

Original articles

Analysis of inexact Krylov subspace methods for approximating the matrix exponential

Khanh N. Dinh, Roger B. Sidje*

Department of Mathematics, The University of Alabama, Tuscaloosa, AL 35487, USA

Received 26 September 2015; received in revised form 24 May 2016; accepted 2 January 2017

Available online xxxx

Abstract

Krylov subspace methods have proved quite effective at approximating the action of a large sparse matrix exponential on a vector. Their numerical robustness and matrix-free nature have enabled them to make inroads into a variety of applications. A case in point is solving the chemical master equation (CME) in the context of modeling biochemical reactions in biological cells. This is a challenging problem that gives rise to an extremely large matrix due to the curse of dimensionality. Inexact Krylov subspace methods that build on truncation techniques have helped solve some CME models that were considered computationally out of reach as recently as a few years ago. However, as models grow, truncating them means using an even smaller fraction of their whole extent, thereby introducing more inexactness. But experimental evidence suggests an apparent success and the aim of this study is to give theoretical insights into the reasons why. Essentially, we show that the truncation can be put in the framework of inexact Krylov methods that relax matrix–vector products and compute them expediently by trading accuracy for speed. This allows us to analyze both the residual (or defect) and the error of the resulting approximations to the matrix exponential from the viewpoint of inexact Krylov methods.

© 2017 International Association for Mathematics and Computers in Simulation (IMACS). Published by Elsevier B.V. All rights reserved.

Keywords: Matrix exponential; Inexact Krylov method; Chemical master equation

1. Introduction

Given a large sparse nonsymmetric matrix $A \in \mathbb{R}^{n \times n}$ and vector $p_0 \in \mathbb{R}^n$, letting $v = p_0$ and taking $m \ll n$ Arnoldi steps with a starting vector $v_1 = v/\|v\|$, where $\|\cdot\|$ means the 2-norm, we obtain an orthonormal basis $V_m = [v_1, \dots, v_m] \in \mathbb{R}^{n \times m}$ of the Krylov subspace $\mathcal{K}_m(A, v) = \text{span}\{v, Av, \dots, A^{m-1}v\}$, and an upper Hessenberg matrix $H_m \in \mathbb{R}^{m \times m}$ that satisfy

$$AV_m = V_{m+1}\bar{H}_m = V_m H_m + h_{m+1,m} v_{m+1} e_m^T, \quad (1a)$$

$$H_m = V_m^T A V_m, \quad (1b)$$

* Corresponding author.

E-mail addresses: kdinh@crimson.ua.edu (K.N. Dinh), roger.b.sidje@ua.edu (R.B. Sidje).

where $\mathbf{e}_m = (0, \dots, 0, 1)^T$, and $\bar{\mathbf{H}}_m \in \mathbb{R}^{(m+1) \times m}$ is \mathbf{H}_m augmented with $h_{m+1,m} \mathbf{e}_m^T$ under its last row. The standard Krylov approximation to the matrix exponential takes the form

$$\exp(\tau \mathbf{A})\mathbf{v} \approx \mathbf{V}_m \exp(\tau \mathbf{H}_m) \beta \mathbf{e}_1, \quad \mathbf{e}_1 = (1, 0, \dots, 0)^T, \quad \beta = \|\mathbf{v}\|. \tag{2}$$

It is well-known that (1) is also the cornerstone for building very efficient Krylov subspace solution techniques for other problems such as eigenvalue problems or linear systems. In the latter, there has been recent interest in transitioning from exact to inexact (or relaxed) matrix–vector products in the Arnoldi process [3,4,15], either out of necessity or deliberately, trading accuracy for speed. It is customary to model these inexact products as

$$\mathbf{A}\mathbf{v}_k \approx (\mathbf{A} + \mathbf{E}_k)\mathbf{v}_k, \tag{3}$$

where \mathbf{E}_k is some error matrix that varies at each invocation, and note that setting $\mathbf{E}_k = \mathbf{0}$ recovers the exact evaluation. To make the difference clear, we refer to the classical method as the exact Arnoldi and it is not meant to imply exact arithmetic. The foremost implication of such a relaxation is that the classical Arnoldi relationship (1) does not hold anymore, but Simoncini and Szyld [15] made the key observation that we end up with

$$(\mathbf{A} + \mathcal{E}_m)\mathbf{V}_m = \mathbf{V}_m \mathbf{H}_m + h_{m+1,m} \mathbf{v}_{m+1} \mathbf{e}_m^T, \quad \mathcal{E}_m = \sum_{k=1}^m \mathbf{E}_k \mathbf{v}_k \mathbf{v}_k^T, \tag{4}$$

which is similar to (1), except that \mathbf{V}_m , which still remains orthonormal, is now a basis of a Krylov subspace obtained by a perturbed \mathbf{A} . When we use the computed \mathbf{V}_m and \mathbf{H}_m from (4) in GMRES for instance, classical error bounds do not apply anymore. However, from theoretical and experimental evidence (such as [14]), the method can withstand cases where the norm of the perturbation \mathcal{E}_m grows quite large.

