

See discussions, stats, and author profiles for this publication at: <https://www.researchgate.net/publication/317062090>

A Structure Preserving Krylov Subspace Method for Large Scale Differential Riccati Equations

Article · May 2017

CITATION

1

READS

58

2 authors:



[Antti Koskela](#)

University of Helsinki

17 PUBLICATIONS 42 CITATIONS

[SEE PROFILE](#)



[Hermann Mena](#)

Yachay Tech

47 PUBLICATIONS 195 CITATIONS

[SEE PROFILE](#)

Some of the authors of this publication are also working on these related projects:



Feedback Control of time-dependent partial differential equations [View project](#)

A STRUCTURE PRESERVING KRYLOV SUBSPACE METHOD FOR LARGE SCALE DIFFERENTIAL RICCATI EQUATIONS

ANTTI KOSKELA* AND HERMANN MENA†

Abstract. We propose a Krylov subspace approximation method for the symmetric differential Riccati equation $\dot{X} = AX + XA^T + Q + XSX$, $X(0) = X_0$. The method is based on projecting the large scale equation onto a Krylov subspace spanned by the matrix A and the low rank factors of X_0 and Q . We prove that the method is structure preserving in a sense that it preserves two important properties of the exact flow, namely the positivity of the exact flow, and also the property of monotonicity under certain practically relevant conditions. We also provide theoretical a priori error analysis which shows a superlinear convergence of the method. This behavior is illustrated in the numerical experiments. Moreover, we carry out a derivation of an efficient a posteriori error estimate as well as discuss multiple time stepping combined with a cut of the rank of the numerical solution.

Key words. Differential Riccati equations, LQR optimal control problems, large scale ordinary differential equations, Krylov subspace methods, matrix exponential, exponential integrators, model order reduction, low rank approximation.

AMS subject classifications. 65F10, 65F60, 65L20, 65M22, 93A15, 93C05

1. Introduction. Large scale differential Riccati equations (DREs) arise in the numerical treatment of optimal control problems governed by partial differential equations. This is the case in particular when solving a linear quadratic regulator problem (LQR), a widely studied problem in control theory. We shortly describe the finite dimensional LQR problem. For more details, we refer to [1, 8]. The differential Riccati equation arises in the finite horizon case, i.e., when a finite time integral cost functional is considered. The functional has then the quadratic form

$$J(x, u) = \int_0^T (x(t)^T C^T C x(t) + u(t)^T u(t)) dt + x(T)^T G x(T), \quad (1.1)$$

where $x \in \mathbb{R}^n$, $C \in \mathbb{R}^{q \times n}$ ($q \ll n$) and $u \in \mathbb{R}^r$ ($r \ll n$). The coefficient matrix G of the penalizing term $x(T)^T G x(T)$ is symmetric, nonnegative and has a low rank. The functional (1.1) is constrained by the system of differential equations

$$\dot{x}(t) = Ax(t) + Bu(t), \quad x(0) = x_0, \quad t \in [0, T] \quad (1.2)$$

where the matrix $A \in \mathbb{R}^{n \times n}$ is sparse and $B \in \mathbb{R}^{n \times r}$. The number of columns of B correspond to the number of controls and matrix C represent an observation matrix. Under suitable conditions (see [1, 8]), the optimal control \tilde{u} minimizing the functional (1.1) is given by

$$\tilde{u}(t) = K(t)\tilde{x}(t), \quad \text{where} \quad K(t) = -B^T X(t), \quad (1.3)$$

$X(t)$ is the unique solution of

$$\dot{X} + A^T X + XA - XBB^T X + C^T C = 0, \quad X(T) = G, \quad (1.4)$$

*Department of Mathematics, Royal Institute of Technology (KTH), Stockholm, SeRC Swedish e-Science Research Center, akoskela@kth.se. The first author was supported by the Dahlquist research fellowship.

†Department of Mathematics, Yachay Tech, Urcuquí, Ecuador, mena@yachaytech.edu.ec.

and $\tilde{x}(t)$ satisfies

$$\dot{\tilde{x}}(t) = (A - BB^T X(T - t))\tilde{x}(t), \quad \tilde{x}(0) = x_0.$$

As a result, the central computational problem becomes that of solving the final value problem (1.4) which, with a change of variables, can be written as a initial value problem.

In most control problems fast and slow modes are present which implies that the associated DRE becomes stiff. This is obviously the case when the problem arises from a spatial discretization of a parabolic partial differential equation. This implies that the use of explicit time integration methods is out of question. Furthermore, if the matrix G is symmetric and nonnegative, then so is the solution $X(t)$ of (1.4), for all $0 \leq t \leq T$ (see e.g. [1]).

We propose a Krylov subspace approximation method for large scale differential Riccati equations of the form (1.4). This method is closely related to projection techniques considered for large scale algebraic Riccati equations [30, 36]. Moreover, a similar projection method for DREs has been recently proposed in [17]. **Our approach differs from that of [17] in the fact that the initial value G of (1.4) is contained in the Krylov subspace allowing nonzero G and multiple time stepping.** Moreover, we carry out a full a priori error analysis of our method and derive a posteriori estimates which in numerical experiments are shown to be efficient.

Essentially, the method is based on projecting the matrices A, Q, S and X_0 on an appropriate Krylov subspace, namely on the *block Krylov subspace* spanned by A and the low rank factors of X_0 and Q . The projected small dimensional system is solved using existing linearization techniques. When using a Padé approximant to solve the the linearization of the small dimensional system, the total approximation can be shown to be structure preserving, meaning that the property of the positivity is preserved, and also under certain conditions the property of monotonicity.

The linearization approach for DREs is a standard method. This allows an efficient integration for dense problems, see e.g. [29]. Another approach, the so called Davison–Maki method [9], uses the fundamental solution of the linearized system. A modified variant, avoiding some numerical instabilities, is proposed in [25]. However, the application of these methods for large scale problems is impossible due to the high dimensionality of the initial value in the linearized differential equation.

The problem of solving large scale DREs has received recently considerable attention. In [4, 5] the authors proposed efficient numerical methods for solving DREs capable of exploiting several of the above described properties: sparsity of A , low rank structure of B , C and G , and the symmetry of the solution. These methods are based on a matrix valued implementation of the BDF and Rosenbrock methods. However, several difficulties arise when approximating the optimal control (1.3) in the large scale setting. One difficulty is to evaluate the state equation $x(t)$ and Riccati equation $X(t)$ in the same mesh. In [39] a refined ADI integration method is proposed which addresses the high storage requirements of large scale DRE integrators. In a recent studies an efficient splitting method [41] and adaptive high-order splitting schemes [42] for large scale DREs have been proposed.

Our Krylov subspace approach is also strongly related to techniques used for approximation of the product of an matrix function and a vector, $f(A)b$. There, instead of computing the matrix $f(A)$ explicitly, an efficient alternative is to project the problem onto the Krylov subspace spanned by A and b . The effectiveness of this approach comes from the fact that generating Krylov subspaces is for the most part

based on operations of the form $b \rightarrow Ab$, which is a cheap operation for sparse A . The approximation of products of the form $f(A)b$ by using Krylov subspace methods has recently been an active topic of research: we mention the work on classical Krylov subspaces [13], [15], [27], [32], extended Krylov subspaces [27], and rational Krylov subspaces [26], [3].

The paper is organized as follows. In Section 2 we describe some preliminaries. Then, in Section 3, the structure preserving method is proposed. In Section 4, the error analysis first for the differential Lyapunov equation (a simplified version of the DRE), and then for the DRE is presented. In Section 5 a posteriori error estimation is described. In Section 6 the rank cut and multiple time stepping are discussed. Numerical examples and conclusions of Sections 7 and 8 conclude the article.

Notation and definitions. Throughout the article $\|\cdot\|$ will denote the Euclidean norm, or its induced matrix norm, i.e., the spectral norm. We say that a matrix A is nonnegative if it is symmetric positive semidefinite, and write $A \geq 0$. For symmetric matrices A and B we write $B \geq A$ if $B - A \geq 0$.

We will repeatedly use the notion of the *logarithmic norm* of a matrix $A \in \mathbb{C}^{n \times n}$. It can be defined via the *field of values* $\mathcal{F}(A)$, which is defined as

$$\mathcal{F}(A) = \{x^*Ax : x \in \mathbb{C}^n, \|x\| = 1\}.$$

Then, the logarithmic norm $\mu(A)$ of A is defined by

$$\mu(A) := \{\max \operatorname{Re} z : z \in \mathcal{F}(A)\}.$$

We will also repeatedly use the exponential like function φ_1 defined by

$$\varphi_1(z) = \frac{e^z - 1}{z} = \sum_{\ell=0}^{\infty} \frac{z^\ell}{(\ell+1)!}.$$

2. Preliminaries. From now on we consider the time invariant symmetric differential Riccati equation (DRE) written in the form

$$\begin{aligned} \dot{X}(t) &= AX(t) + X(t)A^T + Q - X(t)SX(t), \\ X(0) &= X_0, \end{aligned} \tag{2.1}$$

where $t \in [0, T]$ and $A, Q, S, X_0 \in \mathbb{R}^{n \times n}$, $Q^T = Q$, $S^T = S$. Specifically, we focus on the low rank positive semidefinite case, where

$$X_0 = ZZ^T, \quad Q = CC^T, \tag{2.2}$$

for some $Z \in \mathbb{R}^{n \times p}$ and $C \in \mathbb{R}^{n \times q}$, $p, q \ll n$.

