbf{0} \\ \hline \mathbf{0} & 1 \end{array} \right) \cdot \left( \begin{array}{ccc|c} \hat{a}_1^2 & \cdots & \hat{a}_{n-1}^2 & \hat{a}_n^2 \\ \hline 0 & \cdots & 0 & \hat{a}_{n,n}^2 \end{array} \right) = \left( \begin{array}{ccc|c} \tilde{Q}^1 \hat{a}_1^2 & \cdots & \tilde{Q}^1 \hat{a}_{n-1}^2 & \tilde{Q}^1 \hat{a}_n^2 \\ \hline 0 & \cdots & 0 & \hat{a}_{n,n}^2 \end{array} \right) = \hat{H}^2, \quad (25)$$

where  $\hat{a}_i^2$  denotes the  $i$ th column of  $\hat{A}^2$  without the last row. Notice that, by (24),  $\hat{H}^2$  is again of upper Hessenberg form everywhere except in the last column. We will now replace  $\hat{A}^2$  with  $A^2$  in (24) and (25) which leads to  $\hat{Q}^2 A^2 =: \tilde{H}^2$ . As can be seen in (25), the matrices  $\hat{H}^2$  and  $\tilde{H}^2$  only differ in the last column because  $\hat{A}^2$  and  $A^2$  only differ in the last column. These new values have to be computed by applying  $\hat{Q}^2$  to the last column of  $A^2$ . These are the only values which have to be calculated in this transformation and there is a fill-in of at most  $2d - 1$  non-zero values. We illustrate this entire procedure as follows, where  $+$  denotes the fill-in produced by applying  $\hat{Q}^2$  to the last column of  $A^2$ :

$$\begin{aligned}
Q^1 A^1 &= \begin{array}{c} \left[ \begin{array}{cccccccc} \times & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ \times & \times & 0 & 0 & 0 & 0 & 0 & 0 \\ \times & \times & \times & 0 & 0 & 0 & 0 & 0 \\ \times & \times & \times & \times & 0 & 0 & 0 & 0 \\ \times & \times & \times & \times & \times & 0 & 0 & 0 \\ 0 & \times & \times & \times & \times & \times & 0 & 0 \\ 0 & 0 & \times & \times & \times & \times & \times & \times \\ 0 & 0 & 0 & \times & \times & \times & \times & \times \\ 0 & 0 & 0 & 0 & \times & \times & \times & \times \\ 0 & 0 & 0 & 0 & 0 & \times & \times & \times \\ 0 & 0 & 0 & 0 & 0 & 0 & \times & \times \end{array} \right] \\ \downarrow \\ \left[ \begin{array}{cccccccc} \times & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ \times & \times & 0 & 0 & 0 & 0 & 0 & 0 \\ \times & \times & \times & 0 & 0 & 0 & 0 & 0 \\ \times & \times & \times & \times & 0 & 0 & 0 & 0 \\ \times & \times & \times & \times & \times & 0 & 0 & 0 \\ 0 & \times & \times & \times & \times & \times & 0 & 0 \\ 0 & 0 & \times & \times & \times & \times & \times & \times \\ 0 & 0 & 0 & \times & \times & \times & \times & \times \\ 0 & 0 & 0 & 0 & \times & \times & \times & \times \\ 0 & 0 & 0 & 0 & 0 & \times & \times & \times \\ 0 & 0 & 0 & 0 & 0 & 0 & \times & \times \end{array} \right] \end{array} = \begin{array}{c} \left[ \begin{array}{cccccccc} \times & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ \times & \times & \times & \times & \times & 0 & 0 & 0 \\ 0 & \times & \times & \times & \times & \times & \times & \times \\ 0 & 0 & \times & \times & \times & \times & \times & \times \\ 0 & 0 & 0 & \times & \times & \times & \times & \times \\ 0 & 0 & 0 & 0 & \times & \times & \times & \times \\ 0 & 0 & 0 & 0 & 0 & \times & \times & \times \\ 0 & 0 & 0 & 0 & 0 & 0 & \times & \times \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & \times \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{array} \right] = H^1 \\ \downarrow \\ \left[ \begin{array}{cccccccc} \times & \times & \times & \times & 0 & 0 & 0 & 0 \\ \times & \times & \times & \times & \times & 0 & 0 & 0 \\ 0 & \times & \times & \times & \times & \times & 0 & 0 \\ 0 & 0 & 0 & \times & \times & \times & \times & + \\ 0 & 0 & 0 & 0 & \times & \times & \times & + \\ 0 & 0 & 0 & 0 & 0 & \times & \times & + \\ 0 & 0 & 0 & 0 & 0 & 0 & \times & \times \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & \times \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & \times \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & \times \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & \times \end{array} \right] = \tilde{H}^2 \quad (26) \end{array}
\end{aligned}$$