The analysis of Simoncini and Szyld [15] provided insights into inexact GMRES for solving a linear system $\mathbf{A}\mathbf{x} = \mathbf{b}$, but it has so far remained unclear how inexactness affects the Krylov approximation (2). Since we now have (4) instead of (1), we also lose classical error bounds on the matrix exponential (e.g., Gallopoulos and Saad [5], Saad [11], Hochbruck and Lubich [9]). Thus our study fills a gap in the literature by looking at the error in the inexact Krylov counterpart of (2). We additionally offer another related way of assessing the accuracy by investigating the defect or residual [2] from the fact that (2) arises when solving a system of linear ODEs of the form

$$\begin{cases} \mathbf{p}'(t) = \mathbf{A}\mathbf{p}, & t \in [0, t_f] \\ \mathbf{p}(0) = \mathbf{p}_0, & \text{initial condition.} \end{cases} \tag{5}$$

It is worth recalling that, in the exact case, the effectiveness of approximating $\exp(\mathbf{A})\mathbf{v}$ by projecting it onto $\mathcal{K}_m(\mathbf{A}, \mathbf{v})$ hinges on the fact that all polynomials of \mathbf{A} of degree $\leq m - 1$ can be calculated exactly through \mathbf{H}_m , or more precisely,

$$q_{m-1}(\mathbf{A})\mathbf{v} = \mathbf{V}_m q_{m-1}(\mathbf{H}_m) \beta \mathbf{e}_1,$$

where q_{m-1} is any polynomial of degree $\leq m - 1$. By the same reasoning as in the exact case (e.g., Saad [11, Lemma 3.1]), it can be shown from (4) that, for the same polynomial q_{m-1} ,

$$q_{m-1}(\mathbf{A} + \mathcal{E}_m)\mathbf{v} = \mathbf{V}_m q_{m-1}(\mathbf{H}_m) \beta \mathbf{e}_1.$$

The significance of all this is that the inexact Krylov subspace method for $\mathcal{K}_m(\mathbf{A}, \mathbf{v})$ with the relaxation matrices $\mathbf{E}_k, k = 1, \dots, m$, can be seen as the exact Krylov subspace method for $\mathcal{K}_m(\tilde{\mathbf{A}}, \mathbf{v})$ with

$$\tilde{\mathbf{A}} = \mathbf{A} + \mathcal{E}_m = \mathbf{A} + \sum_{k=1}^m \mathbf{E}_k \mathbf{v}_k \mathbf{v}_k^T,$$

which is another simple way to understand the method.

The rest of the paper is organized as follows: Section 2 gives some background on modeling biochemical reactions and on the finite state projection (FSP) algorithm for solving the underlying chemical master equation (CME), which was the challenging application that initially motivated the research presented here. Section 3 analyzes the residual (or defect) of the inexact Krylov method both when the ODE problem is homogeneous or nonhomogeneous, with two different approaches considered for the latter. Section 4 analyzes the error and includes a special treatment that

exploits the structure of the matrix when it arises from stochastic processes such as that involved in the CME. Section 5 reports some numerical experiments. Section 6 finally wraps the presentation with some concluding remarks.

2. Inexact chemical master equation — a motivation

In a biological cell containing different molecular species undergoing various chemical reactions, the state is a vector of integers counting the different species of molecules. Such a discrete formulation is prompted by key regulatory molecules that exist in small numbers, making a continuous formulation (via concentrations) inappropriate. The counter of a species goes up or down when a chemical reaction occurs, depending on whether the reaction produces or consumes that species. Starting from a particular state, the cell will transition to different states as reactions happen. The CME has the form (5) and depicts the evolution of the system's probability distribution, which characterizes the probability of finding the cell in a given state at a given time. The challenge here is that even with simple biochemical models having 4 or 5 reaction channels and a relatively low count of each molecular species, there can be millions of possible states, and the variety of models means that calculations cannot rely on generic simplifications.

The finite state projection (FSP) algorithm of Munsky and Khammash [10] is a model reduction method to cope with the huge size of the CME, and we outline it here because it vividly illustrates how inexactness comes into play. With $J = \{1, \dots, k\}$, and k being the cardinality of J , let

$$A = \left(\begin{array}{c|c} A_J & * \\ \hline * & * \end{array} \right) \in \mathbb{R}^{n \times n},$$

i.e., A_J is a $k \times k$ submatrix of the true operator A . The FSP algorithm takes

$$p(t_f) = \exp(t_f A) p_0 \approx \begin{pmatrix} \exp(t_f A_J) & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{pmatrix} \begin{pmatrix} p_J(0) \\ \mathbf{0} \end{pmatrix}. \quad (6)$$

Note that $p_0 = p(0)$ is also truncated according to J . Munsky and Khammash [10] assessed the loss of the probability mass and gave a theoretical justification of the merit of this approach from a probabilistic point of view. Our study recasts the analysis in the framework of inexact methods.

