



MAX-PLANCK-GESELLSCHAFT

**Max Planck Institute Magdeburg
Preprints**

Peter Benner

Zvonimir Bujanović

**On the solution of large-scale algebraic
Riccati equations by using low-dimensional
invariant subspaces**



MAX-PLANCK-INSTITUT
FÜR DYNAMIK KOMPLEXER
TECHNISCHER SYSTEME
MAGDEBURG

Authors addresses:

Peter Benner
Max Planck Institute for Dynamics of Complex Technical Systems,
Computational Methods in Systems and Control Theory,
D-39106 Magdeburg, Germany
benner@mpi-magdeburg.mpg.de

Zvonimir Bujanović
Max Planck Institute for Dynamics of Complex Technical Systems,
Computational Methods in Systems and Control Theory,
D-39106 Magdeburg, Germany
zbujanov@math.hr

Imprint:

Max Planck Institute for Dynamics of Complex Technical Systems, Magdeburg

Publisher:
Max Planck Institute for
Dynamics of Complex Technical Systems

Address:
Max Planck Institute for
Dynamics of Complex Technical Systems
Sandtorstr. 1
39106 Magdeburg

<http://www.mpi-magdeburg.mpg.de/preprints/>

Abstract

This article discusses an approach to solving large-scale algebraic Riccati equations (AREs) by computing a low-dimensional stable invariant subspace of the associated Hamiltonian matrix. We give conditions on AREs to admit solutions of low numerical rank and show that these can be approximated via Hamiltonian eigenspaces. We discuss strategies on choosing the proper eigenspace that yields a good approximation, and different formulas for building the approximation itself. Similarities of our approach with several other methods for solving AREs are shown: closely related are the projection-type methods that use various Krylov subspaces and the qADI algorithm. The aim of this paper is merely to analyze the possibilities of computing approximate Riccati solutions from low-dimensional subspaces related to the corresponding Hamiltonian matrix and to explain commonalities among existing methods rather than providing a new algorithm.

1 Introduction

Finding the solution of the continuous-time algebraic Riccati equation

$$A^*P + PA + Q - PGP = 0 \quad (1)$$

is of great interest to the control theory community, which uses it in e.g. linear-quadratic optimal regulator problems, H_2 and H_∞ controller design and balancing-related model reduction. Current applications require efficient algorithms in cases where $A \in \mathbb{R}^{n \times n}$ is a large sparse matrix and $Q = C^*C$, $G = BB^*$ are positive semidefinite low-rank matrices, $B \in \mathbb{R}^{n \times m}$, $C \in \mathbb{R}^{p \times n}$ with $m, p \ll n$. In particular, one is interested in obtaining the stabilizing solution P , which is the unique positive semidefinite solution that makes the closed-loop matrix $A - GP$ stable.

There are several competitive methods to tackle this problem, designed to exploit the expected low-rank structure of the solution. These methods include the Newton-ADI (Alternate Direction Implicit) [14, 16, 5, 21] and the various projection-type methods, usually based on approximations using Krylov [11, 13, 22] or rational Krylov [9, 8] subspaces generated by the matrices A^* and A^{-*} and the initial (block-)vector C^* .

In this paper, we follow up on the approach introduced in [3, 10], and further pursued in [1]. There it is suggested to compute a low-dimensional stable invariant subspace of the Hamiltonian matrix $\mathcal{H} = \begin{bmatrix} A & G \\ Q & -A^* \end{bmatrix}$ via a symplectic Lanczos procedure, which is then used to approximate the stabilizing solution of the Riccati equation. In order to justify the existence of such a low-dimensional subspace that yields a good approximation, in Section 2 we first discuss the properties of the input matrices A, G, Q that imply the rapid decay in the singular values of the Riccati solution.

Section 3 addresses the questions on selecting the eigenvalues of \mathcal{H} towards which the Lanczos procedure should be steered, and on constructing an approximation to the Riccati solution once the invariant subspace is available. In particular, we shall see that the former question is especially difficult, even when the entire eigenvalue decomposition of \mathcal{H} is known in advance.

Finally, in Section 4 we relate Krylov methods for computing Hamiltonian eigenspaces to the aforementioned projection-type methods for solving the Riccati equation. This gives some new insights on the latter, in particular on the shift selection in the rational Krylov method. We also show that the approximation built from the Hamiltonian eigenspace is closely tied to another method for large-scale Riccati equations—namely, the qADI method of Wong and Balakrishnan [29].

We briefly recall several properties of Hamiltonian matrices that are needed in the rest of the paper, and introduce some notation. Eigenvalues of real Hamiltonian matrices come in quadruples: if $\lambda \in \mathbb{C}$ is an eigenvalue of \mathcal{H} , then so are $\bar{\lambda}$, $-\lambda$ and $-\bar{\lambda}$. For the purpose of finding the stabilizing solution of the Riccati equation, we are particularly interested in stable eigenvalues. The tuple of stable eigenvalues of \mathcal{H} will be denoted as $\vec{\lambda}^{\mathcal{H}} = (\lambda_1^{\mathcal{H}}, \dots, \lambda_n^{\mathcal{H}})$, and the tuple of all eigenvalues of A as $\vec{\lambda}^A = (\lambda_1^A, \dots, \lambda_n^A)$; ordering of the eigenvalues in a tuple is arbitrary unless specified otherwise. The stable part of the complex plane is denoted with $\mathbb{C}_- = \{z \in \mathbb{C} : \text{Re}(z) < 0\}$. We are assuming that none of the eigenvalues of \mathcal{H} lies on the imaginary axis, which follows from the usual assumption that the pair (A, B) is stabilizable (i.e. $\text{rank}[\xi I - A, B] = n$, for all $\xi \in \mathbb{C} \setminus \mathbb{C}_-$), and that the pair (A, C) is detectable (i.e. (A^*, C^*) is stabilizable), see e.g. [15].

Suppose that the columns of the matrix $\begin{bmatrix} X_k \\ Y_k \end{bmatrix}$ span a k -dimensional stable invariant subspace of \mathcal{H} , i.e.

$$\mathcal{H} \begin{bmatrix} X_k \\ Y_k \end{bmatrix} = \begin{bmatrix} X_k \\ Y_k \end{bmatrix} \Lambda_k,$$

for some $\Lambda_k \in \mathbb{C}^{k \times k}$ with stable spectrum, and $X_k, Y_k \in \mathbb{C}^{n \times k}$. When $k = n$, the stabilizing Riccati solution is given by a simple formula: $P = -Y_n X_n^{-1}$. Furthermore, the eigenvalues of the closed-loop matrix $A - GP$ coincide with the eigenvalues of Λ_n , which are $\lambda_1^{\mathcal{H}}, \dots, \lambda_n^{\mathcal{H}}$. For any $k \leq n$, the stable invariant subspace is isotropic, implying that $X_k^* Y_k = Y_k^* X_k$. The subspace \mathcal{S} is called isotropic if $x^* \mathcal{J} y = 0$, for all $x, y \in \mathcal{S}$, where $\mathcal{J} = \begin{bmatrix} 0 & I_n \\ -I_n & 0 \end{bmatrix}$.

If \tilde{P} is an approximate solution to (1), then the Riccati residual of \tilde{P} will be denoted as

$$\mathcal{R}(\tilde{P}) = A^* \tilde{P} + \tilde{P} A + Q - \tilde{P} G \tilde{P}.$$

Finally, the subspace spanned by the columns of a matrix Z is denoted as $\text{span}\{Z\}$, and the singular values of Z as $\sigma_1(Z) \geq \sigma_2(Z) \geq \dots$. For a matrix $S = \begin{bmatrix} x_1 & x_2 & \dots & x_k \\ y_1 & y_2 & \dots & y_k \end{bmatrix}$, where $x_j, y_j \in \mathbb{C}^n$, and the associated subspace $\mathcal{S} = \text{span}\{\begin{bmatrix} x_1 \\ y_1 \end{bmatrix}, \dots, \begin{bmatrix} x_k \\ y_k \end{bmatrix}\}$, it will be useful to separately study the \mathbb{X} component and the \mathbb{Y} component, defined as

$$\begin{aligned} \mathbb{X}(S) &= [x_1 \ x_2 \ \dots \ x_k], & \mathbb{Y}(S) &= [y_1 \ y_2 \ \dots \ y_k]; \\ \mathbb{X}(\mathcal{S}) &= \text{span}\{x_1, \dots, x_k\}, & \mathbb{Y}(\mathcal{S}) &= \text{span}\{y_1, \dots, y_k\}. \end{aligned}$$

2 Singular value decay of the Riccati solution

Since the matrices involved in the Riccati equation are assumed to be large and sparse, it is important to have the ability of representing the Riccati solution P in a compact form as well. Note that in general, P will be a fully populated, dense matrix even if the coefficients are sparse. Both the Newton-ADI and the projection methods compute the approximation to the solution in a factored form, such as $P \approx ZZ^*$, or $P \approx Y\Delta Y^*$, where $Y, Z \in \mathbb{R}^{n \times k}$, $\Delta \in \mathbb{R}^{k \times k}$, and $k \ll n$. For such an approximation to be effective, one relies on the fact that the matrix P has a low numerical rank, or equivalently, that its singular values are decaying rapidly.

In this section, we are going to justify such an assumption by providing an upper bound on the fractions σ_k/σ_1 , where $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_n$ are the singular values of the matrix P . The technique we are going to use is similar to the one in [25], where the bound for the Lyapunov equation is established. Our technique utilizes the Sylvester-ADI method [28, 27], and the convenient error expression for the approximations it generates. We recall this method in the following definition.

Definition 1. *Given the matrices $\mathcal{A}, \mathcal{B}, \mathcal{C} \in \mathbb{R}^{n \times n}$, two sequences of shifts $(\alpha_j)_j \subset \mathbb{C}$ and $(\beta_j)_j \subset \mathbb{C}$, and an initial approximation P_0 , the Sylvester-ADI method builds a sequence $(P_j)_j \subset \mathbb{C}^{n \times n}$ of approximations to the solution P of the Sylvester equation*

$$\mathcal{A}P + P\mathcal{B} = \mathcal{C}.$$

The sequence is generated by the recurrence

$$\begin{cases} (\mathcal{A} + \beta_j I)P_{j-1/2} = \mathcal{C} - P_{j-1}(\mathcal{B} - \beta_j I), \\ P_j(\mathcal{B} + \alpha_j I) = \mathcal{C} - (\mathcal{A} - \alpha_j I)P_{j-1/2}, \end{cases} \quad \text{for } j = 1, 2, \dots$$

Assuming $P_0 = 0$, it is not difficult to see that the rank of P_k will be equal to $k \cdot \text{rank}(\mathcal{C})$, and that the sequence generated by the Sylvester-ADI method satisfies the following error expression [27]:

$$P - P_k = \left(\prod_{j=0}^k (\mathcal{A} - \alpha_j I)(\mathcal{A} + \beta_j I)^{-1} \right) \cdot P \cdot \left(\prod_{j=0}^k (\mathcal{B} - \beta_j I)(\mathcal{B} + \alpha_j I)^{-1} \right). \quad (2)$$

By rewriting the Riccati equation as a Lyapunov or a Sylvester equation, and by using the ADI iterates as low-rank approximations to the Riccati solution, we can obtain bounds on its singular value decay.

Theorem 2. *Let $\sigma_1 \leq \sigma_2 \leq \dots \leq \sigma_n$ denote the singular values of the matrix P , the solution of (1) with $G = BB^*$ and $Q = C^*C$, where $B \in \mathbb{R}^{n \times m}$, $C \in \mathbb{R}^{p \times n}$. Then, for every k such that $(m+p)k+1 \leq n$ and any selection of shifts $\tau_1, \dots, \tau_k \in \mathbb{C}_-$, the following bound holds:*

$$\frac{\sigma_{(m+p)k+1}}{\sigma_1} \leq \kappa_{A-GP}^2 \cdot \left(\max_{i=1, \dots, n} \prod_{j=1}^k \frac{|\lambda_i^H - \tau_j|}{|\lambda_i^H + \overline{\tau_j}|} \right)^2. \quad (3)$$

If the matrix A is stable, we have a slightly improved bound when $pk + 1 \leq n$:

$$\frac{\sigma_{pk+1}}{\sigma_1} \leq \kappa_{A-GP} \kappa_A \cdot \max_{i=1, \dots, n} \prod_{j=1}^k \frac{|\lambda_i^{\mathcal{H}} - \tau_j|}{|\lambda_i^{\mathcal{H}} + \bar{\tau}_j|} \cdot \max_{i=1, \dots, n} \prod_{j=1}^k \frac{|\lambda_i^A - \tau_j|}{|\lambda_i^A + \bar{\tau}_j|}. \quad (4)$$

Here κ_A and κ_{A-GP} denote the condition numbers of the eigenvector matrices of A and $A - GP$, respectively.

