

Krylov subspace methods for linear systems with tensor product structure*

Daniel Kressner¹ Christine Tobler¹

April 23, 2009

Abstract

The numerical solution of linear systems with certain tensor product structures is considered. Such structures arise, for example, from the finite element discretization of a linear PDE on a d -dimensional hypercube. Linear systems with tensor product structure can be regarded as linear matrix equations for $d = 2$ and appear to be their most natural extension for $d > 2$. A standard Krylov subspace method applied to such a linear system suffers from the curse of dimensionality and has a computational cost that grows exponentially with d . The key to breaking the curse is to note that the solution can often be very well approximated by a vector of low tensor rank. We propose and analyse a new class of methods, so called *tensor Krylov subspace methods*, which exploit this fact and attain a computational cost that grows linearly with d .

1 Introduction

This paper is concerned with certain linear systems that can be written as the sum of d Kronecker products of matrices. More specifically, we consider for $d = 2$,

$$(A_1 \otimes I_{n_2} + I_{n_1} \otimes A_2)x = b_1 \otimes b_2, \quad (1)$$

and for $d = 3$,

$$(A_1 \otimes I_{n_2} \otimes I_{n_3} + I_{n_1} \otimes A_2 \otimes I_{n_3} + I_{n_1} \otimes I_{n_2} \otimes A_3)x = b_1 \otimes b_2 \otimes b_3, \quad (2)$$

where $A_s \in \mathbb{R}^{n_s \times n_s}$, $b_s \in \mathbb{R}^{n_s}$, and I_{n_s} denotes the $n_s \times n_s$ identity matrix. For general $d \in \mathbb{N}$, the linear system takes the form

$$Ax = b, \quad (3)$$

with

$$A = \sum_{s=1}^d I_{n_1} \otimes \cdots \otimes I_{n_{s-1}} \otimes A_s \otimes I_{n_{s+1}} \otimes \cdots \otimes I_{n_d}, \quad (4)$$

$$b = b_1 \otimes \cdots \otimes b_d. \quad (5)$$

¹Seminar for Applied Mathematics, D-MATH, ETH Zurich, Raemistr. 101, CH-8092 Zurich. {kressner,tobler}@math.ethz.ch

*Supported by the SNF research module *Preconditioned methods for large-scale model reduction* within the SNF ProDoc *Efficient Numerical Methods for Partial Differential Equations*.

Classical Krylov subspace methods for solving linear systems, such as conjugate gradient or GMRES, are not well suited for solving (3). To illustrate this, let us consider the case of constant dimensions, $n_s \equiv n$. Then every vector in the Krylov subspace basis has length n^d and a single scalar product requires $2n^d$ operations. The purpose of this paper is to develop Krylov subspace methods having computational costs and memory requirements that scale *linearly*, rather than *exponentially*, in d .

The following model problem shall illustrate the type of applications leading to (3). Consider the partial differential equation

$$-\Delta u = f \quad \text{in } \Omega, \quad u|_{\partial\Omega} = 0, \quad (6)$$

where $\Omega = [0, 1]^d$ is the d -dimensional hypercube. In each space variable y_s , we choose a finite element basis $\mathbb{V}_s = \{v_1^{(s)}, \dots, v_{n_s}^{(s)}\}$, $s = 1, \dots, d$, with $v_i^{(s)}(0) = v_i^{(s)}(1) = 0$. The corresponding $n_s \times n_s$ mass and stiffness matrices for the one-dimensional Laplacian are denoted by M_s and B_s , respectively. For discretizing the variational formulation of the d -dimensional problem (6) we use the tensorized functions

$$v_{i_1, \dots, i_d}(y_1, \dots, y_d) = v_{i_1}(y_1)v_{i_2}(y_2) \cdots v_{i_d}(y_d),$$

yielding the mass and stiffness matrices

$$\mathcal{M} = M_1 \otimes \cdots \otimes M_d, \quad \mathcal{B} = \sum_{s=1}^d M_1 \otimes \cdots \otimes M_{s-1} \otimes B_s \otimes M_{s+1} \otimes \cdots \otimes M_d.$$

Hence, $\mathcal{A} = \mathcal{M}^{-1/2}\mathcal{B}\mathcal{M}^{-1/2}$ is of the form (4) with $A_s = (M_s)^{-1/2}B_s(M_s)^{-1/2}$. If f is separable, $f = f_1(y_1)f_2(y_2) \cdots f_d(y_d)$, then the discretized right hand side takes the form (5). Otherwise, any sufficiently smooth f can be well approximated by a short sum of separable functions, see, e.g., [5, 6, 10], and the solution of the discretized equation can still be obtained from linear systems of the type (3) by superposition.

For $d = 2$, the equation (1) can be reformulated as follows:

$$A_1 X + X A_2^\top = b_1(b_2)^\top, \quad (7)$$

where $x = \text{rowvec}(X)$, with the rowvec operator stacking the rows of a matrix $X \in \mathbb{R}^{n_1 \times n_2}$ into a single column vector $x \in \mathbb{R}^{n_1 \cdot n_2}$. The linear matrix equation (7) is usually called *Sylvester equation*, which has been studied quite intensively, often motivated by applications in systems and control theory. In fact, most results and algorithms presented in this paper are already known for $d = 2$. In particular, several variants of Krylov subspace algorithms for solving (7) have been developed and analysed, see [16, 17, 19, 20, 21]. The novelty of our work is in the extension to $d > 2$; we will point out relevant connections to the case $d = 2$ whenever suitable. A notable exception is the convergence bound for extended Krylov subspace methods we give in Section 6; this result addresses an open question even for the case $d = 2$.

Grasedyck [10] has combined an integral representation of the solution x to (3) with quadrature based on sinc interpolation [23] to show that x can be well approximated by vectors of low tensor rank and to develop a numerical algorithm that scales linearly with d . To the best of our knowledge, this was the first and so far the only algorithm for efficiently approximating x for high dimensions. Somewhat a drawback, the algorithm relies on computing matrix exponentials of scalar multiples of A_s , which might become expensive for larger

matrices. In contrast, the approach proposed in this paper solely relies on matrix-vector multiplications with A_s . If available, matrix-vector products with $(A_s)^{-1}$ can be used to speed up convergence. A variant of Grasedyck's algorithm is still invoked for solving smaller subsystems.

The rest of this paper is organized as follows. Section 2 contains some preliminary results, mainly concerning tensor notation and approximations of low tensor rank to the solution of (3). In Section 3, we will describe the newly proposed tensor Krylov subspace method and discuss some implementation details, such as the efficient computation of the residual. The convergence of this method is analysed in Section 4 for the (symmetric and non-symmetric) positive definite case. Section 5 provides a discussion on solving the compressed systems needed in the course of the tensor Krylov subspace method. In Section 6, we propose an extension of the tensor Krylov subspace method, which is suitable if matrix-vector products not only with A_s but also with $(A_s)^{-1}$ can be performed. Section 7 contains some numerical experiments with academic examples to illustrate the theoretical results obtained in this paper. Finally, some conclusions and possible future research directions are outlined in Section 8.

2 Preliminaries

The following lemma is a consequence of well-known properties of the Kronecker product [15].

Lemma 2.1 *Consider the matrix \mathcal{A} defined in (4). Then $\Lambda(\mathcal{A})$, the set of eigenvalues of \mathcal{A} , is given by all possible sums of eigenvalues of A_1, A_2, \dots, A_d :*

$$\Lambda(\mathcal{A}) = \{\lambda_1 + \lambda_2 + \dots + \lambda_d : \lambda_s \in \Lambda(A_s)\}. \quad (8)$$

The linear system (3) has a unique solution if and only if $\Lambda(\mathcal{A})$ contains no zero eigenvalues, which – by Lemma 2.1 – is equivalent to

$$\lambda_1 + \lambda_2 + \dots + \lambda_d \neq 0, \quad \forall \lambda_s \in \Lambda(A_s). \quad (9)$$

In the case of the Sylvester equation (7), this corresponds to the well-known condition $\Lambda(A_1) \cap \Lambda(-A_2) = \emptyset$.

We recall that a non-symmetric matrix A is called *positive definite* if its symmetric part $(A + A^\top)/2$ is positive definite. By Lemma 2.1, the matrix \mathcal{A} is positive definite if and only if

$$\mu_1 + \mu_2 + \dots + \mu_d > 0, \quad \forall \mu_s \in \Lambda(A_s + A_s^\top)/2. \quad (10)$$

The following lemma recalls an integral representation of the solution x from [10]. The proof is included for completeness, as it already demonstrates the importance of the separability of the exponential.

Lemma 2.2 *If \mathcal{A} is positive definite then the solution of the linear system (3) admits the representation*

$$x = - \int_0^\infty (\exp(-tA_1)b_1 \otimes \dots \otimes \exp(-tA_d)b_d) dt.$$

Proof. Since \mathcal{A} is positive definite, we have the representation

$$\mathcal{A}^{-1} = - \int_0^\infty \exp(-t\mathcal{A}) dt = - \int_0^\infty \prod_{s=1}^d \exp(-t\hat{A}_s) dt, \quad (11)$$

where we used the fact that the terms

$$\hat{A}_s = I \otimes \cdots \otimes I \otimes A_s \otimes I \otimes \cdots \otimes I, \quad (12)$$

contributing to the sum in \mathcal{A} , commute. This yields

$$\begin{aligned} \mathcal{A}^{-1}b &= - \int_0^\infty \prod_{s=1}^d \exp(-t\hat{A}_s)(b_1 \otimes \cdots \otimes b_d) dt \\ &= - \int_0^\infty \prod_{s=1}^d (I \otimes \cdots \otimes I \otimes \exp(-tA_s) \otimes I \otimes \cdots \otimes I)(b_1 \otimes \cdots \otimes b_d) dt \\ &= - \int_0^\infty (\exp(-tA_1)b_1 \otimes \cdots \otimes \exp(-tA_d)b_d) dt, \end{aligned}$$

concluding the proof. \square

Note that Lemma 2.2 still holds if we impose the less restrictive condition that the eigenvalues of \mathcal{A} have positive real part. It is worth emphasizing that the proof of Lemma 2.2 heavily relies on the commutativity of the matrices \hat{A}_s defined in (12), as do all the other developments in this paper.

2.1 Tensor arithmetic and decompositions

This section provides a brief overview of tensor arithmetic concepts needed in the rest of the paper. We refer to the recent survey [3] for more details.

A d -way tensor $v \in \mathbb{R}^{n_1 \times n_2 \times \cdots \times n_d}$ is an element of the tensor product of the vector spaces $\mathbb{R}^{n_1}, \mathbb{R}^{n_2}, \dots, \mathbb{R}^{n_d}$ for fixed integers n_1, \dots, n_d . The coordinates of v (with respect to a choice of bases) form a multi-dimensional array. The element at the multi-index $\mathfrak{J} = (i_1, i_2, \dots, i_n)$ in such an array is denoted by $v_{\mathfrak{J}}$. A tensor can be represented as a vector in $\mathbb{R}^{n_1 n_2 \cdots n_d}$ by simply stacking the elements $v_{\mathfrak{J}}$ in lexicographical order. *In the following, we will identify tensors with their vector representations.* For $d = 2$, this means that a matrix $A \in \mathbb{R}^{n_1 \times n_2}$ is identified with the vector $\text{rowvec}(A) \in \mathbb{R}^{n_1 n_2}$, where rowvec stacks the transposed rows of A on top of each other. This identification of tensors with vectors is unambiguous as soon as the order d and the dimensions n_1, \dots, n_d are fixed.

A tensor $v \in \mathbb{R}^{n_1 \times n_2 \times \cdots \times n_d}$ is of *tensor rank one* if its vector representation can be written as a Kronecker product of d vectors:

$$v = v^{(1)} \otimes v^{(2)} \otimes \cdots \otimes v^{(d)}, \quad v^{(s)} \in \mathbb{R}^{n_s}. \quad (13)$$

Note that this implies that the element at the multi-index $\mathfrak{J} = (i_1, i_2, \dots, i_n)$ takes the form $v_{\mathfrak{J}} = v_{i_1}^{(1)} v_{i_2}^{(2)} \cdots v_{i_d}^{(d)}$.

A tensor is called *supersymmetric* if it is invariant under any permutations of the indices. For matrices, supersymmetry coincides with the usual notion of symmetry. The following lemma generalizes a well-known result on the symmetry of solutions to Lyapunov matrix equations.

