

Arnoldi decomposition, GMRES, and preconditioning for linear discrete ill-posed problems

Silvia Gazzola^a, Silvia Noschese^{b,1}, Paolo Novati^{c,2}, Lothar Reichel^{d,*,3}

^a Department of Mathematical Sciences, University of Bath, Bath BA2 7AY, United Kingdom

^b Dipartimento di Matematica “Guido Castelnuovo”, SAPIENZA Università di Roma, P.le A. Moro, 2, I-00185 Roma, Italy

^c Dipartimento di Matematica e Geoscienze, Università di Trieste, via Valerio, 12/1, I-34127 Trieste, Italy

^d Department of Mathematical Sciences, Kent State University, Kent, OH 44242, USA

ARTICLE INFO

Article history:

Received 14 December 2018

Received in revised form 26 February 2019

Accepted 27 February 2019

Available online 6 March 2019

Keywords:

Linear discrete ill-posed problem

Arnoldi process

GMRES

Truncated iteration

Tikhonov regularization

Truncated singular value decomposition

ABSTRACT

GMRES is one of the most popular iterative methods for the solution of large linear systems of equations that arise from the discretization of linear well-posed problems, such as boundary value problems for elliptic partial differential equations. The method is also applied to the iterative solution of linear systems of equations that are obtained by discretizing linear ill-posed problems, such as many inverse problems. However, GMRES does not always perform well when applied to the latter kind of problems. This paper seeks to shed some light on reasons for the poor performance of GMRES in certain situations, and discusses some remedies based on specific kinds of preconditioning. The standard implementation of GMRES is based on the Arnoldi process, which also can be used to define a solution subspace for Tikhonov or TSVD regularization, giving rise to the Arnoldi–Tikhonov and Arnoldi–TSVD methods, respectively. The performance of the GMRES, the Arnoldi–Tikhonov, and the Arnoldi–TSVD methods is discussed. Numerical examples illustrate properties of these methods.

© 2019 IMACS. Published by Elsevier B.V. All rights reserved.

1. Introduction

This paper considers the solution of linear systems of equations

$$A\mathbf{x} = \mathbf{b}, \quad A \in \mathbb{C}^{m \times m}, \quad \mathbf{x}, \mathbf{b} \in \mathbb{C}^m, \quad (1)$$

with a large matrix A with many “tiny” singular values of different orders of magnitude. In particular, A is severely ill-conditioned and may be rank-deficient. Linear systems of equations (1) with a matrix of this kind are commonly referred to as linear discrete ill-posed problems. They arise, for instance, from the discretization of linear ill-posed problems, such as Fredholm integral equations of the first kind with a smooth kernel.

* Corresponding author.

E-mail addresses: s.gazzola@bath.ac.uk (S. Gazzola), noschese@mat.uniroma1.it (S. Noschese), novati@units.it (P. Novati), reichel@math.kent.edu (L. Reichel).

¹ Research partially supported by a grant from SAPIENZA Università di Roma and by INdAM-GNCS.

² Research partially supported by INdAM-GNCS and FRA-University of Trieste.

³ Research partially supported by National Science Foundation grants DMS-1720259 and DMS-1729509.

In many linear discrete ill-posed problems that arise in science and engineering, the right-hand side vector \mathbf{b} is determined through measurements and is contaminated by a measurement error $\mathbf{e} \in \mathbb{C}^m$, which we sometimes will refer to as “noise.” Thus,

$$\mathbf{b} = \mathbf{b}_{\text{exact}} + \mathbf{e},$$

where $\mathbf{b}_{\text{exact}} \in \mathbb{C}^m$ denotes the unknown error-free right-hand side associated with \mathbf{b} . We will assume that $\mathbf{b}_{\text{exact}}$ is in the range of A , denoted by $\mathcal{R}(A)$, because this facilitates the use of the discrepancy principle to determine a suitable value of a regularization parameter; see below for details. The error-contaminated right-hand side \mathbf{b} is not required to be in $\mathcal{R}(A)$. We remark that the solution methods described also can be used with other techniques for determining the regularization parameter; see, e.g., [34,52] for discussions and illustrations of a variety of methods.

We would like to compute the solution of minimal Euclidean norm, $\mathbf{x}_{\text{exact}}$, of the consistent linear discrete ill-posed problem

$$A\mathbf{x} = \mathbf{b}_{\text{exact}}. \quad (2)$$

Since the right-hand side $\mathbf{b}_{\text{exact}}$ is not known, we seek to determine an approximation of $\mathbf{x}_{\text{exact}}$ by computing an approximate solution of the available linear system of equations (1). We note that, due to the severe ill-conditioning of the matrix A and the error \mathbf{e} in \mathbf{b} , the least-squares solution of minimal Euclidean norm of (1) generally is not a useful approximation of $\mathbf{x}_{\text{exact}}$ due to severe propagation of the error \mathbf{e} into the solution.

A popular approach to determine a meaningful approximation of $\mathbf{x}_{\text{exact}}$ is to apply an iterative method to the solution of (1) and terminate the iterations early enough so that the error in \mathbf{b} is not significantly propagated into the computed approximate solution. The most popular iterative methods for the solution of large linear discrete ill-posed problems are LSQR by Paige and Saunders [19,25,27,51], which is based on partial Golub–Kahan decomposition of A , and GMRES [7,8,22], which is based on partial Arnoldi decomposition of A . Here “GMRES” refers to both the standard GMRES method proposed by Saad and Schultz [57] and to modifications discussed in [17,41].

The LSQR method requires the evaluation of two matrix-vector products in each iteration, one with A and one with its conjugate transpose, which we denote by A^* . GMRES only demands the computation of one matrix-vector product with A per iteration. This makes GMRES attractive to use when it is easy to evaluate matrix-vector products with A but not with A^* . This is, for instance, the case when A approximates a Fredholm integral operator of the first kind and matrix-vector products with A are evaluated by a multipole method. Then A is not explicitly formed and matrix-vector products with A^* are difficult to compute; see, e.g., [24] for a discussion on the multipole method. It may be difficult to evaluate matrix-vector products with A^* also when solving nonlinear problems and A represents a Jacobian matrix, whose entries are not explicitly computed; see [12] for a discussion on such a solution method.

The fact that GMRES does not require the evaluation of matrix-vector products with A^* may result in that, for many linear discrete ill-posed problems (1), this method requires fewer matrix-vector product evaluations than LSQR to determine a desired approximate solution, see, e.g., [4,5,8] for illustrations, as well as [6] for related examples. However, there also are linear discrete ill-posed problems (1), whose solution with LSQR requires fewer matrix-vector product evaluations than GMRES, or for which LSQR furnishes a more accurate approximation of $\mathbf{x}_{\text{exact}}$ than GMRES; see below for illustrations, as well as [30]. Reasons for poor performance of GMRES include:

- (1) Low-dimensional solution subspaces used by GMRES may be poorly suited to represent $\mathbf{x}_{\text{exact}}$. It is often not possible to rectify this problem by carrying out many iterations, since this typically results in severe propagation of the error \mathbf{e} in \mathbf{b} into the iterates determined by GMRES.
- (2) The desired solution $\mathbf{x}_{\text{exact}}$ may be approximated accurately in solution subspaces generated by GMRES, but the method determines iterates that furnish poor approximations of $\mathbf{x}_{\text{exact}}$.
- (3) The GMRES iterates suffer from contamination of propagated error due to the fact that the initial vector in the Arnoldi decomposition used for the solution of (1) is a normalization of the error-contaminated vector \mathbf{b} when, as is commonly done, the initial approximate solution is chosen to be $\mathbf{x}_0 = \mathbf{0}$.

It is the purpose of the present paper to discuss the above mentioned shortcomings of GMRES, illustrate situations when they occur, and provide some remedies. Section 2 recalls the Arnoldi process and GMRES, and shows that the solution subspaces used by GMRES may be inappropriate. Also LSQR is briefly discussed, and distances to relevant classes of matrices are introduced. In Section 3, we define the set of generalized Hermitian matrices and the set of generalized Hermitian positive semidefinite matrices. The distance of the matrix A in (1) to these sets sheds light on how quickly GMRES applied to the solution of the linear system of equations (1) will converge. Section 4 describes “preconditioning techniques.” The “preconditioners” discussed do not necessarily reduce the condition number, and they are not guaranteed to reduce the number of iterations. Instead, they are designed to make the matrix of the preconditioned linear system of equations closer to the set of generalized Hermitian positive semidefinite matrices. This often results in that the computed solution is a more accurate approximation of $\mathbf{x}_{\text{exact}}$ than approximate solutions of the unpreconditioned linear system (1). In Section 5 we consider the situation when GMRES applied to the solution of (1) yields poor approximations of $\mathbf{x}_{\text{exact}}$, but the solution subspace generated by the Arnoldi process contains an accurate approximation of $\mathbf{x}_{\text{exact}}$. We propose to carry out sufficiently

many steps of the Arnoldi process and determine an approximation of $\mathbf{x}_{\text{exact}}$ by Tikhonov regularization or truncated singular value decomposition in the solution subspace so generated. Both regularization methods allow the use of a solution subspace of larger dimension than GMRES. A few computed examples that illustrate the discussion of the previous sections are presented in Section 6, and Section 7 contains concluding remarks.

2. GMRES and LSQR for linear discrete ill-posed problems

GMRES is a popular iterative method for the solution of large linear systems of equations with a square nonsymmetric matrix (1) that arise from the discretization of well-posed problems; see, e.g., Saad [56]. The k th iterate, \mathbf{x}_k , determined by GMRES, when applied to the solution of (1) with initial iterate $\mathbf{x}_0 = \mathbf{0}$, satisfies

$$\|\mathbf{A}\mathbf{x}_k - \mathbf{b}\| = \min_{\mathbf{x} \in \mathbb{K}_k(\mathbf{A}, \mathbf{b})} \|\mathbf{A}\mathbf{x} - \mathbf{b}\|, \quad \mathbf{x}_k \in \mathbb{K}_k(\mathbf{A}, \mathbf{b}), \quad (3)$$

where

$$\mathbb{K}_k(\mathbf{A}, \mathbf{b}) = \text{span}\{\mathbf{b}, \mathbf{A}\mathbf{b}, \dots, \mathbf{A}^{k-1}\mathbf{b}\}$$

is a Krylov subspace and $\|\cdot\|$ denotes the Euclidean vector norm. The choice $\mathbf{x}_0 = \mathbf{0}$ is quite common, and will be tacitly assumed throughout this paper. Also, we will assume that $1 \leq k \ll m$ is sufficiently small so that $\dim(\mathbb{K}_k(\mathbf{A}, \mathbf{b})) = k$, which guarantees that the iterate \mathbf{x}_k is uniquely defined. The standard implementation of GMRES [56,57] is based on the Arnoldi process, given in Algorithm 1 with the modified Gram–Schmidt implementation.

Algorithm 1 (The Arnoldi process).

```

0. Input  $\mathbf{A} \in \mathbb{C}^{m \times m}$ ,  $\mathbf{b} \in \mathbb{C}^m \setminus \{\mathbf{0}\}$ 
1.  $\mathbf{v}_1 := \mathbf{b}/\|\mathbf{b}\|$ ;
2. for  $j = 1, 2, \dots, k$  do
3.    $\mathbf{w} := \mathbf{A}\mathbf{v}_j$ ;
4.   for  $i = 1, 2, \dots, j$  do
5.      $h_{i,j} := \mathbf{v}_i^* \mathbf{w}$ ;  $\mathbf{w} := \mathbf{w} - \mathbf{v}_i h_{i,j}$ ;
6.   end for
7.    $h_{j+1,j} := \|\mathbf{w}\|$ ;  $\mathbf{v}_{j+1} := \mathbf{w}/h_{j+1,j}$ ;
8. end for
```

Algorithm 1 generates orthonormal vectors $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_{k+1}$, the first k of which form a basis for $\mathbb{K}_k(\mathbf{A}, \mathbf{b})$. Define the matrices $V_j = [\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_j]$ for $j \in \{k, k+1\}$. The scalars $h_{i,j}$ determined by the algorithm define an upper Hessenberg matrix $H_{k+1,k} = [h_{i,j}] \in \mathbb{C}^{(k+1) \times k}$. Using these matrices, the recursion formulas for the Arnoldi process can be expressed as a *partial Arnoldi decomposition*,

$$\mathbf{A}V_k = V_{k+1}H_{k+1,k}. \quad (4)$$

The above relation is applied to compute the GMRES iterate \mathbf{x}_k as follows: Express (3) as

$$\min_{\mathbf{x} \in \mathbb{K}_k(\mathbf{A}, \mathbf{b})} \|\mathbf{A}\mathbf{x} - \mathbf{b}\| = \min_{\mathbf{y} \in \mathbb{C}^k} \|\mathbf{A}V_k\mathbf{y} - \mathbf{b}\| = \min_{\mathbf{y} \in \mathbb{C}^k} \|H_{k+1,k}\mathbf{y} - \mathbf{e}_1\|\|\mathbf{b}\|, \quad (5)$$

where the orthonormality of the columns of V_{k+1} and the fact that $\mathbf{b} = \|\mathbf{b}\|V_{k+1}\mathbf{e}_1$ have been exploited. Throughout this paper $\mathbf{e}_j = [0, \dots, 0, 1, 0, \dots, 0]^*$ denotes the j th axis vector. The small minimization problem on the right-hand side of (5) can be solved conveniently by QR factorization of $H_{k+1,k}$; see [56]. Denote its solution by \mathbf{y}_k . Then $\mathbf{x}_k = V_k\mathbf{y}_k$ solves (3) and $\mathbf{r}_k = \mathbf{b} - \mathbf{A}\mathbf{x}_k$ is the associated residual error. Since $\mathbb{K}_{k-1}(\mathbf{A}, \mathbf{b}) \subset \mathbb{K}_k(\mathbf{A}, \mathbf{b})$, we have $\|\mathbf{r}_k\| \leq \|\mathbf{r}_{k-1}\|$; generically this inequality is strict. Note that $\|\mathbf{r}_k\| = \|\mathbf{e}_1\|\|\mathbf{b}\| - H_{k+1,k}\mathbf{y}_k\|$, so that the norm of the residual error can be monitored using projected quantities, which are inexpensive to compute. We remark that a reorthogonalization procedure can be considered with Algorithm 1, by running an additional Gram–Schmidt step for the vector \mathbf{w} after step 6 has been performed: this has the effect of assuring the columns of V_{k+1} a better numerical orthogonality.