With J now an arbitrary subset of $\{1, \dots, n\}$, and the corresponding submatrix A_J padded with zeros as necessary, the 'truncation' can be formalized as

$$A = A_J + R_J,$$

where R_J is the error matrix from the truncation. Then (6) is in turn further approximated by an inexact Krylov method for the matrix exponential, which we saw previously means that the Arnoldi process for approximating the solution uses inexact (or relaxed) matrix–vector products of the form

$$A_J v \approx (A_J + S_J) v,$$

where S_J models some error in the product. Combining with the truncation, we get

$$A v \approx (A - R_J + S_J) v,$$

so that $E = -R_J + S_J$ captures both error terms. In particular, if $S_J = 0$, only the truncation error is in effect, whereas if $R_J = 0$, there will only be the error from the relaxed matrix–vector product. From now on, our theoretical analysis will simply assume that the true matrix A interacts through the inexact evaluation $A v_k \approx (A + E_k) v_k$ without regard to the real source of the error.

3. Bounds on the residual

Given the differential equation (5), consider an approximation $p_m(t) \approx p(t)$, then in the terminology of ODEs the 'residual' or 'defect' is defined as

$$r_m(t) = p_m'(t) - A p_m(t). \quad (7)$$

This definition reminds what happens when approximating the solution to a linear system $\mathbf{Ax} = \mathbf{b}$. Take GMRES, which uses the approximation $\mathbf{x}_m = \mathbf{x}_0 + \mathbf{V}_m \mathbf{y}_m$, where \mathbf{x}_0 is an initial guess, $\mathbf{y}_m = \bar{\mathbf{H}}_m^\dagger \beta \mathbf{e}_1$, with \mathbf{V}_m and $\bar{\mathbf{H}}_m$ arising from the Arnoldi process for $\mathcal{K}_m(\mathbf{A}, \mathbf{r}_0)$, and denoting $r_0 = \mathbf{b} - \mathbf{Ax}_0$, $\beta = \|\mathbf{r}_0\|$, and $\bar{\mathbf{H}}_m^\dagger$ the pseudo-inverse of $\bar{\mathbf{H}}_m$. Applying (4), the inexact GMRES method ends up with (the norm of) the *computed* residual

$$\tilde{\mathbf{r}}_m = \mathbf{r}_0 - \mathbf{V}_{m+1} \bar{\mathbf{H}}_m \mathbf{y}_m, \tag{8}$$

while the *true* residual satisfies

$$\mathbf{r}_m = \mathbf{b} - \mathbf{A}(\mathbf{x}_0 + \mathbf{V}_m \mathbf{y}_m) = \tilde{\mathbf{r}}_m + \mathcal{E}_m \mathbf{V}_m \mathbf{y}_m. \tag{9}$$

From (8) and (9), and the fact that $\mathcal{E}_m \mathbf{V}_m = [\mathbf{E}_1 \mathbf{v}_1, \dots, \mathbf{E}_m \mathbf{v}_m]$ because \mathbf{V}_m is orthonormal, there is an unknown *residual gap* that satisfies

$$\delta_m = \|\mathbf{r}_m - \tilde{\mathbf{r}}_m\| = \|\mathcal{E}_m \mathbf{V}_m \mathbf{y}_m\| = \|[\mathbf{E}_1 \mathbf{v}_1, \dots, \mathbf{E}_m \mathbf{v}_m] \mathbf{y}_m\|. \tag{10}$$

The inexact solver still reports $\|\tilde{\mathbf{r}}_m\|$ as the estimate of the residual, but the residual gap (10) is not obvious, and so the reliability of the final solution is not guaranteed. Simoncini and Szyld [15] obtained the bound (see also [6,14]):

$$\delta_m \leq \sum_{k=1}^m |\eta_k^{(m)}| \cdot \|\mathbf{E}_k \mathbf{v}_k\| \leq \sum_{k=1}^m |\eta_k^{(m)}| \cdot \|\mathbf{E}_k\|, \quad \mathbf{y}_m = \bar{\mathbf{H}}_m^\dagger \beta \mathbf{e}_1 = (\eta_1^{(m)}, \eta_2^{(m)}, \dots, \eta_m^{(m)})^T.$$