Lemma 2.3 *Consider the linear system $\mathcal{A}x = b$, where \mathcal{A} takes the form (4) with constant coefficients $A = A_1 = \cdots = A_d$. If \mathcal{A} is invertible and b represents a supersymmetric tensor then the solution x is also the representation of a supersymmetric tensor.*

Proof. Assume that x is the solution of $\mathcal{A}x = b$ and that x is not supersymmetric. Then there is a permutation $\pi : \{1, \dots, d\} \rightarrow \{1, \dots, d\}$ such that the vector x_π , the representation of the tensor with elements $x_{\pi(\mathfrak{J})}$ for every multi-index \mathfrak{J} , is different from x . The structure of \mathcal{A} implies $\mathcal{A}x_\pi = b_\pi$. Since b is supersymmetric, $\mathcal{A}x_\pi = b_\pi = b$ contradicting the unique solvability of $\mathcal{A}x = b$. \square

The *CANDECOMP/PARAFAC (CP) decomposition* represents a tensor as a sum of rank one tensors. In vector language, this means

$$v = \sum_{i=1}^k v_i^{(1)} \otimes v_i^{(2)} \otimes \dots \otimes v_i^{(d)}, \quad v_i^{(s)} \in \mathbb{R}^{n_s}. \quad (14)$$

If v admits a representation (14) then we say that v has *tensor rank at most k* . In our context we do not need the concept of exact tensor rank, which is much more subtle than for matrices. Defining the $n_s \times k$ matrices $V_s = [v_1^{(s)}, v_2^{(s)}, \dots, v_k^{(s)}]$, a more compact way of writing (14) is

$$v = \llbracket V_1, V_2, \dots, V_d \rrbracket$$

The *Tucker decomposition* is another popular tensor decomposition. For an integer tuple $\mathfrak{K} = (k_1, \dots, k_d)$ it takes the form

$$v = \sum_{\mathfrak{J} \leq \mathfrak{K}} c_{\mathfrak{J}} v_{i_1}^{(1)} \otimes v_{i_2}^{(2)} \otimes \dots \otimes v_{i_d}^{(d)} =: \llbracket c; V_1, V_2, \dots, V_d \rrbracket, \quad (15)$$

where the sum is taken over all multi-indices $\mathfrak{J} = (i_1, \dots, i_d)$ that are elementwise not larger than \mathfrak{K} . It is worth emphasizing that $V_s = [v_1^{(s)}, v_2^{(s)}, \dots, v_{k_s}^{(s)}]$ is now an $n_s \times k_s$ matrix, i.e., the number of columns of V_s may vary with s . The $k_1 \times \dots \times k_d$ tensor formed from the elements $c_{\mathfrak{J}}$ is called the *core tensor*. Note that the CP decomposition (14) is a special case of (15) for constant $k_s \equiv k$ and an “identity” core tensor that is zero except for $c_{i,i,\dots,i} = 1$ for $i = 1, \dots, k$.

Alternatively, the Tucker decomposition can be written as

$$v = (V_1 \otimes \dots \otimes V_d) c, \quad (16)$$

where c is to be understood as the vector representation of the core tensor in (15). This also reveals that v is in the subspace spanned by $V_1 \otimes V_2 \otimes \dots \otimes V_d$.

Notation 2.4 We write the multi-dimensional Kronecker product as

$$\bigotimes_{s=1}^d v^{(s)} := v^{(1)} \otimes \dots \otimes v^{(d)}.$$

Note that the order in which the index s is evaluated is important, as the Kronecker product does not commute.

2.2 Low tensor rank approximations

Solving (3) for larger d requires to work with a data sparse representation of x . For this purpose, x will be approximated by a low rank tensor. The following theorem provides a fundamental connection between approximations of x by low rank tensors and separable approximations to the reciprocal of a sum of d variables.

Theorem 2.5 Consider the linear system (3) with coefficient matrices $A_s \in \mathbb{R}^{n_s \times n_s}$, assume that the system is solvable and that the matrices A_s are diagonalizable: $P_s^{-1}A_sP_s = \Lambda_s$ with invertible P_s and diagonal Λ_s . Let $f_i^{(s)} : \Omega_s \rightarrow \mathbb{R}$, $i = 1, \dots, k$, be analytic functions such that Ω_s contains the eigenvalues of A_s and

$$\left\| \frac{1}{y_1 + y_2 + \dots + y_d} - \sum_{i=1}^k f_i^{(1)}(y_1) f_i^{(2)}(y_2) \cdots f_i^{(d)}(y_d) \right\|_{\infty} \leq \epsilon(k), \quad (17)$$

where $\|\cdot\|_{\infty}$ denotes the supremum norm on $\Omega_1 \times \dots \times \Omega_d$. Then there is a rank- k tensor x_k such that

$$\|x - x_k\|_2 \leq \kappa \epsilon(k) \|b\|_2,$$

where $\kappa = \kappa_2(P_1)\kappa_2(P_2)\cdots\kappa_2(P_d)$ and $\kappa_2(\cdot)$ denotes the 2-norm condition number of a matrix.

Proof. By a similarity transformation with the matrix $\mathcal{P} = P_1 \otimes \dots \otimes P_d$, the linear system (3) is transformed into

$$\left(\sum_{s=1}^d I_{n_1} \otimes \dots \otimes I_{n_{s-1}} \otimes \Lambda_s \otimes I_{n_{s+1}} \otimes \dots \otimes I_{n_d} \right) \tilde{x} = \tilde{b}.$$

with $\tilde{x} = \mathcal{P}^{-1}x$ and $\tilde{b} = \mathcal{P}^{-1}b$. This is a diagonal linear system and the entry of the solution \tilde{x} at the multi-index $\mathfrak{J} = (i_1, \dots, i_d)$ is given by

$$\tilde{x}_{\mathfrak{J}} = \frac{\tilde{b}_{\mathfrak{J}}}{\lambda_{i_1}^{(1)} + \lambda_{i_2}^{(2)} + \dots + \lambda_{i_d}^{(d)}},$$

where $\lambda_{i_1}^{(1)} + \dots + \lambda_{i_d}^{(d)} \neq 0$ from Lemma 2.1. Similarly, if we define the rank- k tensor

$$x_k = \sum_{j=1}^k f_j^{(1)}(A_1)b_1 \otimes f_j^{(2)}(A_2)b_2 \otimes \dots \otimes f_j^{(d)}(A_d)b_d, \quad (18)$$

the entry of the correspondingly transformed tensor $\tilde{x}_k = \mathcal{P}^{-1}x_k$ at \mathfrak{J} is given by

$$\tilde{b}_{\mathfrak{J}} f_j^{(1)}(\lambda_{i_1}^{(1)}) f_j^{(2)}(\lambda_{i_2}^{(2)}) \cdots f_j^{(d)}(\lambda_{i_d}^{(d)}).$$

Hence, with $\mathfrak{K} = (k, \dots, k)$,

$$\begin{aligned} \|\tilde{x} - \tilde{x}_k\|_2^2 &= \sum_{\mathfrak{J} \leq \mathfrak{K}} |\tilde{b}_{\mathfrak{J}}|^2 \left| \frac{1}{\lambda_{i_1}^{(1)} + \lambda_{i_2}^{(2)} + \dots + \lambda_{i_d}^{(d)}} - \sum_{j=1}^k f_j^{(1)}(\lambda_{i_1}^{(1)}) \cdots f_j^{(d)}(\lambda_{i_d}^{(d)}) \right|^2 \\ &\leq \epsilon(k)^2 \sum_{\mathfrak{J} \leq \mathfrak{K}} |\tilde{b}_{\mathfrak{J}}|^2 = \epsilon(k)^2 \|\tilde{b}\|_2^2. \end{aligned}$$

Combining this bound with $\|x - x_k\|_2 = \|\mathcal{P}(\tilde{x} - \tilde{x}_k)\|_2 \leq \|\mathcal{P}\|_2 \|\tilde{x} - \tilde{x}_k\|_2$, $\|\tilde{b}\|_2 \leq \|\mathcal{P}^{-1}\|_2 \|b\|_2$ and $\kappa = \kappa(\mathcal{P}) = \|\mathcal{P}\|_2 \|\mathcal{P}^{-1}\|_2$ yields the statement of the theorem. \square

Theorem 2.5 provides an upper bound on the error for the best approximation of x by a rank- k tensor. To be practically useful, we still need to address the approximation problem (17). The following technical lemma by Braess and Hackbusch [7, Sec. 2] will turn out to be very helpful for this purpose.

Lemma 2.6 ([7]) *Let $s_k(y) = \sum_{i=1}^k \omega_i \exp(-\alpha_i y)$ with $\alpha_i, \omega_i \in \mathbb{R}$. Then there is a choice of $\alpha_i > 0, \omega_i > 0$ (depending on k and $R > 1$) such that*

$$\sup_{y \in [1, R]} \left| \frac{1}{y} - s_k(y) \right| \leq 16 \exp\left(\frac{-k\pi^2}{\log(8R)}\right).$$

Corollary 2.7 *Consider the linear system $\mathcal{A}x = b$ with \mathcal{A} and b of the form (4)–(5). If \mathcal{A} is symmetric positive definite then there exists an approximation x_k of tensor rank at most k , such that*

$$\|x - x_k\|_2 \leq \frac{16}{\lambda_{\min}(\mathcal{A})} \exp\left(\frac{-k\pi^2}{\log(8\kappa(\mathcal{A}))}\right) \|b\|_2. \quad (19)$$

Proof. As \mathcal{A} is symmetric positive definite, all coefficient matrices A_s are symmetric and $\mathcal{A}x = b$ has a unique solution. Therefore, Theorem 2.5 can be applied with P_s orthogonal (hence, $\kappa = 1$). The result follows directly from this theorem combined with the following observation. Applying Lemma 2.6, we use the substitution $y = \frac{y_1 + \dots + y_d}{\lambda_{\min}(\mathcal{A})}$, $y \in [1, \frac{\lambda_{\max}(\mathcal{A})}{\lambda_{\min}(\mathcal{A})}] =: [1, R]$ and obtain

$$\lambda_{\min}(\mathcal{A}) \left| \frac{1}{y_1 + \dots + y_d} - \sum_{i=1}^k \frac{\omega_i}{\lambda_{\min}(\mathcal{A})} e^{-\alpha_i y_1 / \lambda_{\min}(\mathcal{A})} \dots e^{-\alpha_i y_d / \lambda_{\min}(\mathcal{A})} \right| \leq 16 \exp\left(\frac{-k\pi^2}{\log(8\kappa(\mathcal{A}))}\right),$$

yielding a bound on the quantity $\epsilon(k)$ defined in (17). \square

Corollary 19 shows that the solution x can be well approximated by a low-rank tensor, provided that \mathcal{A} is symmetric positive definite. In comparison, the bound in [11] yields

$$\|x - x_{2k+1}\|_2 \leq \frac{C_{st}}{\lambda_{\min}(\mathcal{A})} \exp(-\sqrt{k}) \|b\|_2, \quad (20)$$

where C_{st} is independent of \mathcal{A} and k .¹ Experimentally, we found $C_{st} \approx 1.5$. It is important to emphasize that the convergence rate predicted by (20) does not depend on the condition number of \mathcal{A} . However, this comes at the expense of having \sqrt{k} instead of k in the exponent. It is therefore of interest to compare (20) with the bound of Corollary 2.7 for different $\kappa(\mathcal{A})$, see Figure 1. It turns out that the bound of Corollary 2.7 is often significantly better, except for very large values of $\kappa(\mathcal{A})$ and small k .

Remark 2.8 *Note that Lemma 2.2 also suggests an algorithm for calculating x_k :*

$$x_k = \sum_{j=1}^k \tilde{\omega}_j \bigotimes_{s=1}^d \exp(-\tilde{\alpha}_j A_s) b_s, \quad (21)$$

with $\tilde{\alpha}_j = \alpha_j / \lambda_{\min}(\mathcal{A})$, $\tilde{\omega}_j = \omega_j / \lambda_{\min}(\mathcal{A})$, and α_j, ω_j as in Lemma 2.6. The coefficients α_j, ω_j only depend on k and $R = \kappa(\mathcal{A}) > 1$. This method is discussed in somewhat more detail in Section 5.

¹Note that there appears to be a typo in the bound of Lemma 5 in [10]. As [10] refers to the results of [11], we have quoted the bound from the latter.

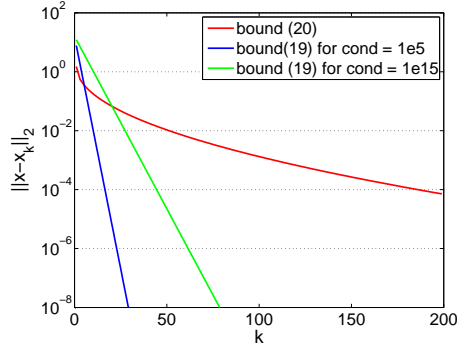


Figure 1: Comparison of the convergence bounds (19) and (20), assuming $\lambda_{\min}(\mathcal{A}) = 1$

3 The tensor Krylov subspace method

In the following, we will develop numerical algorithms for approximating the solution x to the linear system (3). Note that Section 2.2, in particular Remark 2.8, already provides a rather effective method for computing low tensor rank approximations. The computational effort grows linearly with d and the convergence rate depends very mildly, at most logarithmically on the conditioning of \mathcal{A} . However, the involvement of matrix exponentials appears to be a major drawback of this method. The expressions $\exp(-\tilde{\alpha}_j A_s) b_s$ must be computed rather exactly to guarantee a good accuracy of x , which may be regarded expensive compared to, say, a simple matrix-vector multiplication. An approach based on matrix exponentials is certainly feasible in the case of small-sized coefficients; for example if each A_s corresponds to the discretization of a one-dimensional problem, but will become computationally unattractive for larger coefficients A_s . For this purpose, we will propose a method which only requires matrix-vector multiplications with A_s .