Proof. For the first bound, rewrite the Riccati equation (1) as a Lyapunov equation

$$(A - GP)^*P + P(A - GP) = -Q - PGP,$$

and assume that $A - GP$ and PGP are known. Using the notation of Definition 1, consider the Sylvester-ADI method with $\mathcal{A} = (A - GP)^*$, $\mathcal{B} = A - GP$, and $\mathcal{C} = -Q - PGP$. Let the shifts used by the method be $\alpha_j = \bar{\tau}_j$ and $\beta_j = \tau_j$. Note that the matrix \mathcal{C} has rank less than or equal to $m + p$, and thus the iterate P_k is of rank at most $(m + p)k$. Finally, assume that the matrix $A - GP$ is diagonalizable with the eigenvalue decomposition $A - GP = X_{A-GP} \Lambda_{\mathcal{H}} X_{A-GP}^{-1}$, where $\Lambda_{\mathcal{H}} = \text{diag}(\lambda_1^{\mathcal{H}}, \dots, \lambda_n^{\mathcal{H}})$. Inserting this into the error expression (2) yields

$$\begin{aligned} P - P_k &= X_{A-GP}^{-*} \left(\prod_{j=0}^k (\Lambda_{\mathcal{H}}^* - \bar{\tau}_j I)(\Lambda_{\mathcal{H}}^* + \tau_j I)^{-1} \right) X_{A-GP}^* \cdot P \\ &\quad X_{A-GP} \left(\prod_{j=0}^k (\Lambda_{\mathcal{H}} - \tau_j I)(\Lambda_{\mathcal{H}} + \bar{\tau}_j I)^{-1} \right) X_{A-GP}^{-1}. \end{aligned}$$

Taking norms, we have

$$\sigma_{(m+p)k+1} \leq \|P - P_k\| \leq \kappa_{A-GP}^2 \cdot \sigma_1 \cdot \left(\max_{i=1, \dots, n} \prod_{j=1}^k \frac{|\lambda_i^{\mathcal{H}} - \tau_j|}{|\lambda_i^{\mathcal{H}} + \bar{\tau}_j|} \right)^2.$$

The second bound is obtained similarly: this time we rewrite the Riccati equation as the Sylvester equation

$$A^*P + P(A - GP) = -Q, \quad (5)$$

and apply the Sylvester-ADI with $\mathcal{A} = A^*$, $\mathcal{B} = A - GP$, $\mathcal{C} = -Q$, and with the same selection of shifts as before. Now the matrix \mathcal{C} has rank p , and thus $\text{rank } P_k \leq pk$. Using the eigenvalue decomposition $A = X_A \Lambda_A X_A^{-1}$, where $\Lambda_A = \text{diag}(\lambda_1^A, \dots, \lambda_n^A)$, the expression (2) reduces to

$$\begin{aligned} P - P_k &= X_A^{-*} \left(\prod_{j=0}^k (\Lambda_A^* - \bar{\tau}_j I)(\Lambda_A^* + \tau_j I)^{-1} \right) X_A^* \cdot P \\ &\quad X_{A-GP} \left(\prod_{j=0}^k (\Lambda_{\mathcal{H}} - \tau_j I)(\Lambda_{\mathcal{H}} + \bar{\tau}_j I)^{-1} \right) X_{A-GP}^{-1}. \end{aligned}$$

The bound now follows easily by taking norms. \square

Several remarks are in order. First, note that each of the fractions

$$\frac{|\lambda_i^{\mathcal{H}} - \tau_j|}{|\lambda_i^{\mathcal{H}} + \bar{\tau}_j|}$$

is less than 1, since both $\lambda_i^{\mathcal{H}}$ and τ_j are in the left half-plane. If A is stable, the same holds for the fractions

$$\frac{|\lambda_i^A - \tau_j|}{|\lambda_i^A + \bar{\tau}_j|}$$

as well. Thus even when using a single shift $\tau_j = \tau$, $j = 1, 2, \dots$, we have an exponential decay in the singular values:

$$\frac{\sigma_{(m+p)k+1}}{\sigma_1} \leq \kappa_{A-GP}^2 \cdot \rho^{2k}, \quad \rho = \max_{i=1, \dots, n} \frac{|\lambda_i^{\mathcal{H}} - \tau|}{|\lambda_i^{\mathcal{H}} + \bar{\tau}|} < 1.$$

Of course, the decay rate obtained this way may be very slow¹. An optimal decay rate is obtained by solving the minimax problem

$$\min_{\tau_1, \dots, \tau_k \in \mathbb{C}_-} \max_{i=1, \dots, n} \prod_{j=1}^k \frac{|\lambda_i^{\mathcal{H}} - \tau_j|}{|\lambda_i^{\mathcal{H}} + \bar{\tau}_j|}, \quad (6)$$

or the variant with two eigenvalue tuples,

$$\min_{\tau_1, \dots, \tau_k \in \mathbb{C}_-} \left(\max_{i=1, \dots, n} \prod_{j=1}^k \frac{|\lambda_i^{\mathcal{H}} - \tau_j|}{|\lambda_i^{\mathcal{H}} + \bar{\tau}_j|} \cdot \max_{i=1, \dots, n} \prod_{j=1}^k \frac{|\lambda_i^A - \tau_j|}{|\lambda_i^A + \bar{\tau}_j|} \right).$$

These problems are instances of the well-known task of determining the optimal ADI-shifts, and the solution is not known in the general case when the eigenvalues $\lambda_i^{\mathcal{H}}$ and λ_i^A are placed arbitrarily in the complex plane. Various heuristics can be deployed; see e.g. [19]. One possible relaxation of (6) is to allow only shifts which coincide with the eigenvalues of \mathcal{H} . In that case, the shifts may be constructed in the following way: if $\tau_1, \dots, \tau_{k-1}$ are already chosen amongst the eigenvalues of \mathcal{H} , the next one is selected as

$$\tau_k = \arg \max_{\lambda \in \{\lambda_1^{\mathcal{H}}, \dots, \lambda_n^{\mathcal{H}}\}} \prod_{j=1}^{k-1} \frac{|\lambda - \tau_j|}{|\lambda + \bar{\tau}_j|}.$$

The obtained sequence of shifts establishes an ordering of the eigenvalues of \mathcal{H} which we are going to name the ‘‘ADI-minimax’’ ordering. (For definiteness, let τ_1 be the eigenvalue closest to the origin.)

¹Also, note that when the eigenvalue $\lambda_i^{\mathcal{H}}$ is close to the imaginary axis, then $|\lambda_i^{\mathcal{H}} - \tau_j|/|\lambda_i^{\mathcal{H}} + \bar{\tau}_j| \approx 1$ unless τ_j is very close to $\lambda_i^{\mathcal{H}}$. Thus in the case where \mathcal{H} has many such eigenvalues, we cannot expect the decay to be fast—a new shift is needed to resolve each of them. This effect is also observed in Example 3.

As the second remark, we notice a relation between certain Cauchy matrices and the bounds (3) and (4). For the tuples $\vec{\alpha} = (\alpha_1, \dots, \alpha_n)$ and $\vec{\beta} = (\beta_1, \dots, \beta_n)$ of complex numbers, let

$$\mathcal{C}(\vec{\alpha}, \vec{\beta}) = \begin{bmatrix} \frac{1}{\vec{\alpha}_1 + \vec{\beta}_1} & \frac{1}{\vec{\alpha}_1 + \vec{\beta}_2} & \cdots & \frac{1}{\vec{\alpha}_1 + \vec{\beta}_n} \\ \frac{1}{\vec{\alpha}_2 + \vec{\beta}_1} & \frac{1}{\vec{\alpha}_2 + \vec{\beta}_2} & \cdots & \frac{1}{\vec{\alpha}_2 + \vec{\beta}_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{1}{\vec{\alpha}_n + \vec{\beta}_1} & \frac{1}{\vec{\alpha}_n + \vec{\beta}_2} & \cdots & \frac{1}{\vec{\alpha}_n + \vec{\beta}_n} \end{bmatrix}$$

denote the associated Cauchy matrix. Suppose for the moment that $p = 1$, i.e. that $Q = C^*C$ where C^* is a one-column matrix. Using the eigenvalue decompositions of A and $A - GP$, the equation (5) can be written as

$$\begin{aligned} \Lambda_A^*(X_A^* P X_{A-GP}) + (X_A^* P X_{A-GP}) \Lambda_{\mathcal{H}} &= -X_A^* Q X_{A-GP} \\ &= -(X_A^* C^*)(X_{A-GP} C^*)^*. \end{aligned}$$

Let $D_A = \text{diag}(X_A^* C^*)$, $D_{A-GP} = \text{diag}(X_{A-GP} C^*)$ denote the diagonal matrices whose diagonal elements are the entries in the vectors $X_A^* C^*$ and $X_{A-GP} C^*$ respectively. Then

$$\Lambda_A^* \tilde{P} + \tilde{P} \Lambda_{\mathcal{H}} = -e e^*, \quad (7)$$

where we have denoted $\tilde{P} = D_A^{-1} X_A^* P X_{A-GP} D_{A-GP}^{-*}$, and $e \in \mathbb{R}^n$ is the vector consisting of all ones. It is easy to see that the solution to this Sylvester equation is $\tilde{P} = -\mathcal{C}(\vec{\lambda}^A, \vec{\lambda}^{\mathcal{H}})$, and that from

$$P = -X_A^{-*} D_A \cdot \mathcal{C}(\vec{\lambda}^A, \vec{\lambda}^{\mathcal{H}}) \cdot D_{A-GP}^* X_{A-GP}^{-1}$$

we have

$$\sigma_k \leq \kappa_A \kappa_{A-GP} \|C\|^2 \cdot \sigma_k(\mathcal{C}(\vec{\lambda}^A, \vec{\lambda}^{\mathcal{H}})). \quad (8)$$

Here we have used $\|D_A\| = \|\text{diag}(X_A^* C^*)\| = \|X_A^* C^*\|_{\infty} \leq \|X_A^* C^*\| \leq \|X_A\| \|C\|$, and the well-known property $\sigma_k(\mathcal{ABC}) \leq \|A\| \sigma_k(\mathcal{B}) \|C\|$, see e.g. [6]. The equation (8) shows how the decay in the singular values of P is bounded by the decay in the singular values of the Cauchy matrix. By applying the Sylvester-ADI iteration directly to (7), the maximum from (4) reappears:

$$\frac{\sigma_{k+1}(\mathcal{C}(\vec{\lambda}^A, \vec{\lambda}^{\mathcal{H}}))}{\sigma_1(\mathcal{C}(\vec{\lambda}^A, \vec{\lambda}^{\mathcal{H}}))} \leq \max_{i=1, \dots, n} \prod_{j=1}^k \frac{|\lambda_i^{\mathcal{H}} - \tau_j|}{|\lambda_i^{\mathcal{H}} + \bar{\tau}_j|} \cdot \max_{i=1, \dots, n} \prod_{j=1}^k \frac{|\lambda_i^A - \tau_j|}{|\lambda_i^A + \bar{\tau}_j|}.$$

Similarly, we note that another Cauchy matrix $\hat{P} = -\mathcal{C}(\vec{\lambda}^{\mathcal{H}}, \vec{\lambda}^{\mathcal{H}})$ is the solution of the Lyapunov equation

$$\Lambda_{\mathcal{H}}^* \hat{P} + \hat{P} \Lambda_{\mathcal{H}} = -e e^*, \quad (9)$$

and that the Sylvester-ADI method yields the bound for its singular values that has already appeared in (3):

$$\frac{\sigma_{k+1}(\mathcal{C}(\vec{\lambda}^{\mathcal{H}}, \vec{\lambda}^{\mathcal{H}}))}{\sigma_1(\mathcal{C}(\vec{\lambda}^{\mathcal{H}}, \vec{\lambda}^{\mathcal{H}}))} \leq \left(\max_{i=1, \dots, n} \prod_{j=1}^k \frac{|\lambda_i^{\mathcal{H}} - \tau_j|}{|\lambda_i^{\mathcal{H}} + \bar{\tau}_j|} \right)^2. \quad (10)$$

There is an interesting observation that ties the ADI iterates and the Cauchy matrices even closer. Suppose that $-\mathcal{C}(\vec{\lambda}^{\mathcal{H}}, \vec{\lambda}^{\mathcal{H}}) = L\Delta L^*$ is the LDL-factorization of the (positive definite) matrix $-\mathcal{C}(\vec{\lambda}^{\mathcal{H}}, \vec{\lambda}^{\mathcal{H}})$. One can show that the k -th ADI iterate \hat{P}_k for the equation (9) is equal to

$$\hat{P}_k = L \begin{bmatrix} \Delta(1:k, 1:k) & \\ & 0 \end{bmatrix} L^*,$$

where the shifts used are the eigenvalues ordered as they appear in the tuple $\vec{\lambda}^{\mathcal{H}}$; we omit the proof for brevity. This fact also explains why the bounds in [2] and [25] have “striking resemblance [...] even though the approaches that lead to the results are very different” [25]: the approximations obtained using the Cholesky factorization and the ADI are in fact the same. Also note that our “ADI-minimax” eigenvalue ordering is a very slight modification of the “Cholesky ordering” defined in [2].