Returning to the differential equation (5) where the residual of an approximation is defined by (7), and using (2) based on the exact Arnoldi process, we get the Krylov approximation $\mathbf{p}_m(\tau) = \mathbf{V}_m \exp(\tau \mathbf{H}_m) \beta \mathbf{e}_1$ that leads to

$$\begin{aligned} \mathbf{p}'_m(\tau) - \mathbf{A} \mathbf{p}_m(\tau) &= \mathbf{p}'_m(\tau) - \mathbf{A} \mathbf{V}_m \exp(\tau \mathbf{H}_m) \beta \mathbf{e}_1 \\ &= \mathbf{V}_m \exp(\tau \mathbf{H}_m) \mathbf{H}_m \beta \mathbf{e}_1 - [\mathbf{V}_m \mathbf{H}_m + h_{m+1,m} \mathbf{v}_{m+1} \mathbf{e}_m^T] \exp(\tau \mathbf{H}_m) \beta \mathbf{e}_1 \\ &= -\beta h_{m+1,m} \left(\mathbf{e}_m^T \exp(\tau \mathbf{H}_m) \mathbf{e}_1 \right) \mathbf{v}_{m+1}. \end{aligned}$$

Expokit [12] uses variants of this estimate (with some scaling) to monitor the accuracy. We define the *computed* residual as

$$\tilde{\mathbf{r}}_m(\tau) = -\beta h_{m+1,m} \left(\mathbf{e}_m^T \exp(\tau \mathbf{H}_m) \mathbf{e}_1 \right) \mathbf{v}_{m+1},$$

because this would still be the economical quantity (or a related one from it) used to estimate the *true* residual. As we shall see in the following section, with the inexact Arnoldi process (4), the *true* residual in the ODE problem becomes

$$\mathbf{r}_m(\tau) = \tilde{\mathbf{r}}_m(\tau) + \mathcal{E}_m \mathbf{V}_m \exp(\tau \mathbf{H}_m) \beta \mathbf{e}_1,$$

which is reminiscent of (9), but with $\mathbf{y}_m(\tau) = \exp(\tau \mathbf{H}_m) \beta \mathbf{e}_1$ instead of the least squares solution $\mathbf{y}_m = \bar{\mathbf{H}}_m^\dagger \beta \mathbf{e}_1$. Simoncini and Szyld [15] refined their bounds by exploiting properties satisfied by the components of \mathbf{y}_m through Givens rotations in the case of GMRES. Our analysis will derive bounds without assuming those properties since Givens rotations are not involved in the case of the matrix exponential.

3.1. Homogeneous case

The Krylov technique for solving (5) is typically done by using the integration scheme

$$\begin{cases} \mathbf{p}(0) = \mathbf{p}_0 \\ \mathbf{p}(t_{k+1}) = \exp(\tau_k \mathbf{A}) \mathbf{p}(t_k), \end{cases} \tag{11}$$

with some strategy for choosing the stepsizes $\tau_k = t_{k+1} - t_k$. The problem remains how to effectively approximate $\exp(\tau \mathbf{A}) \mathbf{v}$ given τ and \mathbf{v} .

Now, using (2) based on the inexact Arnoldi process (4), the true residual of the resulting Krylov approximation $\mathbf{p}_m(\tau) = \mathbf{V}_m \exp(\tau \mathbf{H}_m) \beta \mathbf{e}_1$ satisfies

$$\begin{aligned} \mathbf{r}_m &= \mathbf{p}'_m - \mathbf{A}\mathbf{p}_m \\ &= (\mathbf{V}_m \exp(\tau \mathbf{H}_m) \beta \mathbf{e}_1)' - \mathbf{A}\mathbf{V}_m \exp(\tau \mathbf{H}_m) \beta \mathbf{e}_1 \\ &= \mathbf{V}_m \exp(\tau \mathbf{H}_m) \mathbf{H}_m \beta \mathbf{e}_1 - (\mathbf{V}_m \mathbf{H}_m + h_{m+1,m} \mathbf{v}_{m+1} \mathbf{e}_m^T - \mathcal{E}_m \mathbf{V}_m) \exp(\tau \mathbf{H}_m) \beta \mathbf{e}_1 \\ &= -h_{m+1,m} (\mathbf{e}_m^T \exp(\tau \mathbf{H}_m) \beta \mathbf{e}_1) \mathbf{v}_{m+1} + \mathcal{E}_m \mathbf{V}_m \exp(\tau \mathbf{H}_m) \beta \mathbf{e}_1 \\ &= \tilde{\mathbf{r}}_m + \mathcal{E}_m \mathbf{V}_m \exp(\tau \mathbf{H}_m) \beta \mathbf{e}_1, \end{aligned}$$

where $\tilde{\mathbf{r}}_m$ is the computed residual defined earlier. The dependency on the time τ will not be made explicit unless there is a risk of ambiguities. The quantity

$$\delta_m^{\text{res}} = \|\mathbf{r}_m - \tilde{\mathbf{r}}_m\| = \|\mathcal{E}_m \mathbf{V}_m \exp(\tau \mathbf{H}_m) \beta \mathbf{e}_1\| \tag{12}$$

is the residual gap between \mathbf{r}_m and $\tilde{\mathbf{r}}_m$. It depends on the relaxation matrices \mathbf{E}_k and therefore cannot be computed in a straightforward manner. From the derivation above, we can bound the norm of the true residual as

$$\|\mathbf{r}_m\| \leq \|\tilde{\mathbf{r}}_m\| + \delta_m^{\text{res}}. \tag{13}$$