3.1 Tensorized Krylov subspaces

We let

$$\mathcal{K}_{k_s}(A_s, b_s) = \text{span}\{b_s, A_s b_s, \dots, A_s^{k_s-1} b_s\}, \quad s = 1, \dots, d,$$

denote the *Krylov subspace* obtained from $k_s - 1$ successive matrix-vector products of A_s with b_s . In view of the PDE (6), each $\mathcal{K}_{k_s}(A_s, b_s)$ could be seen as a subspace corresponding to one coordinate y_s of the domain. To obtain a subspace for all d coordinates, we tensorize and take the linear hull.

Definition 3.1 *Let \mathcal{A}, b be as in equations (3)–(4) and consider a multi-index $\mathfrak{K} = (k_1, \dots, k_d)$ with $k_s \in \mathbb{N}$. Then*

$$\mathcal{K}_{\mathfrak{K}}^{\otimes}(\mathcal{A}, b) := \text{span}(\mathcal{K}_{k_1}(A_1, b_1) \otimes \dots \otimes \mathcal{K}_{k_d}(A_d, b_d))$$

*is called the **tensorized Krylov subspace** associated with \mathcal{A} and b .*

Equivalently, the tensorized Krylov subspace can be defined as

$$\mathcal{K}_{\mathfrak{K}}^{\otimes}(\mathcal{A}, b) = \text{span}\{A_1^{i_1-1} b_1 \otimes \dots \otimes A_d^{i_d-1} b_d : \mathcal{I} \leq \mathfrak{K}\}. \quad (22)$$

A more computationally oriented definition is obtained as follows. Define d matrices U_s , $s = 1, \dots, d$, such that the columns of each U_s span $\mathcal{K}_{k_s}(A_s, b_s)$. Then $\mathcal{K}_{\mathfrak{K}}^{\otimes}(\mathcal{A}, b)$ is spanned by the columns of $\mathcal{U} = U_1 \otimes \dots \otimes U_d$. Combined with equations (15)–(16), this shows that $\mathcal{K}_{\mathfrak{K}}^{\otimes}(\mathcal{A}, b)$ is spanned by the Tucker decompositions $\llbracket c; U_1, U_2, \dots, U_d \rrbracket$ for all possible core tensors c .

In the following, we discuss an extension of the well-known relation between Krylov subspaces and matrix polynomials. Given a multi-index \mathfrak{K} we call $p : \mathbb{R}^d \rightarrow \mathbb{R}$ a *multivariate polynomial* of degree less than \mathfrak{K} if p is a polynomial of degree at most $k_s - 1$ in the s th variable. The space of all such multivariate polynomials is denoted by $\Pi_{\mathfrak{K}}^{\otimes}$. Each $p \in \Pi_{\mathfrak{K}}^{\otimes}$ can be written as

$$p(y_1, \dots, y_d) = \sum_{\mathfrak{L} \leq \mathfrak{K}} c_{\mathfrak{L}} y_1^{l_1-1} y_2^{l_2-1} \dots y_d^{l_d-1}, \quad c_{\mathfrak{L}} \in \mathbb{R},$$

where the sum is taken over multi-indices \mathfrak{L} that satisfy $1 \leq l_s \leq k_s$. The *evaluation* of p at the matrix \mathcal{A} , defined in (4) and represented by the matrix tuple (A_1, \dots, A_d) , is defined as

$$p(A_1, \dots, A_d) := \sum_{\mathfrak{L} \leq \mathfrak{K}} c_{\mathfrak{L}} A_1^{l_1-1} \otimes A_2^{l_2-1} \otimes \dots \otimes A_d^{l_d-1}. \quad (23)$$

Lemma 3.2 *With the notation introduced above,*

$$\mathcal{K}_{\mathfrak{K}}^{\otimes}(\mathcal{A}, b) = \text{span}\{p(A_1, \dots, A_d)b : p \in \Pi_{\mathfrak{K}}^{\otimes}\}.$$

Proof. Let $g = p(A_1, \dots, A_d)b$ for some $p \in \Pi_{\mathfrak{K}}^{\otimes}$. By definition (23), this is equivalent to

$$g = \sum_{\mathfrak{L} \leq \mathfrak{K}} c_{\mathfrak{L}} A_1^{l_1-1} b_1 \otimes \dots \otimes A_d^{l_d-1} b_d, \quad (24)$$

i.e., g is a linear combination of elements from $\mathcal{K}_{k_1}(A_1, b_1) \otimes \dots \otimes \mathcal{K}_{k_d}(A_d, b_d)$ and therefore – by definition – $g \in \mathcal{K}_{\mathfrak{K}}^{\otimes}(\mathcal{A}, b)$.

For the other direction, we note that (22) implies that any $g \in \mathcal{K}_{\mathfrak{K}}^{\otimes}(\mathcal{A}, b)$ can be written in the form (24), which concludes the proof. \square

Remark 3.3 Lemma 3.2 reveals an important difference between standard and tensorized Krylov subspace. For $k_0 \in \mathbb{N}$, the standard Krylov subspace satisfies

$$\mathcal{K}_{k_0}(\mathcal{A}, b) = \{p(\mathcal{A})b : p \in \Pi_{k_0}\},$$

where Π_{k_0} denotes the space of all *univariate* polynomials of degree at most k_0 . For a given $p \in \Pi_{k_0}$, we define the multivariate polynomial

$$p(y_1, y_2, \dots, y_d) := p(y_1 + y_2 + \dots + y_d). \quad (25)$$

By direct computation $p(A_1, \dots, A_d)b = p(\mathcal{A})b$, which – together with Lemma 3.2 – shows

$$\mathcal{K}_{k_0}(\mathcal{A}, b) \subset \mathcal{K}_{\mathfrak{K}}^{\otimes}(\mathcal{A}, b)$$

for $\mathfrak{K} = (k_0, \dots, k_0)$. On the other hand, it is obvious that not every multivariate polynomial takes the particular form (25) and hence $\mathcal{K}_{k_0}(\mathcal{A}, b) \neq \mathcal{K}_{\mathfrak{K}}^{\otimes}(\mathcal{A}, b)$ for $d > 1$ and nontrivial choices of \mathcal{A}, b . To summarize: *Tensorized Krylov subspaces are richer than standard Krylov subspaces.*

3.2 Basic Algorithm

In this section, we present the basics of the newly proposed tensor Krylov subspace algorithm. This algorithm approximates the solution x of the linear system (3) by an element from $\mathcal{K}_{\mathfrak{R}}^{\otimes}(\mathcal{A}, b)$.

To start with, we require a basis of $\mathcal{K}_{\mathfrak{R}}^{\otimes}(\mathcal{A}, b)$. For this purpose, the standard Arnoldi method is used to compute matrices $U_s \in \mathbb{R}^{n_s \times k_s}$ such that the columns of each U_s form an orthonormal basis of the Krylov subspace $\mathcal{K}_{k_s}(A_s, b_s)$. A brief description of the Arnoldi

Algorithm 1 Arnoldi method

Input: Matrix $A_s \in \mathbb{R}^{n_s \times n_s}$, Vector $b_s \in \mathbb{R}^{n_s}$, $k_s \in \mathbb{N}$.

Output: Matrix $U_s \in \mathbb{R}^{n_s \times k_s}$ containing an orthonormal basis of $\mathcal{K}_{k_s}(A_s, b_s)$.

```

 $u_1^{(s)} \leftarrow b_s / \|b_s\|_2$ 
for  $j = 1, \dots, k_s$  do
   $w \leftarrow A_s u_j^{(s)}$ 
   $h_{1:j,j}^{(s)} \leftarrow [u_1^{(s)}, \dots, u_j^{(s)}]^\top w$ 
   $\tilde{u}_{j+1}^{(s)} \leftarrow w - [u_1^{(s)}, \dots, u_j^{(s)}] h_{1:j,j}^{(s)}$ 
   $h_{j+1,j}^{(s)} = \|\tilde{u}_{j+1}^{(s)}\|_2$ 
   $u_{j+1}^{(s)} = \tilde{u}_{j+1}^{(s)} / \|\tilde{u}_{j+1}^{(s)}\|_2$ 
end for
Set  $U_s = [u_1^{(s)}, \dots, u_{k_s}^{(s)}]$ .

```

method is provided in Algorithm 1; more algorithmic details can be found, e.g., in [24]. We assume that a suitable reorthogonalization strategy is performed such that the columns of U_s are also numerically orthonormal.

Upon successful completion of Algorithm 1, one obtains the so called *Arnold decomposition*

$$A_s U_s = U_s H_s + h_{k_s+1, k_s}^{(s)} u_{k_s+1}^{(s)} e_{k_s}^\top, \quad (26)$$

where the upper Hessenberg matrix H_s collects the coefficients $h_{ij}^{(s)}$:

$$H_s = \begin{bmatrix} h_{11}^{(s)} & h_{12}^{(s)} & \cdots & h_{1, k_s}^{(s)} \\ h_{21}^{(s)} & \ddots & \ddots & \vdots \\ & \ddots & \ddots & h_{k_s-1, k_s}^{(s)} \\ & & h_{k_s, k_s-1}^{(s)} & h_{k_s, k_s}^{(s)} \end{bmatrix} = U_s^\top A_s U_s. \quad (27)$$

Note that if A_s is symmetric then H_s inherits this symmetry and becomes a tridiagonal matrix.

As discussed above, the tensor Krylov subspace $\mathcal{K}_{\mathfrak{R}}^{\otimes}(\mathcal{A}, b)$ is spanned by the columns of $\mathcal{U} = U_1 \otimes \dots \otimes U_d$. To extract an approximation to the solution of the linear system (3) from the tensor Krylov subspace, we define $x_{\mathfrak{R}} = \mathcal{U}y$ where y solves the *compressed linear system*

$$\mathcal{H}y = \tilde{b}, \quad \text{with } \mathcal{H} = \mathcal{U}^\top \mathcal{A} \mathcal{U} \text{ and } \tilde{b} = \mathcal{U}^\top b. \quad (28)$$

The solvability of (28) will be discussed in more detail in Section 4. Note that \mathcal{H} and \tilde{b} take the form

$$\begin{aligned}\mathcal{H} &= \mathcal{U}^\top \mathcal{A} \mathcal{U} = \bigotimes_{s=1}^d U_s^\top \left(\sum_{s=1}^d I_{n_1} \otimes \cdots \otimes I_{n_{s-1}} \otimes A_s \otimes I_{n_{s+1}} \otimes \cdots \otimes I_{n_d} \right) \bigotimes_{s=1}^d U_s \\ &= \sum_{s=1}^d I_{k_1} \otimes \cdots \otimes I_{k_{s-1}} \otimes H_s \otimes I_{k_{s+1}} \otimes \cdots \otimes I_{k_d}, \\ \tilde{b} &= \mathcal{U}^\top \bigotimes_{s=1}^d b_s = \bigotimes_{s=1}^d U_s^\top b_s = \prod_{s=1}^d \|b_s\|_2 \bigotimes_{s=1}^d e_1.\end{aligned}$$

It is important to note that the compressed system $\mathcal{H}y = \tilde{b}$ inherits the Kronecker product structure from the original original linear system. Solution methods applicable to the original system can therefore also be applied to the compressed system, see also Section 5. Note that the computational effort for building up and storing the bases of the tensorized Krylov subspace grows only linearly with the number of dimensions. Assuming that the cost for solving the compressed system admits the same growth, we therefore obtain a numerical method with an overall cost that scales linearly with d .

For small dimensions d , it might be feasible to store an explicit representation of the solution y to (28). In this case, the approximation $x_{\mathfrak{R}}$ is represented by the Tucker decomposition

$$x_{\mathfrak{R}} = \mathcal{U}y = \llbracket y; U_1, U_2, \dots, U_d \rrbracket$$

with core tensor y . If y itself is represented by a Tucker decomposition, then $x_{\mathfrak{R}}$ admits again a Tucker decomposition with the same core tensor as y . For high dimensions d , such a representation is not admissible and we will discuss in Section 5 how to represent (or rather approximate) y by a CP decomposition,

$$y = \sum_{i=1}^t y_i^{(1)} \otimes \cdots \otimes y_i^{(d)}. \quad (29)$$

Then $x_{\mathfrak{R}}$ is also represented by a CP decomposition:

$$x_{\mathfrak{R}} \approx \mathcal{U} \sum_{i=1}^t y_i^{(1)} \otimes \cdots \otimes y_i^{(d)} = \sum_{i=1}^t U_1 y_i^{(1)} \otimes \cdots \otimes U_d y_i^{(d)},$$

Algorithm 2 summarizes the proposed tensor Krylov subspace method for solving (3).

Algorithm 2 Tensor Krylov subspace method

Input: Coefficients $A_s \in \mathbb{R}^{n_s \times n_s}$ and $b_s \in \mathbb{R}^{n_s}$ of the linear system (3).

Output: Approximation $x_{\mathfrak{R}} = \mathcal{U}y$ to the solution x of (3).

for $s = 1$ to d **do**

 Apply Algorithm 1 to A_s and b_s to compute U_s, H_s .

end for

 Compute/approximate solution y to the compressed equation $\mathcal{H}y = \tilde{b}$ in (28).

3.3 Computation of the residual

To monitor the convergence of Algorithm 2, one can compute the norm of the residual $r_{\mathfrak{R}} = b - \mathcal{A}x_{\mathfrak{R}}$. The following lemma extends known results for Lyapunov and Sylvester equations [17] to compute this norm in a cheaper way when $x_{\mathfrak{R}}$ is obtained by a Krylov subspace method.