Finally, we comment on the size of κ_{A-GP} in (3) and (4). Suppose now that $m = 1$, i.e. that $G = BB^*$, where B is a one-column matrix. Then the row-vector $K = B^*P$ is the unique solution to the following eigenvalue assignment problem:

$$\text{Given } A \text{ and } B, \text{ find } K \text{ such that } \text{eig}(A - BK) = \{\lambda_1^{\mathcal{H}}, \dots, \lambda_n^{\mathcal{H}}\}.$$

Mehrmann and Xu [18] give a bound on the condition number of the closed-loop matrix $A - BK$, which in our case coincides with $A - GP$. In our notation, the bound is

$$\kappa_{A-GP} \leq \sqrt{n} \cdot \|X_A \text{diag}(X_A^{-1}B)\| \cdot \|(X_A \text{diag}(X_A^{-1}B))^{-1}\| \cdot \|\tilde{\mathcal{C}}\|_F \|\tilde{\mathcal{C}}^{-1}\|_F. \quad (11)$$

The matrix $\tilde{\mathcal{C}}$ is associated to yet another Cauchy matrix, $\tilde{\mathcal{C}} = \mathcal{C}(-\vec{\lambda}^A, \vec{\lambda}^{\mathcal{H}}) \cdot D$, where D is the diagonal scaling matrix that makes the columns of $\tilde{\mathcal{C}}$ have unit norm. The following formula [18] holds for the condition number of this matrix in the Frobenius norm:

$$\|\tilde{\mathcal{C}}\|_F \|\tilde{\mathcal{C}}^{-1}\|_F = \sqrt{n \sum_{i,j=1}^n \frac{\sum_{\ell=1}^n \prod_{k=1, k \neq \ell}^n |\lambda_i^A - \lambda_k^{\mathcal{H}}|^2}{\prod_{k=1, k \neq i}^n |\lambda_i^{\mathcal{H}} - \lambda_k^{\mathcal{H}}|^2} \cdot \frac{\prod_{k=1, k \neq j}^n |\lambda_j^A - \lambda_k^{\mathcal{H}}|^2}{\prod_{k=1, k \neq j}^n |\lambda_j^A - \lambda_k^A|^2}}}. \quad (12)$$

While the Cauchy matrices $\mathcal{C}(\vec{\lambda}^{\mathcal{H}}, \vec{\lambda}^{\mathcal{H}})$ and $\mathcal{C}(\vec{\lambda}^A, \vec{\lambda}^{\mathcal{H}})$, each generated by two stable tuples, exhibit rapid singular value decay, the matrix $\mathcal{C}(-\vec{\lambda}^A, \vec{\lambda}^{\mathcal{H}})$ shows a very slow decay. Moreover, the matrix $\tilde{\mathcal{C}}$ has the smallest condition number (up to the factor \sqrt{n}) among all matrices obtained by scaling the columns of the matrix $\mathcal{C}(-\vec{\lambda}^A, \vec{\lambda}^{\mathcal{H}})$ [24]. However, it seems that getting an upper limit for (12) is not an easy task. Informally, one could provide a bound for this expression if the eigenvalues of \mathcal{H} are not far away from the eigenvalues of A . If we can reorder the eigenvalues of A and \mathcal{H} so that $|\lambda_j^A - \lambda_j^{\mathcal{H}}| \ll |\lambda_j^A - \lambda_k^{\mathcal{H}}|$ for a majority of indices j and all $k \neq j$, then $|\lambda_j^A - \lambda_k^{\mathcal{H}}| \approx |\lambda_j^A - \lambda_k^A| \approx |\lambda_j^{\mathcal{H}} - \lambda_k^{\mathcal{H}}|$, and the fractions under the square root are close to 1. Such an assumption is not unrealistic: the Hamiltonian matrix \mathcal{H} can

be considered as a low-rank perturbation of the matrix \mathcal{A} defined below, which has eigenvalues equal to $\vec{\lambda}^A$ plus its mirrors along the imaginary axis:

$$\mathcal{H} = \underbrace{\begin{bmatrix} A & \\ & -A^* \end{bmatrix}}_{\mathcal{A}} + \begin{bmatrix} & BB^* \\ C^*C & \end{bmatrix}. \quad (13)$$

Analysis in the simplest, symmetric case indicates that most eigenvalues of \mathcal{H} cannot escape far away from those of \mathcal{A} . In that case, $p = m$ and $C^* = B$, and we can rewrite (13) as

$$\mathcal{H} = \underbrace{\mathcal{A} - \begin{bmatrix} B & 0 \\ 0 & B \end{bmatrix} \begin{bmatrix} B & 0 \\ 0 & B \end{bmatrix}^*}_{\tilde{\mathcal{A}}} + \begin{bmatrix} B \\ B \end{bmatrix} \begin{bmatrix} B \\ B \end{bmatrix}^*.$$

The eigenvalues λ_j^A , $\lambda_j^{\tilde{\mathcal{A}}}$ and $\lambda_j^{\mathcal{H}}$ of \mathcal{A} , $\tilde{\mathcal{A}}$ and \mathcal{H} , respectively, are real; suppose that they are ordered from the smallest to the largest. The matrix $\tilde{\mathcal{A}}$ is a perturbation of the symmetric matrix \mathcal{A} by a positive semidefinite matrix of rank $2m$, thus [26],

$$\lambda_{j-2m}^A \leq \lambda_j^{\tilde{\mathcal{A}}} \leq \lambda_j^A,$$

for all $j = 2m + 1, \dots, 2n$. Similarly, the matrix \mathcal{H} is a perturbation of the matrix $\tilde{\mathcal{A}}$ by a positive semidefinite matrix of rank m , and

$$\lambda_j^{\tilde{\mathcal{A}}} \leq \lambda_j^{\mathcal{H}} \leq \lambda_{j+m}^{\tilde{\mathcal{A}}},$$

for all $j = 1, \dots, n - m$. Put together,

$$\lambda_{j-2m}^A \leq \lambda_j^{\mathcal{H}} \leq \lambda_{j+m}^A,$$

and therefore only $3m$ eigenvalues of \mathcal{H} can escape the interval $[\lambda_1^A, \lambda_{2n}^A]$. All the other eigenvalues of \mathcal{H} are interlaced with the eigenvalues of \mathcal{A} , and each can appear in a narrow window between the two nearby eigenvalues of \mathcal{A} .

In the non-symmetric case, this argument does not pass, and one can construct examples where many eigenvalues of \mathcal{H} move far away from those of A and $-A^*$, for carefully chosen perturbations B, C with large norms. However, behavior similar as for the symmetric matrix \mathcal{H} is still observed for generic perturbations.

We conclude this section with an example that illustrates our discussion.

Example 3. Consider the matrix $A \in \mathbb{R}^{100 \times 100}$ obtained by semi-discretization of the partial differential equation

$$\frac{\partial u}{\partial t} = \Delta u + 10u_x + 100u_y$$

on the unit square, and vectors $B, C^* \in \mathbb{R}^{100}$ with entries normally distributed in $[-10, 10]$. The eigenvalues of A and of the Hamiltonian matrix $\mathcal{H} = \begin{bmatrix} A & BB^* \\ C^*C & -A^* \end{bmatrix}$ are

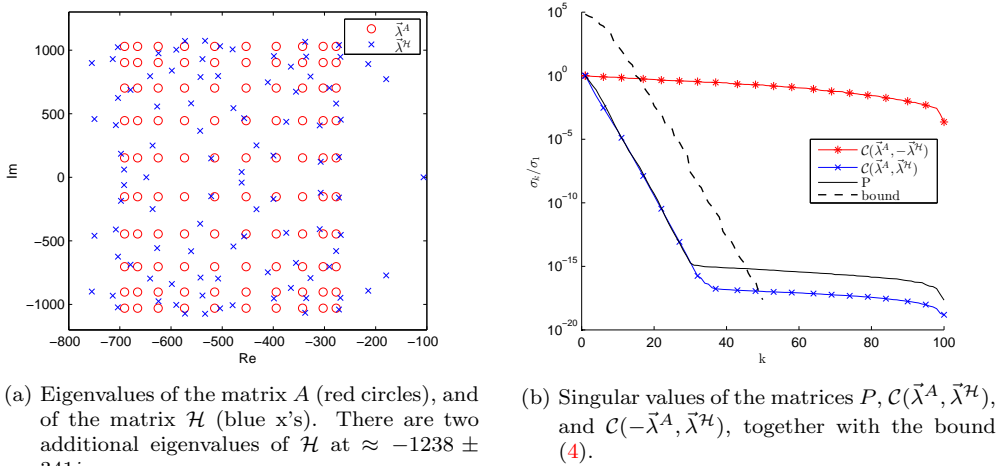


Figure 1: The relation of the singular value decay of the Riccati solution P and the singular values of the Cauchy matrix $\mathcal{C}(\vec{\lambda}^A, \vec{\lambda}^{\mathcal{H}})$.

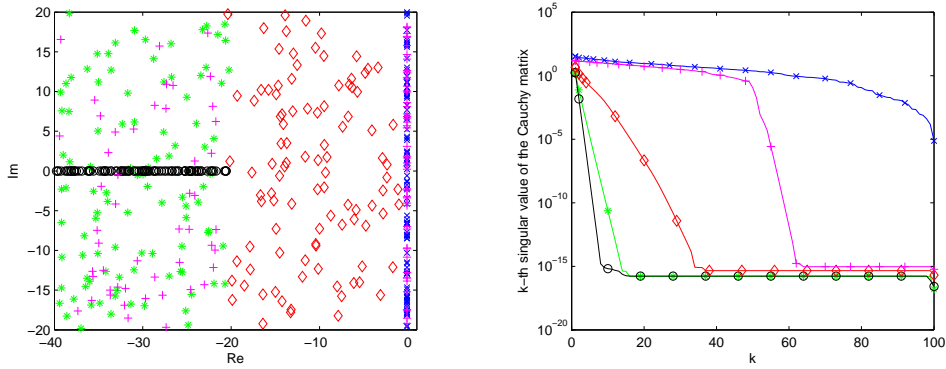
shown in Figure 1a. As we commented above, most eigenvalues of \mathcal{H} are only slight perturbations of the eigenvalues of A . Figure 1b shows the rapid decay in the singular values of the solution P of the Riccati equation associated with \mathcal{H} —already the 20th singular value is approximately 10^{10} times smaller than the first one. It is interesting to observe that the decay in singular values of the Cauchy matrix $\mathcal{C}(\vec{\lambda}^A, \vec{\lambda}^{\mathcal{H}})$ closely follows the one of P ; the same holds true for the matrix $\mathcal{C}(\vec{\lambda}^{\mathcal{H}}, \vec{\lambda}^{\mathcal{H}})$ which is not shown in the figure.

The slope of the line for the bound (4) is very similar to the slope in the decay of the singular values of P . Shifts τ_j have been chosen here as the eigenvalues of \mathcal{H} ordered in the “ADI-minimax” ordering.

In this example, the bound (11) for κ_{A-GP} is quite pessimistic. While the condition number of the matrix $A - GP$ is $\kappa_{A-GP} \approx 104.1$, the condition numbers of the two matrices $X_A \text{diag}(X_A^{-1}B)$ and \tilde{C} are 1987.7 and 3753.5, respectively. Nevertheless, in Figure 1b we see that the singular values of $\mathcal{C}(-\vec{\lambda}^A, \vec{\lambda}^{\mathcal{H}})$ (and then of \tilde{C} as well) decline drastically slower than the ones of $\mathcal{C}(\vec{\lambda}^A, \vec{\lambda}^{\mathcal{H}})$.

Further demonstration of the rapid decay for the Cauchy matrices generated by stable vectors is shown in Figure 2. On the left side, several different selections for the eigenvalue tuple $\vec{\lambda}^{\mathcal{H}}$ are depicted: black circles correspond to 100 randomly chosen points in the interval $[-40, -20]$, green asterisks to 100 points in the rectangle $[-40, -20] \times [-20i, -20i]$, red diamonds to 100 points in the rectangle $[-21, -1] \times [-20i, 20i]$, blue x's to 100 points on the line $-0.1 + [-20i, 20i]$, and purple pluses to 50 points in the rectangle $[-40, -20] \times [-20i, -20i]$ and 50 points on the line $-0.1 + [-20i, 20i]$.

For each selection of the tuple, on the right figure we plot the singular value decay of the matrix $\mathcal{C}(\vec{\lambda}^{\mathcal{H}}, \vec{\lambda}^{\mathcal{H}})$. The decay is clearly faster when the eigenvalues are further away from the imaginary axis, which can also be justified by the bound (10).



(a) Various placements of $\vec{\lambda}^{\mathcal{H}}$; each color specifies another tuple of 100 eigenvalues. (b) Singular values of $\mathcal{C}(\vec{\lambda}^{\mathcal{H}}, \vec{\lambda}^{\mathcal{H}})$; to each line corresponds the placement $\vec{\lambda}^{\mathcal{H}}$ of the same color in the left figure.

Figure 2: Singular value decay of the Cauchy matrix $\mathcal{C}(\vec{\lambda}^{\mathcal{H}}, \vec{\lambda}^{\mathcal{H}})$ and its dependency on the placement of the eigenvalue tuple $\vec{\lambda}^{\mathcal{H}}$ in the complex plane.