2.1. Linearization. We recall a fact that will be needed later on (see e.g. [1, Thm. 3.1.1.]).

LEMMA 1 (Associated linear system). *The DRE (2.1) is equivalent to solving the linear system of differential equations*

$$\frac{d}{dt} \begin{bmatrix} U(t) \\ V(t) \end{bmatrix} = \begin{bmatrix} A & Q \\ S & -A^T \end{bmatrix} \begin{bmatrix} U(t) \\ V(t) \end{bmatrix}, \quad \begin{bmatrix} U(0) \\ V(0) \end{bmatrix} = \begin{bmatrix} I \\ X_0 \end{bmatrix} \tag{2.3}$$

where $U(t), V(t) \in \mathbb{R}^{n \times n}$. If the solution $X(t)$ exists on the interval $[0, T]$, then the solution of (2.3) exists, $U(t)$ is invertible on $[0, T]$, and

$$X(t) = V(t)U^{-1}(t).$$

Note also that the matrix $H = \begin{bmatrix} A & Q \\ S & -A^T \end{bmatrix}$ is Hamiltonian, i.e., it holds that

$$(JH)^T = JH, \quad \text{where} \quad J = \begin{bmatrix} 0 & I \\ -I & 0 \end{bmatrix}. \quad (2.4)$$

This linearization approach for the differential Riccati equation is a standard method for solving finite dimensional DREs, and leads to efficient integration methods for dense problems, see e.g. [29].

2.2. Integral representation of the exact solution. For the exact solution of (2.1) we have the following integral representation (see also [28, Thm. 8]).

THEOREM 2 (Exact solution of the DRE). *The exact solution of the DRE (2.1) is given by*

$$\begin{aligned} X(t) = e^{tA} X_0 e^{tA^T} &+ \int_0^t e^{(t-s)A} Q e^{(t-s)A^T} ds \\ &- \int_0^t e^{(t-s)A} X(s) S X(s) e^{(t-s)A^T} ds. \end{aligned} \quad (2.5)$$

Proof. Can be verified by elementary differentiation. \square

2.3. Positivity and monotonicity of the exact flow. We recall two important properties of the symmetric DRE, namely the positivity of the exact solution (see e.g. [11, Prop. 1.1]) and the monotonicity of the solution with relative to the initial data (see e.g. [12, Thm. 2]). By these we mean the following.

THEOREM 3 (Positivity and monotonicity of the solution). *For the solution $X(t)$ of the symmetric DRE (2.1) it holds:*

1. (Positivity) $X(t)$ is symmetric positive semidefinite and it exists for all $t > 0$.
2. (Monotonicity) Consider two symmetric DREs of the (2.1) corresponding to the linearized systems of the form (2.3) with the coefficient matrices

$$H = \begin{bmatrix} A & Q \\ S & -A^T \end{bmatrix} \quad \text{and} \quad \tilde{H} = \begin{bmatrix} \tilde{A} & \tilde{Q} \\ \tilde{S} & -\tilde{A}^T \end{bmatrix}$$

and let J be the skew-symmetric matrix (2.4). Then, if $\tilde{H}J \leq HJ$, and if $0 \leq X_0 \leq \tilde{X}_0$, then for every $t \geq 0$: $X(t) \leq \tilde{X}(t)$.

Proof. See the proofs of [11, Prop. 1.1] and [12, Thm. 2]. \square

We will later show that our proposed numerical method preserves the properties of Theorem 3.

2.4. Bound for the solution. Using the positivity property of $X(t)$ (Thm. 3) we obtain the following bound for the norm of the solution. This will be repeatedly needed in the analysis of the proposed method.

LEMMA 4 (Bound for the exact solution). *For the solution $X(t)$ of (2.1) it holds*

$$\|X(t)\| \leq e^{2t\mu(A)} \|X_0\| + t\varphi_1(2t\mu(A)) \|Q\|. \quad (2.6)$$

Proof. Since X_0 , Q and $X(t)$ are all symmetric positive definite, we see that the first two terms on the right hand side of (2.5) are symmetric positive semidefinite and

the third term is symmetric negative semidefinite. Moreover, since $X(t)$ is symmetric positive definite by Theorem 3, and for a symmetric positive definite matrix A it holds that $\|A\| = \max_{\|x\|=1} x^*Ax$, we see that

$$\|X(t)\| \leq \|e^{tA}X_0e^{tA^T} + \int_0^t e^{(t-s)A}Qe^{(t-s)A^T}ds\|.$$

Using the well-known bound $\|e^{tA}\| \leq e^{t\mu(A)}$ (see e.g. [14]), the fact that $\mu(A^T) = \mu(A)$ and that $\varphi_1(z) = \int_0^t e^{(t-s)z}ds$, the claim follows. \square

From Lemma 4 we immediately get the following corollary.

COROLLARY 5. *The solution $X(t)$ satisfies*

$$\max_{s \in [0, t]} \|X(s)\| \leq \max\{1, e^{2t\mu(A)}\} \|X_0\| + t \max\{1, \varphi_1(2t\mu(A))\} \|Q\|.$$

3. A structure preserving projection method using block Krylov subspaces. In this Section we derive the projection method and show its structure preserving properties. Given an orthogonal basis matrix V , we first show how to carry out the projection of the system (2.1) using the Galerkin condition. The choice of the matrix V is then motivated by a moment-matching approach. This means that we Taylor expand the exact solution $X(t)$ with respect to time variable t and inspect which subspaces are sufficient to approximate the coefficient matrices of the expansion.

Alternatively, the fact that V needs to contain certain Krylov subspaces in order to obtain a polynomial approximation of $X(t)$ can be seen from the point of view of Krylov subspace approximation of the matrix exponential. This is strongly related to the approach taken by Saad already in [34] for the algebraic Lyapunov equation. This approach will give some auxiliary results that are needed in the convergence analysis of the method.

3.1. Block Krylov subspace approximation of the matrix exponential.

Block Krylov subspace methods are based on the idea of projecting a high dimensional problem involving a matrix $A \in \mathbb{R}^{n \times n}$ and a block matrix $B \in \mathbb{R}^{n \times \ell}$ onto a lower dimensional subspace, a block Krylov subspace $\mathcal{K}_k(A, B)$, which is defined by

$$\mathcal{K}_k(A, B) = \text{span}\{B, AB, A^2B, \dots, A^{k-1}B\}.$$

Usually, an orthogonal basis matrix V_k for $\mathcal{K}_k(A, B)$ is generated using an Arnoldi type iteration, and this matrix is then used for the projections. There exist several Arnoldi type methods to produce an orthogonal basis matrix for $\mathcal{K}_k(A, B)$, and we choose here the method given in [33] which is listed algorithmically as follows.

1. **Input:** $A \in \mathbb{R}^{n \times n}$, $B \in \mathbb{R}^{n \times \ell}$ and number of iterations k .
2. **Start.** Compute QR decomposition: $B = U_1R_1$.
3. **Iterate.** for $j = 1, \dots, k$ compute:

$$H_{ij} = U_i^*AU_j, \quad i = 1, \dots, j,$$

$$W_j = AU_j - \sum_{i=1}^j U_i H_{ij},$$

$$W_j = U_{j+1}H_{j+1,j} \quad (QR \text{ decomposition of } W_j).$$

As usual, the orthogonalization can be carried out at step 3 in a modified Gram–Schmidt manner and reorthogonalization can be performed if needed.

This algorithm gives an orthogonal basis matrix $V_k = [U_1 \ \dots \ U_k] \in \mathbb{C}^{n \times k\ell}$ for the block Krylov subspace $\mathcal{K}_k(A, B)$ and the projected block Hessenberg matrix

$$H_k = V_k^* A V_k. \quad (3.1)$$

This means that the $\ell \times \ell$ (i, j) -block of H_k is given by H_{ij} in the above algorithm. Moreover, the following Arnoldi relation holds:

$$A V_k = V_k H_k + H_{k+1, k} U_{k+1} E_k^T, \quad (3.2)$$

where $E_k = [0 \ I_\ell]^T \in \mathbb{R}^{k\ell \times \ell}$.

If A has its field of values on a line, e.g., is Hermitian or skew-Hermitian, then there exists $\theta \in \mathbb{R}$ such that $e^{i\theta} A$ is Hermitian. By (3.1) this implies that H_k is block tridiagonal. Then, the orthogonalization in the above algorithm has to be done only against two previous blocks, and the second line of step 3 can be replaced by

$$W_j = A U_j - \sum_{i=\max\{1, j-2\}}^j U_i H_{ij}.$$

The resulting algorithm is called the *block Lanczos iteration*.

Since $\mathcal{K}_k(A, B) \subset R(V_k)$, we have the following result for the block Krylov approximation of matrix polynomials.

LEMMA 6. *For all polynomials p_{k-1} of degree $\leq k-1$ the following holds:*

$$p_{k-1}(A)B = V_k V_k^* p_{k-1}(A)B = V_k p_{k-1}(H_k) V_k^* B.$$

Proof. The proof goes analogously to the proof of [32, Lemma 3.1], where B is a vector. \square

Lemma 6 motivates to carry out the approximation of the product of the matrix exponential and a block matrix as

$$\exp(A)B \approx V_k \exp(H_k) V_k^T B = V_k \exp(H_k) E_1 \|B\|. \quad (3.3)$$

For a vector B , the approximation (3.3) was considered already in [13] and [15], and for the case of a block matrix B it has been considered also in [31].

Since the columns of V_k are orthonormal, we have

$$\mathcal{F}(H_k) = \mathcal{F}(V_k^* A V_k) \subset \mathcal{F}(A).$$

and from this it follows that $\mu(H_k) \leq \mu(A)$. Clearly, it also holds that $\|H_k\| \leq \|A\|$. From these bounds and from Lemma 6 we get the following result.

LEMMA 7. *For the approximation (3.3) holds the following bound:*

$$\|\exp(tA)B - V_k \exp(tH_k) V_k^T B\| \leq 2 \max\{1, e^{t\mu(A)}\} \frac{\|tA\|^k}{k!} \|B\|. \quad (3.4)$$

Proof. The proof goes analogously to the proof of [15, Thm 2.1], where B is a vector. \square

Lemma 7 will be needed repeatedly in the error analysis of our proposed method.

3.2. Derivation of the method. Our strategy is to approximate the symmetric positive semidefinite solution $X(t)$ by a low rank matrix

$$X_k(t) = V_k Y_k(t) V_k^T, \quad (3.5)$$

V_k orthogonal. We first describe how do obtain the numerical approximation (3.5) by using the Galerkin condition, given an orthonormal matrix V_k . An analogous approach for the algebraic Riccati equation can be found in [36].