Lemma 3.4 *Let $x_{\mathfrak{R}}$ be computed by Algorithm 2. Then the residual $r_{\mathfrak{R}} = \mathcal{A}x_{\mathfrak{R}} - b$ satisfies*

$$\|r_{\mathfrak{R}}\|_2^2 = \sum_{s=1}^d |h_{k_s+1, k_s}^{(s)}|^2 \sum_{\substack{\mathfrak{L} \leq \mathfrak{R} \\ l_s = k_s}} |y_{\mathfrak{L}}|^2.$$

Proof. For $\beta \in \{0, 1\}^d$, we define $\mathcal{U}_{\beta} = U_{\beta,1} \otimes U_{\beta,2} \otimes \cdots \otimes U_{\beta,d}$, where $U_{\beta,s} = U_s$ if $\beta(s) = 0$ and the columns of $U_{\beta,s}$ form a basis of $\text{span}(U_s)^\perp$ otherwise. Note that $\text{span}(\mathcal{U}_{\beta}) \perp \text{span}(\mathcal{U}_{\gamma})$ unless $\beta = \gamma$. Hence $\mathbb{R}^{n_1 n_2 \cdots n_d}$ can be written as the orthogonal sum of all $\text{span}(\mathcal{U}_{\beta})$ and therefore

$$\|r_{\mathfrak{R}}\|_2^2 = \sum_{\beta \in \{0,1\}^d} \|\mathcal{U}_{\beta}^\top r\|_2^2 \quad (30)$$

For $\beta \equiv 0$, $\mathcal{U}_{\beta} = \mathcal{U}$ and

$$\mathcal{U}^\top r_{\mathfrak{R}} = \mathcal{U}^\top \mathcal{A} \mathcal{U} y - \mathcal{U}^\top b = \mathcal{H} y - \tilde{b} = 0.$$

For general β ,

$$\mathcal{U}_{\beta}^\top \mathcal{A} \mathcal{U} = \sum_{s=1}^d U_{\beta,1}^\top U_1 \otimes \cdots \otimes U_{\beta,s-1}^\top U_{s-1} \otimes U_{\beta,s}^\top A_s U_s \otimes U_{\beta,s+1}^\top U_{s+1} \otimes \cdots \otimes U_{\beta,d}^\top U_d.$$

Note that $U_{\beta,j}^\top U_j = 0$ if $\beta(j) = 1$. Hence, $\mathcal{U}_{\beta}^\top \mathcal{A} \mathcal{U} = 0$ if β contains more than one entry 1. Combined with the fact that $U_{\beta}^\top b = 0$ unless $\beta \equiv 0$, this implies that only terms corresponding to β with exactly one entry 1 contribute to the sum (30). Let us consider such a $\beta_s = (0, \dots, 0, 1, 0, \dots, 0)$ with a single entry 1 at the s th position. Then

$$\begin{aligned} \|\mathcal{U}_{\beta_s}^\top r_{\mathfrak{R}}\|_2^2 &= \|(I \otimes \cdots \otimes I \otimes U_{\beta_s}^\top A_s U_s \otimes I \otimes \cdots \otimes I) y\|_2^2 \\ &= \|(I \otimes \cdots \otimes I \otimes h_{k_s+1, k_s}^{(s)} U_{\beta_s}^\top u_{k_s+1}^{(s)} e_{k_s}^\top \otimes I \otimes \cdots \otimes I) y\|_2^2 \\ &= |h_{k_s+1, k_s}^{(s)}|^2 \|(I \otimes \cdots \otimes I \otimes e_{k_s}^\top \otimes I \otimes \cdots \otimes I) y\|_2^2 \\ &= |h_{k_s+1, k_s}^{(s)}|^2 \sum_{\substack{\mathfrak{L} \leq \mathfrak{R} \\ l_s = k_s}} |y_{\mathfrak{L}}|^2, \end{aligned} \quad (31)$$

where we used the Arnoldi decomposition (26) and $\|U_{\beta_s}^\top u_{k_s+1}^{(s)}\|_2 = 1$. This completes the proof as $\|r\|_2^2$ is obtained by summing up the terms (31) for $s = 1, \dots, d$. \square

Although the expression provided by Lemma 3.4 reduces the cost for computing the residual norm significantly, it still scales exponentially with d simply because almost all elements of y need to be accessed. This exponential growth can be avoided if y is represented in a data-sparse format. Consider, for example, a CP decomposition (29) of y . Then

$$\sum_{\substack{\mathfrak{L} \leq \mathfrak{R} \\ l_s = k_s}} |y_{\mathfrak{L}}|^2 = \left\| \sum_{i=1}^t e_{k_s}^\top y_i^{(s)} \bigotimes_{j \neq s} y_i^{(j)} \right\|_2^2,$$

where the cost for computing the latter expression scales linearly with d , assuming that t remains constant as d grows. Hence, the overall cost for evaluating $\|r_{\mathfrak{R}}\|_2^2$ scales quadratically with d .

4 Convergence analysis

In the following, we will develop a convergence analysis for the tensor Krylov subspace method in special cases. It is clear that this can only be performed in a meaningful way if the unique solvability of the compressed system (28) is guaranteed. The following lemma is an extension of the usual positive definiteness condition in the convergence analysis of FOM methods for standard linear systems, see [22] and the references therein.

Lemma 4.1 *Given the equation (3), suppose that the eigenvalues of the symmetric parts $(A_s + A_s^\top)/2$ are contained in intervals $[\alpha_s, \beta_s]$. Let $\mathcal{U} = U_1 \otimes \cdots \otimes U_d$ where the columns of each $U_s \in \mathbb{R}^{n_s \times k_s}$ form an orthonormal basis. Then the compressed matrix $\mathcal{U}^\top \mathcal{A} \mathcal{U}$ is invertible if*

$$\left[\sum_{s=1}^d \alpha_s, \sum_{s=1}^d \beta_s \right] \cap \{0\} = \emptyset. \quad (32)$$

Proof. The Cauchy interlacing theorem implies that the eigenvalues of the compressed symmetric parts $U_s^\top A_s U_s + U_s^\top A_s^\top U_s = U_s^\top (A_s + A_s^\top) U_s$ are also contained in $[\alpha_s, \beta_s]$. Combined with Lemma 2.1, this shows that any eigenvalue of $\mathcal{U}^\top \mathcal{A} \mathcal{U} + \mathcal{U}^\top \mathcal{A}^\top \mathcal{U}$ can be written as $\mu = \sum_{s=1}^d \mu_s$, $\mu_s \in [\alpha_s, \beta_s]$ and therefore $\mu \in [\sum_{s=1}^d \alpha_s, \sum_{s=1}^d \beta_s]$. Then (32) implies that the set of all such μ is either negative or positive. This shows that the symmetric part of $U_s^\top A_s U_s$ is positive or negative definite, which concludes the proof. \square

It is common to call a non-symmetric matrix to be positive/negative definite if its symmetric part is positive/negative definite. By Lemma 2.1, the condition (32) is equivalent to the definiteness of \mathcal{A} . For this condition to be satisfied it is sufficient but not necessary² that all coefficients A_s are either positive definite or negative definite. In particular, the compressed system is solvable in the special case when all A_s are symmetric positive definite.

For the development of our convergence analysis it is central to note that the *residual* $r_{\mathfrak{R}} = \mathcal{A}x_{\mathfrak{R}} - b$, with $x_{\mathfrak{R}}$ produced by Algorithm 2, satisfies

$$\mathcal{U}^\top r_{\mathfrak{R}} = \mathcal{U}^\top \mathcal{A}x_{\mathfrak{R}} - \mathcal{U}^\top b = \mathcal{U}^\top \mathcal{A} \mathcal{U} y - \tilde{b} = \mathcal{H}y - \tilde{b} = 0.$$

In other words, the following *Galerkin condition* holds:

$$\mathcal{A}x_{\mathfrak{R}} - b \perp \mathcal{K}_{\mathfrak{R}}^\otimes(\mathcal{A}, b). \quad (33)$$

4.1 The symmetric positive definite case

We first consider the case that \mathcal{A} is *symmetric* positive definite. This allows us to introduce the weighted Euclidean norm

$$\|c\|_{\mathcal{A}} := \|\mathcal{A}^{1/2}c\|_2, \quad c \in \mathbb{R}^{n_1 n_2 \cdots n_s},$$

and relate the Galerkin condition (33) to a linear least-squares problem.

²On the other hand, it is easy to see that one can always shift the matrices A_s in a way such that each A_s inherits the definiteness of \mathcal{A} but \mathcal{A} itself is not altered.

Proposition 4.2 *Let x denote solution of (3), where \mathcal{A} is symmetric positive definite. Then the Galerkin condition (33) for an approximation $x_{\mathfrak{R}}$ implies*

$$x_{\mathfrak{R}} = \arg \min_{\tilde{x} \in \mathcal{K}_{\mathfrak{R}}^{\otimes}(\mathcal{A}, b)} \|\tilde{x} - x\|_{\mathcal{A}}.$$

Proof. This result is well known and we include the proof only for the sake of completeness. The condition (33) can be written as

$$0 = \mathcal{U}^{\top}(\mathcal{A}\mathcal{U}y - b) = (\mathcal{A}^{1/2}\mathcal{U})^{\top}(\mathcal{A}^{1/2}\mathcal{U}y - \mathcal{A}^{-1/2}b),$$

which corresponds to the normal equations of the linear least-squares problem

$$y = \arg \min_{\tilde{y} \in \mathbb{R}^{\mathfrak{R}}} \|\mathcal{A}^{1/2}\mathcal{U}\tilde{y} - \mathcal{A}^{-1/2}b\|_2.$$

This can be written as

$$x_{\mathfrak{R}} = \arg \min_{\tilde{x} \in \mathcal{K}_{\mathfrak{R}}^{\otimes}(\mathcal{A}, b)} \|\mathcal{A}^{1/2}\tilde{x} - \mathcal{A}^{-1/2}b\|_2 = \arg \min_{\tilde{x} \in \mathcal{K}_{\mathfrak{R}}^{\otimes}(\mathcal{A}, b)} \|\tilde{x} - x\|_{\mathcal{A}}.$$

□

An immediate consequence of Proposition 4.2 is that the error in the \mathcal{A} -norm of $x_{\mathfrak{R}}$ decreases monotonically as any of the individual dimensions k_s in the tensorized Krylov subspace $\mathcal{K}_{\mathfrak{R}}^{\otimes}(\mathcal{A}, b)$ grows.

The following theorem turns Proposition 4.2 into a multivariate polynomial approximation problem. Let us recall from Section 3.1 that $\Pi_{\mathfrak{R}}^{\otimes}$ denotes the space of all multivariate polynomials of degree at most $k_s - 1$ in the s th variable.

Theorem 4.3 *Under the assumptions of Proposition 4.2,*

$$\|x_{\mathfrak{R}} - x\|_{\mathcal{A}} \leq \sqrt{\|\mathcal{A}\|_2} \|b\|_2 \min_{p \in \Pi_{\mathfrak{R}}^{\otimes}} \max_{\lambda_s \in \Lambda(A_s)} \left| p(\lambda_1, \dots, \lambda_d) - \frac{1}{\lambda_1 + \dots + \lambda_d} \right|.$$

Proof. By Proposition 4.2, $\|x_{\mathfrak{R}} - x\|_{\mathcal{A}} = \min_{\tilde{x} \in \mathcal{K}_{\mathfrak{R}}^{\otimes}(\mathcal{A}, b)} \|\tilde{x} - x\|_{\mathcal{A}}$. By Lemma 3.2, $\tilde{x} \in \mathcal{K}_{\mathfrak{R}}^{\otimes}(\mathcal{A}, b)$ if and only if there exists $p \in \Pi_{\mathfrak{R}}^{\otimes}$ such that $\tilde{x} = p(A_1, \dots, A_d)b$. Hence,

$$\begin{aligned} \|\tilde{x} - x\|_{\mathcal{A}} &\leq \sqrt{\|\mathcal{A}\|_2} \|\tilde{x} - x\|_2 = \sqrt{\|\mathcal{A}\|_2} \min_{p \in \Pi_{\mathfrak{R}}^{\otimes}} \|p(A_1, \dots, A_d)b - \mathcal{A}^{-1}b\|_2 \\ &\leq \sqrt{\|\mathcal{A}\|_2} \|b\|_2 \min_{p \in \Pi_{\mathfrak{R}}^{\otimes}} \|p(A_1, \dots, A_d) - \mathcal{A}^{-1}\|_2. \end{aligned}$$

Let Q_s be an orthogonal matrix such that $Q_s^{\top}A_sQ_s = \Lambda_s$ is a diagonal matrix containing the eigenvalues of A_s on the diagonal. With $\mathcal{Q} = Q_1 \otimes \dots \otimes Q_d$ it is easy to see that

$$\mathcal{Q}^{\top}p(A_1, \dots, A_d)\mathcal{Q} = p(\Lambda_1, \dots, \Lambda_d), \quad \mathcal{Q}^{\top}\mathcal{A}\mathcal{Q} = \sum_{s=1}^d I \otimes \dots \otimes I \otimes \Lambda_s \otimes I \otimes \dots \otimes I,$$

which are both diagonal matrices. Therefore,

$$\begin{aligned} \|p(A_1, \dots, A_d) - \mathcal{A}^{-1}\|_2 &= \|p(\Lambda_1, \dots, \Lambda_d) - \left(\sum_{s=1}^d I \otimes \dots \otimes I \otimes \Lambda_s \otimes I \otimes \dots \otimes I\right)^{-1}\|_2 \\ &= \max_{\lambda_s \in \Lambda(A_s)} \left| p(\lambda_1, \dots, \lambda_d) - \frac{1}{\lambda_1 + \dots + \lambda_d} \right|, \end{aligned}$$

which concludes the proof. \square

We follow the standard technique of relaxing the min-max problem of Theorem 4.3 such that the maximum is taken over the intervals $[\alpha_s, \beta_s] := [\lambda_{\min}(A_s), \lambda_{\max}(A_s)]$ instead of the discrete sets of eigenvalues. This can only increase the bound and we therefore obtain

$$\|x_{\mathfrak{K}} - x\|_{\mathcal{A}} \leq \sqrt{\|\mathcal{A}\|_2} \|b\|_2 E_{\Omega}(\mathfrak{K}),$$

where

$$E_{\Omega}(\mathfrak{K}) := \min_{p \in \Pi_{\mathfrak{K}}^{\otimes d}} \left\| p(\lambda_1, \dots, \lambda_d) - \frac{1}{\lambda_1 + \dots + \lambda_d} \right\|_{\Omega}, \quad (34)$$

with $\|\cdot\|_{\Omega}$ defined as the supremum norm on $\Omega := [\alpha_1, \beta_1] \times \dots \times [\alpha_d, \beta_d]$.