Assume that a fairly accurate bound $\delta > 0$ for the norm of the noise \mathbf{e} in \mathbf{b} is available,

$$\|\mathbf{e}\| \leq \delta, \quad (6)$$

and let $\tau \geq 1$ be a user-chosen parameter that is independent of δ . The *discrepancy principle* prescribes that the iterations of GMRES applied to the solution of (1) be terminated as soon as an iterate \mathbf{x}_k has been determined such that the associated residual error \mathbf{r}_k satisfies

$$\|\mathbf{r}_k\| \leq \tau\delta. \quad (7)$$

The purpose of this stopping criterion is to terminate the iterations before the iterates \mathbf{x}_k are severely contaminated by propagated error that stems from the error \mathbf{e} in \mathbf{b} . Note that the residual $\mathbf{r}_{\text{exact}} = \mathbf{b} - A\mathbf{x}_{\text{exact}}$ satisfies the inequality (7). This follows from (6) and the consistency of (2). Also iterations with LSQR are commonly terminated with the discrepancy principle; see, e.g., [7,19,25] for discussions on the use of the discrepancy principle for terminating iterations with GMRES and LSQR. Other approaches to determine k are discussed in the literature and can be applied instead; see, e.g., [34,52]. We will use the discrepancy principle in the computed examples, because its properties are well understood.

The LSQR method [51] is an implementation of the conjugate gradient method applied to the normal equations,

$$A^*A\mathbf{x} = A^*\mathbf{b}, \quad (8)$$

with a Hermitian positive semidefinite matrix. LSQR circumvents the explicit formation of A^*A . When using the initial iterate $\mathbf{x}_0 = \mathbf{0}$, LSQR determines approximate solutions of (1) in a sequence of nested Krylov subspaces $\mathbb{K}_k(A^*A, A^*\mathbf{b})$, $k = 1, 2, \dots$. The k th iterate, \mathbf{x}_k , computed by LSQR satisfies

$$\|A\mathbf{x}_k - \mathbf{b}\| = \min_{\mathbf{x} \in \mathbb{K}_k(A^*A, A^*\mathbf{b})} \|A\mathbf{x} - \mathbf{b}\|, \quad \mathbf{x}_k \in \mathbb{K}_k(A^*A, A^*\mathbf{b});$$

see [2,51] for further details on LSQR.

When δ in (6) is fairly large, only a few iterations can be carried out by GMRES or LSQR before (7) is satisfied. In particular, an accurate approximation of $\mathbf{x}_{\text{exact}}$ can then be determined by GMRES only if $\mathbf{x}_{\text{exact}}$ can be approximated well in a low-dimensional Krylov subspace $\mathbb{K}_k(A, \mathbf{b})$. Moreover, it has been observed that GMRES based on the Arnoldi process applied to A with initial vector \mathbf{b} may determine iterates \mathbf{x}_k that are contaminated by more propagated error than iterates generated by LSQR; see [30]. A reason for this is that the first column of the matrix V_k in the Arnoldi decomposition (4) (i.e., the first basis vector for the GMRES solution) is a normalization of the error-contaminated vector \mathbf{b} , and the error \mathbf{e} in \mathbf{b} is propagated to all columns of V_k by the Arnoldi process.

A remedy for the latter difficulty is to use a modification of the Arnoldi decomposition,

$$A\hat{V}_k = V_{k+j}H_{k+j,k}, \quad (9)$$

with $j \geq 2$. The columns of $\hat{V}_k \in \mathbb{C}^{m \times k}$ form an orthonormal basis for the Krylov subspace $\mathbb{K}_k(A, A^{(j-1)}\mathbf{b})$, in which we are looking for an approximate solution. Moreover, the columns of $V_{k+j} \in \mathbb{C}^{m \times (k+j)}$ form an orthonormal basis for $\mathbb{K}_{k+j}(A, \mathbf{b})$, and all entries of the matrix $H_{k+j,k} \in \mathbb{C}^{(k+j) \times k}$ below the j th subdiagonal vanish; see [17] for details. The special case when $j = 2$ is discussed in [41]. When $j = 1$, the decomposition (9) simplifies to (4). A reason why applying the decomposition (9) may be beneficial is that, in our typical applications, the matrix A is a low-pass filter. Therefore, the high-frequency error in the vector $\hat{V}_k\mathbf{e}_1 = A^{(j-1)}\mathbf{b}/\|A^{(j-1)}\mathbf{b}\|$ is damped.

The following examples illustrate that GMRES may perform poorly also when there is no error in \mathbf{b} , in the sense that GMRES may require many iterations to solve the system of equations or not be able to compute a solution at all. While the coefficient matrices of these examples are artificial, related examples can be found in Hansen's *Regularization Tools* [28] (see, e.g., the test problem `heat`) and also arise in image restoration when the available image has been contaminated by motion blur; see [14, Section 4].

Example 2.1. Let A in (1) be the downshift matrix of upper Hessenberg form,

$$A = \begin{bmatrix} 0 & & \cdots & 0 \\ 1 & 0 & & 0 \\ 0 & 1 & 0 & \cdots & 0 \\ \vdots & 0 & \ddots & \ddots & 0 \\ 0 & \cdots & 0 & 1 & 0 \end{bmatrix} \in \mathbb{C}^{m \times m}, \quad (10)$$

and let $\mathbf{b} = \mathbf{e}_2$. The minimal-norm solution of the linear system of equations (1) is $\mathbf{x}_{\text{exact}} = \mathbf{e}_1$. Since $\mathbb{K}_k(A, \mathbf{b}) = \text{span}\{\mathbf{e}_2, \mathbf{e}_3, \dots, \mathbf{e}_{k+1}\}$, it follows that the solution of (3) is $\mathbf{x}_k = \mathbf{0}$ for $1 \leq k < m$. These solutions are poor approximations of $\mathbf{x}_{\text{exact}}$. GMRES breaks down at step m due to division by zero in Algorithm 1. This depends on that the matrix A is singular. Thus, when m is large GMRES produces poor approximations of $\mathbf{x}_{\text{exact}}$ for many iterations before breakdown. While breakdown of GMRES can be handled, see [53], the lack of convergence of the iterates towards $\mathbf{x}_{\text{exact}}$ for many steps remains. The poor performance of GMRES in this example stems from the facts that A is a shift operator and the desired solution $\mathbf{x}_{\text{exact}}$ has few nonvanishing entries. We remark that the minimal-norm solution $\mathbf{x}_{\text{exact}} = \mathbf{e}_1$ of (1) lives in $\mathbb{K}_1(A^*A, A^*\mathbf{b})$ and LSQR determines this solution in one step. \square

Example 2.1 illustrates that replacing a linear discrete ill-posed problem (1) with a non-Hermitian coefficient matrix A by a linear discrete ill-posed problem (8) having a Hermitian positive semidefinite matrix A^*A may be beneficial. The matrix A^* in (8) may be considered a preconditioner for the linear system of equations (1). To shed some light on the

possible benefit of this kind of replacement, with the aim of developing suitable preconditioners different from A^* , we will discuss the distance of a square matrix A to the set of Hermitian matrices \mathbb{H} , the set of anti-Hermitian (skew-Hermitian) matrices \mathbb{A} , the set of normal matrices \mathbb{N} , the set of Hermitian positive semidefinite matrices \mathbb{H}_+ , and the set of Hermitian negative semidefinite matrices \mathbb{H}_- . We are interested in the distance to the set of normal matrices, because it is known that GMRES may converge slowly when the matrix A in (1) is far from \mathbb{N} . Specifically, the rate of convergence of GMRES may be slow when A has a spectral factorization with a very ill-conditioned eigenvector matrix; see [39, Theorem 3] and [40] for discussions. Note that when A belongs to the classes \mathbb{N} , \mathbb{H} , \mathbb{A} , or \mathbb{H}_+ , the Arnoldi process and GMRES can be simplified; see, e.g., Eisenstat [18], Huckle [32], Paige and Saunders [50], and Saad [56, Section 6.8].

We remark that the dependence of the convergence behavior of GMRES on the eigenvalues and eigenvectors of A is complicated; see, e.g., Greenbaum et al. [23] and Du et al. [15] for discussions. In particular, one can construct matrices A with a specified distribution of eigenvalues in \mathbb{C} such that GMRES applied to the solution of (2) converges arbitrarily slowly. The convergence also depends on the right-hand side. It is therefore not obvious that replacing the matrix A in (1) by a matrix that is closer to the sets \mathbb{H} , \mathbb{H}_+ , or \mathbb{H}_- by choosing a suitable preconditioner and applying GMRES to the preconditioned linear system of equations so obtained will result in faster convergence. The distribution of eigenvalues is important for the rate of convergence when the eigenvector matrix is not too ill-conditioned; see, e.g., [40] and [56, Section 6.11] for discussions. This is the case for many linear discrete ill-posed problems. It is therefore meaningful to consider preconditioners that change the distribution of the eigenvalues of the system matrix. Further desirable properties for preconditioners for linear discrete ill-posed problems will be commented on below.

We measure distances between a matrix A and the sets \mathbb{H} , \mathbb{A} , \mathbb{N} , and \mathbb{H}_\pm in the Frobenius norm, which for a matrix M is defined as $\|M\|_F = (\text{Trace}(M^*M))^{1/2}$. The following proposition considers the matrix of Example 2.1.

Proposition 2.1. *Let the matrix $A \in \mathbb{C}^{m \times m}$ be defined by (10). The relative distances in the Frobenius norm to the sets of the Hermitian and anti-Hermitian matrices are*

$$\frac{\text{dist}_F(A, \mathbb{H})}{\|A\|_F} = \frac{1}{\sqrt{2}} \quad (11)$$

and

$$\frac{\text{dist}_F(A, \mathbb{A})}{\|A\|_F} = \frac{1}{\sqrt{2}}, \quad (12)$$

respectively. Moreover,

$$\frac{\text{dist}_F(A, \mathbb{N})}{\|A\|_F} \leq \frac{1}{\sqrt{m}}, \quad (13)$$

$$\frac{\text{dist}_F(A, \mathbb{H}_+)}{\|A\|_F} = \frac{\sqrt{3}}{2} \quad (14)$$

and

$$\frac{\text{dist}_F(A, \mathbb{H}_-)}{\|A\|_F} = \frac{\sqrt{3}}{2}. \quad (15)$$

Proof. The distance (11) is shown in [43, Section 5], and (12) can be shown similarly. Thus, the matrix A is equidistant to the sets \mathbb{H} and \mathbb{A} . The upper bound (13) for the distance to the set of normal matrices is achieved for a circulant matrix; see [43, Section 9]. The distance to \mathbb{H}_+ is given by

$$\text{dist}_F(A, \mathbb{H}_+) = \left(\sum_{\lambda_i(A_{\mathbb{H}}) < 0} \lambda_i^2(A_{\mathbb{H}}) + \|A_{\mathbb{A}}\|_F^2 \right)^{1/2}; \quad (16)$$

see Higham [31, Theorem 2.1]. Here $A_{\mathbb{H}} = (A + A^*)/2$ and $A_{\mathbb{A}} = (A - A^*)/2$ denote the Hermitian and skew-Hermitian parts of A , respectively, and $\lambda_1(A_{\mathbb{H}}), \dots, \lambda_m(A_{\mathbb{H}})$ are the eigenvalues of $A_{\mathbb{H}}$. We note that the distance in the Frobenius norm to the set \mathbb{H}_+ is the same as the distance to the set of Hermitian positive definite matrices. The eigenvalues of $A_{\mathbb{H}}$ are known to be

$$\lambda_j(A_{\mathbb{H}}) = \cos \frac{\pi j}{m+1}, \quad j = 1, 2, \dots, m; \quad (17)$$

see, e.g., [44, Section 2]. The expression (14) now follows from $\|A\|_F^2 = m - 1$, $\|A_{\mathbb{A}}\|^2 = (m - 1)/2$, and the fact that the sum in (16) evaluates to $(m - 1)/4$. Finally, (15) follows from

$$\text{dist}_F(A, \mathbb{H}_-) = \left(\sum_{\lambda_i(A_{\mathbb{H}}) > 0} \lambda_i^2(A_{\mathbb{H}}) + \|A_{\mathbb{A}}\|_F^2 \right)^{1/2}$$

and the fact that the eigenvalues (17) are allocated symmetrically with respect to the origin. \square

Proposition 2.1 shows the matrix (10) to be close to a normal matrix, and Example 2.1 illustrates that closeness to normality is not sufficient for GMRES to give an accurate approximation of the solution within a few iterations. Indeed, we can modify the matrix (10) to obtain a normal matrix and, as the following example shows, GMRES requires many iterations to solve the resulting linear system of equations when the order m of the matrix is large.