When all matrix–vector products are exact, i.e., all $\mathbf{E}_k = \mathbf{0}$, then $\mathcal{E}_m = \mathbf{0}$ and $\delta_m^{\text{res}} = 0$ as expected. More generally with $\mathbf{y}_m = \exp(\tau \mathbf{H}_m) \beta \mathbf{e}_1$, the residual gap can be written as

$$\delta_m^{\text{res}} = \|\mathcal{E}_m \mathbf{V}_m \mathbf{y}_m\| = \|[\mathbf{E}_1 \mathbf{v}_1, \dots, \mathbf{E}_m \mathbf{v}_m] \mathbf{y}_m\|,$$

and similarly to [15, Prop. 4.1], if we write $\mathbf{y}_m = (\eta_1^{(m)}, \dots, \eta_m^{(m)})^T$, the following upper bound on the residual gap holds:

$$\delta_m^{\text{res}} \leq \sum_{k=1}^m |\eta_k^{(m)}| \cdot \|\mathbf{E}_k\|, \tag{14}$$

and furthermore we have the following result:

Proposition 3.1. Given $\epsilon^{\text{res}} > 0$, if $\|\mathbf{E}_k\| \leq \frac{\epsilon^{\text{res}}}{m\beta \|\exp(\tau \mathbf{H}_m)\|}$, $k = 1, \dots, m$, then we have

$$\delta_m^{\text{res}} \leq \epsilon^{\text{res}} \tag{15}$$

and therefore $\|\mathbf{r}_m\| \leq \|\tilde{\mathbf{r}}_m\| + \epsilon^{\text{res}}$.

Proof. We have $|\eta_k^{(m)}| \leq \|\mathbf{y}_m\| = \|\exp(\tau \mathbf{H}_m) \beta \mathbf{e}_1\| \leq \beta \|\exp(\tau \mathbf{H}_m)\|$, and so combining the condition on $\|\mathbf{E}_k\|$ with the inequality (14) gives (15). The bound on $\|\mathbf{r}_m\|$ then follows naturally from the relation $\|\mathbf{r}_m\| \leq \|\tilde{\mathbf{r}}_m\| + \delta_m^{\text{res}}$. \square

Remark 3.2. As Giraud et al. [6] pointed out in the context of inexact GMRES, Proposition 3.1 does not allow us to anticipate the relaxation matrices \mathbf{E}_k in advance because of the dependence on \mathbf{H}_m . It can however be used in a postmortem manner to check if the error condition is satisfied.

3.2. Nonhomogeneous case

3.2.1. Bounding the residual gap via the φ function

It is well known (see, e.g., Expokit [12]) that the numerical solution to the system of nonhomogeneous ODEs

$$\begin{cases} \mathbf{p}'(t) = \mathbf{A}\mathbf{p}(t) + \mathbf{b} \\ \mathbf{p}(0) = \mathbf{p}_0 \end{cases} \tag{16}$$

with constant $\mathbf{A} \in \mathbb{R}^{n \times n}$ and $\mathbf{b} \in \mathbb{R}^n$, can be found using the integration scheme

$$\begin{cases} \mathbf{p}(0) = \mathbf{p}_0 \\ \mathbf{p}(t_{k+1}) = \tau_k \varphi(\tau_k \mathbf{A}) [\mathbf{A}\mathbf{p}(t_k) + \mathbf{b}] + \mathbf{p}(t_k) \end{cases} \tag{17}$$

where $\varphi(\tau\mathbf{A}) = \sum_{i=0}^{\infty} \frac{(\tau\mathbf{A})^i}{(i+1)!}$ [8, Chap. 2.1, Chap. 10.7]. This integration scheme circumvents using the representation of the analytical solution of (16), $\mathbf{p}(t) = \exp(t\mathbf{A})\mathbf{p}_0 + t\varphi(t\mathbf{A})\mathbf{b}$, which would need both $\mathcal{K}_m(\mathbf{A}, \mathbf{p}_0)$ and $\mathcal{K}_m(\mathbf{A}, \mathbf{b})$ instead of only one Krylov subspace as implied in (17). Using this scheme, we can derive a result similar to Proposition 3.1, but this can be avoided by the augmented approach shown below.