We anticipate that in the next step we would like to apply the same shift again. However, since  $\hat{Q}^2$  acts on the first row, the requirement (23) does not hold. If we were to naively shift all values of  $\hat{Q}^2$  again to the top left by one entry and add  $e_m^T$  in the last row and column, the resulting matrix  $\hat{Q}^3$  would not be unitary. Figuratively speaking we would cut one Givens-rotation in half, since there can be no rotation acting on the “zero”th row.<sup>7</sup>

Therefore we have to remove the Givens rotation acting on the first row in  $\hat{Q}^2$  in the left-most descending sequence, which is marked as gray in (26). This can be done by applying the inverse rotations. The rotations in this sequence do not commute, since they are ordered consecutively. Thus we remove the entire sequence and add it again only this time starting in the second and ending in the  $(n - 1)$ st row. This again costs  $\mathcal{O}(nd)$ .

Starting with (26) we remove the left-most descending sequence, which results in a fill-in in the

<sup>7</sup>It is of course possible to allow Givens rotations acting on non-consecutive rows. However these rotations are difficult to remove leading to an ever increasing number of rotations.





**Algorithm 1:** QH-Shift Algorithm

**Input:** A sequence of  $d$ -banded consecutive matrices  $\{A^k\}_{k=1,\dots,k_{\max}}$  as in (18)  
**Output:** A sequence of upper triangular matrices  $\{R^k\}_{k=1,\dots,k_{\max}}$  with bandwidth  $d$

- 1 First step: Compute QH factorization  $Q^1 A^1 = H^1$  using  $2d - 1$  sequences of Givens rotations as in (22);
- 2 Compute QR factorization  $G^1 Q^1 A^1 = R^1$  using one more sequence of Givens rotations;
- 3 **for**  $k = 2, \dots, k_{\max}$  **do**
- 4     Shift factorization  $Q^{k-1} A^{k-1} = H^{k-1} \rightarrow \hat{Q}^k A^k = \tilde{H}^k$  as in (26);
- 5     Move rotations acting on the first row to the left-most sequence as in Theorem 3.2;
- 6     Remove the last rotation acting on the first row by replacing the left-most sequence as in (27);
- 7     Bring  $\hat{H}^k$  to Hessenberg form by removing  $2d - 1$  entries in the last column as in (28);
- 8     Compute QR factorization  $G^k Q^k A^k = R^k$ ;
- 9 **end**

The reordering of Givens rotations applied in (29) is described in the following theorem:

**Theorem 3.2** *Let  $l, s, m \in \mathbb{N}$ ,  $l + s \leq m$  and let*

$$Q = (G_l G_{l-1} \cdots G_1)(G_{l+1} G_l \cdots G_1) \cdots (G_{l+s-1} G_{l+s-2} \cdots G_1) \in \mathbb{C}^{m \times m}$$

*be a product of  $s$  descending sequences of Givens rotations, each starting in the first row and decreasing in length (from left to right).*