Denote the right hand of the DRE (2.1) as

$$F(X) = AX + XA^T + Q - XSX.$$

Then, the Galerkin condition

$$V_k^T (\dot{X}_k(t) - F(X_k(t))) V_k = 0$$

directly gives the projected differential equation

$$\dot{Y}_k(t) = A_k Y_k(t) + Y_k(t) A_k^T + Q_k + Y_k(t) S_k Y_k(t), \quad Y(0) = V_k^T X_0 V_k, \quad (3.6)$$

where

$$A_k = V_k^T A V_k, \quad Q_k = V_k^T Q V_k, \quad S_k = V_k^T S V_k.$$

Our choice of the subspace V_k is guided by the following lemma.

LEMMA 8. *Let $X(t)$ satisfy the DRE (2.1) and let $k \geq 1$. Then, there exist matrices $Z_k \in \mathbb{R}^{k(p+q) \times k(p+q)}$ and $W_k \in \mathbb{R}^{n \times k(p+q)}$ such that*

$$X^{(k-1)}(0) = W_k Z_k W_k^T, \quad (3.7)$$

and $R(W_k) \subset \mathcal{K}_k(A, [Z \ C])$.

Proof. The proof goes by induction. The claim clearly holds for $k = 1$ since $X(0) = ZZ^T$. Suppose the claim holds for $k = \ell$, $\ell \in \mathbb{N}_+$, i.e., $X^{(\ell-1)}(0) = W_\ell Z_\ell W_\ell^T$ for some matrices $Z_\ell \in \mathbb{R}^{\ell(p+q) \times \ell(p+q)}$ and $W_\ell \in \mathbb{R}^{n \times \ell(p+q)}$ such that $R(W_\ell) \subset \mathcal{K}_\ell(A, [Z \ C])$. Then, from the DRE (2.1) we infer that

$$X^{(\ell)}(0) = AW_\ell Z_\ell W_\ell^T + W_\ell Z_\ell (AW_\ell)^T + CC^T + W_\ell Z_\ell W_\ell^T S W_\ell Z_\ell W_\ell^T. \quad (3.8)$$

The claim clearly holds for all the four terms on the right hand side of (3.8) for $k = \ell + 1$, and therefore also for their linear combination $X^{(\ell)}(0)$. \square

Lemma 8 directly gives the following corollary.

COROLLARY 9. *Let $X(t)$ satisfy the DRE (2.1). Let $k \geq 1$ and consider the Taylor polynomial*

$$X_k(t) = \sum_{\ell=0}^{k-1} \frac{X^{(\ell)}(0)}{\ell!} t^\ell.$$

Then, there exist matrices $Z_k \in \mathbb{R}^{k(p+q) \times k(p+q)}$ and $W_k \in \mathbb{R}^{n \times k(p+q)}$ such that

$$X_k(t) = W_k Z_k W_k^T, \quad (3.9)$$

and $R(W_k) \subset \mathcal{K}_k(A, [Z \ C])$.

Algorithm 1: Krylov subspace approximation of the DRE (2.1)

Input : Time step size h , Krylov subspace size k , matrices $A, S \in \mathbb{R}^{n \times n}$,
 $Z \in \mathbb{R}^{n \times p}$ and $C \in \mathbb{R}^{n \times q}$, $p, q \ll n$.

- 1 Carry out k steps of the block Arnoldi iteration to obtain
 - the orthogonal basis matrix V_k of $\mathcal{K}_k(A, [Z \ C])$
 - the block Hessenberg matrix $A_k = V_k^T A V_k$
 - the matrices $C_k = V_k^T C$ and $Z_k = V_k^T Z$
- 2 Compute $S_k = V_k^T S V_k$.
- 3 Compute the solution $Y_k(t)$ of the small dimensional system

$$\begin{aligned} \dot{Y}_k(t) &= A_k Y_k(t) + Y_k(t) A_k^T + C_k C_k^T - Y_k(t) S_k Y_k(t), \\ Y(0) &= Z_k Z_k^T. \end{aligned} \quad (3.10)$$

- 4 Approximate $X(t) \approx X_k(t) = V_k Y_k(t) V_k^T$.
-

3.3. The method. Motivated by Corollary 9, we approximate $X(t)$ in the block Krylov subspace $\mathcal{K}_k(A, [Z \ C])$. To this end, an orthogonal basis matrix $V_k \in \mathbb{R}^{n \times k(p+q)}$ for $\mathcal{K}_k(A, [Z \ C])$ is generated using the block Arnoldi iteration. Then, led by the Galerkin condition (equation (3.6)), we carry out the approximation as listed in Algorithm 1.

3.4. Computing the small dimensional system. To solve the small dimensional system (3.10) we use the *modified Davison–Maki method* which is presented in [25]. This method can be described as follows.

By Lemma 1, the solution of the projected system (3.10) is given by

$$Y_k(t) = W_k(t) U_k(t)^{-1}, \text{ where } \begin{bmatrix} U_k(t) \\ W_k(t) \end{bmatrix} = \exp \left(t \begin{bmatrix} A_k^T & C_k C_k^T \\ S_k & -A_k \end{bmatrix} \right) \begin{bmatrix} Z_k Z_k^T \\ I_k \end{bmatrix}. \quad (3.11)$$

Instead of directly evaluating $Y_k(t)$ by (3.11), which was the idea of the original Davison–Maki method [9], we perform substepping in order to avoid numerical instabilities arising from the inversion of the matrix $U_k(t)$ in (3.11). This is exactly the modified Davison–Maki method, and with m substeps it can be described by the following pseudocode.

1. **Input:** Hamiltonian matrix $\begin{bmatrix} A_k^T & C_k C_k^T \\ S_k & -A_k \end{bmatrix}$, $Y_0 = Z_k Z_k^T$,
time $t > 0$, substep size $\Delta t = t/m$, $m \in \mathbb{Z}_+$.
2. **Set:** $Y_k = Y_0$.
3. **Iterate.** for $j = 1, \dots, m$:

$$\begin{bmatrix} U_k \\ W_k \end{bmatrix} = \exp \left(\Delta t \begin{bmatrix} A_k^T & C_k C_k^T \\ S_k & -A_k \end{bmatrix} \right) \begin{bmatrix} Y_k \\ I_k \end{bmatrix}, \quad Y_k = W_k U_k^{-1}.$$

For computing the matrix exponential in Step 3, we use the 13th order diagonal Padé approximant which is implemented in Matlab as 'expm' command [20].

3.5. Structure preserving properties of the approximation. We next inspect the two properties stated in Theorem 3. We show that the proposed projection method preserves the property of the positivity of the exact flow, and it also preserves the property of monotonicity under the condition that the matrix V_k used for the projection stays constant when the initial data for the DRE is changed.

THEOREM 10. *The numerical approximation given by Algorithm 1 preserves the property of positivity stated in Theorem 3.*

Proof. The projected coefficient matrices S_k , $C_k C_k^T$ and the initial value $Z_k Z_k^T$ of the small system (3.10) are clearly all symmetric nonnegative. Thus the small system (3.10) is a symmetric DRE. By Theorem 10 of [11], an application of a symplectic Runge–Kutta scheme with positive weights b_i (see [12] for details) gives as a result a symmetric nonnegative solution. As the s th order diagonal Padé approximant equals the stability function of the s -stage Gauss–Legendre method (see e.g. [24, p. 46]), all the substeps of the modified Davison–Maki method (Subsection 3.4) give symmetric nonnegative matrices and as a result $Y_k(t)$ is symmetric nonnegative as well as $X_k(t) = V_k Y_k(t) V_k^T$. \square

THEOREM 11. *The numerical approximation given by Algorithm 1 preserves the property of monotonicity in the following sense. Consider two DREs corresponding to linearizations with the coefficient matrices*

$$H = \begin{bmatrix} A & Q \\ S & -A^T \end{bmatrix} \quad \text{and} \quad \tilde{H} = \begin{bmatrix} \tilde{A} & \tilde{Q} \\ \tilde{S} & -\tilde{A}^T \end{bmatrix} \quad (3.12)$$

such that

$$\tilde{H}J \leq HJ, \quad 0 \leq X_0 \leq \tilde{X}_0. \quad (3.13)$$

Suppose both systems are projected using the same orthogonal matrix $V_k \in \mathbb{R}^{n \times k}$, giving as a result small k -dimensional systems of the form (3.10) for the matrices $Y_k(t)$ and $\tilde{Y}_k(t)$. Then, for the matrices $X_k(t) = V_k Y_k(t) V_k^T$ and $\tilde{X}_k(t) = V_k \tilde{Y}_k(t) V_k^T$ we have

$$X_k(t) \leq \tilde{X}_k(t).$$

Proof. Consider the projected systems of the form (3.10) corresponding to $Y_k(t)$ and $\tilde{Y}_k(t)$ with the projected coefficient matrices A_k , Q_k and S_k , and \tilde{A}_k , \tilde{Q}_k and \tilde{S}_k , respectively. Consider also the corresponding linearizations of the form (3.12) with the Hamiltonian matrices H_k and \tilde{H}_k .

By the reasoning of the proof of Theorem 10, the projected systems corresponding to $Y_k(t)$ and $\tilde{Y}_k(t)$ are symmetric DREs. Moreover, from (3.13) it clearly follows that $\tilde{H}_k J_k \leq H_k J_k$, where $J_k = \begin{bmatrix} 0 & I \\ -I & 0 \end{bmatrix} \in \mathbb{R}^{2k \times 2k}$, and also that $0 \leq Y_0 \leq \tilde{Y}_0$. By Theorem 6 of [12], an application of a symplectic Runge–Kutta scheme with positive weights b_i (see [12] for details) preserves the monotonicity. Thus the Padé approximants of the substeps of the modified Davison–Maki method (Subsection 3.4) preserve the monotonicity. Thus, $Y_k(t) \leq \tilde{Y}_k(t)$ and as a consequence $X_k(t) \leq \tilde{X}_k(t)$. \square

Remark 12. As the basis matrix V_k given by Algorithm 1 is independent of the matrix $S = BB^T$ in the DRE (2.1), where B is the control matrix in the original linear system (1.2), we see that Algorithm 1 preserves monotonicity under modifications of B . This is important as the property of monotonicity is exactly the ability to modify the costs and the optimal control through modifications in the input data (see [12]).