There are several ways to approach the multivariate polynomial approximation problem (34). Inspired by work in [6], we have first followed an approach based on interpolation by tensor Chebyshev polynomials in [25]. Unfortunately, the Lebesgue constants needed to be taken care of in this approach grow exponentially with d , leading to rather loose bounds for high dimensions. For example, if $k_1 = \dots = k_d =: k$ a factor proportional to $\left(\frac{2}{\pi} \ln(k)\right)^{d-1}$ is introduced by the Lebesgue constants. This factor can be avoided when following a completely different approach, essentially mimicking the proof of [21] on a scalar level for arbitrary dimensions. Lemma A.1 in the Appendix contains the approximation result obtained in this manner. Combining Theorem 4.3 with Lemma A.1 yields the following convergence bound.

Corollary 4.4 *Under the assumptions of Proposition 4.2,*

$$\|x_{\mathfrak{K}} - x\|_{\mathcal{A}} \leq \frac{\sqrt{\|\mathcal{A}\|_2} \|b\|_2}{\lambda_{\min}(\mathcal{A})} \sum_{s=1}^d \frac{\sqrt{\kappa_s} + 1}{\sqrt{\kappa_s}} \cdot \left(\frac{\sqrt{\kappa_s} - 1}{\sqrt{\kappa_s} + 1}\right)^{k_s},$$

where $\kappa_s = 1 + \frac{\beta_s - \alpha_s}{\lambda_{\min}(\mathcal{A})}$.

For $d = 2$, a bound similar to the bound of Corollary 4.4 has been obtained in [21, Proposition 3.1]. In fact, the bound of [21] has a somewhat smaller constant by avoiding the detour via multivariate polynomial approximations. In principle, the proof techniques [21] could also be extended to $d > 2$ but our approach has the advantage of also admitting convergence bounds for the extended Krylov subspace method, see Section 6.

Remark 4.5 It is instructive to compare the error bound of Corollary 4.4 with the well-known error bound of the classical CG method applied to the linear system (3):

$$\|x_k - x\|_{\mathcal{A}} \leq 2 \|b\|_{\mathcal{A}} \left(\frac{\sqrt{\kappa(\mathcal{A})} - 1}{\sqrt{\kappa(\mathcal{A})} + 1}\right)^k \leq 2 \sqrt{\|\mathcal{A}\|_2} \|b\|_2 \left(\frac{\sqrt{\kappa(\mathcal{A})} - 1}{\sqrt{\kappa(\mathcal{A})} + 1}\right)^k, \quad (35)$$

where $\kappa(\mathcal{A}) = \|\mathcal{A}\|_2 \|\mathcal{A}^{-1}\|_2 = \frac{\beta_1 + \dots + \beta_d}{\alpha_1 + \dots + \alpha_d}$.

To simplify the discussion, assume $\alpha_1 = \dots = \alpha_d =: \alpha$ and $\beta_1 = \dots = \beta_d =: \beta$, in which case it is reasonable to choose $\mathfrak{K} = (k, \dots, k)$. Then the bound of Corollary 4.4 simplifies to

$$\|x_{\mathfrak{K}} - x\|_{\mathcal{A}} \leq d \frac{\sqrt{\kappa} + 1}{\sqrt{\kappa}} \cdot \frac{\sqrt{\|\mathcal{A}\|_2} \|b\|_2}{\lambda_{\min}(\mathcal{A})} \left(\frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1} \right)^k, \quad (36)$$

where $\kappa = 1 + \frac{\beta - \alpha}{d\alpha} = \frac{d-1}{d} + \frac{\kappa(\mathcal{A})}{d}$. For larger $\kappa(\mathcal{A})$ this means that the effective condition number that determines the convergence rate in (36) is divided by d in comparison to (35). This indicates that the tensor Krylov subspace method can be expected to require $1/\sqrt{d}$ times the iterations needed by classical CG, see also Remark 3.3. More surprisingly, the convergence of the tensor Krylov subspace method *improves* with increasing number of dimensions d , assuming that the condition number of \mathcal{A} remains constant. This benefit from higher dimensions was already noted for $d = 2$ in [21] and is confirmed by the numerical experiments in Section 7.

Remark 4.6 To avoid unnecessary work, it is of interest to choose $\mathfrak{K} = (k_1, \dots, k_d)$ such that the summands in the convergence bound of Corollary 4.4 are balanced, i.e., the terms $\left(\frac{\sqrt{\kappa_s} - 1}{\sqrt{\kappa_s} + 1} \right)^{k_s}$ are nearly constant for all s . Taking the logarithm yields

$$k_s \log \left(\frac{\sqrt{\kappa_s} - 1}{\sqrt{\kappa_s} + 1} \right) \approx -2 \frac{k_s}{\sqrt{\kappa_s}}$$

for large κ_s . Hence it is reasonable to choose $\mathfrak{K} = (k_1, \dots, k_d)$ such that $\frac{k_s}{\sqrt{\kappa_s}}$ is nearly constant across all dimensions s .

4.2 The non-symmetric positive definite case

To obtain convergence bounds in the case that \mathcal{A} is *non-symmetric* positive definite, we follow the proof technique by Simoncini and Druskin [21] for Lyapunov equations. To simplify the presentation it is assumed that $\|b_s\|_2 = 1$ throughout the rest of this section.

Lemma 4.7 *Let x denote the solution of (3), where \mathcal{A} is positive definite. For the approximation $x_{\mathfrak{K}}$ obtained by Algorithm 2 it holds that*

$$\|x_{\mathfrak{K}} - x\|_2 \leq \sum_{s=1}^d \int_0^{\infty} e^{-\hat{\alpha}_s t} \|x_{k_s}^{(s)} - x^{(s)}\|_2 dt,$$

where $x^{(s)} = e^{-tA_s} b_s$, $x_{k_s}^{(s)} = U_s e^{-tH_s} e_1$, and $\hat{\alpha}_s = \sum_{j \neq s} \alpha_j$ with $\alpha_j = \lambda_{\min}(A_j + A_j^{\top})/2$.

Proof. By Lemma 2.2, both x and $x_{\mathfrak{K}}$ admit integral representations:

$$x = \int_0^{\infty} x^{(1)} \otimes \dots \otimes x^{(d)} dt, \quad x_{\mathfrak{K}} = \int_0^{\infty} x_{k_1}^{(1)} \otimes \dots \otimes x_{k_d}^{(d)} dt.$$

Note that $\|x^{(s)}\|_2 \leq e^{-\alpha_s t}$ as well as $\|x_{k_s}^{(s)}\|_2 \leq e^{-\alpha_s t}$. We have

$$\begin{aligned}
\|x_{\mathfrak{R}} - x\|_2 &\leq \int_0^\infty \left\| \bigotimes_{j=1}^d x_{k_j}^{(j)} - \bigotimes_{j=1}^d x^{(j)} \right\|_2 dt \\
&= \int_0^\infty \left\| \sum_{s=1}^d \left(\bigotimes_{j=1}^{s-1} x_{k_j}^{(j)} \right) \otimes (x_{k_s}^{(s)} - x^{(s)}) \otimes \left(\bigotimes_{j=s+1}^d x^{(j)} \right) \right\|_2 dt \\
&\leq \sum_{s=1}^d \int_0^\infty \left(\prod_{j=1}^{s-1} \|x_{k_j}^{(j)}\|_2 \right) \|x_{k_s}^{(s)} - x^{(s)}\|_2 \left(\prod_{j=s+1}^d \|x^{(j)}\|_2 \right) dt \\
&\leq \sum_{s=1}^d \int_0^\infty e^{-\hat{\alpha}_s t} \|x_{k_s}^{(s)} - x^{(s)}\|_2 dt,
\end{aligned}$$

which concludes the proof. \square

Note that the term $\|x_{k_s}^{(s)} - x^{(s)}\|_2$ appearing in the bound of Lemma 4.7 corresponds to the approximation error of the usual Krylov subspace approximation to the matrix exponential $e^{-tA_s} b_s$. Any reasonably good bound on this approximation error could be inserted to yield a bound on $\|x_{\mathfrak{R}} - x\|_2$. In the following, we demonstrate this procedure for the case that the fields of values $F(A_s) = \{w^\top A_s w : w \in \mathbb{C}^{n_s}, \|w\|_2 = 1\}$ for $s = 1, \dots, d$ are contained in ellipses.

Theorem 4.8 *Additionally to the assumptions of Lemma 4.7, suppose that the field of values of each $A_s \in \mathbb{R}^{n_s \times n_s}$ is contained in an ellipse of center $(c_s, 0)$, foci $(c_s \pm f_s, 0)$ and semi-axes $a_s^{(1)}$ and $a_s^{(2)}$. Then*

$$\|x_{\mathfrak{R}} - x\|_2 \leq \sum_{s=1}^d \frac{4}{\sqrt{(\hat{\alpha}_s + c_s)^2 - f_s^2}} \frac{\rho_s}{\rho_s - 1} \rho_s^{-k_s},$$

where

$$r_s = \frac{a_s^{(1)} + a_s^{(2)}}{2}, \quad \rho_s = \frac{c_s + \hat{\alpha}_s}{2r_s} + \frac{1}{2r_s} \sqrt{(c_s + \hat{\alpha}_s)^2 - f_s^2}$$

for $s = 1, \dots, d$.

Proof. By the proof of Proposition 4.1 in [21],

$$\|x_{k_s}^{(s)} - x^{(s)}\|_2 \leq 4 \sum_{j=k_s}^{\infty} e^{-c_s t} I_j(f_s t) \left(\frac{2r_s}{f_s} \right)^j,$$

where I_j denotes the j th modified Bessel function, see also the proof of Lemma A.1. Combined with Lemma 4.7, this yields

$$\begin{aligned}
\|x_{\mathfrak{R}} - x\|_2 &\leq 4 \sum_{s=1}^d \sum_{j=k_s}^{\infty} \int_0^\infty e^{-(\hat{\alpha}_s + c_s)t} I_j(f_s t) \left(\frac{2r_s}{f_s} \right)^j dt \\
&= \sum_{s=1}^d \frac{4}{\sqrt{(\hat{\alpha}_s + c_s)^2 - f_s^2}} \sum_{j=k_s}^{\infty} \rho_s^{-j} \\
&= \sum_{s=1}^d \frac{4}{\sqrt{(\hat{\alpha}_s + c_s)^2 - f_s^2}} \frac{\rho_s}{\rho_s - 1} \rho_s^{-k_s}.
\end{aligned}$$

□

In a manner analogous to the proof of Theorem 4.8 the other results from [21] for non-symmetric positive definite matrices, dealing for example with a field of values contained in a wedge-shaped set, could be extended to arbitrary dimensions.