*Then  $Q$  can be described as a product of  $s$  sequences of Givens rotations of the form*

$$Q = (G_{l+s-1} G_{l+s-2} \cdots G_1)(G_{l+1} G_l \cdots G_2)(G_{l+2} G_{l+1} \cdots G_2) \cdots (G_{l+s-1} G_{l+s-2} \cdots G_2)$$

**Proof.** We start by noting the so-called *shift-through lemma* of [29] (see Lemma 9.38 there), which states that a product of 3 Givens rotations of the form  $G_2 G_1 G_2$  can be transformed into 3 Givens rotations of the form  $G_1 G_2 G_1$  and vice versa, i.e.

$$Q = \begin{matrix} \curvearrowright & \curvearrowright & \\ \curvearrowright & \curvearrowright & \\ \curvearrowright & & \end{matrix} = \begin{matrix} \curvearrowright & & \\ \curvearrowright & \curvearrowright & \\ \curvearrowright & & \end{matrix}.$$

Of course the values  $c, s$  of all rotations involved change. Note that this only holds for products without intermediate rotations, e.g. it could not be applied to

$$Q = \begin{matrix} \curvearrowright & \curvearrowright & \\ \curvearrowright & \curvearrowright & \\ \curvearrowright & & \end{matrix}.$$

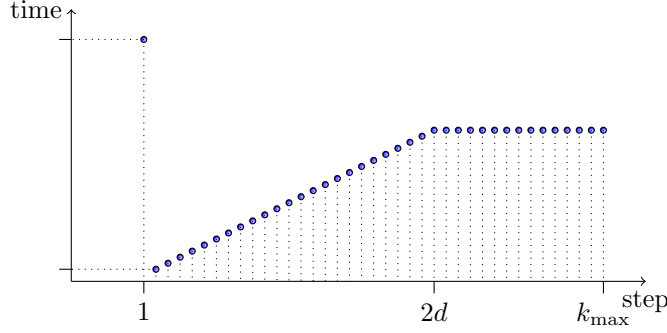


Figure 3.1: QH-Shift Algorithm without reset: Time per step

As described in [3] a repeated application of this lemma to 2 descending sequences of Givens rotations  $(G_l G_{l-1} \cdots G_1)(G_{l+1} G_l \cdots G_1)$  leads to the *shift-through lemma of higher length*:

$$(G_l \cdots G_1)(G_{l+1} \cdots G_1) = (G_{l+1} \cdots G_1)(G_{l+1} \cdots G_2)$$

Figuratively speaking we have moved the rotation from the top right to the lower left. If the right sequence is of higher length than the left sequence, the additional Givens rotations will be added to the left sequence in the end, i.e.

$$(G_l G_{l-1} \cdots G_1)(G_{l+t} G_{l+t-1} \cdots G_1) = (G_{l+t} G_{l+t-1} \cdots G_1)(G_{l+1} G_l \cdots G_2)$$

The theorem follows directly from applying the shift-through-lemma of higher length pairwise  $s - 1$  times to the descending sequences from the right to the left. ■

**Restarted QH-Shift.** As stated before, the number of descending Givens sequences to which Theorem 3.2 is applied in Algorithm 1 grows with each step. Therefore the algorithm is fastest in the second step and then slows down until it reaches step  $2d$ , as is illustrated in Figure 3.1. Depending on the time required for the initial QH factorization in step 1 and the time required in the following steps, it may be more efficient to “restart” the method after step number  $r$  (for some  $r \in \{2, \dots, 2d\}$ ) in order to take advantage of the cheap steps with number  $2, \dots, r$ . This is illustrated in Figure 3.2. The time parameters required to determine the optimal point  $r$  for a restart can be estimated during runtime (see e.g. Figure 4.2 below).