Example 2.2. Let A be a circulant matrix obtained by setting the $(1, m)$ -entry of the matrix (10) to one, and let the right-hand side \mathbf{b} be the same as in Example 2.1. Then the solution is $\mathbf{x}_{\text{exact}} = \mathbf{e}_1$. Similarly as in Example 2.1, GMRES yields the iterates $\mathbf{x}_k = \mathbf{0}$ for $1 \leq k < m$. The solution is not achieved until the iterate \mathbf{x}_m is computed. This depends on the distribution of the eigenvalues of A . A related example is presented by Nachtigal et al. [39]. We remark that the matrix A^*A is the identity, so the first iterate determined by LSQR with initial iterate $\mathbf{x}_0 = \mathbf{0}$ is $\mathbf{x}_{\text{exact}}$. Thus, LSQR performs much better than GMRES also for this example. \square

The construction of preconditioners for linear discrete ill-posed problems (1) with an error-contaminated right-hand side \mathbf{b} is delicate, because we do not want the preconditioner to give rise to severe propagation of the error \mathbf{e} in \mathbf{b} into the computed iterates; see Hanke et al. [26] for an insightful discussion on the construction of preconditioners for linear discrete ill-posed problems. Despite this difficulty and the dependence of convergence not only on the eigenvalue distribution of the system matrix, we have found that suitable preconditioners that make the system matrix be close to the set \mathbb{H}_+ and have eigenvalues that cluster in a small region in the complex plane, typically yield fairly rapid convergence. It is beneficial if the field of values of the preconditioned matrix does not contain the origin. Since the convergence of GMRES applied to the solution of (1) is invariant under multiplication of the matrix A by a complex rotation $e^{i\varphi}$, where $i = \sqrt{-1}$ and $-\pi < \varphi \leq \pi$, it suffices that the preconditioned matrix is close to a “rotated” Hermitian positive semidefinite matrix $e^{i\varphi}N$, where $N \in \mathbb{H}_+$. In the following we refer to the set of “rotated” Hermitian matrices as the set of *generalized Hermitian matrices*, and denote it by \mathbb{G} . It contains normal matrices $A \in \mathbb{C}^{m \times m}$, whose eigenvalues are collinear; see below. The set of “rotated” Hermitian positive semidefinite matrices is denoted by \mathbb{G}_+ and referred to as the set of *generalized Hermitian positive semidefinite matrices*.

3. Generalized Hermitian and Hermitian positive semidefinite matrices

In the following, we show some properties of generalized Hermitian and generalized Hermitian positive semidefinite matrices.

Proposition 3.1. *The matrix $A \in \mathbb{C}^{m \times m}$ is generalized Hermitian if and only if there exist $\varphi \in (-\pi, \pi]$ and $\alpha \in \mathbb{C}$ such that*

$$A = e^{i\varphi}B + \alpha I,$$

where $B \in \mathbb{C}^{m \times m}$ is an Hermitian matrix and $I \in \mathbb{C}^{m \times m}$ denotes the identity.

Proof. Let $A \in \mathbb{C}^{m \times m}$ be a generalized Hermitian matrix. Then there is a unitary matrix $U \in \mathbb{C}^{m \times m}$ such that $A = U\Lambda U^*$, where $\Lambda = \text{diag}[\lambda_1, \dots, \lambda_m] \in \mathbb{C}^{m \times m}$, and A has collinear eigenvalues, i.e., there exist $\varphi \in (-\pi, \pi]$ and $\alpha \in \mathbb{C}$ such that

$$\Lambda = e^{i\varphi}D + \alpha I,$$

where $D = \text{diag}[d_1, \dots, d_m] \in \mathbb{R}^{m \times m}$, so that $\lambda_i = e^{i\varphi}d_i + \alpha$ for $1 \leq i \leq m$. Thus, the matrix

$$B := e^{-i\varphi}(A - \alpha I) = e^{-i\varphi}(U\Lambda U^* - \alpha I) = U(e^{-i\varphi}(\Lambda - \alpha I))U^* = UDU^*$$

is Hermitian.

Conversely, if $B = e^{-i\varphi}(A - \alpha I)$ is Hermitian, then $B = UDU^*$, where $U \in \mathbb{C}^{m \times m}$ is unitary and $D \in \mathbb{R}^{m \times m}$ is diagonal. Hence, $A = e^{i\varphi}B + \alpha I = U(e^{i\varphi}D + \alpha I)U^*$ has collinear eigenvalues and is unitarily diagonalizable. \square

Proposition 3.2. *If the matrix $Z = [z_{i,j}] \in \mathbb{C}^{m \times m}$ is generalized Hermitian, then there exist $\theta \in (-\pi, \pi]$ and $\gamma \in \mathbb{R}$ such that*

$$z_{i,j} = \begin{cases} \bar{z}_{j,i} e^{i\theta}, & \text{if } i \neq j, \\ \bar{z}_{i,i} e^{i\theta} + \gamma e^{i\frac{\theta+\pi}{2}}, & \text{if } i = j, \end{cases}$$

where the bar denotes complex conjugation.

Proof. It follows from Proposition 3.1 that there exist an angle $\phi \in (-\pi, \pi]$ and a scalar $\beta \in \mathbb{C}$ such that $e^{i\phi}Z + \beta I$ is Hermitian, i.e.,

$$e^{i\phi}Z + \beta I = e^{-i\phi}Z^* + \bar{\beta}I.$$

Thus,

$$Z = e^{-2i\phi}Z^* - 2e^{i\frac{\pi-2\phi}{2}}\text{Im}(\beta)I.$$

Setting $\theta = -2\phi$ and $\gamma = -2\text{Im}(\beta)$ concludes the proof. \square

Proposition 3.3. Let $A = [a_{i,j}] \in \mathbb{C}^{m \times m}$. If

$$m \text{Trace}(A^2) \neq \text{Trace}(A), \quad (18)$$

then the unique closest generalized Hermitian matrix $\hat{A} = [\hat{a}_{i,j}] \in \mathbb{C}^{m \times m}$ to A in the Frobenius norm is given by

$$\hat{a}_{i,j} = \begin{cases} \frac{1}{2}(a_{i,j} + \bar{a}_{j,i}e^{i\hat{\theta}}), & \text{if } i \neq j, \\ \frac{1}{2}(a_{i,i} + \bar{a}_{i,i}e^{i\hat{\theta}} + \hat{\gamma}e^{i\frac{\hat{\theta}+\pi}{2}}), & \text{if } i = j, \end{cases} \quad (19)$$

where

$$\hat{\theta} = \arg(\text{Trace}(A^2) - \frac{1}{m}(\text{Trace}(A))^2), \quad \hat{\gamma} = \frac{2}{m} \text{Im}(e^{-i\frac{\hat{\theta}}{2}} \text{Trace}(A)).$$

Moreover, the distance of A to the set \mathbb{G} of generalized Hermitian matrices is given by

$$\text{dist}_F(A, \mathbb{G}) = \sqrt{\frac{\|A\|_F^2}{2} - \frac{1}{2} \text{Re} \left(e^{-i\hat{\theta}} \sum_{i,j=1}^m a_{i,j} a_{j,i} \right) - \frac{1}{m} \left(\text{Im} \left(e^{-i\frac{\hat{\theta}}{2}} \sum_{i=1}^m a_{i,i} \right) \right)^2}.$$

If (18) is violated, then there are infinitely many matrices $\hat{A}(\theta) = [\hat{a}_{i,j}(\theta)] \in \mathbb{C}^{m \times m}$, depending on an arbitrary angle θ , at the same minimal distance from A , whose entries are given by

$$\hat{a}_{i,j}(\theta) = \begin{cases} \frac{1}{2}(a_{i,j} + \bar{a}_{j,i}e^{i\theta}), & \text{if } i \neq j, \\ \frac{1}{2}(a_{i,i} + \bar{a}_{i,i}e^{i\theta} + \hat{\gamma}e^{i\frac{\theta+\pi}{2}}), & \text{if } i = j. \end{cases}$$

Proof. The entries of the generalized Hermitian matrix $Z(\theta, \gamma) = [z_{i,j}(\theta, \gamma)] \in \mathbb{C}^{m \times m}$, that minimizes the distance of A in the Frobenius norm from the set \mathbb{G} for the given angle θ and real γ , can be determined by minimizing $\|A - Z\|_F^2$, where the matrix $Z \in \mathbb{C}^{m \times m}$ is subject to the equality constraints of Proposition 3.2. Indeed, since $z_{i,j} = \bar{z}_{j,i}e^{i\theta}$, for $i > j$, the squared distance of A from Z reads

$$\|A - Z\|_F^2 = \sum_{\substack{i,j=1 \\ i < j}}^m (|z_{j,i} - a_{j,i}|^2 + |\bar{z}_{j,i}e^{i\theta} - a_{i,j}|^2) + \sum_{i=1}^m |z_{i,i} - a_{i,i}|^2. \quad (20)$$

Each term of the first sum in (20) can be written as $|z_{j,i} - a_{j,i}|^2 + |z_{j,i} - \bar{a}_{i,j}e^{i\theta}|^2$. Therefore, the sum is minimized by setting $z_{j,i} = \frac{1}{2}(a_{j,i} + \bar{a}_{i,j}e^{i\theta})$ for any $i < j$. Analogously, in the second sum in (20), since $z_{i,i} = \bar{z}_{i,i}e^{i\theta} + \gamma e^{i\frac{\theta+\pi}{2}}$, one minimizes each term by setting $z_{i,i} = \frac{1}{2}(a_{i,i} + \bar{a}_{i,i}e^{i\theta} + \gamma e^{i\frac{\theta+\pi}{2}})$. We conclude that the entries of $Z(\theta, \gamma)$ are given by

$$z_{i,j}(\theta, \gamma) = \begin{cases} \frac{1}{2}(a_{i,j} + \bar{a}_{j,i}e^{i\theta}), & \text{if } i \neq j, \\ \frac{1}{2}(a_{i,i} + \bar{a}_{i,i}e^{i\theta} + \gamma e^{i\frac{\theta+\pi}{2}}), & \text{if } i = j. \end{cases}$$

Substituting these values into $\|A - Z\|_F$ yields

$$\begin{aligned} d(\theta, \gamma) = \|A - Z(\theta, \gamma)\|_F^2 &= \frac{1}{4} \sum_{\substack{i,j=1 \\ i \neq j}}^m |a_{i,j} - \bar{a}_{j,i}e^{i\theta}|^2 + \frac{1}{4} \sum_{i=1}^m |a_{i,i} - (\bar{a}_{i,i}e^{i\theta} + \gamma e^{i\frac{\theta+\pi}{2}})|^2 \\ &= \frac{\|A\|_F^2}{2} + \frac{m}{4}\gamma^2 - \frac{1}{2} \text{Re} \left(e^{-i\theta} \sum_{i,j=1}^m a_{i,j} a_{j,i} \right) - \gamma \text{Im} \left(e^{-i\frac{\theta}{2}} \sum_{i=1}^m a_{i,i} \right). \end{aligned}$$

The desired values of θ and γ are determined by minimizing $d(\theta, \gamma)$. It follows that $\partial d(\theta, \gamma)/\partial \gamma = 0$ if and only if

$$\gamma = \hat{\gamma}(\theta) = \frac{2}{m} \operatorname{Im} \left(e^{-i\frac{\theta}{2}} \sum_{i=1}^m a_{i,i} \right).$$

Thus, we obtain

$$d(\theta, \hat{\gamma}(\theta)) = \frac{\|A\|_F^2}{2} - \frac{1}{2} \operatorname{Re} \left(e^{-i\theta} \sum_{i,j=1}^m a_{i,j} a_{j,i} \right) - \frac{1}{m} \left(\operatorname{Im} \left(e^{-i\frac{\theta}{2}} \sum_{i=1}^m a_{i,i} \right) \right)^2.$$

It follows that $d'(\theta, \hat{\gamma}(\theta)) = 0$ if and only if

$$\left(\operatorname{Re}(w_1) - \frac{1}{m} \operatorname{Re}(w_2^2) \right) \sin \theta = \left(\operatorname{Im}(w_1) - \frac{1}{m} \operatorname{Im}(w_2^2) \right) \cos \theta,$$

where $w_1 = \sum_{i,j=1}^m a_{i,j} a_{j,i}$ and $w_2 = \sum_{i=1}^m a_{i,i}$. Thus, if $m w_1 \neq w_2^2$, one has

$$\hat{\theta} = \arg(w_1 - \frac{1}{m} w_2^2).$$

This concludes the proof. \square

Corollary 3.1. Let the matrix $A = [a_{i,j}] \in \mathbb{C}^{m \times m}$ have trace zero. If

$$\sum_{i,j=1}^m a_{i,j} a_{j,i} \neq 0, \tag{21}$$

then the unique closest generalized Hermitian matrix $\hat{A} = [\hat{a}_{i,j}] \in \mathbb{C}^{m \times m}$ to A in the Frobenius norm is given by

$$\hat{a}_{i,j} = \frac{1}{2} (a_{i,j} + \bar{a}_{j,i} e^{i\hat{\theta}}),$$

where

$$\hat{\theta} = \arg \left(\sum_{i,j=1}^m a_{i,j} a_{j,i} \right).$$

Moreover,

$$\operatorname{dist}_F(A, \mathbb{G}) = \sqrt{\frac{\|A\|_F^2 - |\sum_{i,j=1}^m a_{i,j} a_{j,i}|}{2}}.$$

If (21) is violated, then there are infinitely many matrices $\hat{A}(\theta) = [\hat{a}_{i,j}(\theta)] \in \mathbb{C}^{m \times m}$, depending on an arbitrary angle θ , at the same minimal distance from A , namely

$$\hat{a}_{i,j}(\theta) = \frac{a_{i,j} + \bar{a}_{j,i} e^{i\theta}}{2}, \quad 1 \leq i, j \leq m.$$

Proof. The result follows by observing that the optimal values of $\hat{\theta}$ and $\hat{\gamma}$ determined by Proposition 3.3 are given by $\hat{\theta} = \arg(\operatorname{Trace}(A^2))$ and $\hat{\gamma} = 0$. \square

We refer to a generalized Hermitian matrix $A \in \mathbb{C}^{m \times m}$, whose eigenvalues for suitable $\varphi \in (-\pi, \pi]$ and $\alpha \in \mathbb{C}$ satisfy

$$\lambda_i = \rho_i e^{i\varphi} + \alpha, \quad \text{with } \rho_i \geq 0, \quad 1 \leq i \leq m,$$

as a generalized Hermitian positive semidefinite matrix.