4. A priori error analysis. We first consider the approximation of the DRE without the quadratic term $-X S X$, i.e., we consider the differential Lyapunov equation. This clarifies the presentation as the derived bounds will be later needed when considering the approximation of the differential Riccati equation. We note, however, that the bounds for the Lyapunov equation are applicable outside of the scope of the optimal control problems, e.g., when considering time integration of a two dimensional inhomogeneous heat equation.

4.1. Error analysis for the Lyapunov equation. Consider the symmetric Lyapunov differential equation with low rank initial data and constant low rank inhomogeneity,

$$\begin{aligned}\dot{X}(t) &= AX(t) + X(t)A^T + CC^T, \\ X(0) &= ZZ^T,\end{aligned}\tag{4.1}$$

where $Z \in \mathbb{R}^{n \times p}$ and $C \in \mathbb{R}^{n \times q}$, $p, q \ll n$. Then, the approximation is given by $X_k(t) = V_k Y_k(t) V_k^T$, where $Y_k(t)$ is a solution of the projected system (3.10) with $S = 0$. For the error of this approximation we obtain the following bound.

THEOREM 13. *Let $A \in \mathbb{R}^{n \times n}$, $Z \in \mathbb{R}^{n \times p}$, $C \in \mathbb{R}^{n \times q}$, and let $X(t)$ be the solution of (4.1). Let $V_k \in \mathbb{R}^{n \times m(q+p)}$ be an orthogonal basis of the block Krylov subspace $\mathcal{K}_k(A, [Z \ C])$. Let $Y_k(t)$ be the solution of the projected system (3.10) with $S = 0$, and let $X_k(t) = V_k Y_k(t) V_k^T$. Then,*

$$\|X(t) - X_k(t)\| \leq 4 \max\{1, e^{2t\mu(A)}\} \|A\|^k \left(\frac{t^k}{k!} \|X_0\| + \frac{t^{k+1}}{(k+1)!} \|Q\| \right).$$

Proof. Using the integral representation of Theorem 2 for both $X(t)$ and $Y_k(t)$, we see that

$$X(t) - X_k(t) = Err_{1,k}(t) + Err_{2,k}(t),$$

where

$$Err_{1,k}(t) = e^{tA} ZZ^T e^{tA^T} - V_k e^{tH_k} V_k^T ZZ^T V_k e^{tH_k^T} V_k^T \tag{4.2}$$

and

$$\begin{aligned}Err_{2,k}(t) &= \int_0^t e^{(t-s)A} CC^T e^{(t-s)A^T} ds \\ &\quad - \int_0^t V_k e^{(t-s)H_k} V_k^T CC^T V_k e^{(t-s)H_k^T} V_k^T ds.\end{aligned}\tag{4.3}$$

Adding and subtracting $e^{tA} ZZ^T V_k e^{tH_k^T} V_k^T$ to the right hand side of (4.2) gives

$$\begin{aligned}Err_1(t) &= e^{tA} Z (e^{tA} Z - V_k e^{tH_k} V_k^T Z)^T \\ &\quad - (V_k e^{tH_k} V_k^T Z - e^{tA} Z) V_k e^{tH_k^T} V_k^T ZZ^T.\end{aligned}$$

Using Lemma 7 to bound the norm of $\exp(tA)Z - V_k \exp(tH_k) V_k^T Z$, and using the fact that $\mu(H_k) \leq \mu(A)$, gives

$$\|Err_1(t)\| \leq 4 \max(1, e^{2t\mu(A)}) \frac{\|tA\|^k}{k!} \|X_0\|.$$

Then, similarly, adding and subtracting the term $\int_0^t e^{(t-s)A} C C^T V_k e^{(t-s)H_k^T} V_k^T ds$ to (4.3) and applying Lemma 7 shows that

$$\begin{aligned} \|Err_2(t)\| &\leq 4\|Q\| \int_0^t \max\{1, e^{2(t-s)\mu(A)}\} \frac{\|(t-s)A\|^k}{k!} ds \\ &\leq 4\|Q\| \max\{1, e^{2t\mu(A)}\} \|A\|^k \frac{t^{k+1}}{(k+1)!} \end{aligned}$$

which shows the claim. \square

4.2. Refined error bounds for the Lyapunov equation. Although Theorem 13 shows the superlinear convergence speed for the approximation of the Lyapunov equation (4.1), sharper bounds can be obtained, e.g., by using the bounds given in [21]. As an example we state the following. If A is symmetric negative semidefinite with its spectrum inside the interval $[-4\rho, 0]$, and V_k is an orthonormal basis matrix for the block Krylov subspace $\mathcal{K}_k(A, B)$, we have (see [21, Thm. 2]) for the error $\varepsilon_k := \|e^{tA}B - V_k e^{tH_k} V_k^* B\|$ the bound

$$\begin{aligned} \varepsilon_k &\leq 10 e^{-k^2/(5\rho t)} \|B\|, \quad \sqrt{4\rho t} \leq k \leq 2\rho t, \\ \varepsilon_k &\leq 10 (\rho t)^{-1} e^{-\rho t} \left(\frac{e\rho t}{k}\right)^k \quad k \geq 2\rho t. \end{aligned} \tag{4.4}$$

Using (4.4) and following the proof of Theorem 13, we get the following bound for the case of a symmetric negative semidefinite A .

THEOREM 14. *Let $A \in \mathbb{R}^{n \times n}$, $Z \in \mathbb{R}^{n \times p}$ and $C \in \mathbb{R}^{n \times q}$ define the Lyapunov equation 4.1. Let $V_k \in \mathbb{R}^{n \times m(q+p)}$ be an orthogonal basis matrix of the subspace $\mathcal{K}_k(A, [Z \ C])$. Let $Y_k(t)$ be the solution of the projected (using V_k) system (3.10) with $S = 0$, and let $X_k(t) = V_k Y_k(t) V_k^T$. Then, for the error $\varepsilon_k := \|X(t) - X_k(t)\|$ it holds that*

$$\begin{aligned} \varepsilon_k &\leq 20 e^{-k^2/(5\rho t)} (\|X_0\| + t\|Q\|), \quad \sqrt{4\rho t} \leq k \leq 2\rho t, \\ \varepsilon_k &\leq 20 (\rho t)^{-1} e^{-\rho t} \left(\frac{e\rho t}{k}\right)^k (\|X_0\| + t\|Q\|), \quad k \geq 2\rho t, \end{aligned} \tag{4.5}$$

The bound (4.5) can be illustrated with the following simple numerical example. Let $A \in \mathbb{R}^{400 \times 400}$ be the tridiagonal matrix $10^2 \cdot \text{diag}(1, -2, 1)$, $t = 0.05$, and let $Z \in \mathbb{R}^{400}$ and $C \in \mathbb{R}^{400}$ be random vectors. Figure 5.1 shows the convergence of the algorithm vs. the a priori bound (4.5).

4.3. Error for the approximation of the Riccati equation. Here, we state our main theorem which shows the superlinear convergence property of Algorithm 1 when applied to the DRE (2.1). Its proof, which is essentially based on Lemma 7 and Grönwall's lemma, is lengthy and is left to the appendix.

First, however, we state a bound for the norm of the numerical solution $X_k(t)$ which will be needed in the proof of the main theorem.

LEMMA 15. *Suppose, $X_0 = ZZ^T$, $Q = CC^T$ and that S is symmetric nonnegative. Then, $X_k(t)$ is symmetric nonnegative, and satisfies the bound*

$$\|X_k(t)\| \leq e^{2t\mu(A)} \|X_0\| + t\varphi_1(2t\mu(A)) \|Q\|.$$

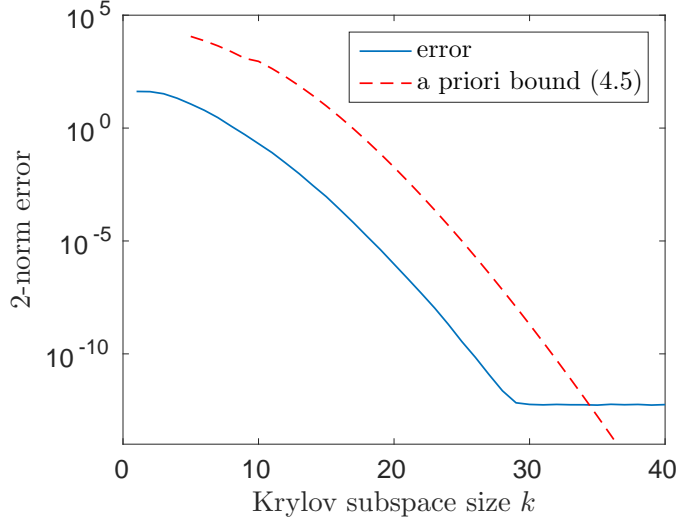


FIGURE 4.1. Convergence of the approximation vs. the a priori bound given in the equation (4.5).

Proof. As ZZ^T , CC^T and S are symmetric and nonnegative, we see from (3.10) that so are the orthogonally projected matrices $Z_k Z_k^T$, $C_k C_k^T$ and S_k . Thus, the projected system is a symmetric DRE. Applying Lemma 4 to the projected system, and using the facts that $\mu(H_k) \leq \mu(A)$, $\|Q_k\| \leq \|Q\|$ and $\|V_k Y_0 V_k^T\| \leq \|X_0\|$, shows the claim. \square

From Lemma 15 we immediately get the following bound.