3.2.2. *Bounding the residual gap via an augmented matrix exponential*

An indirect way to solve (16) takes root in the analytical solution $\mathbf{p}(t) = \exp(t\mathbf{A})\mathbf{p}_0 + t\varphi(t\mathbf{A})\mathbf{b}$, and has proved convenient in other circumstances such as [1,13]. Define the augmented matrix

$$\mathbf{A}^+ = \begin{pmatrix} \mathbf{A} & \mathbf{b} \\ \mathbf{0} & 0 \end{pmatrix} \in \mathbb{R}^{(n+1) \times (n+1)},$$

then we have

$$\exp(t\mathbf{A}^+) = \begin{pmatrix} \exp(t\mathbf{A}) & t\varphi(t\mathbf{A})\mathbf{b} \\ \mathbf{0} & 1 \end{pmatrix},$$

so that the solution can be fetched as $\mathbf{p}(t) = [\exp(t\mathbf{A}^+)\mathbf{p}_0^+]_{1:n}$ with $\mathbf{p}_0^+ = \begin{pmatrix} p_0 \\ 1 \end{pmatrix}$.

The problem now amounts to getting $\exp(t\mathbf{A}^+)\mathbf{p}_0^+$. Transforming the problem back to the form $\exp(t\mathbf{A}^+)\mathbf{p}_0^+$ not only inherits the analysis done in the homogeneous case in an elegant way, but also enables seamless code re-use. Furthermore, the inexactness in the matrix–vector product with \mathbf{A}^+ is only triggered from \mathbf{A} through

$$\mathbf{A}^+ + \mathbf{E}_k^+ = \begin{pmatrix} \mathbf{A} + \mathbf{E}_k & \mathbf{b} \\ \mathbf{0} & 0 \end{pmatrix},$$

and so results can nicely be recast in terms of \mathbf{E}_k . For this reason, we will not dwell any further on the nonhomogeneous case in the rest of our presentation.

4. **Bounds on the error**

4.1. *General upper bound on the error*

As pointed out in the introduction, the inexact method for $\mathcal{K}_m(\mathbf{A}, \mathbf{v})$ with the perturbation matrices $\mathbf{E}_k, k = 1, \dots, m$, can be seen as the exact method for $\mathcal{K}_m(\tilde{\mathbf{A}}, \mathbf{v})$ with

$$\tilde{\mathbf{A}} = \mathbf{A} + \mathcal{E}_m = \mathbf{A} + \sum_{k=1}^m \mathbf{E}_k \mathbf{v}_k \mathbf{v}_k^T.$$

Therefore if we define

$$\begin{aligned} \tilde{\epsilon}_m &= \|\exp(\tau\tilde{\mathbf{A}})\mathbf{v} - \mathbf{V}_m \exp(\tau\mathbf{H}_m)\beta\mathbf{e}_1\| \\ &= \|\exp(\tau(\mathbf{A} + \mathcal{E}_m))\mathbf{v} - \mathbf{V}_m \exp(\tau\mathbf{H}_m)\beta\mathbf{e}_1\|, \end{aligned}$$

then bounds on $\tilde{\epsilon}_m$ have been given in the literature of exact Krylov methods [5,11,9]. However, our main focus is not so much on the bounds on $\tilde{\epsilon}_m$, but on the *true* error

$$\epsilon_m = \|\exp(\tau\mathbf{A})\mathbf{v} - \mathbf{V}_m \exp(\tau\mathbf{H}_m)\beta\mathbf{e}_1\|.$$

The relationship between ϵ_m and $\tilde{\epsilon}_m$ is straightforward from the triangle inequality

$$\epsilon_m \leq \|\exp(\tau\mathbf{A})\mathbf{v} - \exp(\tau\tilde{\mathbf{A}})\mathbf{v}\| + \|\exp(\tau\tilde{\mathbf{A}})\mathbf{v} - \mathbf{V}_m \exp(\tau\mathbf{H}_m)\beta\mathbf{e}_1\| = \tilde{\epsilon}_m + \delta_m^{err},$$

where we define the *error gap*

$$\delta_m^{err} = \|\exp(\tau\mathbf{A})\mathbf{v} - \exp(\tau\tilde{\mathbf{A}})\mathbf{v}\|. \tag{18}$$

With upper bounds on $\tilde{\epsilon}_m$ ready in hand, it remains to get good bounds on δ_m^{err} , which turns out to be a matrix perturbation analysis that we discuss next.

4.2. Bounding the error gap

In [8, Chap. 10.2], Higham obtained

$$\exp(\tau\tilde{A}) = \exp(\tau A) + \int_0^\tau \exp((\tau - s)A)\mathcal{E}_m \exp(sA)ds + \mathcal{O}(\|\tau\mathcal{E}_m\|^2). \tag{19}$$

Post-multiplying each side by \mathbf{v} , we get

$$\exp(\tau\tilde{A})\mathbf{v} = \exp(\tau A)\mathbf{v} + \int_0^\tau \exp((\tau - s)A)\mathcal{E}_m \exp(sA)\mathbf{v}ds + \mathcal{O}(\|\tau\mathcal{E}_m\|^2).$$

Hence

$$\delta_m^{err} \leq \int_0^\tau \|\exp((\tau - s)A)\| \|\mathcal{E}_m\| \|\exp(sA)\| \|\mathbf{v}\| ds + \mathcal{O}(\|\tau\mathcal{E}_m\|^2),$$

which leads directly to the following statement.