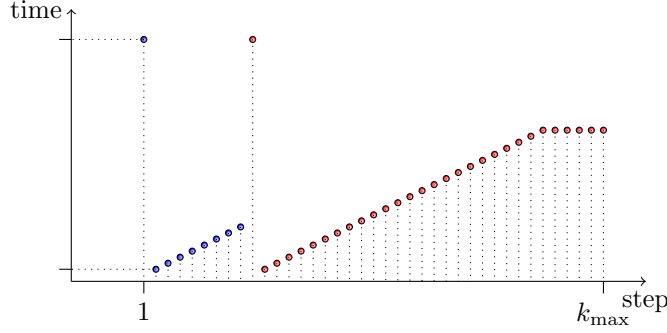


Figure 3.2: QH-Shift Algorithm with one restart after  $r = 9$  steps.

## 4 Applications

### 4.1 Laurent Operators with local impurities

We start with bi-infinite matrices with constant diagonals, also known as Laurent operators. Let  $a$  be a continuous function on the complex unit circle  $\mathbb{T}$ , and denote by  $(a_j)_{j \in \mathbb{Z}}$  the sequence of Fourier coefficients of  $a$  so that

$$a(t) = \sum_{j \in \mathbb{Z}} a_j t^j, \quad t = e^{i\theta} \in \mathbb{T}. \quad (30)$$

We denote the corresponding bounded linear operator on  $\ell^2 := \ell^2(\mathbb{Z})$  (see [8]), as well as the infinite matrix

$$\begin{pmatrix} \ddots & \ddots & \ddots & \ddots & \ddots & \ddots & \ddots \\ \ddots & a_0 & a_{-1} & a_{-2} & a_{-3} & a_{-4} & \ddots \\ \ddots & a_1 & a_0 & a_{-1} & a_{-2} & a_{-3} & \ddots \\ \ddots & a_2 & a_1 & a_0 & a_{-1} & a_{-2} & \ddots \\ \ddots & a_3 & a_2 & a_1 & a_0 & a_{-1} & \ddots \\ \ddots & a_4 & a_3 & a_2 & a_1 & a_0 & \ddots \\ \ddots & \ddots & \ddots & \ddots & \ddots & \ddots & \ddots \end{pmatrix},$$

by  $L(a)$ . Via the Fourier transform, the convolution operator  $L(a)$  corresponds to multiplication on  $L^2(\mathbb{T})$  by the function  $a$  from (30), which is referred to as the symbol of the operator  $L(a)$ . In particular, the spectrum of  $L(a)$  is the image of  $\mathbb{T}$  under the function  $a$ . Moreover, Laurent operators are normal and therefore  $\text{spec}_\varepsilon(L(a)) = \text{spec}(L(a)) + \varepsilon\mathbb{D}$  can be explicitly computed. We lose these properties when we add so-called local impurities, meaning operators  $E : \ell^2 \rightarrow \ell^2$  with a finitely supported matrix. In condensed matter physics this corresponds to a local impurity in an otherwise periodic crystal structure. The resulting operator  $L(a) + E$  is in general non-normal, and its spectrum and pseudospectra are difficult to approximate (see e.g. [6, 7]). So let us apply our algorithm.

In this example we consider exponentially decreasing Fourier coefficients  $a_k$  as  $k \rightarrow \pm\infty$ , so that  $L(a)$  is in the Wiener algebra  $\mathcal{W}$  (see Section 2.1). More precisely, when approximating  $L(a)$  by operators with finite bandwidth, say  $L_d(a)$  with bandwidth  $d \in \mathbb{N}$ , we can explicitly give upper

bounds on the approximation error as in (11). Here this means

$$\|L(a) - L_d(a)\| \leq \llbracket L(a) - L(b) \rrbracket = \sum_{|k|>d} |a_k| \leq \eta_d \quad (31)$$

for some error  $\eta_d > 0$ . For simplicity, we choose  $E$  to be of bandwidth  $\leq d$ , although impurities with larger support could be treated analogously with an appropriate error estimation in (31). The resulting lower and upper bounds on  $\text{spec}_\varepsilon(L(a) + E)$  for a given bandwidth  $d$  and cut-size  $N$  can be summed up, using (8) and (17), as follows:

$$\bigcup_{c=0}^{d-1} \Gamma_{\varepsilon-\eta_d}^{n,c}(L_d(a) + E) \subset \text{spec}_\varepsilon(L(a) + E) \subset \bigcap_{c=0}^{d-1} \Gamma_{\varepsilon+\eta_d+\delta_N}^{n,c}(L_d(A) + E), \quad (32)$$

where  $\varepsilon > 0$  and  $\delta_N \in \mathcal{O}(\frac{1}{N})$  denotes the approximation error introduced in (16). In order to compute the spectral inclusion sets  $\Gamma_\varepsilon^{n,c}(L_d(a) + E)$  from (14), we can make use of the fact that we only need to consider a finite number (growing linearly with  $n$ ) of positions  $k$ , since  $L(a) + E$  is constant along its diagonals as we move away from the support of  $E$ .

We rewrite the pseudospectral sub- and supersets from (32) as

$$\{\lambda \in \mathbb{C} : F_l(\lambda) < \varepsilon - \eta_d\} \text{ and} \quad (33)$$

$$\{\lambda \in \mathbb{C} : F_u(\lambda) < \varepsilon + \eta_d + \delta_N\}, \quad (34)$$

respectively, where

$$F_l(\lambda) := \min_{c=0,\dots,d-1} \left( \min_{k \in c+d\mathbb{Z}} \left( \min(\nu((A - \lambda I)|_{\ell^2(J_k^n)}), \nu((A - \lambda I)^*|_{\ell^2(J_k^n)})) \right) \right)$$

$$F_u(\lambda) := \max_{c=0,\dots,d-1} \left( \min_{k \in c+d\mathbb{Z}} \left( \min(\nu((A - \lambda I)|_{\ell^2(J_k^n)}), \nu((A - \lambda I)^*|_{\ell^2(J_k^n)})) \right) \right)$$

with  $A := L_d(a) + E$ . If we are only interested in  $\text{spec}_\varepsilon(A)$  for a few values  $\varepsilon > 0$ , then we can approximate connected components of these sets by determining the boundary of each set. To that end, we use continuation methods to determine all  $\lambda \in \mathbb{C}$  which satisfy

$$F_l(\lambda) - (\varepsilon - \eta_d) = 0, \quad (35)$$

$$F_u(\lambda) - (\varepsilon + \eta_d + \delta_n) = 0, \quad (36)$$

respectively. There have been several approaches to computing the boundary curves of pseudospectral of finite square matrices using gradient-based methods, see e.g. [10, 4]. However, since  $F_u$  and  $F_l$  involve several nested minima and maxima of smallest singular values of rectangular matrices, they are non-smooth, so that these methods are not well suited to our case.

In [22] a piecewise linear (PL)-continuation method was used to approximate pseudospectral boundaries of matrices, and this algorithm can be easily modified to be applied to  $F_u$  and  $F_l$  as well. The idea is to triangulate the complex plane and determine all triangles in which the signs of  $F_u$  (respectively  $F_l$ ) are not equal on all three vertices. The functions have therefore to be evaluated on these vertices only, and each function evaluation consists of two applications of the QH-shift method (one for  $(A - \lambda I)$  and one for  $(A - \lambda I)^*$ ). We note the difficulty of finding a starting point, an initial triangle, which can be solved using a coarse grid or a bisection based method (see [22]). The boundaries can be estimated prior using the coarse upper bound  $\text{spec}_\varepsilon(L(a) + E) \subset \text{spec}(L(a)) + (\varepsilon + \|E\|)\mathbb{D}$ .