Proposition 3.4. The matrix $A \in \mathbb{C}^{m \times m}$ is generalized Hermitian positive semidefinite if and only if there are constants $\varphi \in (-\pi, \pi]$ and $\alpha \in \mathbb{C}$ such that

$$A = e^{i\varphi} B + \alpha I,$$

where the matrix $B \in \mathbb{C}^{m \times m}$ is Hermitian positive semidefinite.

Proof. The proposition follows from the proof of Proposition 3.1, where we use the fact that the diagonal entries of the diagonal matrix D are nonnegative. \square

We are interested in measuring the distance between A and the set \mathbb{G}_+ of generalized Hermitian positive semidefinite matrices in the Frobenius norm. We deduce from (19) that, if (18) holds, then the unique closest generalized Hermitian matrix is of the form

$$\widehat{A} = \frac{A + e^{i\hat{\theta}} A^* + \widehat{\gamma} e^{i\frac{\hat{\theta}+\pi}{2}} I}{2} = e^{i\frac{\hat{\theta}}{2}} \widetilde{A} + \frac{\widehat{\gamma}}{2} e^{i\frac{\hat{\theta}+\pi}{2}} I, \quad (22)$$

where \widetilde{A} denotes the Hermitian part of $e^{-i\frac{\hat{\theta}}{2}} A$. The identity (22) shows that the unique closest generalized Hermitian positive semidefinite matrix to A can be written as

$$\widehat{A}_+ := e^{i\frac{\hat{\theta}}{2}} \widetilde{A}_+ + \frac{\widehat{\gamma}}{2} e^{i\frac{\hat{\theta}+\pi}{2}} I,$$

where \widetilde{A}_+ denotes the Hermitian positive semidefinite matrix closest to $e^{-i\frac{\hat{\theta}}{2}} A$. The construction of \widetilde{A}_+ can be easily obtained following [31]. Thus, the distance

$$\text{dist}_F(A, \mathbb{G}_+) = \|A - \widehat{A}_+\|_F \geq \text{dist}_F(A, \mathbb{G})$$

can be computed similarly as (16), taking into account the squared sum of the negative eigenvalues of \widetilde{A} .

4. Some preconditioning techniques

Preconditioning is a popular technique to improve the rate of convergence of GMRES when applied to the solution of many linear systems of equations, including those obtained by the discretization of well-posed problems; see, e.g., [37,56] for discussions and references. This technique replaces a linear system of equations (1) by a left-preconditioned system

$$MAx = Mb \quad (23)$$

or by a right-preconditioned system

$$AMy = b, \quad x := My, \quad (24)$$

and applies GMRES to the solution of one of these preconditioned systems. The matrix $M \in \mathbb{C}^{m \times m}$ is referred to as a preconditioner. In the well-posed setting, M typically is chosen so that the iterates generated by GMRES when applied to (23) or (24) converge to the solution faster than the iterates determined by GMRES applied to the original (unpreconditioned) linear system of equations (1). One would like M to have a structure that allows rapid evaluation of matrix-vector products My , $y \in \mathbb{C}^m$. Left- and right-preconditioners may be applied simultaneously.

Preconditioning also can be applied to the solution of linear discrete ill-posed problems (1); see, e.g., [14,16,26,29,47,54]. The aim of the preconditioner M in this context is to determine a solution subspace $\mathbb{K}_k(MA, Mb)$ for problem (23), or a solution subspace $M\mathbb{K}_k(AM, b)$ for problem (24), that contain accurate approximations of x_{exact} already when their dimension k is small. Moreover, we would like to choose M so that the error e in b is not severely amplified and propagated into the computed iterates when solving (23) or (24). We seek to achieve these goals by choosing particular preconditioners M such that the matrices MA or AM are close to the sets \mathbb{H}_+ or \mathbb{G}_+ . We will comment on the distance of these matrices to the sets \mathbb{H} and \mathbb{A} . Right-preconditioning generally is more useful than left-preconditioning, because the GMRES residual norm for the system (24) can be cheaply evaluated by computing the residual norm of a low-dimensional system of equations, similarly as for unpreconditioned systems; cf. (5) for the latter. Being able to inexpensively compute the norm of the residual error is a favorable feature when a stopping criterion based on the residual norm is used, such as the discrepancy principle. Henceforth, we focus on right-preconditioning. We describe several novel approaches to construct a preconditioner that can be effective in a variety of situations.

When the matrix A is a shift operator, GMRES may not be able to deliver an accurate approximation of x_{exact} within a few iterations (this is the case of Example 2.1). To remedy this difficulty, we propose to approximate A by a circulant matrix C_A . We may, for instance, determine C_A as the solution of the matrix nearness problem discussed in [10,11,42],

$$\min_{C \in \mathbb{C}^{m \times m} \text{ circulant}} \|C - A\|_F, \quad (25)$$

and use the preconditioner

$$M = C_A^{-1}.$$

The minimization problem (25) can be solved easily by using the spectral factorization

$$C_A = W D_A W^*, \quad (26)$$

where the matrix $D_A \in \mathbb{C}^{m \times m}$ is diagonal and $W \in \mathbb{C}^{m \times m}$ is a unitary fast Fourier transform (FFT) matrix; see [13] for details. Hence,

$$\|C_A - A\|_F = \|D_A - W^* A W\|_F,$$

and it follows that D_A is made up of the diagonal entries of $W^* A W$. The computation of the matrix D_A , with the aid of the FFT, requires $\mathcal{O}(m^2 \log_2(m))$ arithmetic floating point operations (flops); see [10,11,42] for details. Alternatively, a circulant preconditioner may be computed as the solution of the matrix nearness problem

$$\min_{C \in \mathbb{C}^{m \times m} \text{ circulant}} \|I - C^{-1} A\|_F. \quad (27)$$

This minimization problem is discussed in [16,58,59]. The solution is given by $C_{AA^*} C_{A^*}^{-1}$; see [59]. The flop count for solving (27), by using the FFT, also is $\mathcal{O}(m^2 \log_2(m))$; see [10,42,59].

A cheaper way to determine a circulant preconditioner (26) is to let $\mathbf{x} \in \mathbb{C}^m$ be a vector with normally distributed random entries with zero mean, define $\mathbf{y} := A\mathbf{x}$, and then determine the diagonal matrix D_A in (26) by requiring that $\mathbf{y} = C_A \mathbf{x}$. This gives

$$D_A = \text{diag}[(W^* \mathbf{y}) / (W^* \mathbf{x})], \quad (28)$$

where the vector division is component-wise. The computation of D_A in this way only requires the evaluation of two fast Fourier transforms and m scalar divisions, which only demands $\mathcal{O}(m \log_2(m))$ flops. We remark that further approaches to construct circulant preconditioners are discussed in the literature; see [10,42]. Moreover, $e^{i\theta}$ -circulants, which allow an angle θ as an auxiliary parameter, can be effective preconditioners: they generalize the preconditioners (25) and (27), and also can be constructed with $\mathcal{O}(m^2 \log_2(m))$ flops; see [45,47]. Having determined the preconditioner M , we apply the Arnoldi process to the matrix AM with initial vector \mathbf{b} . The evaluation of each matrix-vector product with a circulant or an $e^{i\theta}$ -circulant matrix M requires only $\mathcal{O}(m \log_2(m))$ flops when using the FFT. Iterations are carried out until the discrepancy principle is satisfied. Let \mathbf{y}_k be the approximate solution of (24) so obtained. Then $\mathbf{x}_k = M\mathbf{y}_k$ is an approximation of $\mathbf{x}_{\text{exact}}$.

A generic approach to determine a preconditioner M that makes AM closer to the set \mathbb{H}_+ than A is to carry out k_p steps of the Arnoldi process applied to A with initial vector \mathbf{b} . Assuming that no breakdown occurs, this yields a decomposition of the form (4) with k replaced by k_p , and we define the approximation

$$A_{k_p} := V_{k_p+1} H_{k_p+1, k_p} V_{k_p}^* \quad (29)$$

of A . If A_{k_p} contains information about the dominant singular values of the matrix A only, then A_{k_p} is a regularized approximation of A . This property is illustrated numerically in [21] for severely ill-conditioned matrices. Moreover, in a continuous setting and under the assumption that A is a Hilbert–Schmidt operator of infinite rank [55, Chapter 2], it is shown in [48] that the SVD of A can be approximated by computing an Arnoldi decomposition of A . This property is inherited in the discrete setting of the present paper, whenever a suitable discretization of a Hilbert–Schmidt operator is used.

The approximation (29) suggests the simple preconditioner

$$M := A_{k_p}^*. \quad (30)$$

The rank of this preconditioner is at most k_p and, therefore, GMRES applied to the solution of (24) will break down within k_p steps; see, e.g., [3,53] for discussions on GMRES applied to linear systems of equations with a singular matrix. We would like to choose k_p large enough so that GMRES applied to (24) yields a sufficiently accurate approximation of $\mathbf{x}_{\text{exact}}$ within k_p steps. The following proposition sheds light on some properties of the matrix AM when M is defined by (30).

Proposition 4.1. Assume that k_p steps of the Arnoldi process applied to A with initial vector \mathbf{b} can be carried out without breakdown, and let the preconditioner M be defined by (30). Then AM is Hermitian positive semidefinite with rank at most k_p , and $\mathcal{R}(AM) \subset \mathcal{R}(V_{k_p+1})$.

Proof. From (29) and the decomposition (4), with k replaced by k_p , it is immediate to verify that

$$\begin{aligned} AM &= AA_{k_p}^* = AV_{k_p} H_{k_p+1, k_p}^* V_{k_p+1}^* \\ &= V_{k_p+1} H_{k_p+1, k_p} H_{k_p+1, k_p}^* V_{k_p+1}^* = C_{k_p+1, k_p} C_{k_p+1, k_p}^*, \end{aligned}$$

where

$$C_{k_p+1, k_p} = V_{k_p+1} H_{k_p+1, k_p} \in \mathbb{C}^{m \times k_p} \quad (31)$$

is a matrix of rank at most k_p . Finally, for any $\mathbf{z} \in \mathbb{C}^m$, we have

$$AM\mathbf{z} = V_{k_p+1} H_{k_p+1, k_p} H_{k_p+1, k_p}^* V_{k_p+1}^* \mathbf{z}.$$

This shows that $AM\mathbf{z} \in \mathcal{R}(V_{k_p+1})$. \square

Since $AA_{k_p}^*$ is singular, problem (24) should be solved in the least-squares sense, i.e., instead of solving (24) one should compute

$$\mathbf{y} = \arg \min_{\hat{\mathbf{y}} \in \mathbb{C}^m} \|C_{k_p+1, k_p} C_{k_p+1, k_p}^* \hat{\mathbf{y}} - \mathbf{b}\|, \quad \mathbf{x} = A_{k_p}^* \mathbf{y}, \quad (32)$$

where C_{k_p+1, k_p} is defined by (31). It follows from the definition (29) of A_{k_p} , and the fact that $\mathcal{R}(V_{k_p}) = \mathbb{K}_{k_p}(A, \mathbf{b})$, that the solution \mathbf{x} of (32) belongs to $\mathbb{K}_{k_p}(A, \mathbf{b})$. A regularized solution of the minimization problem (32) can be determined in several ways. For instance, one can apply a few steps of the Arnoldi process (Algorithm 1) to compute an approximate solution of the least-squares problem (32), i.e., one applies the Arnoldi process to the matrix $C_{k_p+1, k_p} C_{k_p+1, k_p}^*$ with initial vector $\mathbf{v}_1 = \mathbf{b}/\|\mathbf{b}\|$. We note that the latter application of the Arnoldi process does not require additional matrix-vector product evaluations with the matrix A . Alternatively, we may determine a regularized solution of (32) by using Tikhonov regularization or by truncated singular value decomposition (TSVD) of the matrix C_{k_p+1, k_p} . We will discuss the latter regularization techniques in detail in Section 5. Computational experiments, some of which are reported in Section 6, indicate that it is often possible to determine a meaningful approximation of $\mathbf{x}_{\text{exact}}$ by computing a regularized solution of (32) even when GMRES applied to the original problem (1) yields a poor approximation of $\mathbf{x}_{\text{exact}}$.

The approximation (29) of A also can be used to define the preconditioner

$$M := A_{k_p}^* + (I - V_{k_p} V_{k_p}^*) = V_{k_p} H_{k_p+1, k_p}^* V_{k_p+1}^* + (I - V_{k_p} V_{k_p}^*). \quad (33)$$

The number of steps k_p should be chosen so that the matrix AM is fairly close to the set \mathbb{H}_+ . The preconditioned coefficient matrix defined by this preconditioner,

$$\begin{aligned} AM &= AA_{k_p}^* + A(I - V_{k_p} V_{k_p}^*) \\ &= V_{k_p+1} H_{k_p+1, k_p} H_{k_p+1, k_p}^* V_{k_p+1}^* + A(I - V_{k_p} V_{k_p}^*), \end{aligned} \quad (34)$$

is non-Hermitian. A few steps of the Arnoldi process (Algorithm 1) can be applied to the matrix (34) to determine a regularized solution of (24). However, differently from the situation when using the preconditioner (30), this requires additional matrix-vector product evaluations with A . Regularization of (24) when the preconditioner is defined by (33) can again be achieved by applying Tikhonov or TSVD regularization. An analogue of Proposition 4.1 does not hold for the preconditioner M defined by (33). Instead, we can show the following result.

Proposition 4.2. Assume that $k_p + j$ steps of the Arnoldi process applied to A with initial vector \mathbf{b} can be carried out without breakdown, and let the preconditioner M be defined by (33). Then the iterate \mathbf{y}_j determined at the j th step of GMRES applied to the preconditioned system (24) with initial approximate solution $\mathbf{y}_0 = \mathbf{0}$ belongs to the Krylov subspace $\mathbb{K}_{k_p+j}(A, \mathbf{b})$.