Theorem 4.1. For any arbitrary A , \mathbf{v} and $\tilde{A} = A + \mathcal{E}_m$ from the inexact Krylov subspace method, we have

$$\delta_m^{err} = \|\exp(\tau A)\mathbf{v} - \exp(\tau\tilde{A})\mathbf{v}\| \leq \beta h_A^2 \|\tau\mathcal{E}_m\| + \mathcal{O}(\|\tau\mathcal{E}_m\|^2),$$

where $\beta = \|\mathbf{v}\|$ and $h_A = \max_{s \in [0, \tau]} \|\exp(sA)\|$ is the so-called ‘hump’ on $[0, \tau]$.

When A originates from a Markov chain as is the case in the CME, (i.e., with nonnegative off-diagonal elements, negative diagonal elements and zero column sums) it is known that $\|\exp(sA)\|_1 = 1, \forall s \geq 0$ (see for example [8, section 2.3]). We can then draw from (19) that

$$\|\exp(\tau A)\mathbf{v} - \exp(\tau\tilde{A})\mathbf{v}\|_1 \leq \int_0^\tau \|\exp((\tau - s)A)\|_1 \|\mathcal{E}_m\|_1 \|\exp(sA)\|_1 \|\mathbf{v}\|_1 ds + \mathcal{O}(\|\tau\mathcal{E}_m\|_1^2).$$

and using in addition the fact that a probability vector has $\|\mathbf{v}\|_1 = 1$, we get the following simplified result.

Theorem 4.2. Let A be the transition rate matrix of a Markov chain, then given a probability vector \mathbf{v} and the perturbed $\tilde{A} = A + \mathcal{E}_m$ from the inexact Krylov subspace method, we have

$$\|\exp(\tau A)\mathbf{v} - \exp(\tau\tilde{A})\mathbf{v}\|_1 \leq \|\tau\mathcal{E}_m\|_1 + \mathcal{O}(\|\tau\mathcal{E}_m\|_1^2).$$

4.3. Series expansion of the error

In [11], Saad derived the following series expansion of the error produced by the exact Krylov subspace method

$$\exp(\tau A)\mathbf{v} - V_m \exp(\tau H_m)\beta \mathbf{e}_1 = \tau h_{m+1,m} \sum_{k=1}^\infty \mathbf{e}_m^T \varphi_k(\tau H_m)\beta \mathbf{e}_1 (\tau A)^{k-1} \mathbf{v}_{m+1}, \tag{20}$$

of which the first term in the series was argued to be a good estimate of the error when the stepsize τ is small enough. Here, the functions φ_k are defined as

$$\varphi_k(x) = \sum_{i=0}^\infty \frac{x^i}{(i+k)!}, \quad \varphi_k(0) = \frac{1}{k!},$$

which implies that

$$0 \leq \varphi_k(x) \leq \frac{\varphi_{k-1}(x)}{k} \leq \dots \leq \frac{e^x}{k!}, \quad \text{if } x \geq 0.$$

In the context of the inexact Krylov subspace method, we can use (4) and substitute A for $\tilde{A} = A + \mathcal{E}_m$ in (20), but that will bring up the unwieldy issue of the error gap again. Instead, we generalize the expansion in the following statement that reveals how the terms involving \mathcal{E}_m break out in a strikingly neat way.

Please cite this article in press as: K.N. Dinh, R.B. Sidje, Analysis of inexact Krylov subspace methods for approximating the matrix exponential, Math. Comput. Simulation (2017), <http://dx.doi.org/10.1016/j.matcom.2017.01.002>

Theorem 4.3. *The (true) error in the inexact Krylov subspace method for the matrix exponential has the series expansion*

$$\begin{aligned} & \exp(\tau\mathbf{A})\mathbf{v} - \mathbf{V}_m \exp(\tau\mathbf{H}_m)\beta\mathbf{e}_1 \\ &= \sum_{k=1}^{\infty} (\tau\mathbf{A})^{k-1} \left[\tau h_{m+1,m} \left(\mathbf{e}_m^T \varphi_k(\tau\mathbf{H}_m)\beta\mathbf{e}_1 \right) \mathbf{v}_{m+1} - \tau \mathcal{E}_m \mathbf{V}_m \varphi_k(\tau\mathbf{H}_m)\beta\mathbf{e}_1 \right]. \end{aligned} \tag{21}$$

Proof. As in the proof of [11, Theorem 5.1], define the error in approximating $\varphi_k(\mathbf{A})\mathbf{v}$ by projecting it onto \mathbf{V}_m as $s_m^k = \varphi_k(\mathbf{A})\mathbf{v} - \mathbf{V}_m \varphi_k(\mathbf{H}_m)\beta\mathbf{e}_1$.