COROLLARY 16. *The numerical solution $X_k(t)$ satisfies*

$$\max_{s \in [0, t]} \|X_k(s)\| \leq \max\{1, e^{2t\mu(A)}\} \|X_0\| + t \max\{1, \varphi_1(2t\mu(A))\} \|Q\|.$$

We are now ready to state our main theorem.

THEOREM 17. *Let $A \in \mathbb{R}^{n \times n}$, $Z \in \mathbb{R}^{n \times p}$, $C \in \mathbb{R}^{n \times q}$ and $S \in \mathbb{R}^{n \times n}$ defined the DRE (2.1). Let $X_k(t)$ be the numerical solution given by Algorithm 1. Then, the following bound holds:*

$$\|X(t) - X_k(t)\| \leq C(t) \|A\|^k \left(\frac{t^k}{k!} \|X_0\| + \frac{t^{k+1}}{(k+1)!} \|Q\| \right),$$

where

$$C(t) = 4(1 + 2\|S\|\alpha(t) \max\{1, e^{t\mu(A)}\} C_2(t)) e^{t\|S\|\alpha(t)},$$

$$C_2(t) = 1 + t\|S\|\alpha(t) \varphi_1(t\|S\|\alpha(t) \max\{1, e^{t\mu(A)}\}),$$

and

$$\alpha(t) = \max\{1, e^{2t\mu(A)}\} \|X_0\| + t \max\{1, \varphi_1(2t\mu(A))\} \|Q\|.$$

5. A posteriori error estimation. We consider next a posteriori error estimation of the method by using ideas presented in [7].

Denote the original DRE (2.1) as

$$\dot{X}(t) = F(X(t)), \quad X(0) = X_0.$$

Using the residual matrix $R_k(t) = F(X_k(t)) - \dot{X}_k(t)$ we derive computable error estimates. These derivations are based on the following lemma.

LEMMA 18. *The error $\mathcal{E}_k(t) := X(t) - X_k(t)$ satisfies the equation*

$$\begin{aligned} \mathcal{E}_k(t) = & \int_0^t e^{(t-s)A} R_k(s) e^{(t-s)A^T} ds \\ & - \int_0^t e^{(t-s)A} \left(\mathcal{E}_k(t) S X(s) + X_k(s) S \mathcal{E}_k(s) \right) e^{(t-s)A^T} ds, \end{aligned} \quad (5.1)$$

where

$$R_k(t) = V_{k+1} H_{k+1,k} E_k^T Y_k(t) V_k^T + V_k Y_k(t) E_k H_{k+1,k}^T V_{k+1}^T. \quad (5.2)$$

Proof. We see that the error $\mathcal{E}_k(t)$ satisfies the ODE

$$\begin{aligned} \dot{\mathcal{E}}_k(t) &= \dot{X}(t) - \dot{X}_k(t) = F(X(t)) - F(X_k(t)) + R_k(t) \\ &= A(X(t) - X_k(t)) + (X(t) - X_k(t)) A^T \\ &\quad - X(t) S X(t) + X_k(t) S X_k(t) + R_k(t) \\ &= A(X(t) - X_k(t)) + (X(t) - X_k(t)) A^T \\ &\quad - (X(t) - X_k(t)) S X(t) - X_k(t) S (X(t) - X_k(t)) + R_k(t) \\ &= A \mathcal{E}_k(t) + \mathcal{E}_k(t) A^T - \mathcal{E}_k(t) S X(t) - X_k(t) S \mathcal{E}_k(t) + R_k(t) \end{aligned} \quad (5.3)$$

with the initial value $\mathcal{E}_k(0) = 0$. Applying the variation-of-constants formula to (5.3) gives (5.1).

Next we show the representation (5.2). Since

$$F(X_k(t)) = A V_k Y_k(t) V_k^T + V_k Y_k(t) V_k^T A^T + Q - V_k Y_k(t) V_k^T S V_k Y_k(t) V_k^T$$

and

$$\dot{X}_k(t) = V_k H_k Y_k(t) V_k^T + V_k Y_k(t) H_k^T V_k^T + V_k Q_k V_k^T - V_k Y_k(t) V_k^T S V_k Y_k(t) V_k^T$$

we see that

$$\begin{aligned} R_k(t) &= F(X_k(t)) - \dot{X}_k(t) \\ &= (A V_k - V_k H_k) Y_k(t) V_k^T + V_k Y_k(t) (A V_k - V_k H_k)^T + Q - V_k Q_k V_k^T \\ &= (A V_k - V_k H_k) Y_k(t) V_k^T + V_k Y_k(t) (A V_k - V_k H_k)^T, \end{aligned} \quad (5.4)$$

since $V_k Q_k V_k^T = V_k V_k^T C C^T V_k V_k^T = C C^T = Q$ as $C \in \mathcal{R}(V_k)$. Substituting the Arnoldi relation $A V_k - V_k H_k = V_{k+1} H_{k+1,k} E_k^T$ into (5.4) gives the representation (5.2). \square

To derive an a posteriori estimate, we neglect the second term in equation (5.1) and approximate $e^{(t-s)A} \approx I$ in the integrand of the first term. This leads to the approximation

$$\begin{aligned} \mathcal{E}_k(t) \approx \int_0^t R_k(s) \, ds &= U_{k+1} H_{k+1,k} E_k^T \left(\int_0^t Y_k(s) \, ds \right) V_k^T \\ &+ V_k \left(\int_0^t Y_k(s) \, ds \right)^T E_k H_{k+1,k}^T U_{k+1}^T. \end{aligned} \quad (5.5)$$

Since $U_{k+1}^T V_k = 0$, it holds that

$$\begin{aligned} \left\| \int_0^t R_k(s) \, ds \right\| &= \left\| U_{k+1} H_{k+1,k} E_k^T \left(\int_0^t Y_k(s) \, ds \right) V_k^T \right\| \\ &\leq \left\| H_{k+1,k} E_k^T \left(\int_0^t Y_k(s) \, ds \right) \right\|. \end{aligned} \quad (5.6)$$

The integral $\int_0^t Y_k(s) \, ds$ can be estimated by simply summing

$$\int_0^t Y_k(s) \, ds \approx \sum_{\ell=1}^m \Delta t Y_k(\ell \Delta t),$$

where $\Delta t = t/m$ and where the intermediate values $Y_k(\ell \Delta t)$ can be obtained from the summing and squaring phase of Algorithm 1 (Subsection 3.4). From (5.5) and (5.6) we arrive to a computationally efficient a posteriori estimate

$$est_k(t) := \left\| H_{k+1,k} E_k^T \sum_{\ell=1}^m \Delta t Y_k(\ell \Delta t) \right\|. \quad (5.7)$$

To illustrate the efficiency of this estimate consider the following simple example. let $A \in \mathbb{R}^{400 \times 400}$ be the tridiagonal matrix $10^2 \cdot \text{diag}(1, -2, 1)$, $t = 0.1$, and let $Z \in \mathbb{R}^{400}$ and $C \in \mathbb{R}^{400}$ be random vectors. Figure 5.1 shows the error $\|X(t) - X_k(t)\|$ vs. the estimate (5.7).

6. Rank cut and Multiple time stepping. When performing multiple time stepping, the memory consumption may become a problem as the rank of the numerical solution usually grows at each time step. For example, if $\text{rank} \begin{bmatrix} X_0 & Q \end{bmatrix} = m$, then after k iterations using Algorithm 1 the numerical solution $X_k(t) = V_k Y_k(t) V_k^T$ will have a rank km so that memory for $\mathcal{O}(kmn)$ entries is needed. As a remedy to this, a rank cut can be carried out after each time step. This can be done, for example, as follows.

Let $X \in \mathbb{R}^{n \times n}$ and let $\sigma_1 \leq \sigma_2 \leq \dots \leq \sigma_n$ denote the singular values of a matrix X , and u_i, v_i the corresponding left and right singular vectors. Then we consider for cut of the rank the projector

$$P_\epsilon(X) = \sum_{\sigma_i > \epsilon} \sigma_i u_i v_i^T.$$

Note that P_ϵ gives an orthogonal projection onto \mathcal{M}_ϵ , the manifold of matrices with singular values greater than ϵ . This means, in particular, that

$$\|X - P_\epsilon(X)\| \leq \epsilon.$$

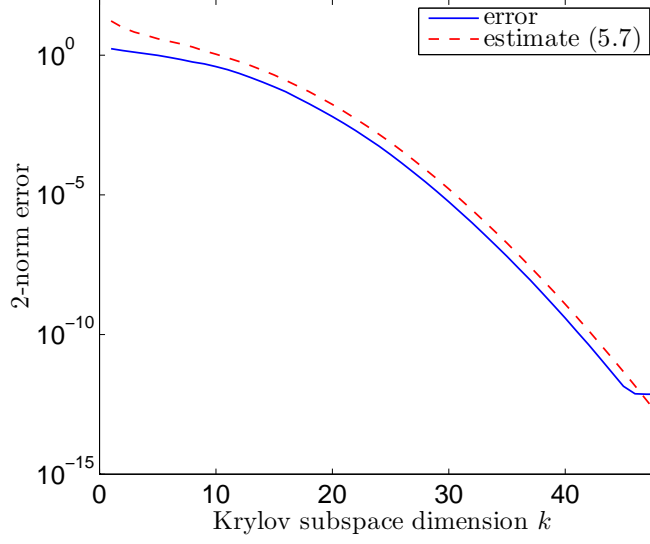


FIGURE 5.1. Convergence of the approximation vs. the a posteriori estimate (5.7).

This projector is applied efficiently on the numerical solution $X_k(t)$ since obviously

$$P_\varepsilon(X_k(t)) = P_\varepsilon(V_k Y_k(t) V_k^T) = V_k P_\varepsilon(Y_k(t)) V_k^T, \quad (6.1)$$

since V_k is orthogonal, and similarly for $P_\varepsilon(X_k(t))$.

Algorithm 2 gives the pseudocode for a simple implementation of the rank cut: at each time step the solution is approximated by the block Krylov subspace depicted by Algorithm 1, after which the rank of the numerical solution is reduced by (6.1).