Proof. We show the proposition by induction. It is immediate to verify that

$$\mathbf{y}_1 \in \mathbb{K}_1(AM, \mathbf{b}) = \text{span}\{\mathbf{b}\} = \mathbb{K}_1(A, \mathbf{b}) \subset \mathbb{K}_{k_p}(A, \mathbf{b}) \subset \mathbb{K}_{k_p+1}(A, \mathbf{b}).$$

Assume that $\mathbf{y}_i \in \mathbb{K}_{k_p+i}(A, \mathbf{b})$. Then, since

$$\mathbf{y}_{i+1} \in \mathbb{K}_{i+1}(AM, \mathbf{b}) \subset \text{span}\{\mathbf{b}, AM\mathbb{K}_{k_p+i}(A, \mathbf{b})\},$$

\mathbf{y}_{i+1} is a linear combination of vectors of this subspace, i.e.,

$$\mathbf{y}_{i+1} = s_1 \mathbf{b} + AM V_{k_p+i} \mathbf{s}_{k_p+i} = s_1 \mathbf{b} + V_{k_p+1} \mathbf{s}_{k_p+1} + V_{k_p+i+1} \mathbf{s}_{k_p+i+1},$$

where $s_1 \in \mathbb{C}$, $\mathbf{s}_{k_p+1} \in \mathbb{C}^{k_p+1}$, $\mathbf{s}_{k_p+i} \in \mathbb{C}^{k_p+i}$, and $\mathbf{s}_{k_p+i+1} \in \mathbb{C}^{k_p+i+1}$. Here we have used the definition (33) of M and the Arnoldi decomposition (4), with k replaced by k_p . Hence, $\mathbf{y}_{i+1} \in \mathcal{R}(V_{k_p+i+1}) = \mathbb{K}_{k_p+i+1}(A, \mathbf{b})$. \square

Assume that the conditions of Proposition 4.2 hold, and let $\mathbf{y}_j = V_{k_p+j} \mathbf{s}_{k_p+j}$ with $\mathbf{s}_{k_p+j} \in \mathbb{C}^{k_p+j}$. Then the corresponding approximate solution \mathbf{x}_j of (24) satisfies

$$\mathbf{x}_j = M\mathbf{y}_j = M V_{k_p+j} \mathbf{s}_{k_p+j} \in \mathcal{R}(V_{k_p+j}) = \mathbb{K}_{k_p+j}(A, \mathbf{b}).$$

Hence, application of the GMRES method with the right-preconditioner (33) determines an approximate solution in the (unpreconditioned) Krylov subspace $\mathbb{K}_{k_p+j}(A, \mathbf{b})$.

We conclude this section by considering two more preconditioners that are related to (30) and (33). They are not designed with the aim of making the preconditioned matrix close to the sets \mathbb{H} or \mathbb{H}_+ . Assume, as above, that the Arnoldi process does not break down during the first k_p steps. Then the matrix A_{k_p} defined by (29) can be computed, and one may use

$$M := A_{k_p} \quad (35)$$

as a preconditioner. Similarly to (30), this preconditioner has rank at most k_p and, assuming that A_{k_p} only contains information about the k_p dominant singular values of A , M may be regarded as a regularized approximation of A . Note that, by exploiting the Arnoldi decomposition (4) with k replaced by $k_p + 1$, one obtains the expression

$$AM = V_{k_p+2} H_{k_p+2, k_p+1} H_{k_p+1, k_p} V_{k_p}^*. \quad (36)$$

We note that when applying a few (at most k_p) steps of GMRES to compute an approximate solution of the preconditioned system (24), no additional matrix-vector product evaluations with the matrix A are necessary, in addition to the $k_p + 1$ matrix-vector product evaluations required to determine the right-hand side of (36). The iterate \mathbf{x}_j determined at the j th step of GMRES applied to the preconditioned system (24) belongs to $\mathcal{R}(V_{k_p+2}) = \mathbb{K}_{k_p+2}(A, \mathbf{b})$.

The preconditioner

$$M := A_{k_p} + (I - V_{k_p} V_{k_p}^*) = V_{k_p+1} H_{k_p+1, k_p} V_{k_p}^* + (I - V_{k_p} V_{k_p}^*) \quad (37)$$

is analogous to (33). This preconditioner also was considered in [36] in the framework of the solution of a sequence of slowly-varying linear systems of equations. Similarly to (35), the preconditioner (37) is not designed to make the preconditioned matrix AM close to the set \mathbb{H}_+ . By using the Arnoldi decomposition (4) with k replaced by $k_p + 1$, we obtain

$$AM = AA_{k_p} + A(I - V_{k_p} V_{k_p}^*) = V_{k_p+2} H_{k_p+2, k_p+1} H_{k_p+1, k_p} V_{k_p}^* + A(I - V_{k_p} V_{k_p}^*).$$

It is evident that, even though $k_p + 1$ steps of the Arnoldi process have been carried out to define M , additional matrix-vector products with A are required when applying the Arnoldi process to the preconditioned system (24). Using the same arguments as in Proposition 4.2, one can show that, if $k_p + j$ steps of the Arnoldi process applied to A with initial vector \mathbf{b} can be carried out without breakdown, then the iterate \mathbf{y}_j determined at the j th iteration of GMRES applied to the preconditioned system (24) and the corresponding approximate solution $\mathbf{x}_j = M\mathbf{y}_j$ of (1) belong to $\mathbb{K}_{k_p+j}(A, \mathbf{b})$. We note that Tikhonov or TSVD regularization can be applied when solving the preconditioned system (24) with either one of the preconditioners (35) or (37).

5. Solving the preconditioned problems

As already suggested in the previous section, instead of using GMRES to solve the preconditioned system (24) with one of the preconditioners described, one may wish to apply additional regularization in order to determine an approximate solution of (1) of higher quality. In the following we discuss application of Tikhonov and TSVD regularization. We refer to the solution methods so obtained as the Arnoldi-Tikhonov and Arnoldi-TSVD methods, respectively. Due to the additional regularization, both these methods allow the use of a solution subspace of larger dimension than preconditioned GMRES without additional regularization. This may result in computed approximations of $\mathbf{x}_{\text{exact}}$ of higher quality.

The Arnoldi-Tikhonov method for (24) determines an approximate solution \mathbf{x}_μ of (1) by first computing the solution \mathbf{y}_μ of the Tikhonov minimization problem

$$\min_{\mathbf{y} \in \mathbb{K}_k(AM, \mathbf{b})} \{\|AM\mathbf{y} - \mathbf{b}\|^2 + \mu \|\mathbf{y}\|^2\}, \quad (38)$$

where $\mu > 0$ is a regularization parameter to be specified, and then evaluates the approximation $\mathbf{x}_\mu = M\mathbf{y}_\mu$ of $\mathbf{x}_{\text{exact}}$. The minimization problem (38) has a unique solution for any $\mu > 0$. Application of k steps of the Arnoldi process to the matrix AM with initial vector \mathbf{b} gives the Arnoldi decomposition

$$AMV_k = V_{k+1} H_{k+1, k}, \quad (39)$$

which is analogous to (4). Using (39), the minimization problem (38) can be expressed as the reduced Tikhonov minimization problem

$$\min_{\mathbf{z} \in \mathbb{C}^k} \{\|H_{k+1, k} \mathbf{z} - \|\mathbf{b}\| \mathbf{e}_1\|^2 + \mu \|\mathbf{z}\|^2\}, \quad (40)$$

whose minimizer \mathbf{z}_μ gives the approximate solution $\mathbf{y}_\mu := V_k \mathbf{z}_\mu$ of (38), so that $\mathbf{x}_\mu := M\mathbf{y}_\mu$ is an approximate solution of (1).

The Arnoldi-TSVD method seeks to determine an approximate solution of (24) by using a truncated singular value decomposition of the (small) matrix $H_{k+1, k}$ in (39). Let $\mathbf{y} := V_k \mathbf{z}$. Then, using (39), we obtain

$$\min_{\mathbf{y} \in \mathbb{K}_k(AM, \mathbf{b})} \|AM\mathbf{y} - \mathbf{b}\| = \min_{\mathbf{z} \in \mathbb{C}^k} \|H_{k+1,k}\mathbf{z} - \|\mathbf{b}\|\mathbf{e}_1\|. \quad (41)$$

Let $H_{k+1,k} = U_{k+1}\Sigma_k W_k^*$ be the singular value decomposition. Thus, the matrices $U_{k+1} \in \mathbb{C}^{(k+1) \times (k+1)}$ and $W_k \in \mathbb{R}^{k \times k}$ are unitary, and

$$\Sigma_k = \text{diag}[\sigma_1^{(k)}, \sigma_2^{(k)}, \dots, \sigma_k^{(k)}] \in \mathbb{R}^{(k+1) \times k}$$

is diagonal (and rectangular), with nonnegative diagonal entries ordered according to $\sigma_1^{(k)} \geq \sigma_2^{(k)} \geq \dots \geq \sigma_k^{(k)} \geq 0$. Define the diagonal matrix

$$\Sigma_k^{(j)} = \text{diag}[\sigma_1^{(k)}, \dots, \sigma_j^{(k)}, 0, \dots, 0] \in \mathbb{R}^{(k+1) \times k}$$

by setting the $k-j$ last diagonal entries of Σ_k to zero, where we assume that j is small enough so that $\sigma_j^{(k)} > 0$. Introduce the associated rank- j matrix $H_{k+1,k}^{(j)} := U_{k+1}\Sigma_k^{(j)}W_k^*$. Let $\mathbf{z}^{(j)}$ denote the minimal norm solution of

$$\min_{\mathbf{z} \in \mathbb{C}^k} \|H_{k+1,k}^{(j)}\mathbf{z} - \|\mathbf{b}\|\mathbf{e}_1\|. \quad (42)$$

Problem (42) describes the truncated singular value decomposition (TSVD) method applied to the solution of the reduced minimization problem in the right-hand side of (41); see, e.g., [19,27] for further details on the TSVD method. Once the solution $\mathbf{z}^{(j)}$ of (42) is computed, we get the approximate solution $\mathbf{y}^{(j)} := V_k\mathbf{z}^{(j)}$ of (41), from which we obtain the approximate solution $\mathbf{x}^{(j)} := M\mathbf{y}^{(j)}$ of (1). A modified TSVD method described in [46] also can be used.

All the solution schemes discussed in this section are inherently multi-parameter regularization methods, i.e., their success depends of the appropriate tuning of more than one regularization parameter. In the remainder of this section we will discuss reliable strategies to effectively choose these parameters. First of all, when one of the preconditioners (30), (33), (35), or (37) is used, an initial number of Arnoldi iterations, k_p , has to be carried out. Since we would like the preconditioners M to be suitable regularized approximations of the matrix A , a natural way to determine k_p is to monitor the expansion of the Krylov subspace $\mathbb{K}_{k_p}(A, \mathbf{b})$. The subdiagonal elements $h_{i+1,i}$, $i = 1, 2, \dots, k$, of the Hessenberg matrix $H_{k+1,k} = [h_{i,j}] \in \mathbb{C}^{(k+1) \times k}$ in (4) are helpful in this respect; see [22,49]. We terminate the initial Arnoldi process as soon as an index k_p such that

$$h_{k_p+1,k_p} < \tau'_1 \quad \text{and} \quad \frac{|h_{k_p+1,k_p} - h_{k_p,k_p-1}|}{h_{k_p,k_p-1}} > \tau''_1, \quad (43)$$

for certain user-specified parameters τ'_1 and τ''_1 , is found. By choosing τ'_1 small, we require some stabilization to take place while generating the Krylov subspace $\mathbb{K}_k(A, \mathbf{b})$; simultaneously, by setting τ''_1 close to 1, we require the subdiagonal entries of $H_{k+1,k}$ to stabilize. In terms of regularization, this criterion is partially justified by the bound

$$\prod_{j=1}^{k_p} h_{j+1,j} \leq \prod_{j=1}^{k_p} \sigma_j,$$

see [38], which states that, on geometric average, the sequence $\{h_{j+1,j}\}_{j \geq 1}$ decreases faster than the singular values. Numerical experiments reported in [22] indicate that the quantity $\|A - V_{k+1}H_{k+1,k}V_k^*\|$ decreases to zero as k increases with about the same rate as the singular values of A . More precisely, even though no theoretical results are available at present, one can experimentally verify that typically

$$\|A - V_{k+1}H_{k+1,k}V_k^*\| \simeq \sigma_{k+1}^{(k+1)},$$

where $\sigma_{k+1}^{(k+1)}$ is the $(k+1)$ st singular value of $H_{k+2,k+1}$ ordered in decreasing order. Here $\|\cdot\|$ denotes the spectral norm of the matrix.

Note that there is no guarantee that the above estimate is tight: Firstly, we would have equality only if the matrices $V_{k+1}U_{k+1}$ and V_kW_k coincide with the matrices made up by the left and right singular vectors, respectively, of the TSVD of the matrix A . If this is not the case, then we may have $\|A - V_{k+1}H_{k+1,k}V_k^*\| \gg \sigma_{k+1}^{(k+1)}$. Secondly, one cannot guarantee that $\sigma_{k+1}^{(k+1)} \geq \sigma_{k+1}$. Nevertheless, experimentally it appears reliable to terminate the Arnoldi iterations when the product $p_\sigma^{(k)} := \sigma_1^{(k)}\sigma_{k+1}^{(k+1)}$ is sufficiently small, i.e., one should stop as soon as

$$p_\sigma^{(k_p)} := \sigma_1^{(k_p)}\sigma_{k_p+1}^{(k_p+1)} < \tau_2, \quad (44)$$

where τ_2 is a user-specified threshold.