We applied this method to a Laurent operator defined by the coefficients

$$a_k := \begin{cases} \left(\frac{1}{2}\right)^k + 1.1 \left(\frac{1}{2i}\right)^k, & k > 0 \\ 3.1, & k = 0 \\ \left(\frac{i}{2}\right)^{-k}, & k < 0 \end{cases} \quad (37)$$

and an impurity  $E$  which is a scaled and shifted  $10 \times 10$  Grcar-matrix. We refer to  $L(a)$  as “the fish”, motivated by the shape of its spectrum. The approximation error (31) can be estimated as  $\eta_d \leq \frac{1}{2^{d-2}}$ . The results can be seen in Figure 4.1. In addition to the upper and lower bounds on  $\text{spec}_\varepsilon(L(a) + E)$ , we see that the superset does not contain the origin, implying that this particular operator is invertible. Furthermore we compare the speed of the QH-shift method, the restarted QH-shift method and the classic QR-decomposition, including the SVD, for a single  $\lambda \in \mathbb{C}$ , as can be seen in Figure 4.2.

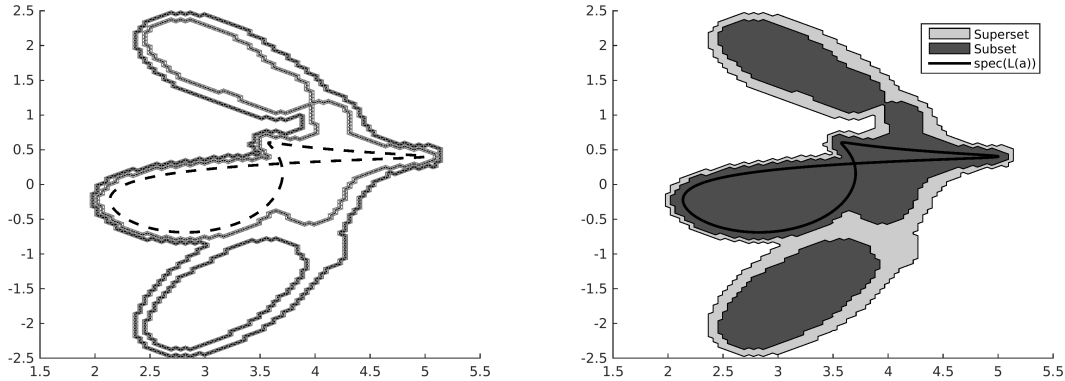


Figure 4.1: Boundary sets for the impure Laurent operator with symbol (37) and a scaled  $10 \times 10$  Grcar-matrix as local impurity. Left: Triangulation of solution curves of (35) and (36) and  $\text{spec}(L(a))$  superimposed as a dotted line; Right: Visualization of the resulting sub- and superset of  $\text{spec}_\varepsilon(L(a) + E)$ . Dimensions are  $d = 15$ ,  $N = 200$  and  $\varepsilon = 0.1$ , and the approximation errors are  $\eta_d \leq 1.2 \cdot 10^{-4}$  and  $\delta_N \leq 0.0512$ . We used a total of 1486 equilateral triangles of side-length 0.05 for the computations.

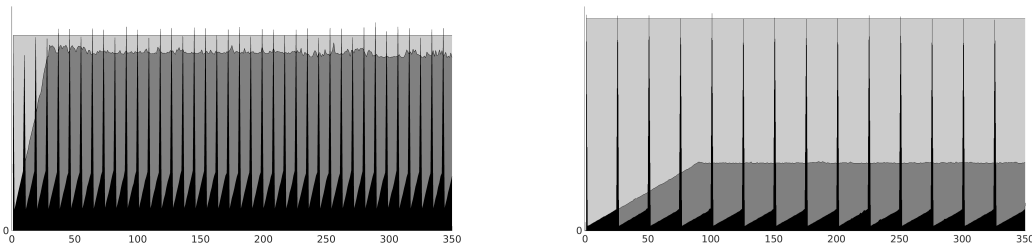


Figure 4.2: Two plots of the computing time required for each step in Algorithm 1 using the QH-Shift method (dark grey), the restarted QH-shift method (black) and, for comparison, the classical full QR decomposition using Givens rotations (light grey). The total computational cost corresponds to the total black, dark or light grey area. Comparison of the left ( $d = 15$ ) and right ( $d = 45$ ) image confirms that our algorithm pays off with increasing bandwidth. The cut-size is  $N = 50$  in both cases, so that  $n = 15 \cdot 50 = 750$  and  $n = 45 \cdot 50 = 2250$ , respectively.