5 Solving the compressed system

The tensor Krylov subspace method, see Algorithm 2, requires the solution of the compressed system

$$\mathcal{H}y = \tilde{b},$$

having the same Kronecker product structure as the original system (3), with the coefficients H_s in upper Hessenberg form. As mentioned in Section 3.2, this system might be solved explicitly by a direct method for small dimensions but for higher dimensions this will quickly become prohibitively expensive. The method already suggested in Remark 2.8 provides a viable alternative. An approximation y_t to the exact solution y is calculated as

$$y_t = \sum_{j=1}^t \frac{\omega_j}{\lambda_{\min}(\mathcal{H})} \bigotimes_{s=1}^d \exp\left(-\frac{\alpha_j}{\lambda_{\min}(\mathcal{H})} H_s\right) \tilde{b}_s,$$

with

$$1/\lambda \approx \sum_{j=1}^t \omega_j e^{-\alpha_j \lambda}, \quad \lambda \in \Lambda(\mathcal{H}).$$

The success of this approach depends of course crucially on the choice of the coefficients $\alpha_j > 0, \omega_j > 0$. Ideally, we would like to solve the min-max problem

$$\min_{\alpha_j, \omega_j \in \mathbb{R}^+} \max_{\lambda \in \Omega} \left| 1/\lambda - \sum_{j=1}^t \omega_j e^{-\alpha_j \lambda} \right|,$$

where $\Omega \subseteq \mathbb{C}$ contains all eigenvalues of \mathcal{H} scaled by some factor $1/\rho$. Section 5.1 discusses the case of symmetric positive definite \mathcal{H} , in which case $\rho = \lambda_{\min}(\mathcal{H})$ and $\Omega = [1, \kappa(\mathcal{H})]$. For this purpose, a finite upper bound on the condition number $\kappa(\mathcal{H})$ must be available, which can be determined from the eigenvalues of H_s and applying Lemma 2.1. For non-symmetric \mathcal{H} (or in case no upper bound on $\kappa(\mathcal{H})$ can be computed) we have $\Omega = \{z \in \mathbb{C} : \Re(z) \geq 1\}$, assuming that \mathcal{H} has only eigenvalues in the right-half complex plane and $\rho = \min\{\Re(\lambda) : \lambda \in \Lambda(\mathcal{H})\}$. This case is discussed in Section 5.2. A comparison of the coefficients α_j resulting in each case can be found in Figure 2.

For both choices of coefficients, the choice of t needs to be determined in advance. We choose t such that the convergence bounds given below are not larger than a tolerance provided by the user. As this tolerance determines the overall quality of the approximation to the large linear system (3), it will usually be chosen rather small, say 10^{-9} .

5.1 Symmetric \mathcal{H} and condition number known

As already mentioned in Section 2, the existence of coefficients $\alpha_j > 0, \omega_j > 0$ satisfying

$$\left| \frac{1}{\lambda} - \sum_{j=1}^t \omega_j e^{-\alpha_j \lambda} \right| \leq 16 \exp\left(\frac{-t\pi^2}{\log(8R)}\right) \quad \forall \lambda \in [1, R]$$

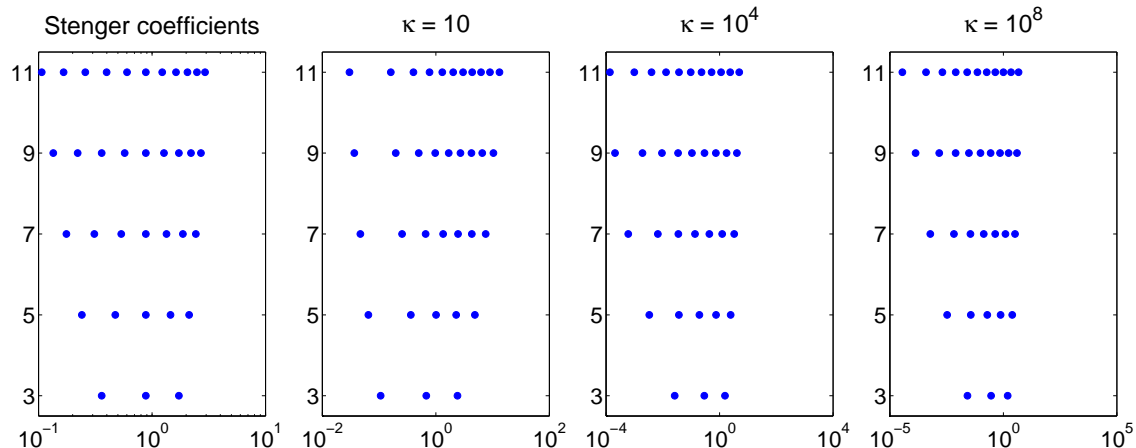


Figure 2: Coefficients α_j proposed in Section 5.2 (leftmost plot) and Section 5.1 (3 rightmost plots, for $\kappa(\mathcal{H}) = 10, 10^4, 10^8$), for $t = 3, 5, \dots, 11$ coefficients.

is proved in [13]. By Corollary 2.7,

$$\|y - y_t\|_2 \leq \frac{16}{\lambda_{\min}(\mathcal{H})} \exp\left(\frac{-t\pi^2}{\log(8\kappa(\mathcal{H}))}\right) \|\tilde{b}\|_2.$$

Unfortunately, there is apparently no simple explicit formula for determining such coefficients $\alpha_j > 0, \omega_j > 0$. The optimal choice of coefficients can be calculated by a Remez-like algorithm for many different values of t and R . This has been performed by Hackbusch [12] and the resulting coefficients were used in this paper.

5.2 Condition number unknown and/or non-symmetric \mathcal{H}

If Ω is a subset of the right-half complex plane, the following bound from [11] applies:

$$\left| \frac{1}{\lambda} - \sum_{j=-t}^t \omega_j e^{-\alpha_j \lambda} \right| \leq C_{st} \exp(|\Im(\lambda)|/\pi) \exp(-\sqrt{t}), \quad \forall \lambda \in \mathbb{C}, \Re(\lambda) \geq 1,$$

where C_{st} does not depend on t or λ . This yields

$$\|y - y_{2t+1}\|_2 \leq C_{st} \|\mathcal{H}^{-1}\|_2 \exp(\mu/\pi) \exp(-\sqrt{t}) \|\tilde{b}\|_2,$$

where $\mu = \max\{|\Im(\lambda)| : \lambda \in \Lambda(\mathcal{H})\}$. In this case, there are explicit formulas for suitable coefficients:

$$\begin{aligned} \alpha_j &= \log\left(\exp(jt^{-1/2}) + \sqrt{1 + \exp(2jt^{-1/2})}\right), \\ \omega_j &= (t + t \exp(-2jt^{-1/2}))^{-1/2}, \text{ for } j = -t, \dots, t. \end{aligned}$$

6 An extended tensor Krylov subspace method

In many cases of practical interest, the convergence of the tensor Krylov subspace method can be significantly accelerated if also matrix-vector products with A_s^{-1} for all or some s are

available, for example by means of sparse direct factorizations. In the following we propose an *extended tensor Krylov subspace method*, in the spirit of the closely related extended Krylov subspace methods for matrix functions [8] and linear matrix equations [20, 14]. The convergence of this method for approximating matrix functions has been recently discussed in [4, 18].

In contrast to the tensor Krylov subspace method described in Section 3, some (or all) of the matrices U_s now represent orthonormal bases for the *extended* Krylov subspace

$$\tilde{\mathcal{K}}_{k_s}(A_s, b_s) := \text{span}(\mathcal{K}_{k_s}(A_s, b_s) \cup \mathcal{K}_{k_s+1}(A_s^{-1}, b_s)),$$

assuming of course that A_s is invertible. Generically, the dimension of the extended Krylov subspace is $2k_s$ and hence $U_s \in \mathbb{R}^{n_s \times 2k_s}$. Arnoldi-type algorithms for computing U_s are discussed, for example, in [18, 20]. The rest of the extended tensor Krylov subspace method is identical with Algorithm 2.

For simplicity, we assume that all A_s are symmetric positive definite and extended Krylov subspaces are used for all $s = 1, \dots, d$. Then

$$\begin{aligned} \tilde{\mathcal{K}}_{k_s}(A_s, b_s) &= \text{span}\{A_s^{-k_s}b_s, \dots, A_s^{-1}b_s, b_s, A_s b_s, \dots, A_s^{k_s-1}b_s\} \\ &= \text{span}\{\ell(A_s)b_s : \ell \in \mathbb{L}_{k_s}\}, \end{aligned}$$

where \mathbb{L}_{k_s} denotes the linear space of Laurent polynomials $\ell(y) = \sum_{j=-k_s}^{k_s-1} c_j y^j$. Similarly for the *tensorized extended Krylov subspace* it holds that

$$\begin{aligned} \tilde{\mathcal{K}}_{\mathfrak{K}}^{\otimes}(\mathcal{A}, b) &:= \text{span}(\tilde{\mathcal{K}}_{k_1}(A_1, b_1) \otimes \dots \otimes \tilde{\mathcal{K}}_{k_d}(A_d, b_d)) \\ &= \text{span}\{\ell(A_1, \dots, A_d)b : \ell \in \mathbb{L}_{\mathfrak{K}}^{\otimes}\}. \end{aligned} \quad (37)$$

Here, $L_{k_s}^{\otimes}$ denotes the space of all multivariate Laurent polynomials

$$\ell(y_1, \dots, y_d) = \sum_{-\mathfrak{K} \leq \mathfrak{L} < \mathfrak{K}} c_{\mathfrak{L}} y_1^{l_1} y_2^{l_2} \dots y_d^{l_d}, \quad c_{\mathfrak{L}} \in \mathbb{R},$$

where $-\mathfrak{K} \leq \mathfrak{L} < \mathfrak{K}$ is to be understood as $-k_s \leq l_s \leq k_s - 1$ for $s = 1, \dots, d$. The evaluation of ℓ at a matrix tuple (A_1, \dots, A_d) is then – analogous to (23) – defined as

$$\ell(A_1, \dots, A_d) = \sum_{-\mathfrak{K} \leq \mathfrak{L} < \mathfrak{K}} c_{\mathfrak{L}} A_1^{l_1} \otimes A_2^{l_2} \otimes \dots \otimes A_d^{l_d}, \quad c_{\mathfrak{L}} \in \mathbb{R}.$$

The identity (37) can be shown in a similar way as Lemma 3.2.

Much of the convergence analysis of Section 4 can be extended in a straightforward way. In particular, the convergence bound of Theorem 4.3 becomes

$$\begin{aligned} \|x_{\mathfrak{K}} - x\|_{\mathcal{A}} &\leq \sqrt{\|\mathcal{A}\|_2} \|b\|_2 \min_{\ell \in \mathbb{L}_{\mathfrak{K}}^{\otimes}} \max_{\lambda_s \in \Lambda(A_s)} \left| \ell(\lambda_1, \dots, \lambda_d) - \frac{1}{\lambda_1 + \dots + \lambda_d} \right| \\ &\leq \sqrt{\|\mathcal{A}\|_2} \|b\|_2 \min_{\ell \in \mathbb{L}_{\mathfrak{K}}^{\otimes}} \max_{\lambda_s \in [\alpha_s, \beta_s]} \left| \ell(\lambda_1, \dots, \lambda_d) - \frac{1}{\lambda_1 + \dots + \lambda_d} \right|, \end{aligned} \quad (38)$$

where $[\alpha_s, \beta_s] = [\lambda_{\min}(A_s), \lambda_{\max}(A_s)]$. To proceed further, we need to address the multivariate Laurent polynomial approximation problem

$$\tilde{E}_{\Omega} := \min_{\ell \in \mathbb{L}_{\mathfrak{K}}^{\otimes}} \left\| \ell(\lambda_1, \dots, \lambda_d) - \frac{1}{\lambda_1 + \dots + \lambda_d} \right\|_{\Omega}, \quad (39)$$

where $\|\cdot\|_{\Omega}$ denotes the supremum norm on $\Omega := [\alpha_1, \beta_1] \times \dots \times [\alpha_d, \beta_d]$.

Lemma 6.1 Consider \tilde{E}_Ω defined in (39) for constant $\alpha_s \equiv \alpha$, $\beta_s \equiv \beta$, and $\mathfrak{K} = (k, \dots, k)$. Then

$$\tilde{E}_\Omega \leq \left(1 + \frac{\beta}{\alpha}\right) \frac{1}{\tilde{\alpha}} \frac{\sqrt{\tilde{\kappa}} + 1}{\sqrt{\tilde{\kappa}}} \cdot \left(\frac{\sqrt{\tilde{\kappa}} - 1}{\sqrt{\tilde{\kappa}} + 1}\right)^k,$$

where $\tilde{\alpha} = d(d-1)^{\frac{1-d}{d}} \alpha^{\frac{d-1}{d}} \beta^{\frac{1}{d}}$, $\tilde{\beta} = \alpha + \beta$, and $\tilde{\kappa} = (d-1)/d + \tilde{\beta}/(d\tilde{\alpha})$.