3 Approximations using the Hamiltonian eigenspaces

In the small-scale setting, the solution to the Riccati equation (1) is obtained by computing the stable invariant subspace of the matrix \mathcal{H} : if

$$\begin{bmatrix} A & G \\ Q & -A^* \end{bmatrix} \begin{bmatrix} X \\ Y \end{bmatrix} = \begin{bmatrix} X \\ Y \end{bmatrix} \cdot \Lambda, \quad (14)$$

where the columns of $X, Y \in \mathbb{C}^{n \times n}$ form the basis of the eigenspace, and all the eigenvalues of $\Lambda \in \mathbb{C}^{n \times n}$ lie in the left half-plane, then the stabilizing solution is given by the formula

$$P = -YX^{-1}.$$

Due to computational complexity and memory capacity limitations, for larger values of n it is feasible to compute only a small-dimensional eigenspace of \mathcal{H} . The discussion from Section 2 can be adapted to demonstrate that there indeed exists such a subspace, which can be used in order to obtain a good approximation to the exact solution P . To see that, let us assume that the matrix Λ is diagonal, and that the matrix containing the eigenvectors of \mathcal{H} has been normalized so that each column of the matrix X has a unit norm. Then, the \mathbb{X} component of (14) (i.e. the first n rows) reads as

$$AX + GY = X\Lambda \quad \Rightarrow \quad A - GP = X\Lambda X^{-1},$$

and we see that the matrix X is the eigenvector matrix of $A - GP$. Once again, we assume that this matrix has a modest condition number.

The \mathbb{Y} component of (14) can be rewritten as the Sylvester equation

$$A^*Y + Y\Lambda = QX,$$

where Y is the unknown. By using the Sylvester-ADI technique as described in the previous section, one can conclude that the singular values of the matrix Y follow the

rapid decay of the singular values in the Cauchy matrix $\mathcal{C}(\vec{\lambda}^A, \vec{\lambda}^{\mathcal{H}})$. Thus already for some $k \ll n$, we can have that $\sigma_k(Y) \leq \epsilon$, where $\epsilon \ll 1$ is a given tolerance. To see that computing only k of the eigenvectors will suffice for a good approximation to P , note that there exists a rank-revealing QR-factorization [12, 7] of Y ,

$$Y\Pi = \begin{bmatrix} Q_1 & Q_2 \end{bmatrix} \begin{bmatrix} R_{11} & R_{12} \\ & R_{22} \end{bmatrix},$$

such that $\|R_{22}\| \leq f(n)\sigma_{k+1}(Y)$. Here the orthonormal matrix Q_1 has k columns, R_{11} is an upper triangular matrix of order k , Π is a permutation matrix, and $f(n)$ is a low-degree polynomial. Let

$$\hat{Y}_k = Q_1 \begin{bmatrix} R_{11} & R_{12} \end{bmatrix}, \quad P_k = -\hat{Y}_k X^{-1}.$$

Then the range of the matrix P_k is spanned by certain k columns of the matrix Y —more precisely, these are the first k columns of $Y\Pi$ and we will denote those as Y_k . Furthermore, we have

$$\begin{aligned} \|P - P_k\| &= \|(-Y + \hat{Y}_k)X^{-1}\| \leq \|Y - \hat{Y}_k\| \cdot \|X^{-1}\| \\ &= \left\| \begin{bmatrix} Q_1 & Q_2 \end{bmatrix} \begin{bmatrix} 0 & 0 \\ & R_{22} \end{bmatrix} \right\| \cdot \|X^{-1}\| = \|R_{22}\| \cdot \|X^{-1}\| \\ &\leq f(n)\epsilon \|X^{-1}\| \leq f(n)\kappa_{A-GP} \cdot \epsilon. \end{aligned}$$

In the last inequality we have used the fact that the matrix X has columns of unit norm. Therefore, one can compute only k eigenvectors of the matrix \mathcal{H} to obtain the matrix Y_k . All columns of the matrix P_k , which is a good approximation to P , are simply certain linear combinations of the columns of Y_k .

The construction of the approximation P_k that we have just described still requires knowledge of the entire stable invariant subspace of \mathcal{H} . In order to make the method practical, one has to resolve the following two issues while avoiding the computation of the full eigenvector matrices X and Y :

1. Which are the k stable eigenpairs of \mathcal{H} that should be computed?
2. Once the k eigenpairs of \mathcal{H} are available, how to compute the matrix P_k (or some other good approximation to P)?

We first turn our attention to the second question. Suppose that we have computed a k -dimensional stable invariant subspace,

$$\mathcal{H} \begin{bmatrix} X_k \\ Y_k \end{bmatrix} = \begin{bmatrix} X_k \\ Y_k \end{bmatrix} \Lambda_k. \quad (15)$$

As we have seen, a good approximation to P may be formed by constructing a matrix whose columns belong to the subspace $\text{span}\{Y_k\}$. One such matrix (denoted from now on as P_k^{proj}) is obtained by solving the Riccati equation projected onto $\text{span}\{Y_k\}$. The procedure is described in Algorithm 1.

Algorithm 1: Riccati solution approximation based on the projection to the invariant subspace of the Hamiltonian matrix

Input: $A, G, Q \in \mathbb{R}^{n \times n}$, $X_k, Y_k \in \mathbb{R}^{n \times k}$ that satisfy (15)

Output: Approximation P_k^{proj} to the solution of $A^*P + PA + Q - PGP = 0$

1 Compute an orthonormal basis U_k for the subspace $\text{span}\{Y_k\}$;

2 Obtain the stabilizing solution Δ to the projected Riccati equation of order k :

$$(U_k^*AU_k)^*\Delta + \Delta(U_k^*AU_k) + (U_k^*QU_k) - \Delta(U_k^*GU_k)\Delta = 0;$$

3 Set $P_k^{\text{proj}} = U_k\Delta U_k^*$;

The approximate solution P_k^{proj} is obviously positive semidefinite (assuming that the projected Riccati equation has one), and direct calculation shows that it satisfies the Galerkin condition

$$Y_k^*\mathcal{R}(P_k^{\text{proj}})Y_k = 0.$$

Another option is to avoid solving the projected Riccati equation, and instead compute the approximation (denoted here as P_k^{dir}) directly from X_k and Y_k . The following formula was suggested by Amodei and Buchot [1]:

$$P_k^{\text{dir}} = -Y_k(Y_k^*X_k)^{-1}Y_k^*. \quad (16)$$

Note that the columns of P_k^{dir} also belong to $\text{span}\{Y_k\}$. For simplicity of exposition, we study here only the generic case where the matrix Y_k is of full column rank, which implies that $Y_k^*X_k$ is non-singular [1]. The motivation for introducing (16) was to mimic certain properties of the exact solution $P = -YX^{-1}$, as shown in the following proposition.

Proposition 4 ([1]). *The matrix P_k^{dir} has the following properties:*

- (a) P_k^{dir} is a symmetric positive semidefinite matrix.
- (b) $P_k^{\text{dir}}X_k = -Y_k$ (c.f. $PX = -Y$).
- (c) Restriction of the closed loop matrix $A - GP_k^{\text{dir}}$ to the subspace $\text{span}\{X_k\}$ is stable.
- (d) P_k^{dir} is the unique matrix of rank k such that (a) and (b) hold.
- (e) P_k^{dir} fulfills the following Galerkin condition: $X_k^*\mathcal{R}(P_k^{\text{dir}}) = 0$.

The Riccati residual of the approximation P_k^{dir} can be written in a special form which suggests the role the selected eigenvectors have in providing the approximation. With the help of this form, the following theorem also explains the dynamics of how the approximation P_k^{dir} improves as the number of eigenvectors k increases.

Theorem 5. Let $\tilde{C}^* = (I - Y_k(Y_k^* X_k)^{-1} X_k^*) C^*$, and let $\tilde{A} = A - G P_k^{\text{dir}}$.

(a) The Riccati residual of the approximation P_k^{dir} can be written as

$$\mathcal{R}(P_k^{\text{dir}}) = \tilde{C}^* \tilde{C}. \quad (17)$$

(b) If $(\Lambda_\ell, \begin{bmatrix} X_\ell \\ Y_\ell \end{bmatrix})$ is an eigenpair of the matrix \mathcal{H} , i.e. $\mathcal{H} \begin{bmatrix} X_\ell \\ Y_\ell \end{bmatrix} = \begin{bmatrix} X_\ell \\ Y_\ell \end{bmatrix} \Lambda_\ell$, where $X_\ell, Y_\ell \in \mathbb{C}^{n \times \ell}$, $\Lambda_\ell \in \mathbb{C}^{\ell \times \ell}$ for some ℓ , then

$$(\tilde{\Lambda}_\ell, \begin{bmatrix} \tilde{X}_\ell \\ \tilde{Y}_\ell \end{bmatrix}) = (\Lambda_\ell, \begin{bmatrix} X_\ell \\ Y_\ell + P_k^{\text{dir}} X_\ell \end{bmatrix})$$

is an eigenpair of the matrix

$$\tilde{\mathcal{H}} = \begin{bmatrix} \tilde{A} & G \\ \tilde{C}^* \tilde{C} & -\tilde{A}^* \end{bmatrix}. \quad (18)$$

When Λ_ℓ is stable, then $Y_\ell + P_k^{\text{dir}} X_\ell = (I - Y_k(Y_k^* X_k)^{-1} X_k^*) Y_\ell$.

(c) Let \tilde{P} denote the exact stabilizing solution of the Riccati equation associated with the Hamiltonian matrix (18),

$$\tilde{A}^* \tilde{P} + \tilde{P} \tilde{A} + \tilde{C}^* \tilde{C} - \tilde{P} G \tilde{P} = 0. \quad (19)$$

Then the exact stabilizing solution P of the original Riccati equation (1) is given by

$$P = P_k^{\text{dir}} + \tilde{P}.$$

(d) Let $(\Lambda_\ell, \begin{bmatrix} X_\ell \\ Y_\ell \end{bmatrix})$ denote an eigenpair of the matrix \mathcal{H} , such that the spectrum of Λ_ℓ is stable and disjoint from the spectrum of Λ_k . Furthermore, let $P_{k+\ell}^{\text{dir}}$ denote the approximate solution to the initial Riccati equation (1) computed using the eigenpair $(\begin{bmatrix} \Lambda_k & \\ & \Lambda_\ell \end{bmatrix}, \begin{bmatrix} X_k & X_\ell \\ Y_k & Y_\ell \end{bmatrix})$ of \mathcal{H} , and let $\tilde{P}_\ell^{\text{dir}}$ denote the approximate solution to the residual Riccati equation (19) computed using the eigenpair $(\tilde{\Lambda}_\ell, \begin{bmatrix} \tilde{X}_\ell \\ \tilde{Y}_\ell \end{bmatrix})$ of $\tilde{\mathcal{H}}$ as defined in (b). Then

$$P_{k+\ell}^{\text{dir}} = P_k^{\text{dir}} + \tilde{P}_\ell^{\text{dir}}. \quad (20)$$

Proof. (a) First, note that from $A X_k + G Y_k = X_k \Lambda_k$, we have

$$\begin{aligned} P_k^{\text{dir}} G P_k^{\text{dir}} &= -P_k^{\text{dir}} G Y_k (Y_k^* X_k)^{-1} Y_k^* = -P_k^{\text{dir}} (X_k \Lambda_k - A X_k) (Y_k^* X_k)^{-1} Y_k^* \\ &= P_k^{\text{dir}} A X_k (Y_k^* X_k)^{-1} Y_k^* + Y_k \Lambda_k (Y_k^* X_k)^{-1} Y_k^*, \end{aligned}$$

and that from $QX_k - A^*Y_k = Y_k\Lambda_k$, we have

$$\begin{aligned} A^*P_k^{\text{dir}} &= -A^*Y_k(Y_k^*X_k)^{-1}Y_k^* = (Y_k\Lambda_k - QX_k)(Y_k^*X_k)^{-1}Y_k^* \\ &= Y_k\Lambda_k(Y_k^*X_k)^{-1}Y_k^* - QX_k(Y_k^*X_k)^{-1}Y_k^*. \end{aligned}$$

Since Q and $Y_k^*X_k$ are symmetric matrices (isotropy!), it follows that

$$P_k^{\text{dir}}A = Y_k(Y_k^*X_k)^{-1}\Lambda_k^*Y_k^* - Y_k(Y_k^*X_k)^{-1}X_k^*Q,$$

and

$$\begin{aligned} P_k^{\text{dir}}GP_k^{\text{dir}} &= Y_k(Y_k^*X_k)^{-1}\Lambda_k^*Y_k^* - Y_k(Y_k^*X_k)^{-1}X_k^*QX_k(Y_k^*X_k)^{-1}Y_k^* \\ &\quad + Y_k\Lambda_k(Y_k^*X_k)^{-1}Y_k^*. \end{aligned}$$

Adding it all together, all the terms containing Λ_k cancel out, and we obtain

$$\begin{aligned} \mathcal{R}(P_k^{\text{dir}}) &= A^*P_k^{\text{dir}} + P_k^{\text{dir}}A + Q - P_k^{\text{dir}}GP_k^{\text{dir}} \\ &= (I - Y_k(Y_k^*X_k)^{-1}X_k^*)Q(I - Y_k(Y_k^*X_k)^{-1}X_k^*)^*. \end{aligned}$$