## 4.2 Singular Integral Operators

Let  $c$  and  $e$  be continuous functions on  $\mathbb{T}$  which define Laurent operators  $L(c)$  and  $L(e)$  on  $\ell^2$  via the Fourier isomorphism (see Section 4.1 and [8]). In our numerical example below, we use the “whale” symbol from [8] and our “fish” symbol from (37). Now we interbreed “whale” and “fish”:

Put  $S_{\ell^2} := P - Q$  on  $\ell^2$ , where  $P$  is the orthogonal projection of  $\ell^2(\mathbb{Z})$  onto  $\ell^2(\mathbb{N}_0)$  and  $Q := I - P$ . Then  $S_{\ell^2}$  corresponds to the so-called Cauchy singular integral operator on  $L^2(\mathbb{T})$  (see e.g. [19, p.130f]). After composition and addition with our multiplication operators by  $c$  and  $e$  on  $L^2(\mathbb{T})$ ,



## References

- [1] G. R. BARRENECHEA, L. BOULTON and N. BOUSSAID: Eigenvalue enclosures, [arXiv:1306.5354](https://arxiv.org/abs/1306.5354)
- [2] Z. BAI, J. DEMMEL, J. DONGARRA, A. RUHE and H. VAN DER VORST: *Templates for the Solution of Algebraic Eigenvalue Problems - A Practical Guide*, siam, first edition, 2000.
- [3] M. VAN BAREL, R. VANDEBRIL, P. VAN DOOREN and K. FREDERIX: Implicit double shift qr-algorithm for companion matrices, *Numerische Mathematik*, Springer, vol. 116(2), p.117-212.
- [4] C. BEKAS and E. GALLOPOULOS: Cobra: Parallel path following for computing the matrix pseudospectrum, *Parallel Computing* **27** (2001), 1879-1896
- [5] J. BEN-ARTZI, A. C. HANSEN, O. NEVANLINNA and M. SEIDEL: New barriers in complexity theory: On the Solvability Complexity Index and towers of algorithms, *C. R. Acad. Sci. Paris, Ser. I* **353** (2015), 931-936.
- [6] A. BÖTTCHER, M. EMBREE and M. LINDNER: Spectral approximation of banded Laurent matrices with localized random perturbations, *Integr. Equ. Oper. Theory* **42** (2002), 142-165.
- [7] A. BÖTTCHER, M. EMBREE and V. I. SOKOLOV: Infinite Toeplitz and Laurent matrices with localized impurities, *Linear Algebra Appl.* **343/344** (2002), 101-118.
- [8] A. BÖTTCHER and B. SILBERMANN: *Introduction to Large Truncated Toeplitz Matrices*, Springer, Berlin, Heidelberg, 1999.
- [9] N. BOSNER: Fast Methods for Large Scale Singular Value Decomposition *Doctoral Thesis, University of Zagreb*, 2006.
- [10] M. BRUEHL: A curve tracing algorithm for computing the pseudospectrum, *BIT* **36(3)** (1996), 441-454
- [11] S. N. CHANDLER-WILDE, R. CHONCHAIYA and M. LINDNER: Upper bounds on the spectra and pseudospectra of Jacobi and related operators, *in preparation*
- [12] S. N. CHANDLER-WILDE and M. LINDNER: Coburn's Lemma and the Finite Section Method for Random Jacobi Operators, *Journal of Functional Analysis*, **270** (2016), 802-841.
- [13] R. CHONCHAIYA: Computing the Spectra and Pseudospectra of Non-Self-Adjoint Random Operators Arising in Mathematical Physics, *PhD Thesis, University of Reading, UK*, 2010. <http://www.tuhh.de/~matmli/files/ChonchaiyaThesis.pdf>
- [14] E. B. DAVIES and M. PLUM: Spectral pollution, *IMA J. Numer. Anal.*, **24** (2004), 417-438.
- [15] G.H. GOLUB and C.F. VAN LOAN: *Matrix Computations*, The Johns Hopkins University Press, first edition, 1983
- [16] R. HAGGER, M. LINDNER and M. SEIDEL: Essential pseudospectra and essential norms of band-dominated operators, *Journal of Mathematical Analysis and Applications*, **437** (2016), 255-291.
- [17] A. C. HANSEN: On the solvability complexity index, the  $n$ -pseudospectrum and approximations of spectra of operators, *J. Amer. Math. Soc.*, **24** (2011), 81-124.
- [18] V. G. KURBATOV: *Functional Differential Operators and Equations*, Kluwer Academic Publishers, Dordrecht, Boston, London 1999.