Proof. By straightforward algebraic manipulation,

$$\frac{1}{\lambda_1 + \dots + \lambda_d} = \frac{1 + \gamma \prod_{s=1}^d \lambda_s^{-1}}{\mu_1 + \dots + \mu_d} \quad (40)$$

with $\mu_s = \lambda_s + \gamma \prod_{j \neq s} \lambda_j^{-1}$. Setting $\gamma = \alpha^{d-1} \beta$, we obtain

$$\mu_s \geq \min_{\lambda \in [\alpha, \beta]} \lambda + \gamma \lambda^{-d+1} = d(d-1)^{\frac{1-d}{d}} \alpha^{\frac{d-1}{d}} \beta^{\frac{1}{d}} = \tilde{\alpha}$$

and

$$\mu_s \leq \max_{\lambda \in [\alpha, \beta]} \lambda + \gamma \lambda^{-d+1} = \alpha + \beta = \tilde{\beta}$$

By Lemma A.1 there is a multivariate polynomial $p(\mu_1, \dots, \mu_d) \in \Pi_{\mathfrak{K}}^\otimes$ such that

$$E_{\tilde{\Omega}} := \left\| p(\mu_1, \dots, \mu_d) - \frac{1}{\mu_1 + \dots + \mu_d} \right\|_{\tilde{\Omega}} \leq \frac{1}{\tilde{\alpha}} \frac{\sqrt{\tilde{\kappa}} + 1}{\sqrt{\tilde{\kappa}}} \cdot \left(\frac{\sqrt{\tilde{\kappa}} - 1}{\sqrt{\tilde{\kappa}} + 1}\right)^k,$$

where $\tilde{\Omega} = [\tilde{\alpha}, \tilde{\beta}]^d$ and $\tilde{\kappa}$ is defined as in the statement of the lemma. By (40), the multivariate Laurent polynomial

$$\ell(\lambda_1, \dots, \lambda_d) := \left(1 + \gamma \prod_{s=1}^d \lambda_s^{-1}\right) p(\mu_1, \dots, \mu_d) \in \mathbb{L}_{\mathfrak{K}}^\otimes$$

satisfies

$$\tilde{E}_\Omega \leq \left\| \ell(\lambda_1, \dots, \lambda_d) - \frac{1}{\lambda_1 + \dots + \lambda_d} \right\|_{\Omega} \leq \left\| 1 + \gamma \prod_{s=1}^d \lambda_s^{-1} \right\|_{\Omega} E_{\tilde{\Omega}} = \left(1 + \frac{\beta}{\alpha}\right) E_{\tilde{\Omega}},$$

which concludes the proof. \square

The factor $\tilde{\kappa}$ that determines the asymptotics of the convergence bound in Lemma 6.1 takes the form

$$\tilde{\kappa} \approx \frac{\tilde{\beta}}{d\tilde{\alpha}} = \frac{\alpha + \beta}{d^2(d-1)^{\frac{1-d}{d}} \alpha^{\frac{d-1}{d}} \beta^{\frac{1}{d}}} \approx \frac{1}{d^2} (d-1)^{\frac{d-1}{d}} \kappa(\mathcal{A})^{\frac{d-1}{d}}$$

for larger β/α . In particular for $d = 2$, when solving Sylvester equations, $\tilde{\kappa} \approx \sqrt{\kappa(\mathcal{A})}/4$. This compares favorably with the factor $\kappa = \kappa(\mathcal{A})/2$ that determines the asymptotics of the convergence bound (36) for the (standard) tensor Krylov subspace method. The linear convergence rate of the extended Krylov subspace method for solving Sylvester equations is therefore bounded by $\sqrt{\tilde{\kappa}} \approx \kappa(\mathcal{A})^{1/4}/2$, which was also observed experimentally in [20]. For large d , the bound of Lemma 6.1 becomes less favorable: $\tilde{\kappa} \xrightarrow{d \rightarrow \infty} \kappa(\mathcal{A})/d$. It is not clear to us whether the bound of Lemma 6.1 could be improved to also reflect the significantly better convergence of the extended tensor Krylov subspace method we observed for higher dimensions.

7 Numerical experiments

We have implemented the tensor Krylov subspace and extended tensor Krylov subspace methods in MATLAB, using the Tensor Toolbox [1, 2] for storing tensors in CP decomposition and for multiplying matrices with tensor. For solving the compressed systems, the coefficients described in Section 5 are used. A tolerance of $\varepsilon = 10^{-9}$ is chosen as an upper bound on the accuracy of the solution to the compressed system.

7.1 Symmetric example

As a first example, we consider the Poisson equation in d dimensions:

$$\begin{aligned} -\Delta u &= f & \text{in } \Omega &= [0, 1]^d \\ u &= 0 & \text{on } \Gamma &:= \partial\Omega, \end{aligned}$$

with separable right-hand side $f(y_1, y_2, \dots, y_d) = g(y_1)g(y_2)\cdots g(y_d)$. A standard finite-difference discretization on equidistant nodes leads to the linear system $\mathcal{A}x = b$, with

$$A_s = \frac{1}{h^2} \begin{pmatrix} 2 & -1 & & & \\ -1 & 2 & -1 & & \\ & & \ddots & \ddots & \ddots \\ & & & -1 & 2 & -1 \\ & & & & -1 & 2 \end{pmatrix}, \quad b_s = \begin{pmatrix} g(z_1^{(s)}) \\ \vdots \\ g(z_n^{(s)}) \end{pmatrix}.$$

It is well-known that the minimal and maximal eigenvalues of A_s are given by

$$\lambda_{\min} = \frac{2}{h^2} \left(1 - \cos \left(\frac{\pi}{n_s + 1} \right) \right), \quad \lambda_{\max} = \frac{2}{h^2} \left(1 - \cos \left(\frac{\pi n_s}{n_s + 1} \right) \right),$$

resulting in

$$\kappa = 1 + \frac{\lambda_{\max} - \lambda_{\min}}{d\lambda_{\min}} = 1 + \frac{2 \cos\left(\frac{\pi}{n+1}\right)}{d(1 - \cos\left(\frac{\pi}{n+1}\right))} \approx \frac{4(n+1)^2}{\pi^2 d}. \quad (41)$$

For simplicity, we have used right-hand side vectors b_s composed of uniformly distributed pseudo-random numbers.

The tensor Krylov subspace method was used with multi-index $\mathfrak{K} = (k, \dots, k)$, as the size n_s and condition number are identical for all A_s . All convergence plots display the convergence of the *relative residual* $\frac{\|Ax_k - b\|_2}{\|b\|_2}$. As the solution x_k cannot be stored explicitly, the norm of the residual must be calculated directly from its CP decomposition. For this purpose, an efficient method is proposed in Section 3.3. As the square of the norm is calculated, the residual will only be accurate to a precision of about 10^{-8} .

Figures 3 and 4 show the convergence of the tensor Krylov subspace method for systems of size $n_s = 200$ and $n_s = 1000$, respectively, and various dimensions d . The observed convergence rates correspond reasonably well to the theoretically predicted convergence rates; in particular, the convergence rate improve for higher dimensions. The plots in Figures 3 and 4 are remarkably similar apart from the different scaling of the x -axis. The convergence curves also nicely confirm the fact that $x_k = x$ holds for $k = n_s$ in exact arithmetic.

In Figure 5, we apply the extended Krylov subspaces method to a system of size $n = 200$. At $k = 40$, the method converges within working precision for all dimensions. The observed

convergence is significantly better than the convergence rate $(\frac{\sqrt{\kappa}-1}{\sqrt{\kappa}+1})^k$ predicted by Lemma 6.1. We suspect that this can be attributed to the fact that eigenvalues converge at both ends of the spectrum rather quickly in the course of the extended Krylov subspace method, leading to a rapid decrease of the effective condition number in the course of the iteration. Interestingly, the observed convergence does not improve when going from 5 to 10 dimensions, as predicted by the theoretical convergence rate.

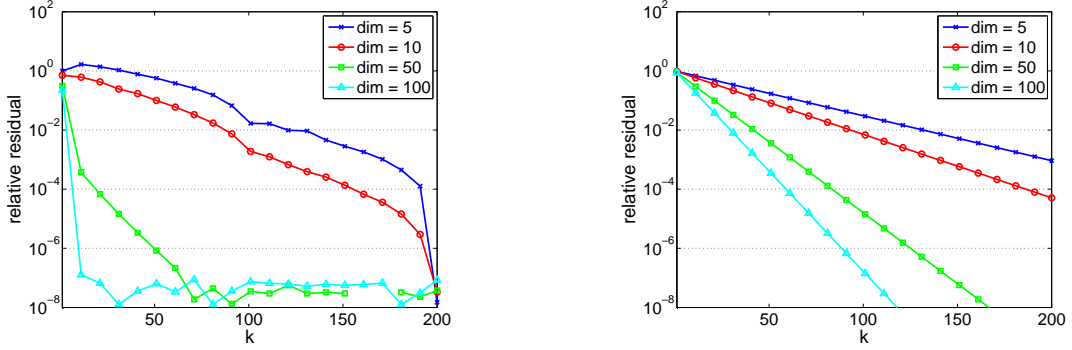


Figure 3: Relative residual of the tensor Krylov subspace method for symmetric example and $n = 200$ (left). Predicted convergence rate $(\frac{\sqrt{\kappa}-1}{\sqrt{\kappa}+1})^k$ (right).

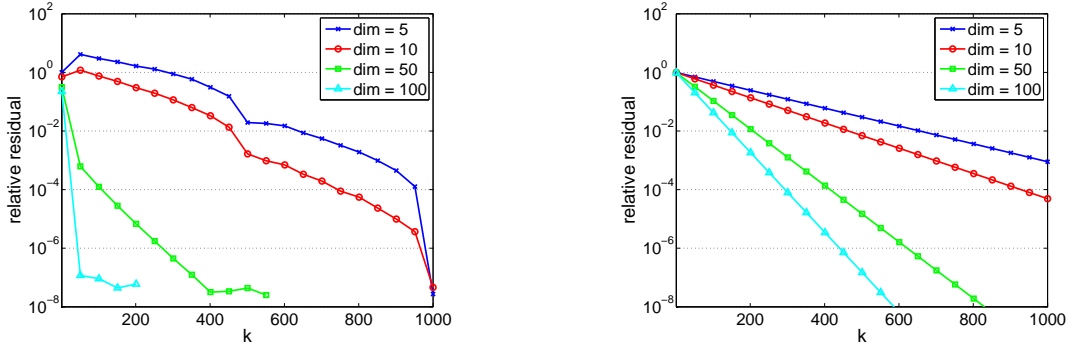


Figure 4: Relative residual of the tensor Krylov subspace method for symmetric example and $n = 1000$ (left). Predicted convergence rate $(\frac{\sqrt{\kappa}-1}{\sqrt{\kappa}+1})^k$ (right).

7.2 Non-symmetric example

We next consider the convection-diffusion equation

$$\begin{aligned} -\Delta u + c^\top \nabla u &= f \quad \text{in } \Omega = [0, 1]^d \\ u &= 0 \quad \text{on } \Gamma := \partial\Omega, \end{aligned}$$

where f is again a separable function. A standard finite-difference discretization on equidistant nodes, combined now with a second order convergent scheme for the convection term,

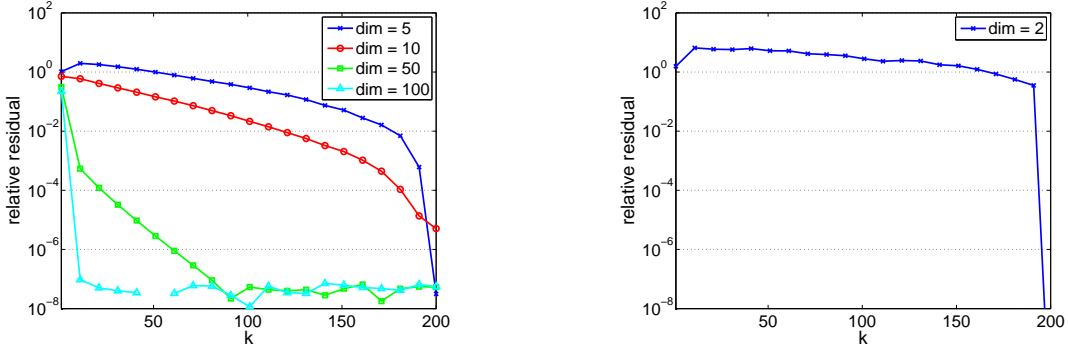


Figure 6: Relative residual for system size 200, with factor $c = 10$ (left) and $c = 100$ (right)

number to its square root. However, it is not clear how the exact inverses A_s^{-1} can be replaced by a preconditioner for A_s ; we observed experimentally that the naive way completely destroys the convergence of the method.

Commutativity The commutativity of the summands of \mathcal{A} is crucial for the theoretical and algorithmic developments of this paper. It is not clear whether the tensor Krylov subspace method can be generalized in a meaningful way to an arbitrary sum $\mathcal{A} = \mathcal{A}_1 + \dots + \mathcal{A}_d$ with $\mathcal{A}_s \mathcal{A}_t = \mathcal{A}_t \mathcal{A}_s$ for all $1 \leq s \leq t \leq d$.

Extension to other Kronecker structures It is important to emphasize that the scope of PDEs that can be addressed by our method is rather restrictive; essentially the domain as well as the coefficients must be separable in all or some of the space variables. Possibly the most promising direction for future research is to explore how a wider scope of high-dimensional PDEs can be addressed, for example when the coefficients are not separable but are represented/approximated by a short sum of separable functions.

Minimal residual methods A method that select the element from the tensorized Krylov subspace which minimizes the norm of the residual could be an attractive alternative to the FOM-like method described in this paper. However, already for $d = 2$ such a MINRES-like method is difficult to realize efficiently [16]. It can be expected that this will be even more difficult for larger d .

9 Acknowledgments

The authors thank Lars Grasedyck for helpful and inspiring discussions on the subject of this paper.

A Appendix

The following lemma contains the approximation result needed in Section 4 for the convergence analysis of the tensor Krylov subspace method.