(b) By simply computing the \mathbb{X} and the \mathbb{Y} component, we show that

$$\begin{bmatrix} \tilde{A} & G \\ \tilde{C}^*\tilde{C} & -\tilde{A}^* \end{bmatrix} \begin{bmatrix} X_\ell \\ Y_\ell + P_k^{\text{dir}}X_\ell \end{bmatrix} = \begin{bmatrix} X_\ell \\ Y_\ell + P_k^{\text{dir}}X_\ell \end{bmatrix} \Lambda_\ell.$$

The \mathbb{X} component is straightforward:

$$\begin{aligned} \tilde{A}X_\ell + G(Y_\ell + P_k^{\text{dir}}X_\ell) &= (A - GP_k^{\text{dir}})X_\ell + GY_\ell + GP_k^{\text{dir}}X_\ell \\ &= AX_\ell + GY_\ell = X_\ell\Lambda_\ell. \end{aligned}$$

The \mathbb{Y} component is a bit more involved: with $\tilde{Q} = \tilde{C}^*\tilde{C}$ we have

$$\begin{aligned} &\tilde{Q}X_\ell - \tilde{A}^*(Y_\ell + P_k^{\text{dir}}X_\ell) \\ &= \tilde{Q}X_\ell - A^*Y_\ell + P_k^{\text{dir}}GY_\ell - A^*P_k^{\text{dir}}X_\ell + P_k^{\text{dir}}GP_k^{\text{dir}}X_\ell \\ &= \tilde{Q}X_\ell + (Y_\ell\Lambda_\ell - QX_\ell) + P_k^{\text{dir}}(X_\ell\Lambda_\ell - AX_\ell) - A^*P_k^{\text{dir}}X_\ell \\ &\quad + P_k^{\text{dir}}GP_k^{\text{dir}}X_\ell \\ &= \tilde{Q}X_\ell + (Y_\ell + P_k^{\text{dir}}X_\ell)\Lambda_\ell - (A^*P_k^{\text{dir}} + P_k^{\text{dir}}A + Q - P_k^{\text{dir}}GP_k^{\text{dir}})X_\ell \\ &= (Y_\ell + P_k^{\text{dir}}X_\ell)\Lambda_\ell, \end{aligned}$$

since $A^*P_k^{\text{dir}} + P_k^{\text{dir}}A + Q - P_k^{\text{dir}}GP_k^{\text{dir}} = \mathcal{R}(P_k^{\text{dir}}) = \tilde{Q}$. When Λ_ℓ is stable, $Y_k^*X_\ell = X_k^*Y_\ell$ since the stable subspace is isotropic. Hence,

$$\begin{aligned} Y_\ell + P_k^{\text{dir}}X_\ell &= Y_\ell - Y_k(Y_k^*X_k)^{-1}Y_k^*X_\ell = Y_\ell - Y_k(Y_k^*X_k)^{-1}X_k^*Y_\ell \\ &= (I - Y_k(Y_k^*X_k)^{-1}X_k^*)Y_\ell. \end{aligned}$$

(c) The claim follows from

$$\begin{aligned} A^*(P_k^{\text{dir}} + \tilde{P}) + (P_k^{\text{dir}} + \tilde{P})A + Q - (P_k^{\text{dir}} + \tilde{P})G(P_k^{\text{dir}} + \tilde{P}) \\ = (A - GP_k^{\text{dir}})^*\tilde{P} + \tilde{P}(A - GP_k^{\text{dir}}) - \tilde{P}G\tilde{P} + \mathcal{R}(P_k^{\text{dir}}), \end{aligned}$$

which is equal to zero since $\mathcal{R}(P_k^{\text{dir}}) = \tilde{Q}$.

(d) Since $P_{k+\ell}^{\text{dir}}$ is the unique positive semidefinite matrix of rank $k + \ell$ that satisfies $P_{k+\ell}^{\text{dir}} \begin{bmatrix} X_k & X_\ell \\ Y_k & Y_\ell \end{bmatrix} = - \begin{bmatrix} Y_k & Y_\ell \end{bmatrix}$, it suffices to show that the same holds true for $P_k^{\text{dir}} + \tilde{P}_\ell^{\text{dir}}$. Firstly, we have

$$(P_k^{\text{dir}} + \tilde{P}_\ell^{\text{dir}})X_k = -Y_k + \tilde{P}_\ell^{\text{dir}}X_k = -Y_k.$$

Here $\tilde{P}_\ell^{\text{dir}}X_k = 0$ since

$$\tilde{Y}_\ell^*X_k = (Y_\ell + P_k^{\text{dir}}X_\ell)^*X_k = Y_\ell^*X_k + X_\ell^*P_k^{\text{dir}}X_k = Y_\ell^*X_k - X_\ell^*Y_k = 0,$$

and the columns of $\begin{bmatrix} X_k & X_\ell \\ Y_k & Y_\ell \end{bmatrix}$ lie in the isotropic stable invariant subspace of \mathcal{H} . Secondly,

$$\begin{aligned} (P_k^{\text{dir}} + \tilde{P}_\ell^{\text{dir}})X_\ell &= P_k^{\text{dir}}X_\ell + \tilde{P}_\ell^{\text{dir}}\tilde{X}_\ell = P_k^{\text{dir}}X_\ell - \tilde{Y}_\ell \\ &= P_k^{\text{dir}}X_\ell - (Y_\ell + P_k^{\text{dir}}X_\ell) = -Y_\ell. \end{aligned}$$

□

From the first claim of Theorem 5, we have that

$$\|\mathcal{R}(P_k^{\text{dir}})\| = \|(I - Y_k(Y_k^*X_k)^{-1}X_k^*)C^*\|^2. \quad (21)$$

Note that the matrix $Y_k(Y_k^*X_k)^{-1}X_k^*$ is the skew projector onto the subspace $\text{span}\{Y_k\}$ along the orthogonal complement of the subspace $\text{span}\{X_k\}$. Thus, using the eigenspace $\text{span}\{\begin{bmatrix} X_k \\ Y_k \end{bmatrix}\}$ for computing the approximation P_k^{dir} has the effect of purging the component that C^* has in $\text{span}\{Y_k\}$. Obviously, finding a Y_k such that all columns of C^* lie inside the subspace $\text{span}\{Y_k\}$ will result in P_k^{dir} being equal to the exact solution P .

The statements (b), (c) and (d) of the theorem show that, by computing more and more eigenpairs of \mathcal{H} , we implicitly drive \tilde{C} and the solution \tilde{P} of the residual Riccati equation (19) towards the zero matrix. Equivalently, the bottom left block in the matrix $\tilde{\mathcal{H}}$ of (18) converges to zero, along with the \mathbb{Y} components of its eigenvectors. Also, notice that the \mathbb{Y} components of the eigenvectors of $\tilde{\mathcal{H}}$ that correspond to the eigenvalues used for computing P_k^{dir} are exactly equal to zero.

Moreover, suppose that we incrementally build P_k^{dir} , by initially using a single eigenpair to obtain P_1^{dir} , and then in each of the $k - 1$ steps adding another eigenpair to obtain $P_2^{\text{dir}}, \dots, P_k^{\text{dir}}$. Theorem 5 shows that

$$0 \leq P_1^{\text{dir}} \leq P_2^{\text{dir}} \leq \dots \leq P_k^{\text{dir}} \leq P = P_k^{\text{dir}} + \tilde{P},$$

where the inequality symbol denotes the Löwner partial ordering. Thus the approximations P_k^{dir} have another desirable property: as k increases, they are monotonically increasing towards the exact solution.

The problem of finding a small subset of eigenvectors of \mathcal{H} that yields a good approximation to P is more difficult—even when the full eigenvalue decomposition of \mathcal{H} is available, and the only task is to filter out the few representative eigenpairs. The formula (21) for the residual norm indicates where the difficulty is when we use P_k^{dir} as the approximation. To obtain a small residual, one has to find eigenvectors such that the skew projection of C^* vanishes. Thus, it seems that strategies for selecting the eigenpairs which are based purely on eigenvalue location have no prospect. For example, computing P_k^{dir} corresponding to the k largest or k smallest eigenvalues of \mathcal{H} may require a large k for the approximation to become good.

Nevertheless, sorting the eigenvalues in the “ADI-minimax” ordering, and expanding the approximation subspace by adding the eigenpairs in that order works very well in most cases, although one can carefully design a Riccati equation for which such a method fails as well, as we will demonstrate in Example 6.

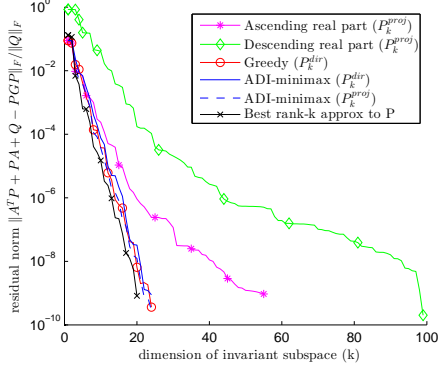
Theorem 5(c) shows that, once P_k^{dir} is computed, one can switch to solving the residual Riccati equation (19). Since \mathbb{Y} components in all stable eigenvectors of \tilde{H} have to be driven to zero eventually, it is reasonable to expand the current invariant subspace with the eigenvector of $\tilde{\mathcal{H}}$ that has the largest \mathbb{Y} component. This approach results in fast convergence. However, an efficient method for finding such an eigenvector is elusive. The following example summarizes the problematics of choosing a proper invariant subspace.

Example 6. Consider the matrix $A \in \mathbb{R}^{100 \times 100}$ obtained by semi-discretization of the PDE

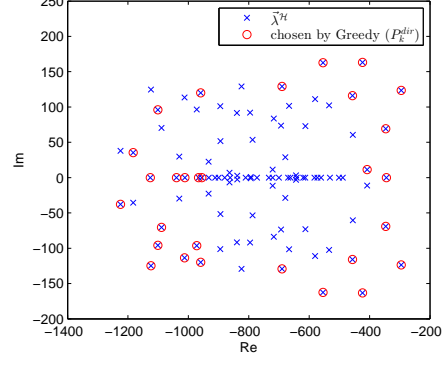
$$\frac{\partial u}{\partial t} = \Delta u + 20u_x - 180u,$$

and matrices $B, C^* \in \mathbb{R}^{100 \times 3}$ with entries chosen randomly from $[0, 10]$. For this small example, we compute the full eigenvalue decomposition of the matrix $\mathcal{H} = \begin{bmatrix} A & BB^* \\ C^*C & -A^* \end{bmatrix}$. By generating various stable subspaces of \mathcal{H} , we compare the effectiveness of approximations P_k^{dir} and P_k^{proj} .

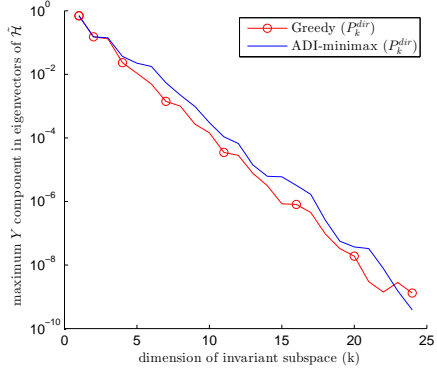
Figure 3a shows the residuals in the Riccati equation for different strategies of expanding the invariant subspace. Suppose, for example, that we sort the eigenvalues of \mathcal{H} by decreasing real part, and then use the first k associated eigenvectors as the basis for the invariant subspace. The green line with diamond marks shows the residuals of the corresponding approximate Riccati solutions P_k^{proj} . Similarly, the purple line with asterisks shows the residuals when the eigenvalues are sorted by increasing real part. In both cases, the final subspace (the one with the residual less than 10^{-9}) has a dimension much larger than necessary. To see that, we computed the exact solution P , and for each k the matrix of rank k that is closest to P in 2-norm. Columns of this matrix span a certain subspace of dimension k , and the Riccati residuals follow the



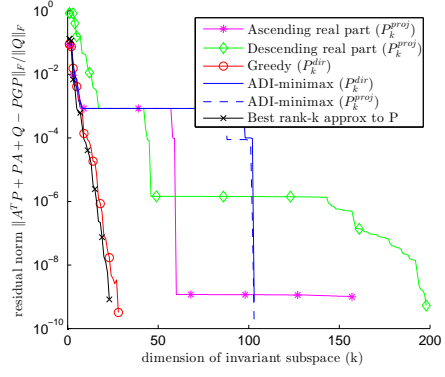
(a) Residuals in the Riccati equation.



(b) Eigenvalues of \mathcal{H} chosen by the greedy method.



(c) Maximum norms of \mathbb{Y} components in eigenvectors of matrices $\tilde{\mathcal{H}}$ in each expansion step.



(d) Stagnation for strategies based solely on eigenvalue location.

Figure 3: Various strategies for selecting the stable invariant subspace of \mathcal{H} .

black line with the x marks. Thus, there exists a subspace of dimension 20 yielding an approximate Riccati solution with the residual of 10^{-9} , while sorting the eigenvalues of \mathcal{H} by the real part requires almost the full stable subspace to reach that accuracy in one case, and a subspace of dimension 55 in the other.