Algorithm 2: Multiple time stepping with rank cutting

Input : Final time T ,
Number of time steps N ,
tolerance tol_K for the Krylov iteration,
rank cut threshold ε

for $i = 1, 2, \dots, N$ **do**

1 Carry out Algorithm 1 with an initial value X_{i-1} using time step size
 $h = T/N$ and a tolerance tol_K to obtain \tilde{X}_i

2 Project \tilde{X}_i onto \mathcal{M}_ε using (6.1): $X_i = P_\varepsilon(\tilde{X}_i)$

end

An analysis of Algorithm 2 is out of the scope of this article, and we will simply illustrate its behavior in the following numerical example.

7. Numerical example. As a numerical example we consider an optimal cooling problem described in [38] (see also Example 2 in [41]). The underlying linear system is of the form

$$\begin{aligned} M\dot{x} &= Ax + Bu, \\ y &= Cx, \end{aligned} \tag{7.1}$$

where the coefficient matrices arise from a finite element discretization of the cross section of a rail. Here the dimension $n = 1357$ and $A, M \in \mathbb{R}^{n \times n}$, $B \in \mathbb{R}^{n \times 7}$ and $C \in \mathbb{R}^{n \times 6}$. This leads to a symmetric DRE of the form (2.1) with the coefficient matrices $\tilde{A} = M^{-1}A$, $Q = C^T C$ and $S = M^{-1}B(M^{-1}B)^T$. Since M and A are sparse, by recomputing a sparse LU decomposition of M the action of the matrix $M^{-1}A$ on a vector can be evaluated cheaply.

We take zero initial value $X_0 = 0$ for the DRE, and integrate up to $T = 10$. Figure 7.1 shows the convergence of Algorithm 1 and the a posteriori error estimate (5.7). For the scaling and squaring part (Subsection 3.4), we set the parameter $m = 10$. Table 7.1 shows the CPU time needed for the block Krylov process and for the scaling squaring part of Algorithm 1, for four different Krylov subspace sizes.

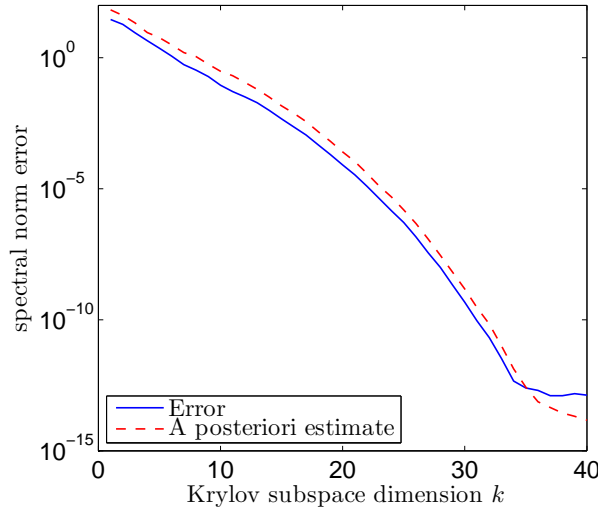


FIGURE 7.1. Convergence of a single step Krylov subspace approximation (Algorithm 1) and its a posteriori estimate (5.7).

k	Krylov iteration	solving small dimensional system
10	0.046	0.037
20	0.16	0.091
30	0.31	0.20
40	0.49	0.44

TABLE 7.1

Timings for the Krylov subspace iteration and for the solving of the dense projected system using the modified Davison–Maki method, when integrating up to $t = 10$ using a single Krylov subspace iteration. Times are in seconds.

Next, we apply Algorithm 2 with $N = 10$ substeps. We set for the Krylov error tolerance $tol_K = \varepsilon$, and use the a posteriori estimate (5.7) as a criterion for stepping the Krylov iteration. Figure 7.2 depicts the final errors at $T = 10$ for 4 different values of ε . As we see the final errors are not far from the tolerances ε used for substeps. Figure 7.3 depicts the growth of the rank in the numerical solution for different tolerances ε . As expected, the cruder the tolerance, the more memory is needed for the numerical approximation.

Interestingly, the substepping approach, i.e. that of Algorithm 2, requires less memory than for a given error tolerance than a single run using Algorithm 1. This is depicted in Table 7.2.

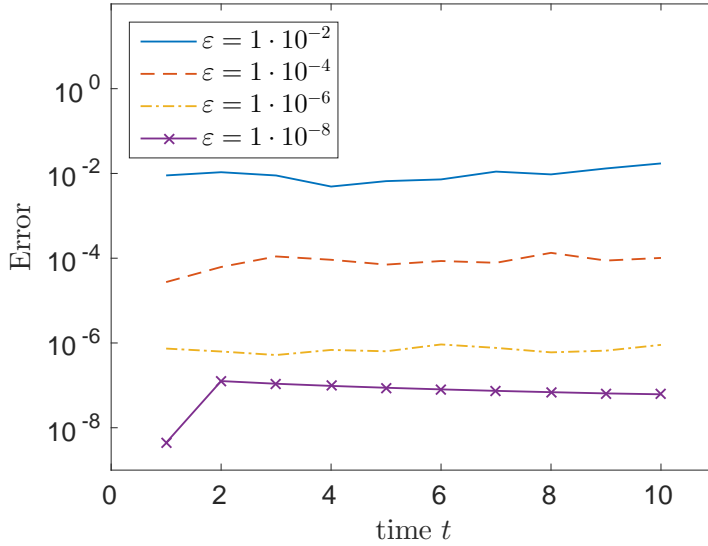


FIGURE 7.2. The error of the numerical solution for different tolerances ε (Algorithm 2).

ε	time stepping	single step iteration
10^{-2}	60	112
10^{-4}	76	147
10^{-6}	112	175
10^{-8}	160	203

TABLE 7.2

Number of columns of the basis matrix V_k along the iteration for the time stepping approach of Algorithm 2 and for one step approximation using Algorithm 1, when an error tolerance ε is required.

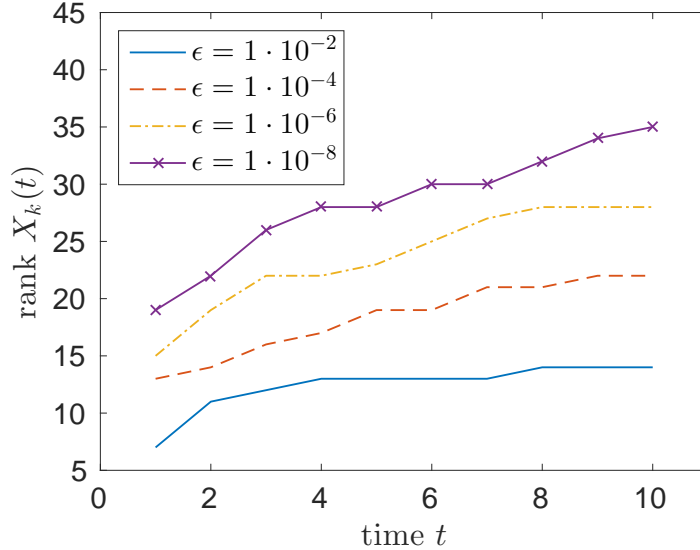


FIGURE 7.3. The rank of the numerical solution for different tolerances ϵ (Algorithm 2).

8. Conclusions and Outlook. We have proposed a Krylov subspace approximation method for large scale differential Riccati equations. We have proven that the method is structure preserving in a sense that it preserves two important properties of the exact flow, namely the property of the positivity and also under certain conditions also the property of the monotonicity. We have also provided theoretical a priori error analysis of the Krylov subspace approximation which shows a superlinear convergence. This behavior was verified in the numerical experiments. In addition, a posteriori error analysis was carried out and the proposed estimate turned out to be very accurate. In order to limit the memory consumption in multiple time stepping, we proposed an efficient algorithm for limiting the rank of the numerical solution. This strategy was shown to be efficient in numerical experiments. However, more numerical and also analytic studies are needed for the rank limitation and multiple time stepping strategies.

We would like to point out that the presented block Krylov subspace method can be straightforwardly extended to the case of unsymmetric differential Riccati equation. A possible extension could also be the nonautonomous case, i.e, the case in which the coefficient matrices are Q , S and A are time dependent (see e.g. [2]). Problems of this kind are crucial in nonlinear control problems in the context of model predictive control.

Acknowledgments. The authors thank Valeria Simoncini for pointing out relevant literature related to the algebraic Riccati equation.

Appendix A. Proof of Theorem 17.

We first state two technical lemmas needed in the proof of the main theorem.