Once the preconditioner M has been determined, other regularization parameters should be suitably chosen: Namely, the number of preconditioned Arnoldi iterations and, in case the Arnoldi–Tikhonov (40) or Arnoldi–TSVD (42) methods are considered, one also has to determine a value for the regularization parameter $\mu > 0$ or truncation parameter $j \in$

\mathbb{N} , respectively. Since choosing the number of Arnoldi iterations is less critical (i.e., one can recover good approximate solutions provided that suitable values for μ or j are set at each iteration), we propose that a maximum allowed number of preconditioned Arnoldi iterations be carried out, and we apply the discrepancy principle to determine the parameters μ or j . Specifically, when using the Arnoldi–Tikhonov method, we choose μ so that the computed solution \mathbf{x}_μ satisfies the discrepancy principle

$$\|\mathbf{A}\mathbf{x}_\mu - \mathbf{b}\| = \tau\delta \quad (45)$$

to avoid severe propagation of the noise \mathbf{e} into \mathbf{x}_μ . We remark that this μ -value can be computed quite rapidly by substituting the Arnoldi decomposition (39) into (45); see [9,22,35] for discussions on unpreconditioned Tikhonov regularization. There also are other approaches to determining the regularization parameter; see, e.g., [34,52].

When applying the Arnoldi-TSVD method, we choose j as small as possible so that the discrepancy principle is satisfied, i.e.,

$$\|H_{k+1,k}^{(j)}\mathbf{z}^{(j)} - \|\mathbf{b}\|\mathbf{e}_1\| \leq \tau\delta, \quad (46)$$

and tacitly assume that $j < k$; otherwise k has to be increased. For most reasonable values of τ and δ , equations (45) and (46) have a unique solution $\mu > 0$ and $j > 0$, respectively.

6. Computed examples

This section illustrates the performance of the preconditioners introduced in Section 4 used with GMRES, or with the Arnoldi–Tikhonov and Arnoldi-TSVD methods described in Section 5. The Arnoldi algorithm is implemented with re-orthogonalization. A first set of experiments considers moderate-scale test problems from [28], and takes into account the preconditioners described in the second part of Section 4 only. A second set of experiments considers realistic large-scale problems arising in the framework of 2D image deblurring, and also includes comparisons with circulant preconditioners. Comparisons with the unpreconditioned counterparts of these methods are presented. All the computations were carried out in MATLAB R2016b on a single processor 2.2 GHz Intel Core i7 computer.

To keep the notations light, we let C_1 , C_2 , and C_3 denote the preconditioners obtained by solving (25), (27), and (28), respectively. Also, we let M_1 , M_2 , M_3 , and M_4 be the preconditioners in (30), (33), (35), and (37), respectively. The unpreconditioned GMRES, Arnoldi–Tikhonov, and Arnoldi-TSVD methods are referred to as “GMRES”, “Tikh”, and “TSVD”, respectively; their preconditioned counterparts are denoted by “GMRES(P_a)”, “Tikh(P_a)”, and “TSVD(P_a)”, where $P \in \{C, M\}$ and $a \in \{1, 2, 3, 4\}$. In the following graphs, specific markers are used for the different preconditioners: ‘o’ denotes C_1 , ‘□’ denotes C_2 , ‘<’ denotes C_3 , ‘◇’ denotes M_4 , and ‘*’ indicates that no preconditioner is used. For some test problems, we report results for LSQR, with associated marker ‘+’. The stopping criteria (43), (44), (45), and (46) are used with the parameters $\tau'_1 = 10^{-4}$, $\tau''_1 = 0.9$, $\tau_2 = 10^{-10}$, and $\tau = 1.01$. Clearly this choice does not work properly for every problem. In general, a reliable choice of these parameters, as well as the choice of the maximum number of Arnoldi iterations to compute the preconditioners M_i , $i = 1, \dots, 4$, is closely related to the decay rate of the singular values and the quality of their approximation. We use the relative reconstruction error norm, defined by $\|\mathbf{x}_{\text{exact}} - \mathbf{x}_k\|/\|\mathbf{x}_{\text{exact}}\|$ or $\|\mathbf{x}_{\text{exact}} - \mathbf{x}_\mu\|/\|\mathbf{x}_{\text{exact}}\|$, as a measure of the reconstruction quality.

6.1. First set of experiments

We consider problems (1) with a real nonsymmetric coefficient matrix of size $m = 200$ and a right-hand side vector that is affected by Gaussian white noise \mathbf{e} , with relative noise level $\|\mathbf{e}\|/\|\mathbf{b}\| = 10^{-2}$. For all the tests, the maximum allowed number of Arnoldi iterations in Algorithm 1 is $k = 60$.

baart. This is a Fredholm integral equation of the first kind [1]. All the methods are tested with and without additional regularization, and with different preconditioners. The standard (unpreconditioned) GMRES method is known to perform well on this test problem. Nonetheless, we can experimentally show that the new preconditioned solvers can outperform standard GMRES in terms of the quality of the computed solutions. In the left frames of Fig. 1, we report the relative error history for different preconditioners (also defined with different parameters k_P) and for different solvers. We can clearly see that, if no additional Tikhonov or TSVD regularization is incorporated (top left frame of Fig. 1), “semi-convergence” appears after only few steps, i.e., the iterates computed during the first few iterations approach $\mathbf{x}_{\text{exact}}$ while, during subsequent iterations, they yield worse approximations of $\mathbf{x}_{\text{exact}}$. Semi-convergence is less evident when the preconditioner M_2 is used. When additional regularization in Tikhonov or TSVD form is incorporated (mid and bottom left frames of Fig. 1), all the preconditioned methods are more stable and exhibit smaller relative errors (when compared to GMRES without Tikhonov or TSVD regularization). For the present test problem, the preconditioners M_3 and M_4 perform the best. Indeed, the reconstructions displayed in the right-hand side of Fig. 1 show that the boundary values of the solution are accurately recovered when M_3 or M_4 are used. Applying the stopping rule (44) to determine the number of Arnoldi steps that define the preconditioner yields $k_P = 9$; the stopping rule (43) gives the same value. We report the behavior of relevant quantities used to set k_P in the top frames of Fig. 3. Note that increasing the number of Arnoldi iterations, k_P , is not always beneficial. Indeed,

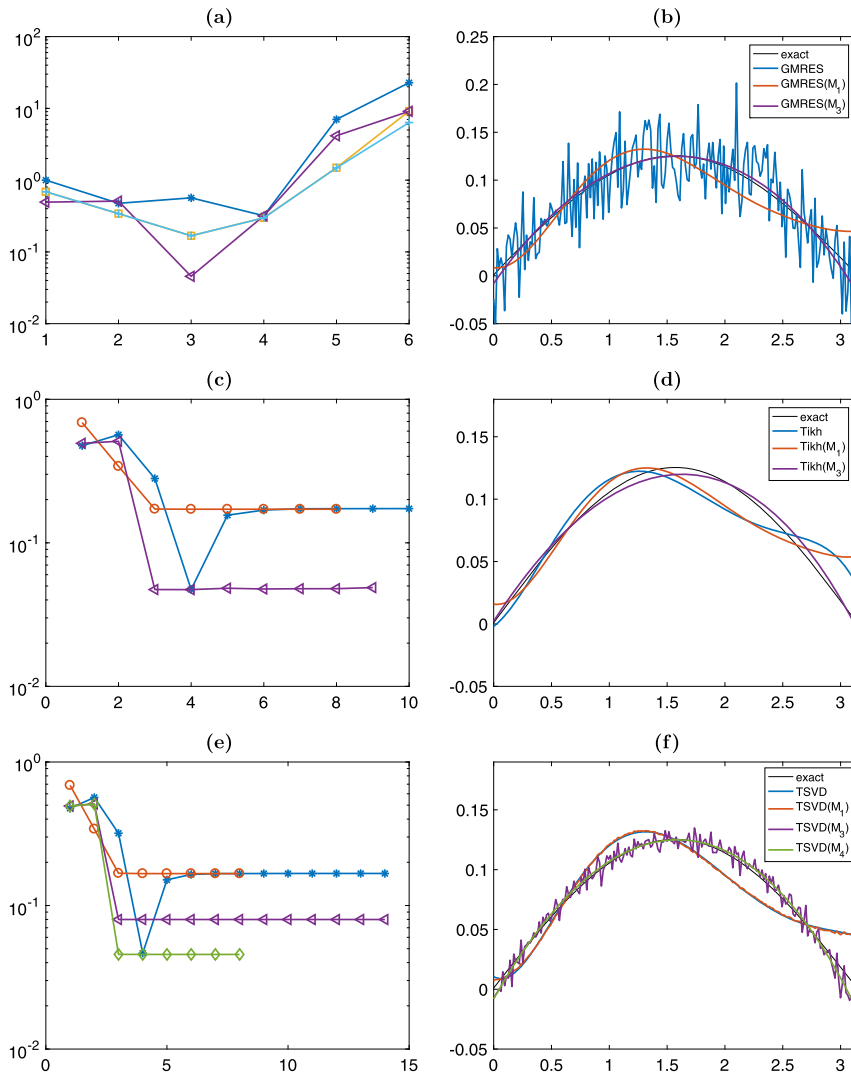


Fig. 1. Test problem baart, with $m = 200$ and $\|\mathbf{e}\|/\|\mathbf{b}\| = 10^{-2}$. (a) Relative error history, without any additional regularization and $k_p = 9$. (b) Best approximations, without any additional regularization and $k_p = 9$. (c) Relative error history, with Arnoldi-Tikhonov and $k_p = 9$. (d) Approximations for $k = 6$, with Arnoldi-Tikhonov and $k_p = 9$. (e) Relative error history, with Arnoldi-TSVD and $k_p = 39$. (f) Approximations for $k = 9$, with Arnoldi-TSVD and $k_p = 39$.

a larger k_p -value may result in a more severe loss of orthogonality in the Arnoldi process (Algorithm 1), even if reorthogonalization is used, so that numerical inaccuracies may affect the computation of all the preconditioners (30)–(37). Moreover, preconditioners (30) and (35) should be rank- k_p regularized approximations of the original matrix A^* and A , respectively; by increasing k_p these approximations become increasingly ill-conditioned and, therefore, less successful in regularizing the problem at hand. The best relative errors attained by each iterative method (considering different choices of solvers and preconditioners) are reported in Table 1, where averages over 30 different realizations of the noise \mathbf{e} in the vector \mathbf{b} are shown.

heat. We consider a discretization of the inverse heat equation formulated as a Volterra integral equation of the first kind. This problem can be regarded as numerically rank-deficient, with numerical rank equal to 195. According to the analysis in [33], GMRES does not converge to the minimum norm solution of (1) for this problem, as the null spaces of A and A^* are different. For this test problem, using the preconditioned methods described in Section 5, with some of the preconditioners derived in Section 4, can make a dramatic difference. When applying stopping rule (44) to determine the number of Arnoldi iterations that define the preconditioners, we obtain $k_p = 23$. The stopping rule (43) yields a similar k_p -value. We report the behavior of relevant quantities used to determine k_p in the bottom frames of Fig. 3. In the left frames of Fig. 2, we report the relative error norm history when different preconditioners (also defined with respect to different parameters k_p) and different solvers are considered. In all these graphs, the unpreconditioned Arnoldi-Tikhonov and Arnoldi-TSVD solutions diverge, with the best approximations being the ones recovered in the first iteration, i.e., the ones belonging to

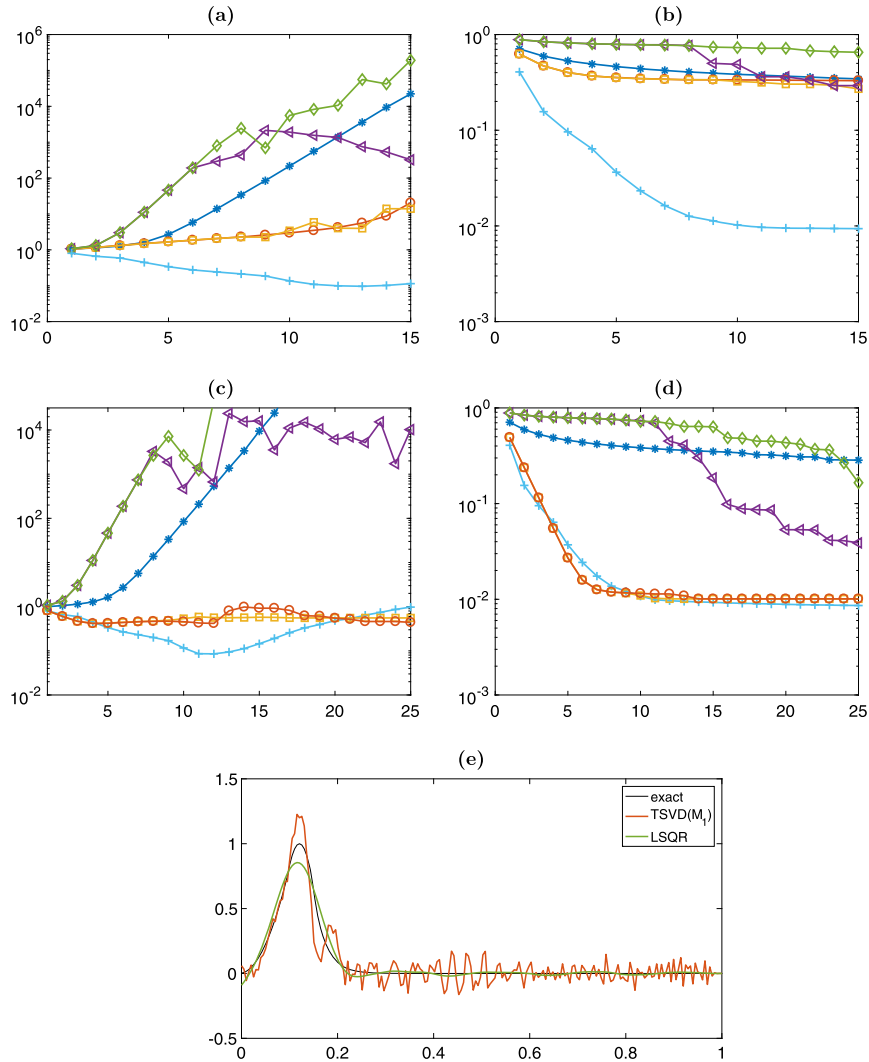


Fig. 2. Test problem heat, with $m = 200$ and $\|e\|/\|b\| = 10^{-2}$. (a) Relative error history, with Arnoldi-TSVD and $k_P = 23$. (b) Relative residual history, with Arnoldi-TSVD and $k_P = 23$. (c) Relative error history, with Arnoldi-Tikhonov and $k_P = 60$. (d) Relative residual history, with Arnoldi-Tikhonov and $k_P = 60$. (e) Best approximations, with Arnoldi-Tikhonov and $k_P = 60$.