On the other hand, using the ADI-minimax ordering of the eigenvalues makes the convergence very fast in this case, with either P_k^{dir} or P_k^{proj} . We must note here that using suboptimal invariant subspaces with P_k^{dir} may lead to numerical problems. The issue is the following: if Y_k does not have a full numerical rank, then the inverse in the formula $P_k^{\text{dir}} = -Y_k(Y_k^* X_k)^{-1} Y_k^*$ is not well-defined. This issue, along with Theorem 5, motivates another possible strategy: in each step we expand the invariant subspace with the eigenvector of $\tilde{\mathcal{H}}$ which has the largest \mathbb{Y} component. This strategy, denoted by “greedy” in the figures, ensures that the matrix Y_k keeps a full column rank, and also tends to provide the strongest decrease in the Riccati residual. Eigenvalues of \mathcal{H} chosen by the greedy strategy are shown in Figure 3b.

Example 7. *In this example, we demonstrate that using only the location of the eigenvalues in the complex plane as a criterion for choosing the invariant subspace is not sufficient. The goal is to construct a Hamiltonian matrix such that many eigenvalues that must be computed for a good approximation are artificially clustered somewhere in the interior of the spectrum (in this case, in the vicinity of the point -700). Let*

$$\mathcal{H}_1 = \left[\begin{array}{cc|cc} A & & BB^* & \\ & 10^{-3}A - 700I & & 10^{-3}BB^* \\ \hline C^*C & & -A^* & \\ & 10^{-3}C^*C & & -(10^{-3}A - 700I)^* \end{array} \right],$$

where A, B, C are as in Example 6. The Riccati equation associated with \mathcal{H}_1 effectively consists of two smaller equations of order 100. Its solution has the form

$$P_1 = \begin{bmatrix} P & \\ & P_\epsilon \end{bmatrix},$$

where $\|P_\epsilon\| \approx 10^{-3} \|P\|$, and P is the exact solution from Example 6. Therefore, in order to capture P_1 with relative precision higher than 10^{-3} , the invariant subspace has to include several eigenvalues of \mathcal{H}_1 that are clustered around -700 . As we see in Figure 3d, all of the methods we described in Example 6 that are based exclusively on eigenvalue location need many iterations until they discover the hidden, but crucial part of the spectrum. This effect manifests itself as the typical stagnation plateaus in the graphs. On the other hand, the greedy method immediately finds the eigenvectors necessary for further progress in convergence. Unlike this artificial example, the ADI-minimax ordering performed quite well in our experimentation with Riccati equations arising in applications and in the typical benchmark problems.

As mentioned before, the ADI-minimax strategy described in the examples requires full knowledge of the spectrum of \mathcal{H} , while the greedy strategy needs all the eigenvectors as well. Both of these are prohibitive for practical purposes. To remedy this issue, one could deploy various heuristics, such as using Ritz values as approximations to the eigenvalues for the ADI-minimax ordering. On the other hand, computing the \mathbb{Y} components of the Ritz vectors for the greedy strategy is more expensive and appears to be less reliable.

Similar techniques are also necessary in other methods for solving the large-scale matrix equations. For example, optimal shifts in ADI-methods are solutions to minimax problems similar to (6), which then have to be solved only approximately in order to be feasible.

4 Relations with the Krylov projection methods and the qADI

In the previous section we have discussed on how to compute the approximation to the solution of the Riccati equation once a certain number of eigenpairs for the Hamiltonian

matrix \mathcal{H} is available, and on the problem of selecting the proper eigenpairs for such a task. However, an interesting interconnection with the already existing methods for the Riccati equation arises when we start taking into consideration the process of computing the eigenpairs itself.

There are several methods available for the partial eigenvalue problem of a large-scale Hamiltonian matrix. In particular, symplectic Lanczos methods [3, 4] compute the exterior eigenvalues and the associated eigenvectors while obeying the Hamiltonian structure. Preservation of the structure is also a key feature of the SHIRA method [17], which computes the eigenvalues closest to some given target $\tau \in \mathbb{C}$. All of these methods approximate the eigenpairs by building up certain Krylov subspaces generated by the matrix $f(\mathcal{H})$, where f is a scalar rational function.

4.1 Krylov subspace methods based on (A^*, C^*)

On the other hand, the common projection methods for the Riccati equation also use Krylov spaces; however, those methods involve the matrix A^* . Indeed, the literature describes methods using the ordinary Krylov subspaces generated by the matrix A^* [13], or the extended ones generated by both A^* and A^{-*} [22, 11], or the rational Krylov subspaces generated with A^* [8, 23]. These methods are summarized in Algorithm 2; we use the following notation for the k -dimensional Krylov subspace generated by some matrix M and an initial (block-)vector v :

$$\mathcal{K}_k(M, v) = \text{span}\{v, Mv, M^2v, \dots, M^{k-1}v\},$$

and the rational Krylov subspace that uses shifts $\vec{\sigma} = (\sigma_1, \dots, \sigma_{k-1}) \in \mathbb{C}^{k-1}$:

$$\mathcal{K}(M, v, \vec{\sigma}) = \text{span}\{v, (M - \sigma_1 I)^{-1}v, (M - \sigma_2 I)^{-1}v, \dots, (M - \sigma_{k-1} I)^{-1}v\}.$$

Algorithm 2: Krylov projection methods for the Riccati equation

Input: $A \in \mathbb{R}^{n \times n}$, $B \in \mathbb{R}^{n \times m}$, $C \in \mathbb{R}^{p \times n}$

Output: Approximation to the solution of $A^*P + PA + C^*C - PBB^*P = 0$

- 1 Build an orthonormal basis U for one of:
 - ordinary Krylov subspace $\mathcal{K}_k(A^*, C^*)$;
 - extended Krylov subspace $\mathcal{K}_\ell(A^*, C^*) + \mathcal{K}_{k-\ell}(A^{-*}, C^*)$; [KPIK]
 - rational Krylov subspace $\mathcal{K}(A^*, C^*, -\vec{\sigma})$; [RKSM]
 - 2 Solve the projected Riccati equation:
$$(U^*AU)^*\Delta + \Delta(U^*AU) + (CU)^*(CU) - \Delta(U^*B)(U^*B)^*\Delta = 0;$$
 - 3 Approximate $P_k^{\text{cry}} = U\Delta U^*$;
-

All variants of the algorithm were first applied to the Lyapunov equation, for which a well-developed theoretical background exists on why precisely the aforementioned subspaces provide a good approximation. Only later was it observed that by simply replacing the projected Lyapunov equation with the projected Riccati equation, one obtains a good approximation for the solution of (1) as well, while keeping the same

projection subspace as in the Lyapunov case. The following theorem gives a new perspective on why that is the case.

Theorem 8. *For any $F \in \mathbb{C}^{n \times p}$, the following equalities hold:*

- (a) $\mathbb{Y}(\mathcal{K}_k(\mathcal{H}, \begin{bmatrix} F \\ C^* \end{bmatrix})) = \mathcal{K}_k(A^*, C^*)$.
- (b) $\mathbb{Y}(\mathcal{K}_k(\mathcal{H}^{-1}, \begin{bmatrix} F \\ C^* \end{bmatrix})) = \mathcal{K}_k(A^{-*}, C^*)$.
- (c) *Suppose that A is stable and $\vec{\sigma} \in \mathbb{C}_-^{k-1}$. If none of $\sigma_1, \dots, \sigma_{k-1}$ is an eigenvalue of \mathcal{H} , then*

$$\mathbb{Y}(\mathcal{K}(\mathcal{H}, \begin{bmatrix} F \\ C^* \end{bmatrix}, \vec{\sigma})) \subseteq \mathcal{K}(A^*, C^*, -\vec{\sigma}),$$

with the equality holding when none of the shifts is an eigenvalue of A and $F = 0$. If σ_j is an eigenvalue of \mathcal{H} , then $\mathcal{K}(A^, C^*, -\vec{\sigma})$ contains the \mathbb{Y} component of the associated eigenvector.*

Proof.

- (a) It suffices to show $\mathbb{Y}(\mathcal{H}^j \begin{bmatrix} F \\ C^* \end{bmatrix}) \in \mathcal{K}_{j+1}(A^*, C^*)$ and $(A^*)^j C^* \in \mathbb{Y}(\mathcal{K}_{j+1}(\mathcal{H}, \begin{bmatrix} F \\ C^* \end{bmatrix}))$, for all $j = 0, \dots, k-1$, which we prove by induction. For $j = 0$, both claims are obvious; suppose they hold true for some $j = \ell$. Then there exist $Z_j, W_j \in \mathbb{C}^{p \times p}$ and $\tilde{F} \in \mathbb{C}^{n \times p}$ such that

$$\begin{aligned} \mathcal{H}^\ell \begin{bmatrix} F \\ C^* \end{bmatrix} &= \begin{bmatrix} \tilde{F} \\ \sum_{j=0}^{\ell} (A^*)^j C^* Z_j \end{bmatrix}, \\ (A^*)^\ell C^* &= \mathbb{Y}(\sum_{j=0}^{\ell} \mathcal{H}^j \begin{bmatrix} F \\ C^* \end{bmatrix} W_j). \end{aligned}$$

Therefore,

$$\mathcal{H}^{\ell+1} \begin{bmatrix} F \\ C^* \end{bmatrix} = \begin{bmatrix} \hat{F} \\ C^* C \tilde{F} - \sum_{j=0}^{\ell} (A^*)^{j+1} C^* Z_j \end{bmatrix} = \begin{bmatrix} \hat{F} \\ \sum_{j=0}^{\ell+1} (A^*)^j C^* \tilde{Z}_j \end{bmatrix},$$

for some $\hat{F} \in \mathbb{C}^{n \times p}$, and with $\tilde{Z}_0 = C \tilde{F}$ and $\tilde{Z}_{j+1} = -Z_j$ for $j = 0, \dots, \ell$. Thus, $\mathbb{Y}(\mathcal{H}^{\ell+1} \begin{bmatrix} F \\ C^* \end{bmatrix}) \in \mathcal{K}_{\ell+1}(A^*, C^*)$.

Similarly,

$$\mathcal{H} \sum_{j=0}^{\ell} \mathcal{H}^j \begin{bmatrix} F \\ C^* \end{bmatrix} W_j = \begin{bmatrix} \acute{F} \\ C^* C \check{F} - (A^*)^{\ell+1} C^* \end{bmatrix},$$

for some $\acute{F}, \check{F} \in \mathbb{C}^{n \times p}$. Therefore, we have that

$$(A^*)^{\ell+1} = \mathbb{Y}(\sum_{j=0}^{\ell} \mathcal{H}^j \begin{bmatrix} F \\ C^* \end{bmatrix} \tilde{W}_j)$$

with $\tilde{W}_0 = C \check{F}$ and $\tilde{W}_{j+1} = -W_j$ for $j = 0, \dots, \ell$, and $(A^*)^{\ell+1} \in \mathbb{Y}(\mathcal{K}_{\ell+1}(\mathcal{H}, \begin{bmatrix} F \\ C^* \end{bmatrix}))$.

(b) The proof uses the same techniques as (a) and (c), so we skip it for the sake of brevity.

(c) Assume that σ_j is not an eigenvalue of \mathcal{H} , and let $(\mathcal{H} - \sigma_j I)^{-1} \begin{bmatrix} F \\ C^* \end{bmatrix} =: \begin{bmatrix} X \\ Y \end{bmatrix}$. Then

$$(\mathcal{H} - \sigma_j I) \begin{bmatrix} X \\ Y \end{bmatrix} = \begin{bmatrix} F \\ C^* \end{bmatrix}, \quad (22)$$

and the \mathbb{Y} component is

$$C^*CX + (-A^* - \sigma_j I)Y = C^*. \quad (23)$$

Since A is stable and $\sigma_j \in \mathbb{C}_-$, the matrix $A^* + \sigma_j I$ is non-singular, and we have

$$Y = (A^* + \sigma_j I)^{-1} C^* \cdot (CX - I_p),$$

where I_p is the identity matrix of order p . Thus $\mathbb{Y}(\mathcal{K}(\mathcal{H}, \begin{bmatrix} F \\ C^* \end{bmatrix}, \vec{\sigma})) \subseteq \mathcal{K}(A^*, C^*, -\vec{\sigma})$.

To show the reverse inclusion, we first prove that $CX - I_p$ is non-singular. Assume the opposite, i.e. that there exists some $z \in \mathbb{C}^p$, $z \neq 0$ such that $CXz = z$. Postmultiplying (23) with z , we have $(A^* + \sigma_j I)Yz = 0$. Again, since $A^* + \sigma_j I$ is non-singular, it follows that $Yz = 0$, and postmultiplying (22) with z yields $(A - \sigma_j I)Xz = 0$. Here we have assumed that $F = 0$. With the additional assumption that σ_j is not in the spectrum of A , it holds that $Xz = 0$ as well². Thus the matrix $\begin{bmatrix} X \\ Y \end{bmatrix}$ is not of full-column rank, which contradicts C^* being of full column rank. Therefore, $CX - I_p$ is non-singular, and

$$(A^* + \sigma_j I)^{-1} C^* = Y(CX - I_p)^{-1} = \mathbb{Y}((\mathcal{H} - \sigma_j I)^{-1} \begin{bmatrix} F \\ C^* \end{bmatrix}) \cdot (CX - I_p)^{-1}.$$

This completes the proof of $\mathbb{Y}(\mathcal{K}(\mathcal{H}, \begin{bmatrix} F \\ C^* \end{bmatrix}, \vec{\sigma})) \supseteq \mathcal{K}(A^*, C^*, -\vec{\sigma})$.