A.1. Auxiliary results. LEMMA 19. *Let $A \in \mathbb{C}^{n \times n}$, $B \in \mathbb{C}^{n \times \ell}$, let V_k be an orthogonal matrix, i.e., $V_k^* V_k = I_k$, such that $\mathcal{K}_k(A, B) \subset R(V_k)$ and let $H_k = V_k^* A V_k$. Then, for all $t, s > 0$ it holds that*

$$\| (e^{tA} - V_k e^{tH_k} V_k^T) e^{sA} B \| \leq 4 \max\{1, e^{(t+s)\mu(A)}\} \frac{\|(t+s)A\|^k}{k!} \|B\|.$$

Proof. Using Lemma ??, we see that

$$\begin{aligned} V_k V_k^T e^{sA} B &= \sum_{\ell=0}^{k-1} \frac{(sA)^\ell}{\ell!} B + V_k V_k^T \sum_{\ell=k}^{\infty} \frac{(sA)^\ell}{\ell!} B \\ &= V_k \sum_{\ell=0}^{k-1} \frac{(sH_k)^\ell}{\ell!} V_k^T B + V_k V_k^T \sum_{\ell=k}^{\infty} \frac{(sA)^\ell}{\ell!} B \\ &= V_k e^{sH_k} V_k^T B - V_k \sum_{\ell=k}^{\infty} \frac{(sH_k)^\ell}{\ell!} V_k^T B + V_k V_k^T \sum_{\ell=k}^{\infty} \frac{(sA)^\ell}{\ell!} B \end{aligned}$$

Therefore,

$$\begin{aligned} (e^{tA} - V_k e^{tH_k} V_k^T) e^{sA} B &= e^{(t+s)A} B - V_k e^{tH_k} V_k^T (V_k V_k^T e^{sA} B) \\ &= e^{(t+s)A} B - V_k e^{(t+s)H_k} V_k^T B \\ &\quad - V_k e^{tH_k} \sum_{\ell=k}^{\infty} \frac{(sH_k)^\ell}{\ell!} V_k^T B + V_k e^{tH_k} V_k^T \sum_{\ell=k}^{\infty} \frac{(sA)^\ell}{\ell!} B \\ &= \sum_{\ell=k}^{\infty} \frac{((t+s)A)^\ell}{\ell!} B - V_k \sum_{\ell=k}^{\infty} \frac{((t+s)H_k)^\ell}{\ell!} V_k^T B \\ &\quad - V_k e^{tH_k} \sum_{\ell=k}^{\infty} \frac{(sH_k)^\ell}{\ell!} V_k^T B + V_k e^{tH_k} V_k^T \sum_{\ell=k}^{\infty} \frac{(sA)^\ell}{\ell!} B \end{aligned} \tag{A.1}$$

Using the bounds $\|e^{tA}\| \leq e^{\mu(A)}$ and $\mu(H_k) \leq \mu(A)$, and the bound (See [15, Lemma B.2])

$$\left\| \sum_{\ell=k}^{\infty} \frac{(tA)^\ell}{\ell!} \right\| \leq \max\{1, e^{\mu(A)}\} \frac{\|tA\|^k}{k!} \quad \text{for all } t \geq 0,$$

on the four terms on the RHS of (A.1), the claim follows. \square

LEMMA 20. *Let $X(s)$ be the solution of the Riccati differential equation (2.1) at time $s > 0$, and let V_k be the orthogonal basis matrix for the Krylov subspace $\mathcal{K}_k(A, [C \ Z])$, and denote $H_k = V_k^T A V_k$. Then, the following bound holds:*

$$\begin{aligned} &\| (e^{(t-s)A} - V_k e^{(t-s)H_k} V_k^T) X(s) \| \\ &\leq 4 C(s) \max\{1, e^{(t+s)\mu(A)}\} \|A\|^k \left(\frac{t^k}{k!} \|X_0\| + \frac{t^{k+1}}{(k+1)!} \|Q\| \right), \end{aligned}$$

where

$$C(s) := 1 + s \|S\| \max_{w \in [0, s]} \|X(w)\| \varphi_1 \left(s \|S\| \max_{w \in [0, s]} \|X(w)\| \max\{1, e^{s\mu(A)}\} \right).$$

Proof. Using the integral representation (2.5) for $X(s)$ we may write

$$(e^{(t-s)A} - V_k e^{(t-s)H_k} V_k^T) X(s) = C_{1,k}(t, s) + C_{2,k}(t, s), \quad (\text{A.2})$$

where

$$\begin{aligned} C_{1,k}(t, s) &= \left[(e^{(t-s)A} - V_k e^{(t-s)H_k} V_k^T) e^{sA} Z \right] Z^T e^{sA^T} \\ &\quad + \int_0^s \left[(e^{(t-s)A} - V_k e^{(t-s)H_k} V_k^T) e^{(s-u)A} C \right] C^T e^{(s-u)A^T} du \end{aligned} \quad (\text{A.3})$$

and

$$C_{2,k}(t, s) = (e^{(t-s)A} - V_k e^{(t-s)H_k} V_k^T) \int_0^s e^{(s-u)A} X(u) S X(u) e^{(s-u)A^T} du. \quad (\text{A.4})$$

By using Lemma 19 on the expressions inside the square brackets in right hand side of (A.3) we obtain the bounds

$$\begin{aligned} &\left\| \left[(e^{(t-s)A} - V_k e^{(t-s)H_k} V_k^T) e^{sA} Z \right] Z^T e^{sA^T} \right\| \\ &\leq 4 \max\{1, e^{(t+s)\mu(A)}\} \|X_0\| \|A\|^k \frac{t^k}{k!} \end{aligned} \quad (\text{A.5})$$

and

$$\begin{aligned} &\left\| \int_0^s \left[(e^{(t-s)A} - V_k e^{(t-s)H_k} V_k^T) e^{(s-u)A} C \right] C^T e^{(s-u)A^T} du \right\| \\ &\leq 4 \|Q\| \int_0^s \frac{((t-u)\|A\|)^k}{k!} \max\{1, e^{(t-u)\mu(A)}\} \max\{1, e^{(s-u)\mu(A)}\} du \\ &\leq 4 \|Q\| \max\{1, e^{(t+s)\mu(A)}\} \|A\|^k \frac{t^{k+1}}{(k+1)!}. \end{aligned}$$

Thus,

$$\|C_{1,k}(t, s)\| \leq 4 \max\{1, e^{(t+s)\mu(A)}\} \|A\|^k \left(\frac{t^k}{k!} \|X_0\| + \frac{t^{k+1}}{(k+1)!} \|Q\| \right). \quad (\text{A.6})$$

From (A.4) we see that

$$\begin{aligned} \|C_{2,k}(t, s)\| &\leq \int_0^s \left\| (e^{(t-s)A} - V_k e^{(t-s)H_k} V_k^T) e^{(s-u)A} X(u) \right\| \\ &\quad \cdot \|S\| \max_{w \in [0, s]} \|X(w)\| e^{(s-u)\mu(A)} du. \end{aligned} \quad (\text{A.7})$$

Next we bound the first factor in the integrand of (A.7). We substitute the integral

representation (2.5) for $X(u)$ to find that

$$\begin{aligned}
& \left\| (e^{(t-s)A} - V_k e^{(t-s)H_k} V_k^T) e^{(s-u)A} X(u) \right\| \\
& \leq \left\| \left[(e^{(t-s)A} - V_k e^{(t-s)H_k} V_k^T) e^{sA} Z \right] Z^T e^{uA^T} \right\| \\
& + \int_0^u \left\| \left[(e^{(t-s)A} - V_k e^{(t-s)H_k} V_k^T) e^{(s-w)A} C \right] C^T e^{(u-w)A^T} \right\| \\
& + \int_0^u \left\| (e^{(t-s)A} - V_k e^{(t-s)H_k} V_k^T) e^{(s-w)A} X(w) \right\| dw \\
& \cdot \|S\| \max_{w \in [0, u]} \|X(w)\| \max\{1, e^{(u-w)\mu(A)}\} dw.
\end{aligned} \tag{A.8}$$

As above when bounding $\|C_{1,k}(t, s)\|$, we use Lemma 19 on the expressions inside the square brackets on right hand side of (A.8), to get the inequality

$$\begin{aligned}
& \left\| (e^{(t-s)A} - V_k e^{(t-s)H_k} V_k^T) e^{(s-u)A} X(u) \right\| \\
& \leq 4\|A\|^k \max\{1, e^{(t+u)\mu(A)}\} \left(\frac{t^k}{k!} \|X_0\| + \frac{t^{k+1}}{(k+1)!} \|Q\| \right) \\
& + \|S\| \max_{w \in [0, u]} \|X(w)\| \max\{1, e^{u\mu(A)}\} \\
& \cdot \int_0^u \left\| (e^{(t-s)A} - V_k e^{(t-s)H_k} V_k^T) e^{(s-w)A} X(w) \right\| dw.
\end{aligned} \tag{A.9}$$

Applying Grönwall's lemma on (A.9), we find that

$$\begin{aligned}
& \left\| (e^{(t-s)A} - V_k e^{(t-s)H_k} V_k^T) e^{(s-u)A} X(u) \right\| \\
& \leq 4\|A\|^k \max\{1, e^{(t+u)\mu(A)}\} \left(\frac{t^k}{k!} \|X_0\| + \frac{t^{k+1}}{(k+1)!} \|Q\| \right) \\
& \cdot e^{u\|S\| \max_{w \in [0, u]} \|X(w)\| \max\{1, e^{u\mu(A)}\}}.
\end{aligned} \tag{A.10}$$

Substituting (A.10) into (A.7), we get

$$\begin{aligned}
& \|C_{2,k}(t, s)\| \leq 4\|S\| \max_{w \in [0, s]} \|X(w)\| \|A\|^k \max\{1, e^{(t+s)\mu(A)}\} \\
& \cdot \left(\frac{t^k}{k!} \|X_0\| + \frac{t^{k+1}}{(k+1)!} \|Q\| \right) s\varphi_1 \left(s\|S\| \max_{w \in [0, s]} \|X(w)\| \max\{1, e^{t\mu(A)}\} \right).
\end{aligned} \tag{A.11}$$

The bounds (A.6) and (A.11) together show the claim. \square

Using Lemmas 19 and 20 we are now ready to prove Theorem 17.