span{b}. Moreover, the approximate solutions computed when using the preconditioners (30)–(37) with $k_P = 23$ do not look much improved. Indeed, while the computed approximate solutions obtained with the preconditioners M_1 and M_2 do not degenerate as quickly with the number of iterations, the computed solutions determined with the preconditioners M_3 and M_4 are worse than those determined by unpreconditioned iterations. However, when the maximum allowed value of k_P (i.e., $k_P = 60$) is chosen, the gain of using a preconditioned approach is evident. While the regularizing preconditioners M_3 and M_4 still perform poorly, the preconditioners M_1 and M_2 , which seek to make the matrix AM Hermitian positive semidefinite by incorporating an approximate regularized version of A^* , allow us to compute an approximate solution, whose quality is close to the one achieved by LSQR. The right frames of Fig. 2 display the history of the corresponding relative residuals (or discrepancies) norms. We can clearly see that the residuals are good indicators of the performance of these methods. Indeed, for $k_P = 23$ all the residuals (except those for LSQR) have a quite large norm and, in particular, the discrepancy principle (7), (45), (46) is far from being satisfied. For $k_P = 60$, the preconditioned Arnoldi-Tikhonov method with the preconditioners M_1 or M_2 eventually satisfies the discrepancy principle. Also, the approximate solution obtained with M_1 reproduces the main features of the exact solution, though some spurious oscillations are present. This is probably due to the tiny value $\mu = 1.2287 \cdot 10^{-8}$ selected for the regularization parameter according to the discrepancy principle (45); spurious oscillations are likely to be damped if a larger value for μ is used. The smallest relative errors attained by each iterative method (considering different choices of solvers and preconditioners) are reported in Table 1, where averages over 30 different realizations of the noise in the vector b are shown.

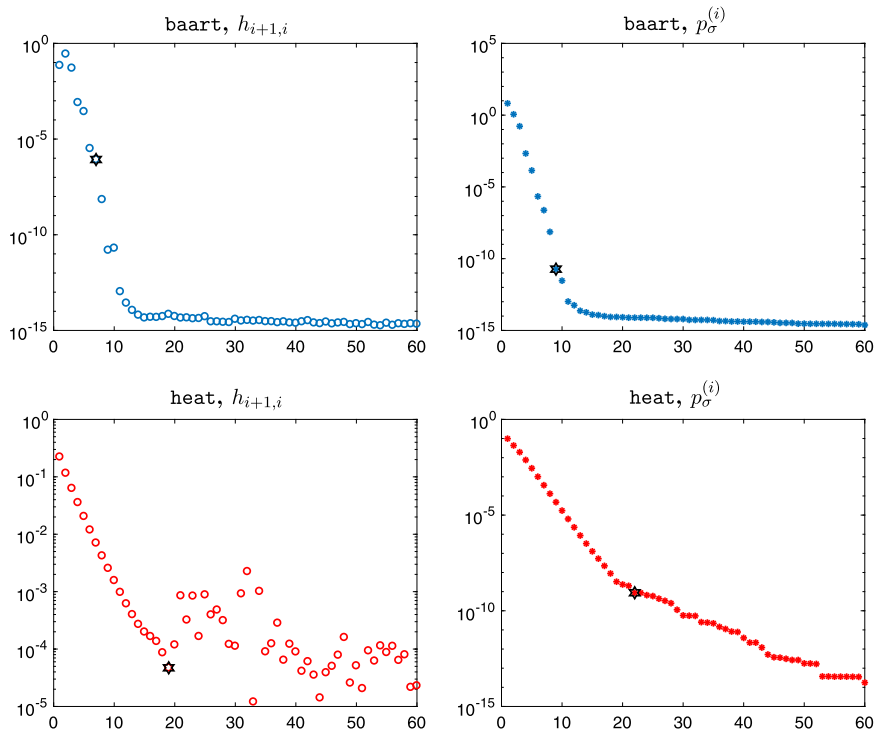


Fig. 3. Illustration of the stopping criteria (43) and (44) for the test problems baart (top row) and heat (bottom row). In the left column, the values of the subdiagonal entries $h_{i+1,i}$ of the Hessenberg matrix $H_{k+1,k}$ are plotted against i ($i = 1, \dots, k$). In the right column, the products $p_{\sigma}^{(i)} := \sigma_1^{(i)} \sigma_{i+1}^{(i+1)}$ of the extremal singular values of $H_{i+1,i}$ are plotted against i .

Table 1

Average values of the best relative errors over 30 runs of the test problems in the first set of experiments, with $\|\mathbf{e}\|/\|\mathbf{b}\| = 10^{-2}$. The smaller parameter k_p satisfies, on average, the stopping rule (44); the larger parameter k_p is obtained adding 30 to the smaller parameter k_p .

baart						
	TSVD		Tikh		none	
	$k_p = 9$	$k_p = 39$	$k_p = 9$	$k_p = 39$	$k_p = 9$	$k_p = 39$
–	4.7202e-02	4.7202e-02	6.7530e-02	6.7530e-02	3.0950e-01	3.0950e-01
M_1	2.2148e-02	1.6744e-01	2.4002e-02	1.7926e-01	1.8452e-02	1.5647e-01
M_2	1.6689e-01	1.2429e-01	1.7733e-01	1.3091e-01	1.5838e-01	1.2517e-01
M_3	4.5578e-02	6.1255e-02	6.6982e-02	6.7486e-02	4.5029e-02	6.1259e-02
M_4	1.7025e-02	4.5678e-02	2.4297e-02	6.8386e-02	1.7027e-02	4.1604e-02
LSQR	1.5787e-01	1.5787e-01	1.5787e-01	1.5787e-01	1.5787e-01	1.5787e-01
heat						
	TSVD		Tikh		none	
	$k_p = 20$	$k_p = 50$	$k_p = 20$	$k_p = 50$	$k_p = 20$	$k_p = 50$
–	6.5870e-01	6.5870e-01	5.6767e-01	5.6767e-01	1.0584e+00	1.0584e+00
M_1	1.0296e+00	3.6071e-01	1.0296e+00	3.6173e-01	1.0296e+00	3.6136e-01
M_2	1.0296e+00	3.6390e-01	1.0119e+00	3.0444e-01	1.0296e+00	3.6390e-01
M_3	1.0747e+00	1.0747e+00	1.0747e+00	1.0747e+00	1.0747e+00	1.0747e+00
M_4	1.0747e+00	1.0747e+00	1.0747e+00	1.0375e+00	1.0747e+00	1.0747e+00
LSQR	9.2105e-02	9.2105e-02	9.2105e-02	9.2105e-02	9.2105e-02	9.2105e-02

6.2. Second set of experiments

We consider 2D image restoration problems, where the available images are contaminated by spatially invariant blur and Gaussian white noise. In this setting, given a point-spread function (PSF) that describes how a single pixel of the exact discrete image $\mathbf{X}_{\text{exact}} \in \mathbb{R}^{N \times N}$ is blurred, the blurring process is modeled as a 2D convolution of the PSF and $\mathbf{X}_{\text{exact}}$. Here and in the following, the PSF is represented as a 2D image. A 2D image restoration problem can be expressed as a linear system of equations (1), where the 1D array \mathbf{b} is obtained by stacking the columns of the 2D blurred and noisy image (so that $m = N^2 = n$), and the square matrix A incorporates the convolution process together with some given boundary conditions. Our experiments consider two different grayscale test images, two different PSFs, and reflective boundary conditions; the

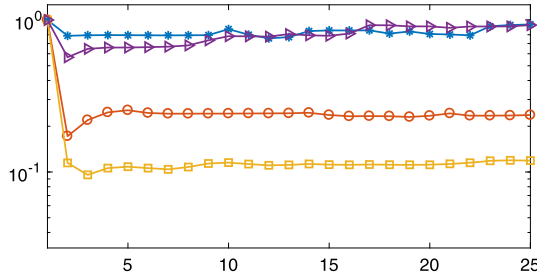


Fig. 4. Image deblurring problem with anisotropic motion blur. Relative error history of GMRES (*) and right-preconditioned GMRES methods (with preconditioners C_1 (○), C_2 (□), and C_3 (▷)).

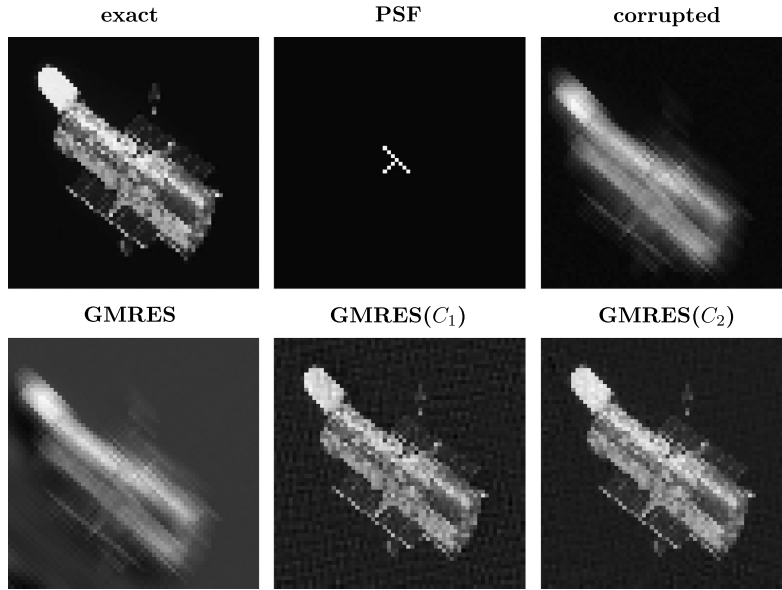


Fig. 5. Image deblurring problem with anisotropic motion blur. The upper row displays the test data. The lower row displays the best reconstructions obtained by: GMRES method ($7.5748e-02$, $k = 12$); GMRES(C_1) method ($1.7202e-01$, $k = 2$); GMRES(C_2) method ($9.5692e-02$, $k = 3$).

exact images are artificially blurred, and noise of several levels is added. These test problems are generated using the routines in *IR Tools* [20]. The maximum allowed number of Arnoldi iterations in Algorithm 1 is set to $k = 100$, and k_p is determined according to (44).

Anisotropic motion blur. For this experiment, a small satellite test image of size 64×64 pixels is taken as the exact image; it is displayed in the top frames of Fig. 5, together with a PSF of size 7×7 pixels modeling motion in two orthogonal directions, and the available corrupted image (with noise level $\|\mathbf{e}\|/\|\mathbf{b}\| = 2 \cdot 10^{-2}$). Thanks to the still moderate size of this test problem, all the preconditioners described in Section 4 can be straightforwardly implemented. GMRES and right-preconditioned GMRES with preconditioners C_1 , C_2 , and C_3 are considered. The preconditioners M_i , $i = 1, \dots, 4$, do not perform well for this restoration problem, even when the maximum number of Arnoldi steps $k = k_p = 100$ is carried out. This is probably due to the fact that the PSF is quite nonsymmetric. For this test problem, $m = n = 4096$ in (1). Fig. 4 displays the relative reconstruction error histories for these solvers. The most effective preconditioner for this problem is C_2 . Moreover, both C_1 and C_2 require only a few iterations to compute an accurate restoration and exhibit a quite stable behavior. We therefore do not consider the Arnoldi-Tikhonov and Arnoldi-TSVD methods for this test problem. Fig. 5 shows the best restorations achieved by each method. Relative errors and the corresponding number of iterations are displayed in the caption.

Isotropic motion blur. The test data for this experiment are displayed in Fig. 6; the exact and corrupted images are of size 256×256 pixels, so that $m = n = 65536$ in (1). We consider a 17×17 PSF modeling diagonal motion blur. The noise level is $5 \cdot 10^{-3}$. Fig. 7 shows the best restorations achieved by each method; relative errors and the corresponding number of iterations are displayed in the caption.

All the methods carry out more iterations than in the previous example. This is due to the smaller amount of noise in the present example, and the larger scale of the problem. Visual inspection of the images in Fig. 7 shows the unpreconditioned Arnoldi-TSVD solution to give a restoration with some motion artifacts, as the restored image displays some shifts in the diagonal directions (i.e., in the direction of the motion blur). These artifacts are less pronounced in the TSVD(M_1) restoration,

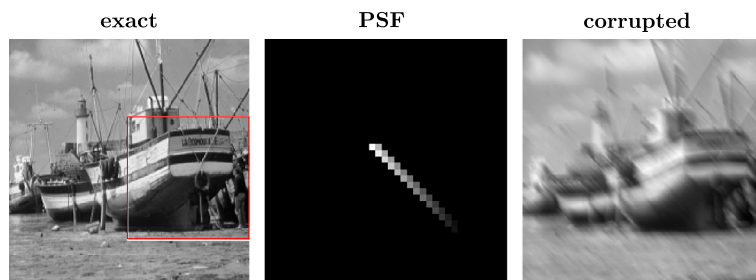


Fig. 6. From left to right: exact image; blow-up (600%) of the diagonal motion PSF; blurred and noisy available image, with $\|\mathbf{e}\|/\|\mathbf{b}\| = 5 \cdot 10^{-3}$.