Finally, assume that σ_j is an eigenvalue of \mathcal{H} . Then

$$\begin{bmatrix} A - \sigma_j I & BB^* \\ C^*C & -A^* - \sigma_j I \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} = 0$$

for the eigenvector $\begin{bmatrix} x \\ y \end{bmatrix}$. Rearranging the \mathbb{Y} component yields

$$y = (A^* + \sigma_j I)^{-1} C^* \cdot Cx \in \mathcal{K}(A^*, C^*, -\vec{\sigma}).$$

□

From the theorem, one can conclude the following: consider the process of computing an invariant subspace of the Hamiltonian matrix \mathcal{H} via any of the Krylov subspaces, with the purpose of solving the Riccati equation. As the computation progresses, there are better and better approximations of the Hamiltonian eigenvectors in the subspace.

²Note that when σ_j is an eigenvalue of A , it can happen that Y is not of full column rank. In such a case, $\mathbb{Y}(\mathcal{K}(\mathcal{H}, \begin{bmatrix} F \\ C^* \end{bmatrix}, \vec{\sigma}))$ is a strict subspace of $\mathcal{K}(A^*, C^*, -\vec{\sigma})$.

In particular, using the \mathbb{Y} components of these eigenvectors, one can eventually build a good approximate Riccati solution such as P_k^{proj} or P_k^{dir} .

On the other hand, consider the associated Krylov process generated by the matrix A^* instead of \mathcal{H} . Theorem 8 states that the subspaces obtained with A^* are the same as the \mathbb{Y} components of the Hamiltonian Krylov subspaces. Therefore, they also contain better and better approximations to the \mathbb{Y} components of the Hamiltonian eigenvectors, which approximately span the range of the Riccati solution.

In fact, we have the following connection between the rational Krylov method and the Hamiltonian eigenspace approximations.

Corollary 9. *Suppose that the matrix C^* has one column, and that the shifts used in the RKSM variant of Algorithm 2 are equal to some k stable eigenvalues of the Hamiltonian matrix \mathcal{H} . Then $P_k^{\text{kry}} = P_k^{\text{proj}}$, where P_k^{proj} is computed by Algorithm 1 in which the invariant subspace is associated to the same eigenvalues. Here we assume a slightly modified definition of the rational Krylov subspace for a k -tuple $\vec{\sigma}$:*

$$\tilde{\mathcal{K}}(M, v, \vec{\sigma}) = \text{span}\{(M - \sigma_1 I)^{-1}v, (M - \sigma_2 I)^{-1}v, \dots, (M - \sigma_k I)^{-1}v\}.$$

Moreover, suppose that we extend Algorithm 1 so that $\begin{bmatrix} X_k \\ Y_k \end{bmatrix}$ spans a Hamiltonian Krylov subspace with the initial vector $\begin{bmatrix} 0 \\ C^* \end{bmatrix}$ instead of an exact invariant subspace. Then $P_k^{\text{proj}} = P_k^{\text{kry}}$, where P_k^{kry} is computed by Algorithm 2 with the same type of Krylov subspace.

To summarize, success of the projection methods that use Krylov subspaces generated by A^* may be explained by the fact that they contain good approximations for \mathbb{Y} components of the eigenvectors of \mathcal{H} .

This relation between rational Krylov subspaces generated by \mathcal{H} and those generated by A^* also motivates an alternative shift selection procedure for RKSM, as the following example demonstrates.

Example 10. *Shift selection in the rational Krylov subspace method for the Lyapunov and the Riccati equation, as well as in the ADI methods, is a very important part of the algorithm, and can have a great influence on the convergence rate. Several heuristic procedures exist for choosing the shifts, in particular,*

- *Penzl shifts [20] are chosen among the Ritz values of the matrix A^* , computed from low-dimensional Krylov subspaces generated by A^* and/or A^{-*} . More precisely, if $\theta_1, \dots, \theta_\ell$ are the available Ritz values, then the shifts $\sigma_1, \sigma_2, \dots$ are computed sequentially as the approximate solutions of the ADI-minimax problem:*

$$\sigma_{j+1} = \arg \max_{\sigma \in \{\theta_1, \dots, \theta_\ell\}} \prod_{i=1}^j \frac{|\sigma - \sigma_i|}{|\sigma + \bar{\sigma}_i|}.$$

- *Adaptive shifts [9] $\sigma_1, \sigma_2, \dots$ are defined as*

$$\sigma_{j+1} = \arg \max_{\sigma \in \delta S_j} \prod_{i=1}^j \frac{|\sigma - \sigma_i|}{|\sigma + \theta_i|}.$$

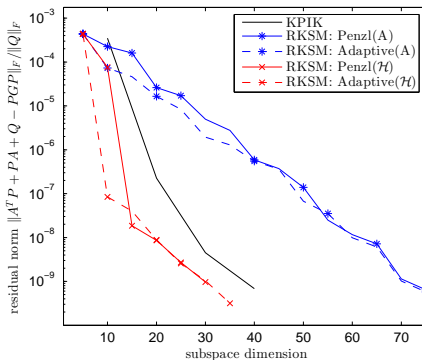


Figure 4: Convergence of the rational Krylov subspace method for the Riccati equation may be faster when shifts are computed using the Hamiltonian matrix \mathcal{H} instead of A^* .

Here, the θ_i approximate the eigenvalues of A^* , and are computed as the Ritz values from the current rational Krylov subspace, and \mathcal{S}_j is the convex hull of the set $\{\theta_1, \dots, \theta_j, \sigma_1, \sigma_2, \sigma_3, \sigma_4\}$. Estimates for the smallest and the largest eigenvalue of A^* in the absolute value, together with their complex conjugates, are taken as the first four shifts.

These procedures are applied for both the Lyapunov and the Riccati equation. However, as we have seen, choosing the shifts so that they approximate the stable eigenvalues of \mathcal{H} instead of the eigenvalues of A^* should be better suited for the Riccati equation. We therefore also consider two alternatives: Penzl and adaptive shifts that use Ritz values and Krylov subspaces generated by the matrix \mathcal{H} .

This is indeed a better choice in certain cases: suppose that the matrix $A \in \mathbb{R}^{3375 \times 3375}$ is obtained by the semi-discretization of

$$\frac{\partial u}{\partial t} = \Delta u$$

on the unit cube, and that the entries of the matrices $B = C^* \in \mathbb{R}^{3375 \times 5}$ are chosen at random from $[-1, 1]$, with normal distribution. Figure 4 shows the convergence of RKSM with all four shift selection procedures. For this Riccati equation, the five leftmost eigenvalues of the Hamiltonian matrix are well separated from all the others, and are also far away from the eigenvalues of A . None of the shifts generated by using the matrix A^* comes anywhere near these eigenvalues. Thus the associated eigenvectors, which contribute significantly to the subspace spanned by the Riccati solution, do not appear in the rational Krylov subspace. On the other hand, these eigenvectors are captured already in the early iterations when \mathcal{H} is used for shift generation, and that yields a much faster convergence. The KPIK method can be understood as a rational Krylov method with shifts 0 and $-\infty$, and thus it is also able to retrieve the leftmost eigenvectors of \mathcal{H} in the subspace relatively fast.

4.2 Hamiltonian eigenspaces and the qADI algorithm

Another method for solving the Riccati equation was introduced by Wong and Balakrishnan [29]. This method, named qADI, is an adaptation of the Lyapunov-ADI method and, unlike the standard Newton-ADI algorithms, it avoids the need to perform the outer Newton steps.

The qADI iterates are defined by the initial approximation $P_0^{\text{adi}} = 0$ and the following recurrence:

$$\begin{aligned} P_{k+1/2}^{\text{adi}}(A + \overline{\sigma_{k+1}}I - GP_k^{\text{adi}}) &= -Q - (A^* - \overline{\sigma_{k+1}}I)P_k^{\text{adi}}, \\ (A^* + \sigma_{k+1}I - P_{k+1/2}^{\text{adi}}G)P_{k+1}^{\text{adi}} &= -Q - P_{k+1/2}^{\text{adi}}(A - \sigma_{k+1}I). \end{aligned}$$

Here the shifts $(\sigma_k)_k \subseteq \mathbb{C}$ are free parameters which one chooses with the goal of accelerating the convergence of the method, similarly as in the Lyapunov-ADI. Our final result shows that the qADI produces exactly the same sequence of approximations as the Hamiltonian subspace approach, if the matrix C^* has one column and the shifts are equal to the eigenvalues of \mathcal{H} .

Theorem 11. *Suppose that $Q = C^*C$, where $C \in \mathbb{C}^{1 \times n}$. Consider the sequence of approximations $(P_k^{\text{dir}})_k$ such that $P_0^{\text{dir}} = 0$ and P_k^{dir} is obtained by computing the Hamiltonian invariant subspaces associated with eigenvalues $\lambda_1, \lambda_2, \dots, \lambda_k$, and the qADI iteration with the shifts equal to these eigenvalues:*

$$P_{k+1/2}^{\text{adi}}(A + \overline{\lambda_{k+1}}I - GP_k^{\text{adi}}) = -Q - (A^* - \overline{\lambda_{k+1}}I)P_k^{\text{adi}}, \quad (24)$$

$$(A^* + \lambda_{k+1}I - P_{k+1/2}^{\text{adi}}G)P_{k+1}^{\text{adi}} = -Q - P_{k+1/2}^{\text{adi}}(A - \lambda_{k+1}I). \quad (25)$$

If the matrices $A + \overline{\lambda_{k+1}}I - GP_k^{\text{adi}}$ and $A^* + \lambda_{k+1}I - P_{k+1/2}^{\text{adi}}G$ are invertible for all k , then $P_k^{\text{adi}} = P_k^{\text{dir}}$.

Proof. We show the claim using induction. When $k = 0$, the claim is trivial. Assume that $P_k^{\text{dir}} = P_k^{\text{adi}}$ for some $k \geq 0$. It suffices to prove that P_{k+1}^{dir} satisfies equation (25). Using notation of Theorem 5 and (20), we have that $P_{k+1}^{\text{dir}} = P_k^{\text{dir}} + \Delta$, with $\Delta = \tilde{P}_1^{\text{dir}}$. Thus (25) is equivalent to:

$$(A^* + \lambda_{k+1}I - P_{k+1/2}^{\text{adi}}G)(P_k^{\text{dir}} + \Delta) = -Q - P_{k+1/2}^{\text{adi}}(A - \lambda_{k+1}I),$$

which we rewrite as

$$\begin{aligned} (A^* + \lambda_{k+1}I)P_k^{\text{dir}} + (\tilde{A}^* + \lambda_{k+1}I)\Delta + (P_k^{\text{adi}} - P_{k+1/2}^{\text{adi}})G\Delta \\ = -Q - P_{k+1/2}^{\text{adi}}(\tilde{A} - \lambda_{k+1}I). \end{aligned} \quad (26)$$

Recall that $\tilde{A} = A - GP_k^{\text{dir}} = A - GP_k^{\text{adi}}$. We replace each of the three terms on the left hand side after making the following observations:

1. By (17), and setting $\tilde{Q} = \tilde{C}^*\tilde{C}$, we have

$$(A^* + \lambda_{k+1}I)P_k^{\text{dir}} = -P_k^{\text{dir}}(\tilde{A} - \lambda_{k+1}I) - Q + \tilde{Q}. \quad (27)$$

2. Note that $\Delta = -\tilde{y}(\tilde{y}^* \tilde{x})^{-1} \tilde{y}^*$ for the eigenvector $\begin{bmatrix} \tilde{x} \\ \tilde{y} \end{bmatrix}$ of $\tilde{\mathcal{H}}$ associated with its eigenvalue λ_{k+1} : $\begin{bmatrix} \tilde{A} - \lambda_{k+1} I & G \\ \tilde{Q} & -\tilde{A}^* - \lambda_{k+1} I \end{bmatrix} \begin{bmatrix} \tilde{x} \\ \tilde{y} \end{bmatrix} = 0$. Thus

$$\begin{aligned} (\tilde{A}^* + \lambda_{k+1} I) \Delta &= -(\tilde{A}^* + \lambda_{k+1} I) \tilde{y}(\tilde{y}^* \tilde{x})^{-1} \tilde{y}^* = -\tilde{Q} \tilde{x}(\tilde{y}^* \tilde{x})^{-1} \tilde{y}^* \\ &= -\tilde{Q} \mathcal{P}, \end{aligned} \quad (28)$$

and

$$G \Delta = -G \tilde{y}(\tilde{y}^* \tilde{x})^{-1} \tilde{y}^* = (\tilde{A} - \lambda_{k+1} I) \tilde{x}(\tilde{y}^* \tilde{x})^{-1} \tilde{y}^* = (\tilde{A} - \lambda_{k+1} I) \mathcal{P}. \quad (29)$$

Here $\mathcal{P} = \tilde{x}(\tilde{y}^* \tilde{x})^{-1} \tilde{y}^*$ denotes the skew projector.