- [19] M. LINDNER: *Infinite Matrices and their Finite Sections: An Introduction to the Limit Operator Method*, Frontiers in Mathematics, Birkhäuser 2006.
- [20] M. LINDNER: Fredholm Theory and Stable Approximation of Band Operators and Their Generalisations, *Habilitation thesis, TU Chemnitz, Germany*, 2009.
- [21] M. LINDNER and M. SEIDEL: An affirmative answer to a core issue on limit operators, *Journal of Functional Analysis*, **267** (2014), 901–917.
- [22] D. MEZHER and B. PHILIPPE: PAT - a Reliable Path-Following Algorithm, *Numerical Algorithms* **29** (2002), 131-152
- [23] V. S. RABINOVICH, S. ROCH and B. SILBERMANN: *Limit Operators and Their Applications in Operator Theory*, Birkhäuser 2004.
- [24] L. REICHEL and L. N. TREFETHEN: Eigenvalues and pseudo-eigenvalues of Toeplitz matrices, *Lin. Alg. Appl.* **162-164** (1992), 153–185.
- [25] M. SEIDEL: On Some Banach Algebra Techniques in Operator Theory, *PhD Thesis, TU Chemnitz*, 2011.
- [26] M. SEIDEL: On  $(N, \varepsilon)$ -pseudospectra of operators on Banach spaces, *Journal of Functional Analysis*, **262** (2012), 4916–4927.
- [27] M. SEIDEL and B. SILBERMANN Finite Sections of Band-Dominated Operators - Norms, Condition Numbers and Pseudospectra, *Operator Theory: Advances and Applications*, **228** (2013), 375–390.
- [28] L. N. TREFETHEN and M. EMBREE: *Spectra and Pseudospectra: the Behavior of Nonnormal Matrices and Operators*, Princeton University Press, Princeton, NJ, 2005.
- [29] R. VANDEBRIL, M. VAN BAREL and N. MASTRONARDI: *Matrix Computations and semiseparable matrices*, Md. Johns Hopkins Univ. Press, volume 1, 2008
- [30] R. VANDEBRIL, M. VAN BAREL and N. MASTRONARDI: *Rotational Figures Latex Package*, <http://people.cs.kuleuven.be/~raf.vandebril/homepage/software/latex.php>, 2010
- [31] T. G. WRIGHT and L. N. TREFETHEN: *Pseudospectra of rectangular matrices*, IMA Journal of Numerical Analysis **22** (2002), 501-519

**Authors' addresses:**

Marko Lindner  
 Torge Schmidt  
 Institut Mathematik  
 TU Hamburg (TUHH)  
 D-21073 Hamburg  
 GERMANY

[lindner@tuhh.de](mailto:lindner@tuhh.de)  
[torge.schmidt@tuhh.de](mailto:torge.schmidt@tuhh.de)