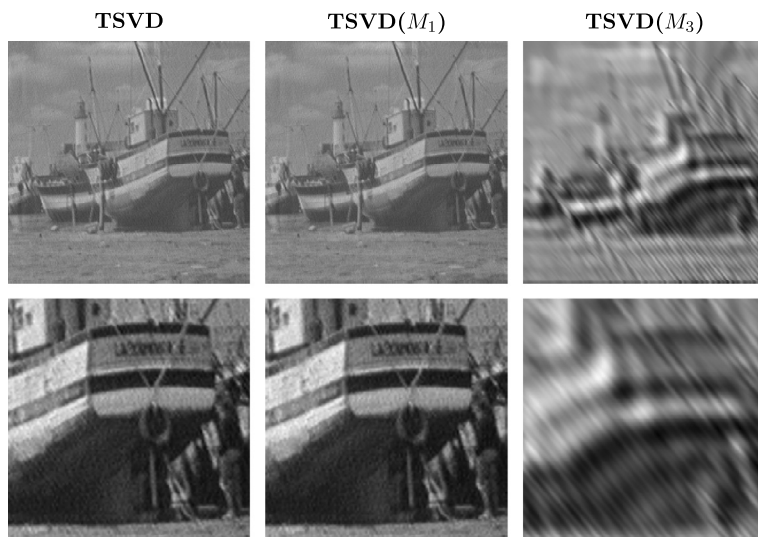


Fig. 7. The lower row displays blow-ups (200%) of the restored images in the upper row. From left to right: unpreconditioned Arnoldi-TSVD method ($1.0481e-01$, $k = 26$); TSVD(M_1) method ($1.0081e-01$, $k_P = 50$, $k = 7$); TSVD(M_3) method ($2.5948e-01$, $k_P = 50$, $k = 35$).

as the preconditioner (30) makes the problem more symmetric. The reconstruction produced by TSVD(M_3) is noticeably worse. Indeed, the preconditioner (35) merely approximates a regularized inverse of A , and this is not desirable when applying the Arnoldi algorithm to very unsymmetric blur. Results obtained when applying Arnoldi–Tikhonov methods are very similar to those achieved with Arnoldi-TSVD methods. We therefore only show the latter.

7. Conclusions

This paper presents an analysis of the GMRES method and the Arnoldi process with applications to the regularization of large-scale linear discrete ill-posed problems. Theoretical properties that involve the distance of the original coefficient matrix to classes of generalized Hermitian matrices are derived. Novel preconditioners based on matrices stemming from the standard Arnoldi decomposition are introduced, and the resulting right-preconditioned linear systems are solved with methods based on the preconditioned Arnoldi algorithm, or the preconditioned Arnoldi–Tikhonov and Arnoldi-TSVD methods. Numerical results on a variety of test problems illustrate that the new preconditioning techniques discussed give approximations of the desired solution $\mathbf{x}_{\text{exact}}$ of higher quality than when no preconditioner is employed.

References

- [1] M.L. Baart, The use of auto-correlation for pseudo-rank determination in noisy ill-conditioned least-squares problems, *IMA J. Numer. Anal.* 2 (1982) 241–247.
- [2] Å. Björck, *Numerical Methods in Matrix Computation*, Springer, New York, 2015.
- [3] P.N. Brown, H.F. Walker, GMRES on (nearly) singular systems, *SIAM J. Matrix Anal. Appl.* 18 (1997) 37–51.
- [4] D. Calvetti, B. Lewis, L. Reichel, Restoration of images with spatially variant blur by the GMRES method, in: F.T. Luk (Ed.), *Proceedings of the Society of Photo-Optical Instrumentation Engineers (SPIE)*, X, in: *Proceedings of the Society of Photo-Optical Instrumentation Engineers (SPIE)*, vol. 4116, The International Society for Optical Engineering, Bellingham, WA, 2000, pp. 364–374.
- [5] D. Calvetti, B. Lewis, L. Reichel, On the choice of subspace for iterative methods for linear discrete ill-posed problems, *Int. J. Appl. Math. Comput. Sci.* 11 (2001) 1069–1092.

- [6] D. Calvetti, B. Lewis, L. Reichel, Krylov subspace iterative methods for nonsymmetric discrete ill-posed problems in image restoration, in: F.T. Luk (Ed.), *Advanced Signal Processing Algorithms, Architectures, and Implementations XI*, in: *Proceedings of the Society of Photo-Optical Instrumentation Engineers (SPIE)*, vol. 4474, The International Society for Optical Engineering, Bellingham, WA, 2001, pp. 224–233.
- [7] D. Calvetti, B. Lewis, L. Reichel, On the regularizing properties of the GMRES method, *Numer. Math.* 91 (2002) 605–625.
- [8] D. Calvetti, B. Lewis, L. Reichel, GMRES, L-curves, and discrete ill-posed problems, *BIT* 42 (2002) 44–65.
- [9] D. Calvetti, S. Morigi, L. Reichel, F. Sgallari, Tikhonov regularization and the L-curve for large discrete ill-posed problems, *J. Comput. Appl. Math.* 123 (2000) 423–446.
- [10] R.H.-F. Chan, X.-Q. Jin, *An Introduction to Iterative Toeplitz Solvers*, SIAM, Philadelphia, 2007.
- [11] T.F. Chan, An optimal circulant preconditioner for Toeplitz systems, *SIAM J. Sci. Stat. Comput.* 9 (1988) 766–771.
- [12] T.F. Chan, K.R. Jackson, Nonlinearly preconditioned Krylov subspace methods for discrete Newton algorithms, *SIAM J. Sci. Stat. Comput.* 5 (1984) 533–542.
- [13] P.J. Davis, *Circulant Matrices*, 2nd ed., Chelsea, New York, 1994.
- [14] M. Donatelli, D. Martin, L. Reichel, Arnoldi methods for image deblurring with anti-reflective boundary conditions, *Appl. Math. Comput.* 253 (2015) 135–150.
- [15] K. Du, J. Duintjer Tebbens, G. Meurant, Any admissible harmonic Ritz value set is possible for GMRES, *Electron. Trans. Numer. Anal.* 47 (2017) 37–56.
- [16] L. Dykes, S. Noschese, L. Reichel, Circulant preconditioners for discrete ill-posed Toeplitz systems, *Numer. Algorithms* 75 (2017) 477–490.
- [17] L. Dykes, L. Reichel, A family of range restricted iterative methods for linear discrete ill-posed problems, in: *Dolomites Research Notes on Approximation*, vol. 6, 2013, pp. 27–36.
- [18] S.C. Eisenstat, Equivalence of Krylov subspace methods for skew-symmetric linear systems, *arXiv:1512.00311*, 2015.
- [19] H.W. Engl, M. Hanke, A. Neubauer, *Regularization of Inverse Problems*, Kluwer, Dordrecht, 1996.
- [20] S. Gazzola, P.C. Hansen, J. Nagy, IR tools: a MATLAB package of iterative regularization methods and large-scale test problems, *Numer. Algorithms* (2018), <https://doi.org/10.1007/s11075-018-0570-7>.
- [21] S. Gazzola, P. Novati, M.R. Russo, Embedded techniques for choosing the parameter in Tikhonov regularization, *Numer. Linear Algebra Appl.* 21 (2014) 796–812.
- [22] S. Gazzola, P. Novati, M.R. Russo, On Krylov projection methods and Tikhonov regularization, *Electron. Trans. Numer. Anal.* 44 (2015) 83–123.
- [23] A. Greenbaum, V. Ptak, Z. Strakoš, Any nonincreasing convergence curve is possible for GMRES, *SIAM J. Matrix Anal. Appl.* 17 (1996) 465–469.
- [24] L. Greengard, V. Rokhlin, A new version of the fast multipole method for the Laplace equation in three dimensions, *Acta Numer.* 6 (1997) 229–269.
- [25] M. Hanke, *Conjugate Gradient Type Methods for Ill-Posed Problems*, Longman, Harlow, 1995.
- [26] M. Hanke, J. Nagy, R. Plemmons, Preconditioned iterative regularization for ill-posed problems, in: L. Reichel, A. Ruttan, R.S. Varga (Eds.), *Numerical Linear Algebra*, de Gruyter, Berlin, 1993, pp. 141–163.
- [27] P.C. Hansen, *Rank-Deficient and Discrete Ill-Posed Problems*, SIAM, Philadelphia, 1998.
- [28] P.C. Hansen, Regularization tools version 4.0 for MATLAB 7.3, *Numer. Algorithms* 46 (2007) 189–194.
- [29] P.C. Hansen, T.K. Jensen, Smoothing-norm preconditioning for regularizing minimum-residual methods, *SIAM J. Matrix Anal.* 29 (2006) 1–14.
- [30] P.C. Hansen, T.K. Jensen, Noise propagation in regularizing iterations for image deblurring, *Electron. Trans. Numer. Anal.* 31 (2008) 204–220.
- [31] N.J. Higham, Computing a nearest symmetric positive semidefinite matrix, *Linear Algebra Appl.* 103 (1988) 103–118.
- [32] T. Huckle, The Arnoldi method for normal matrices, *SIAM J. Matrix Anal. Appl.* 15 (1994) 479–489.
- [33] T.K. Jensen, P.C. Hansen, Iterative regularization with minimal residual methods, *BIT* 47 (2007) 103–120.
- [34] S. Kindermann, Convergence analysis of minimization-based noise level-free parameter choice rules for linear ill-posed problems, *Electron. Trans. Numer. Anal.* 38 (2011) 233–257.
- [35] B. Lewis, L. Reichel, Arnoldi–Tikhonov regularization methods, *J. Comput. Appl. Math.* 226 (2009) 92–102.
- [36] D. Loghini, D. Ruiz, A. Touhami, Adaptive preconditioners for nonlinear systems of equations, *J. Comput. Appl. Math.* 189 (2006) 362–374.
- [37] G. Meurant, *Computer Solution of Large Linear Systems*, Elsevier, Amsterdam, 1999.
- [38] I. Moret, A note on the superlinear convergence of GMRES, *SIAM J. Numer. Anal.* 34 (1997) 513–516.
- [39] N.M. Nachtigal, S.C. Reddy, L.N. Trefethen, How fast are nonsymmetric matrix iterations?, *SIAM J. Matrix Anal. Appl.* 13 (1992) 778–795.
- [40] N.M. Nachtigal, L. Reichel, L.N. Trefethen, A hybrid GMRES algorithm for nonsymmetric linear systems, *SIAM J. Matrix Anal. Appl.* 13 (1992) 796–825.
- [41] A. Neuman, L. Reichel, H. Sadok, Implementations of range restricted iterative methods for linear discrete ill-posed problems, *Linear Algebra Appl.* 436 (2012) 3974–3990.
- [42] M.K. Ng, *Iterative Methods for Toeplitz Systems*, Oxford University Press, Oxford, 2004.
- [43] S. Noschese, L. Pasquini, L. Reichel, The structured distance to normality of an irreducible real tridiagonal matrix, *Electron. Trans. Numer. Anal.* 28 (2007) 65–77.
- [44] S. Noschese, L. Pasquini, L. Reichel, Tridiagonal Toeplitz matrices: properties and novel applications, *Numer. Linear Algebra Appl.* 20 (2013) 302–326.
- [45] S. Noschese, L. Reichel, The structured distance to normality of Toeplitz matrices with application to preconditioning, *Numer. Linear Algebra Appl.* 18 (2011) 429–447.
- [46] S. Noschese, L. Reichel, A modified TSVD method for discrete ill-posed problems, *Numer. Linear Algebra Appl.* 21 (2014) 813–822.
- [47] S. Noschese, L. Reichel, A note on superoptimal generalized circulant preconditioners, *Appl. Numer. Math.* 75 (2014) 188–195.
- [48] P. Novati, Some properties of the Arnoldi based methods for linear ill-posed problems, *SIAM J. Numer. Anal.* 55 (2017) 1437–1455.
- [49] P. Novati, M.R. Russo, A GCV based Arnoldi–Tikhonov regularization method, *BIT* 54 (2014) 501–521.
- [50] C.C. Paige, M.A. Saunders, Solution of sparse indefinite systems of linear equations, *SIAM J. Numer. Anal.* 12 (1975) 617–629.
- [51] C.C. Paige, M.A. Saunders, LSQR: an algorithm for sparse linear equations and sparse least squares, *ACM Trans. Math. Softw.* 8 (1982) 43–71.
- [52] L. Reichel, G. Rodriguez, Old and new parameter choice rules for discrete ill-posed problems, *Numer. Algorithms* 63 (2013) 65–87.
- [53] L. Reichel, Q. Ye, Breakdown-free GMRES for singular systems, *SIAM J. Matrix Anal. Appl.* 26 (2005) 1001–1021.
- [54] L. Reichel, Q. Ye, Simple square smoothing regularization operators, *Electron. Trans. Numer. Anal.* 33 (2009) 63–83.
- [55] J.R. Ringrose, *Compact Non-Self-Adjoint Operators*, Van Nostrand Reinhold, London, 1971.
- [56] Y. Saad, *Iterative Methods for Sparse Linear Systems*, 2nd ed., SIAM, Philadelphia, 2003.
- [57] Y. Saad, M.H. Schultz, GMRES: a generalized minimal residual method for solving nonsymmetric linear systems, *SIAM J. Sci. Stat. Comput.* 7 (1986) 856–869.
- [58] V.V. Strela, E.E. Tyrtshnikov, Which circulant preconditioner is better?, *Math. Compet.* 65 (1996) 137–150.
- [59] E.E. Tyrtshnikov, Optimal and superoptimal circulant preconditioners, *SIAM J. Matrix Anal. Appl.* 13 (1992) 459–473.