Lemma A.1 For the multivariate polynomial approximation error $E_\Omega(\mathfrak{K})$ defined in (34) it holds that

$$E_\Omega(\mathfrak{K}) \leq \frac{1}{\lambda_{\min}(\mathcal{A})} \sum_{s=1}^d \frac{\sqrt{\kappa_s} + 1}{\sqrt{\kappa_s}} \cdot \left(\frac{\sqrt{\kappa_s} - 1}{\sqrt{\kappa_s} + 1} \right)^{k_s}$$

where $\kappa_s = 1 + \frac{\beta_s - \alpha_s}{\lambda_{\min}(\mathcal{A})}$.

Proof. Since \mathcal{A} is assumed to be symmetric positive definite, $\alpha_1 + \dots + \alpha_d > 0$ and hence every $(\lambda_1, \dots, \lambda_d) \in \Omega$ also satisfies $\lambda_1 + \dots + \lambda_d > 0$. This allows us to write

$$\frac{1}{\lambda_1 + \dots + \lambda_d} = \int_0^\infty e^{-(\lambda_1 + \dots + \lambda_d)t} dt = \int_0^\infty \prod_{s=1}^d e^{-\lambda_s t} dt.$$

Following [9, 21], we expand $e^{-\lambda_s t}$ in terms of Chebyshev series:

$$e^{-\lambda_s t} = e^{-\frac{\alpha_s + \beta_s}{2}t} \sum_{j=0}^{\infty} a_j^{(s)}(t) T_j(-g_s(\lambda_s)), \quad (42)$$

where T_j denotes the j th Chebyshev polynomial of the first kind, and the parameter transformations $g_s : [\alpha_s, \beta_s] \rightarrow [-1, 1]$ are defined as

$$g_s(\lambda_s) = \frac{2}{\beta_s - \alpha_s} \lambda_s - \frac{\beta_s + \alpha_s}{\beta_s - \alpha_s}$$

The Chebyshev coefficients $a_j^{(s)}$ take the form

$$a_j^{(s)}(t) = \begin{cases} I_j\left(\frac{\beta_s - \alpha_s}{2}t\right), & j = 0, \\ 2I_j\left(\frac{\beta_s - \alpha_s}{2}t\right), & \text{otherwise,} \end{cases}$$

where I_j denotes the j th modified Bessel function. We define polynomials $p_{k_s}^{(s)}(\lambda_s, t)$ in λ_s by truncating the series (42) after the $(k_s - 1)$ th term. Obviously,

$$\|p_{k_s}^{(s)}(\lambda_s, t) - e^{-\lambda_s t}\|_\Omega \leq e^{-\frac{\alpha_s + \beta_s}{2}t} \sum_{j=k_s}^{\infty} a_j^{(s)}(t). \quad (43)$$

Moreover,

$$\|p_{k_s}^{(s)}(\lambda_s, t)\|_\Omega \leq |p_{k_s}^{(s)}(\alpha_s, t)| = e^{-\frac{\alpha_s + \beta_s}{2}t} \sum_{j=0}^{k_s-1} a_j^{(s)}(t) \leq e^{-\frac{\alpha_s + \beta_s}{2}t} \sum_{j=0}^{\infty} a_j^{(s)}(t) = e^{-\alpha_s t}.$$

Hence,

$$\begin{aligned} \left\| \prod_{s=1}^d e^{-\lambda_s t} - \prod_{s=1}^d p_{k_s}^{(s)}(\lambda_s, t) \right\|_\Omega &= \left\| \sum_{s=1}^d \left[\prod_{j=1}^{s-1} e^{-\lambda_j t} \right] (e^{-\lambda_s t} - p_{k_s}^{(s)}(\lambda(s), t)) \left[\prod_{j=s+1}^d p_{k_j}^{(j)}(\lambda_j, t) \right] \right\|_\Omega \\ &\leq \sum_{s=1}^d \|e^{-\lambda_s t} - p_{k_s}^{(s)}(\lambda_s, t)\|_\Omega \prod_{\substack{j=1 \\ j \neq s}}^d e^{-\alpha_j t} \\ &= \sum_{s=1}^d \|e^{-\lambda_s t} - p_{k_s}^{(s)}(\lambda_s, t)\|_\Omega e^{-(\lambda_{\min}(\mathcal{A}) - \alpha_s)t} \end{aligned} \quad (44)$$

To use these results for our multivariate interpolation problem, we set

$$p(\lambda_1, \dots, \lambda_s) = \sum_{\mathfrak{J} \leq \mathfrak{K}} c_{\mathfrak{J}} p_{i_1}^{(1)}(\lambda_1) p_{i_2}^{(2)}(\lambda_2) \cdots p_{i_d}^{(d)}(\lambda_d),$$

where $\mathfrak{J} = (i_1, \dots, i_d)$ with $1 \leq i_s \leq k_s$, and

$$p_{i_s}^{(s)}(\lambda_s) = T_{i_s-1}(-g_s(\lambda_s)), \quad c_{\mathfrak{J}} = \int_0^\infty \prod_{s=1}^d e^{-\frac{\alpha_s + \beta_s}{2} t} a_{i_s-1}^{(s)}(t) dt.$$

With this choice,

$$\begin{aligned} E_{\Omega}(\mathfrak{K}) &\leq \left\| p(\lambda_1, \dots, \lambda_s) - \frac{1}{\lambda_1 + \dots + \lambda_s} \right\|_{\Omega} \\ &= \left\| \sum_{\mathfrak{J} \leq \mathfrak{K}} c_{\mathfrak{J}} \prod_{s=1}^d p_{i_s}^{(s)}(\lambda_s) - \frac{1}{\lambda_1 + \dots + \lambda_s} \right\|_{\Omega} \\ &= \left\| \int_0^\infty \sum_{\mathfrak{J} \leq \mathfrak{K}} \prod_{s=1}^d e^{-\frac{\alpha_s + \beta_s}{2} t} a_{i_s-1}^{(s)}(t) p_{i_s}^{(s)}(\lambda_s) dt - \int_0^\infty \prod_{s=1}^d e^{-\lambda_s t} dt \right\|_{\Omega} \\ &= \left\| \int_0^\infty \underbrace{\prod_{s=1}^d \sum_{i_s=1}^{k_s} e^{-\frac{\alpha_s + \beta_s}{2} t} a_{i_s-1}^{(s)}(t) p_{i_s}^{(s)}(\lambda_s)}_{=: p_{k_s}^{(s)}(\lambda_s, t)} dt - \int_0^\infty \prod_{s=1}^d e^{-\lambda_s t} dt \right\|_{\Omega}. \\ &\leq \int_0^\infty \left\| \prod_{s=1}^d p_{k_s}^{(s)}(\lambda_s, t) - \prod_{s=1}^d e^{-\lambda_s t} \right\|_{\Omega} dt \\ (44) \Rightarrow &\leq \sum_{s=1}^d \int_0^\infty e^{-(\lambda_{\min}(\mathcal{A}) - \alpha_s)t} \left\| e^{-\lambda_s t} - p_{k_s}^{(s)}(\lambda_s, t) \right\|_{\Omega} dt \\ (43) \Rightarrow &\leq \underbrace{\sum_{s=1}^d \sum_{j=k_s}^{\infty} \int_0^\infty e^{-(\lambda_{\min}(\mathcal{A}) + \frac{\beta_s - \alpha_s}{2} t)} a_j^{(s)}(t) dt}_{=: E_s(k_s)}. \end{aligned} \tag{45}$$

Using $\int_0^\infty e^{-\delta t} I_j(\gamma t) dt = \frac{\gamma^j}{\sqrt{\delta^2 - \gamma^2}(\delta + \sqrt{\delta^2 - \gamma^2})^j}$, which holds for $\delta > \gamma$ [21] and $j \neq 0$, we obtain after some algebraic manipulation

$$\int_0^\infty e^{-(\lambda_{\min}(\mathcal{A}) + \frac{\beta_s - \alpha_s}{2} t)} 2I_j\left(\frac{\beta_s - \alpha_s}{2} t\right) dt = \frac{2}{\lambda_{\min}(\mathcal{A})\sqrt{\kappa_s}} \cdot \left(\frac{\sqrt{\kappa_s} - 1}{\sqrt{\kappa_s} + 1}\right)^j$$

with $\kappa_s = 1 + \frac{\beta_s - \alpha_s}{\lambda_{\min}(\mathcal{A})}$. Hence,

$$E_s(k_s) = \frac{2}{\lambda_{\min}(\mathcal{A})\sqrt{\kappa_s}} \sum_{j=k_s}^{\infty} \left(\frac{\sqrt{\kappa_s} - 1}{\sqrt{\kappa_s} + 1}\right)^j = \frac{\sqrt{\kappa_s} + 1}{\lambda_{\min}(\mathcal{A})\sqrt{\kappa_s}} \cdot \left(\frac{\sqrt{\kappa_s} - 1}{\sqrt{\kappa_s} + 1}\right)^{k_s}$$

which, together with (45), concludes the proof. \square

References

- [1] B. W. Bader and T. G. Kolda. Efficient MATLAB computations with sparse and factored tensors. Technical Report SAND2006-7592, Sandia National Laboratories, Albuquerque, NM and Livermore, CA, December 2006.
- [2] B. W. Bader and T. G. Kolda. Matlab tensor toolbox version 2.2, January 2007. Available from <http://csmr.ca.sandia.gov/~tgkolda/TensorToolbox/>.
- [3] B. W. Bader and T. G. Kolda. Tensor decompositions and applications. *Preprint of article to appear in SIAM Review*, 2009/06.
- [4] B. Beckermann and L. Reichel. Error estimation and evaluation of matrix functions via the Faber transform. Technical report, Université de Lille, 2008.
- [5] G. Beylkin and M. J. Mohlenkamp. Algorithms for numerical analysis in high dimensions. *SIAM J. Sci. Comput.*, 26(6):2133–2159, 2005.
- [6] S. Börm. \mathcal{H}_2 -matrices – an efficient tool for the treatment of dense matrices. Habilitationsschrift, Christian-Albrechts-Universität zu Kiel, 2006.
- [7] D. Braess and W. Hackbusch. Approximation of $1/x$ by exponential sums in $[1, \infty)$. *IMA J. Numer. Anal.*, 25(4):685–697, 2005.
- [8] V. Druskin and L. Knizhnerman. Extended Krylov subspaces: approximation of the matrix square root and related functions. *SIAM J. Matrix Anal. Appl.*, 19(3):755–771, 1998.
- [9] V. L. Druskin and L. A. Knizhnerman. Two polynomial methods for calculating functions of symmetric matrices. *Zh. Vychisl. Mat. i Mat. Fiz.*, 29(12):1763–1775, 1989.
- [10] L. Grasedyck. Existence and computation of low Kronecker-rank approximations for large linear systems of tensor product structure. *Computing*, 72(3-4):247–265, 2004.
- [11] L. Grasedyck, W. Hackbusch, and B. N. Khoromskij. Solution of large scale algebraic matrix Riccati equations by use of hierarchical matrices. *Computing*, 70(2):121–165, 2003.
- [12] W. Hackbusch. Approximation of $1/x$ by exponential sums. Available from http://www.mis.mpg.de/scicomp/EXP_SUM/1_x/tabelle.
- [13] W. Hackbusch. Entwicklungen nach Exponentialsummen, 2008. Technical report, see <http://www.mis.mpg.de/preprints/tr/report-0405.pdf>.
- [14] M. Heyouni. Extended Arnoldi methods for large Sylvester matrix equations. Technical report L.M.P.A., 2008.
- [15] R. A. Horn and C. R. Johnson. *Topics in Matrix Analysis*. Cambridge University Press, Cambridge, 1991.
- [16] D. Y. Hu and L. Reichel. Krylov-subspace methods for the Sylvester equation. *Linear Algebra Appl.*, 172:283–313, 1992.

- [17] I. M. Jaimoukha and E. M. Kasenally. Krylov subspace methods for solving large Lyapunov equations. *SIAM J. Numer. Anal.*, 31:227–251, 1994.
- [18] L. Knizhnerman and V. Simoncini. A new investigation of the extended Krylov subspace method for matrix function evaluations. Technical report, 2008.
- [19] Y. Saad. Numerical solution of large Lyapunov equations. In *Signal processing, scattering and operator theory, and numerical methods (Amsterdam, 1989)*, volume 5 of *Progr. Systems Control Theory*, pages 503–511. Birkhäuser Boston, Boston, MA, 1990.
- [20] V. Simoncini. A new iterative method for solving large-scale Lyapunov matrix equations. *SIAM J. Sci. Comput.*, 29(3):1268–1288, 2007.
- [21] V. Simoncini and V. Druskin. Convergence analysis of projection methods for the numerical solution of large Lyapunov equations. *SIAM Journal on Numerical Analysis*, 47(2):828–843, 2009.
- [22] G. Starke. Field-of-values analysis of preconditioned iterative methods for nonsymmetric elliptic problems. *Numer. Math.*, 78(1):103–117, 1997.
- [23] F. Stenger. *Numerical methods based on sinc and analytic functions*, volume 20 of *Springer Series in Computational Mathematics*. Springer-Verlag, New York, 1993.
- [24] G. W. Stewart. *Matrix Algorithms. Vol. II*. SIAM, Philadelphia, PA, 2001. Eigensystems.
- [25] C. Tobler. Krylov subspace methods for large linear systems with tensor product structure. Master’s thesis, ETH Zürich, September 2008.