After inserting (27), (28), (29) into (26), we need to show that

$$\begin{aligned} -P_k^{\text{dir}}(\tilde{A} - \lambda_{k+1} I) - Q + \tilde{Q} - \tilde{Q} \mathcal{P} + (P_k^{\text{adi}} - P_{k+1/2}^{\text{adi}})(\tilde{A} - \lambda_{k+1} I) \mathcal{P} \\ = -Q - P_{k+1/2}^{\text{adi}}(\tilde{A} - \lambda_{k+1} I), \end{aligned}$$

which, by using $P_k^{\text{adi}} = P_k^{\text{dir}}$, and assuming that $\tilde{A} + \overline{\lambda_{k+1}} I$ is nonsingular, can be factored as

$$\begin{aligned} \left((P_{k+1/2}^{\text{adi}} - P_k^{\text{adi}})(\tilde{A} - \lambda_{k+1} I) (\tilde{A} + \overline{\lambda_{k+1}} I) + \tilde{Q}(\tilde{A} + \overline{\lambda_{k+1}} I) \right) \\ \cdot (\tilde{A} + \overline{\lambda_{k+1}} I)^{-1} (I - \mathcal{P}) = 0. \end{aligned}$$

The expression in the first parentheses can be simplified:

$$\begin{aligned} &(P_{k+1/2}^{\text{adi}} - P_k^{\text{adi}})(\tilde{A} - \lambda_{k+1} I)(\tilde{A} + \overline{\lambda_{k+1}} I) + \tilde{Q}(\tilde{A} + \overline{\lambda_{k+1}} I) \\ &= (P_{k+1/2}^{\text{adi}}(\tilde{A} + \overline{\lambda_{k+1}} I) - P_k^{\text{adi}}(\tilde{A} + \overline{\lambda_{k+1}} I))(\tilde{A} - \lambda_{k+1} I) + \tilde{Q}(\tilde{A} + \overline{\lambda_{k+1}} I) \\ &= (-Q - (A^* - \overline{\lambda_{k+1}} I) P_k^{\text{adi}} - P_k^{\text{adi}}(\tilde{A} + \overline{\lambda_{k+1}} I))(\tilde{A} - \lambda_{k+1} I) + \tilde{Q}(\tilde{A} + \overline{\lambda_{k+1}} I) \\ &= -\tilde{Q}(\tilde{A} - \lambda_{k+1} I) + \tilde{Q}(\tilde{A} + \overline{\lambda_{k+1}} I) \\ &= 2 \operatorname{Re}(\lambda_{k+1}) \tilde{Q}. \end{aligned}$$

For the second equality we have used the defining relation (24) for $P_{k+1/2}^{\text{adi}}$. From the discussion so far, we conclude that the initial claim $P_{k+1}^{\text{dir}} = P_{k+1}^{\text{adi}}$ is equivalent to

$$2 \operatorname{Re}(\lambda_{k+1}) \cdot \tilde{Q} \cdot (\tilde{A} + \overline{\lambda_{k+1}} I)^{-1} (I - \mathcal{P}) = 0. \quad (30)$$

The matrix $\tilde{Q} = \tilde{C}^* \tilde{C}$ is of rank 1, and thus the $n \times n$ matrix on the left-hand side of (30) is also of rank 1. To show that some matrix Z is equal to zero, it suffices to show $z^* Z = 0$ for all vectors z . Note that $z^* \tilde{Q} = 0$ trivially holds for all vectors $z \perp \tilde{C}^*$, and the set of these vectors spans an $(n-1)$ -dimensional subspace in \mathbb{C}^n . Thus it suffices to show that $z^* \cdot 2 \operatorname{Re}(\lambda_{k+1}) \cdot \tilde{Q} \cdot (\tilde{A} + \overline{\lambda_{k+1}} I)^{-1} (I - \mathcal{P}) = 0$ for a single

vector z for which $z^* \tilde{Q} \neq 0$. We can take $z = \tilde{x}$: since $(\tilde{A} + \overline{\lambda_{k+1}}I)^* \tilde{y} = \tilde{Q} \tilde{x}$, we have $\tilde{y}^* = \tilde{x}^* \tilde{Q} (\tilde{A} + \overline{\lambda_{k+1}}I)^{-1}$, and thus

$$\begin{aligned} \tilde{x}^* \cdot 2 \operatorname{Re}(\lambda_{k+1}) \cdot \tilde{Q} \cdot (\tilde{A} + \overline{\lambda_{k+1}}I)^{-1} (I - \mathcal{P}) &= \tilde{y}^* (I - \mathcal{P}) \\ &= \tilde{y}^* (I - \tilde{x}(\tilde{y}^* \tilde{x})^{-1} \tilde{y}^*) = 0. \end{aligned}$$

Note that $\tilde{x}^* \tilde{Q} = 0$ would imply that either $\tilde{A} + \overline{\lambda_{k+1}}I$ is singular (which contradicts the assumption of the theorem), or that $\tilde{y} = 0$ (for which the sequence $(P_k^{\text{dir}})_k$ is not well defined). \square

5 Conclusion

In this paper, we discussed a method for solving the large-scale Riccati equation that is based on computing a low-dimensional invariant subspace of a Hamiltonian matrix. The theoretical potential of this method is promising and backed up by the results on the singular value decay of the Riccati solution, but there is a couple of issues that can render it impractical.

On one hand, there is the problematics of selecting which eigenpairs to compute, which is very similar to shift selection in the ADI and the rational Krylov methods, and which may be resolved in a similar way by various heuristics (such as targeting the eigenvalues of \mathcal{H} that are closest to Penzl or adaptive shifts). On the other hand, the requirement for computing an invariant subspace is more rigid compared to the other methods which can approximate the Riccati solution using arbitrary subspaces. We can relax this requirement, and consider the subspaces obtained by Krylov processes generated with \mathcal{H} . As we have seen, the \mathbb{Y} components of these coincide with Krylov subspaces generated with A^* , and the latter are computed more efficiently since the matrix involved is half the size of \mathcal{H} .

The interconnection between the two Krylov processes also offers an explanation for the success of methods such as RKSM and KPIK, since they can be interpreted as methods that tend to incorporate better and better approximations to \mathbb{Y} components of the eigenvectors of \mathcal{H} . This fact has motivated us to use Ritz values of \mathcal{H} instead of A^* in order to generate shifts for the RKSM. As shift generation with \mathcal{H} using the symplectic Lanczos process is relatively cheap, this combination with the rational Krylov subspace method for (A^*, C^*) seems to be preferable over Krylov subspace methods using \mathcal{H} for solving large-scale AREs.

Moreover, we could also show that the qADI iteration for AREs can be interpreted as a rational Krylov subspace method. This observation may lead to improvements in the qADI method which are currently being investigated.

References

- [1] L. AMODEI AND J.-M. BUCHOT, *An invariant subspace method for large-scale algebraic Riccati equation*, Appl. Numer. Math., 60 (2010), pp. 1067–1082.

- [2] A. ANTOULAS, D. SORENSSEN, AND Y. ZHOU, *On the decay rate of Hankel singular values and related issues*, Sys. Control Lett., 46 (2002), pp. 323–342.
- [3] P. BENNER AND H. FASSBENDER, *An implicitly restarted symplectic Lanczos method for the Hamiltonian eigenvalue problem*, Linear Algebra Appl., 263 (1997), pp. 75–111.
- [4] P. BENNER, H. FASSBENDER, AND M. STOLL, *A Hamiltonian Krylov–Schur-type method based on the symplectic Lanczos process*, Linear Algebra Appl., 435 (2011), pp. 578–600.
- [5] P. BENNER, J.-R. LI, AND T. PENZL, *Numerical solution of large Lyapunov equations, Riccati equations, and linear-quadratic control problems*, Numer. Lin. Alg. Appl., 15 (2008), pp. 755–777.
- [6] R. BHATIA, *Matrix Analysis*, Graduate Texts in Mathematics, Springer New York, 1997.
- [7] S. CHANDRASEKARAN AND I. IPSEN, *On rank-revealing factorisations*, SIAM J. Matrix Anal. Appl., 15 (1994), pp. 592–622.
- [8] V. DRUSKIN, L. KNIZHNERMAN, AND V. SIMONCINI, *Analysis of the rational Krylov subspace and ADI methods for solving the Lyapunov equation*, SIAM J. Numer. Anal., 49 (2011), pp. 1875–1898.
- [9] V. DRUSKIN AND V. SIMONCINI, *Adaptive rational Krylov subspaces for large-scale dynamical systems*, Sys. Control Lett., 60 (2011), pp. 546–560.
- [10] W. FERNG, W.-W. LIN, AND C.-S. WANG, *The shift-inverted J-Lanczos algorithm for the numerical solutions of large sparse algebraic Riccati equations*, Comput. Math. Appl., 33 (1997), pp. 23–40.
- [11] M. HEYOUNI AND K. JBILOU, *An extended block Arnoldi algorithm for large-scale solutions of the continuous-time algebraic Riccati equation*, Electr. Trans. Num. Anal., 33 (2009), pp. 53–62.
- [12] Y. P. HONG AND C.-T. PAN, *Rank-revealing QR factorizations and the singular value decomposition*, Math. Comp., 58 (1992), pp. 213–232.
- [13] K. JBILOU, *Block Krylov subspace methods for large algebraic Riccati equations*, Numer. Algorithms, 34 (2003), pp. 339–353.
- [14] D. KLEINMAN, *On an iterative technique for Riccati equation computations*, IEEE Trans. Automat. Control, AC-13 (1968), pp. 114–115.
- [15] P. LANCASTER AND L. RODMAN, *The Algebraic Riccati Equation*, Oxford University Press, Oxford, 1995.
- [16] J.-R. LI AND J. WHITE, *Low rank solution of Lyapunov equations*, SIAM J. Matrix Anal. Appl., 24 (2002), pp. 260–280.

- [17] V. MEHRMANN AND D. WATKINS, *Structure-preserving methods for computing eigenpairs of large sparse skew-Hamiltonian/Hamiltonian pencils*, SIAM J. Matrix Anal. Appl., 22 (2000), pp. 1905–1925.
- [18] V. MEHRMANN AND H. XU, *Choosing poles so that the single-input pole placement problem is well conditioned*, SIAM J. Matrix Anal. Appl., 19 (1998), pp. 664–681.
- [19] T. PENZL, *A cyclic low rank Smith method for large sparse Lyapunov equations*, SIAM J. Sci. Comput., 21 (2000), pp. 1401–1418.
- [20] ———, *LYAPACK Users Guide*, Tech. Rep. SFB393/00-33, Sonderforschungsbereich 393 *Numerische Simulation auf massiv parallelen Rechnern*, TU Chemnitz, 09107 Chemnitz, Germany, 2000. Available from <http://www.tu-chemnitz.de/sfb393/sfb00pr.html>.
- [21] J. SAAK, *Efficient numerical solution of large scale algebraic matrix equations in PDE control and model order reduction*, PhD thesis, TU Chemnitz, July 2009. Available online from <http://nbn-resolving.de/urn:nbn:de:bsz:chi-200901642>.
- [22] V. SIMONCINI, *A new iterative method for solving large-scale Lyapunov matrix equations*, SIAM J. Sci. Comput., 29 (2007), pp. 1268–1288.
- [23] V. SIMONCINI, D. B. SZYLD, AND M. MONSALVE, *On two numerical methods for the solution of large-scale algebraic Riccati equations*, IMA J. Numer. Anal., 34 (2014), pp. 904–920.
- [24] A. SLUIS, *Condition numbers and equilibration of matrices*, Numer. Math., 14 (1969), pp. 14–23.
- [25] D. SORENSEN AND Y. ZHOU, *Bounds on eigenvalue decay rates and sensitivity of solutions to Lyapunov equations*, Tech. Rep. TR02-07, Dept. of Comp. Appl. Math., Rice University, Houston, TX, June 2002. Available online from <http://www.caam.rice.edu/caam/trs/tr02.html#TR02-07>.
- [26] R. THOMPSON, *The behavior of eigenvalues and singular values under perturbations of restricted rank*, Linear Algebra Appl., 13 (1976), pp. 69–78.
- [27] N. TRUHAR AND R.-C. LI, *On ADI method for Sylvester equations*, Tech. Rep. 2008-02, University of Texas at Arlington, Department of Mathematics, 2008. Available online from http://www.uta.edu/math/preprint/rep2008_02.pdf.
- [28] E. WACHSPRESS, *Iterative solution of the Lyapunov matrix equation*, Appl. Math. Letters, 107 (1988), pp. 87–90.
- [29] N. WONG AND V. BALAKRISHNAN, *Quadratic alternating direction implicit iteration for the fast solution of algebraic Riccati equations*, in Proc. Int. Symposium on Intelligent Signal Processing and Communication Systems, 2005, pp. 373–376.