A.2. Proof of Theorem 17. *Proof.* From the integral representation (2.5) for $X(t)$ and for the solution $Y_k(t)$ of the small dimensional system (3.10), we see that

$$X(t) - X_k(t) = F_{1,k}(t) + F_{2,k}(t), \tag{A.12}$$

where

$$F_{1,k}(t) := e^{tA} X_0 e^{tA^T} - V_k e^{tH_k} V_k^T X_0 V_k e^{tH_k^T} V_k^T \\ + \int_0^t \left(e^{(t-s)A} Q e^{(t-s)A^T} - V_k e^{(t-s)H_k} Q e^{(t-s)H_k^T} V_k^T \right) ds,$$

and

$$F_{2,k}(t) = \int_0^t e^{(t-s)A} X(s) S X(s) e^{(t-s)A^T} ds \\ - \int_0^t V_k e^{(t-s)H_k} V_k^T X_k(s) S X_k(s) V_k e^{(t-s)H_k^T} V_k^T ds. \quad (\text{A.13})$$

Theorem 13 shows that $F_{1,k}(t)$ is bounded as

$$\|F_{1,k}(t)\| \leq 4 \max\{1, e^{2t\mu(A)}\} \|A\|^k \left(\frac{t^k}{k!} \|X_0\| + \frac{t^{k+1}}{(k+1)!} \|Q\| \right). \quad (\text{A.14})$$

We add and subtract the term

$$\int_0^t e^{(t-s)A} X(s) S X_k(s) V_k e^{(t-s)H_k^T} V_k^T ds$$

to (A.13) to obtain

$$F_{2,k}(t) = \int_0^t e^{(t-s)A} X(s) S F_{3,k}(t, s)^T ds \\ + \int_0^t F_{3,k}(t, s) S X_k(s) V_k e^{(t-s)H_k^T} V_k^T ds, \quad (\text{A.15})$$

where

$$F_{3,k}(t, s) = e^{(t-s)A} X(s) - V_k e^{(t-s)H_k} V_k^T X_k(s) \\ = (e^{(t-s)A} - V_k e^{(t-s)H_k} V_k^T) X(s) - V_k e^{(t-s)H_k} V_k^T (X(s) - X_k(s)). \quad (\text{A.16})$$

From (A.15) and (A.16) we see that

$$\|F_{2,k}(t)\| \leq 2 \|S\| \alpha(t) \int_0^t \max\{1, e^{(t-s)\mu(A)}\} \\ \cdot \left(\| (e^{(t-s)A} - V_k e^{(t-s)H_k} V_k^T) X(s) \| + \| X(s) - X_k(s) \| \right) ds. \quad (\text{A.17})$$

where

$$\alpha(t) = \max \left\{ \max_{s \in [0, t]} \|X(s)\|, \max_{s \in [0, t]} \|X_k(s)\| \right\}.$$

The claim follows now from (A.12), (A.14), (A.17), Lemma 19, Grönwall's lemma, Corollary 5 and Corollary 16. \square

REFERENCES

- [1] H. ABOU-KANDIL, G. FREILING, V. IONESCU, AND G. JANK, *Matrix Riccati Equations in Control and Systems Theory*, Birkhäuser, Basel, 2003.
- [2] P. BADER, S. BLANES AND E. PONSODA, *Structure preserving integrators for solving (non-) linear quadratic optimal control problems with applications to describe the flight of a quadrotor*, J. Comput. Appl. Math. 262 (2014), pp. 223–233.
- [3] B. BECKERMANN AND L. REICHEL, *Error estimates and evaluation of matrix functions via the Faber transform*, SIAM J. Numer. Anal., 47 (2009), pp. 3849–3883.
- [4] P. BENNER AND H. MENA, *Rosenbrock methods for solving Riccati differential equations*, IEEE Trans. Autom. Control 58.11 (2013), pp. 2950–2956.
- [5] P. BENNER AND H. MENA, *Numerical solution of the infinite-dimensional LQR problem and the associated riccati differential equations*, J. Numer. Math., De Gruyter (accepted), (2016), DOI: 10.1515/jnma-2016-1039.
- [6] D.A. BINI, B. IANNAZZO AND B. MEINI, *Numerical solution of algebraic Riccati equations*, SIAM, Philadelphia, 2011.
- [7] M.A. BOTCHEV, V. GRIMM AND M. HOCHBRUCK, *Residual, restarting, and Richardson iteration for the matrix exponential*, SIAM J. Sci. Comput. 35.3 (2013), pp. A1376–A1397.
- [8] M.J. CORLESS AND A.E. FRAZHO, *Linear Systems and Control*, Marcel Dekker, New York, 2003.
- [9] E. DAVISON AND M. MAKI, *The numerical solution of the matrix Riccati differential equation*, IEEE Trans. Autom. Control 18.1 (1973), pp. 71–73.
- [10] L. DIECI, *Numerical integration of the differential Riccati equation and some related issues*, SIAM J. Numer. Anal. 29.3 (1992), pp. 781–815.
- [11] L. DIECI AND T. EIROLA, *Positive definiteness in the numerical solution of Riccati differential equations*, Numer. Math. 67.3 (1994), pp. 303–313.
- [12] L. DIECI AND T. EIROLA, *Preserving monotonicity in the numerical solution of Riccati differential equations*, Numer. Math. 74.1 (1996).
- [13] V.L. DRUSKIN AND L.A. KNIZHNERMAN, *Two polynomial methods of calculating functions of symmetric matrices*, USSR Comput. Math. Math. Phys., 29 (1989), pp. 112–121.
- [14] L.N. TREFETHEN AND M. EMBREE, *Spectra and pseudospectra: the behavior of nonnormal matrices and operators*, Princeton University Press, 2005.
- [15] E. GALLOPOULOS AND Y. SAAD, *Efficient solution of parabolic equations by Krylov approximation methods*, SIAM J. Sci. Statist. Comput., 13 (1992), pp. 1236–1264.
- [16] G. GOLUB AND C. VAN LOAN, *Matrix computations*, 3rd edition, The Johns Hopkins University Press, Baltimore, 2012.
- [17] Y. GÜLDOĞAN, M. HACHED, K. JBILOU AND M. KURULAY, *Low rank approximate solutions to large-scale differential matrix Riccati equations*, arXiv preprint arXiv:1612.00499 (12/2016).
- [18] E. HAIRER, S.P. NØRSETT AND G. WANNER, *Solving Ordinary Differential Equations I: Nonstiff Problems*, Vol. 8 of Springer Series in Computational Mathematics, 2nd edition., Springer, 1993.
- [19] N.J. HIGHAM, *Functions of matrices: theory and computation*, SIAM, Philadelphia, (2008).
- [20] N.J. HIGHAM, *The scaling and squaring method for the matrix exponential revisited*, SIAM J. Matrix Anal. Appl., 26 (2005), pp. 1179–1193.
- [21] M. HOCHBRUCK AND C. LUBICH, *On Krylov subspace approximations to the matrix exponential operator*, SIAM J. Numer. Anal., 34 (1997), pp. 1911–1925.
- [22] M. HOCHBRUCK AND A. OSTERMANN, *Exponential integrators*, Acta Numer. 19 (2010), pp. 209–286.
- [23] R.A. HORN AND C.R. JOHNSON, *Topics in matrix analysis*, Cambridge University Press, 1991.
- [24] A. ISELER AND S.P. NØRSETT, *Order Stars: Theory and Applications*, Chapman & Hall, London, 1991.
- [25] C. KENNEY AND R. LEIPNIK, *Numerical integration of the differential matrix Riccati equation*, IEEE Trans. Autom. Control 30.10 (1985), pp. 962–970.
- [26] J. VAN DEN ESHOF AND M. HOCHBRUCK, *Preconditioning Lanczos approximations to the matrix exponential*, SISC 27.4 (2006), pp. 1438–1457.
- [27] L. KNIZHNERMAN AND V. SIMONCINI, *A new investigation of the extended Krylov subspace method for matrix function evaluations*, Numer. Linear Algebra Appl. 17.4 (2010), pp. 615–638.
- [28] V. KUČERA *A review of the matrix Riccati equation*, Kybernetika 9.1 (1973), pp. 42–61.
- [29] A. LAUB, *A Schur method for solving algebraic Riccati equations*, IEEE Trans. Autom. Control 24.6 (1979), pp. 913–921.

- [30] Y. LIN AND V. SIMONCINI, *Minimal residual methods for large scale Lyapunov equations*, Appl. Numer. Math. 72 (2013), pp. 52–71.
- [31] L. LOPEZ AND V. SIMONCINI, *Preserving geometric properties of the exponential matrix by block Krylov subspace methods*, BIT 46.4 (2006), pp.813–830.
- [32] Y. SAAD, *Analysis of some Krylov subspace approximations to the matrix exponential operator*, SIAM J. Numer. Anal., 29 (1992), pp. 209–228.
- [33] Y. SAAD, *Iterative methods for sparse linear systems*, PWS Publishing, Boston, 1996.
- [34] Y. SAAD, *Numerical solution of large Lyapunov equation*, In: Signal Processing, Scattering, Operator Theory and Numerical Methods, Kaashoek MA, van Schuppen JH, Ran ACM (eds). Birkhäuser, Basel, 1990, pp. 503–51.xs
- [35] V. SIMONCINI, *A new iterative method for solving large-scale Lyapunov matrix equations*, SIAM J. Sci. Comput. 29.3 (2007).
- [36] V. SIMONCINI, *Analysis of the rational Krylov subspace projection method for large-scale algebraic Riccati equations*, SIAM J. Matrix Anal. Appl. 37.4 (2016), pp.1655–1674.
- [37] V. SIMONCINI, *Computational methods for linear matrix equations*, SIAM Rev. 58.3 (2016).
- [38] P. BENNER AND J. SAAK, *A semi-discretized heat transfer model for optimal cooling of steel profiles*, In: Dimension reduction of large-scale systems, vol. 45, Lecture Notes in Computational Science and Engineering, P. Benner, V. Mehrmann, and D. Sorensen, Eds. Berlin/Heidelberg, Springer, 2005, pp. 353–356.
- [39] N. LANG, H. MENA AND J. SAAK, *On the benefits of the LDL factorization for large-scale differential matrix equation solvers*, Linear Algebra Appl., Vol. 480 (2015), pp. 44–71.
- [40] M.N. SPIJKER, *Numerical ranges and stability estimates*, Appl. Numer. Math., 13 (1993), pp. 241–249.
- [41] T. STILLFJORD, *Low-rank second-order splitting of large-scale differential Riccati equations*, IEEE Trans. Autom. Control 60.10 (2015), pp. 2791–2796.
- [42] T. STILLFJORD, *Adaptive high-order splitting schemes for large-scale differential Riccati equations*, arXiv preprint arXiv:1612.00677 (12/2016).
- [43] C. VAN LOAN, *Computing integrals involving the matrix exponential*, IEEE Trans. Autom. Control 23.3 (1978), pp. 395–404.