From the definition of φ_k and the fact that $\varphi_k(\mathbf{0})\mathbf{v} = \mathbf{V}_m \varphi_k(\mathbf{0})\beta\mathbf{e}_1$, we directly have

$$\begin{aligned} \varphi_k(\mathbf{A})\mathbf{v} &= \mathbf{A} \varphi_{k+1}(\mathbf{A})\mathbf{v} + \varphi_k(\mathbf{0})\mathbf{v} \\ &= \mathbf{A}[\mathbf{V}_m \varphi_{k+1}(\mathbf{H}_m)\beta\mathbf{e}_1 + s_m^{k+1}] + \varphi_k(\mathbf{0})\mathbf{v} \\ &= \mathbf{V}_m[\varphi_k(\mathbf{0})\beta\mathbf{e}_1 + \mathbf{H}_m \varphi_{k+1}(\mathbf{H}_m)\beta\mathbf{e}_1] + h_{m+1,m} \mathbf{v}_{m+1} \mathbf{e}_m^T \varphi_{k+1}(\mathbf{H}_m)\beta\mathbf{e}_1 \\ &\quad - \mathcal{E}_m \mathbf{V}_m \varphi_{k+1}(\mathbf{H}_m)\beta\mathbf{e}_1 + \mathbf{A} s_m^{k+1} \\ &= \mathbf{V}_m \varphi_k(\mathbf{H}_m)\beta\mathbf{e}_1 + h_{m+1,m} \mathbf{v}_{m+1} \mathbf{e}_m^T \varphi_{k+1}(\mathbf{H}_m)\beta\mathbf{e}_1 - \mathcal{E}_m \mathbf{V}_m \varphi_{k+1}(\mathbf{H}_m)\beta\mathbf{e}_1 + \mathbf{A} s_m^{k+1}, \end{aligned}$$

resulting in another expression for s_m^k through a recurrence

$$s_m^k = h_{m+1,m} \left(\mathbf{e}_m^T \varphi_{k+1}(\mathbf{H}_m)\beta\mathbf{e}_1 \right) \mathbf{v}_{m+1} - \mathcal{E}_m \mathbf{V}_m \varphi_{k+1}(\mathbf{H}_m)\beta\mathbf{e}_1 + \mathbf{A} s_m^{k+1}.$$

Using these expressions of the error terms gives

$$\begin{aligned} \exp(\mathbf{A})\mathbf{v} - \mathbf{V}_m \exp(\mathbf{H}_m)\beta\mathbf{e}_1 &= s_m^0 \\ &= h_{m+1,m} \left(\mathbf{e}_m^T \varphi_1(\mathbf{H}_m)\beta\mathbf{e}_1 \right) \mathbf{v}_{m+1} - \mathcal{E}_m \mathbf{V}_m \varphi_1(\mathbf{H}_m)\beta\mathbf{e}_1 + \mathbf{A} s_m^1 = \dots \\ &= h_{m+1,m} \sum_{k=1}^j \left(\mathbf{e}_m^T \varphi_k(\mathbf{H}_m)\beta\mathbf{e}_1 \right) \mathbf{A}^{k-1} \mathbf{v}_{m+1} - \sum_{k=1}^j \mathbf{A}^{k-1} \mathcal{E}_m \mathbf{V}_m \varphi_k(\mathbf{H}_m)\beta\mathbf{e}_1 + \mathbf{A}^j s_m^j, \end{aligned}$$

in which

$$\|\mathbf{A}^j s_m^j\| \leq \|\mathbf{A}\|^j \|s_m^j\| \leq \|\mathbf{A}\|^j \beta \left(\varphi_j(\|\mathbf{A}\|) + \varphi_j(\|\mathbf{H}_m\|) \right) \leq \beta \left(e^{\|\mathbf{A}\|} + e^{\|\mathbf{H}_m\|} \right) \frac{\|\mathbf{A}\|^j}{j!}$$

converges to 0 as $j \rightarrow \infty$. Taking this into account in the sums above, we get

$$\exp(\mathbf{A})\mathbf{v} - \mathbf{V}_m \exp(\mathbf{H}_m)\beta\mathbf{e}_1 = \sum_{k=1}^{\infty} \mathbf{A}^{k-1} [h_{m+1,m} \left(\mathbf{e}_m^T \varphi_k(\mathbf{H}_m)\beta\mathbf{e}_1 \right) \mathbf{v}_{m+1} - \mathcal{E}_m \mathbf{V}_m \varphi_k(\mathbf{H}_m)\beta\mathbf{e}_1].$$

Finally, if we rescale the inexact Arnoldi relation in (4) with the stepsize τ , we get (21). \square

4.4. Exactness in the case of truncated approximations

The analysis so far has not made any assumption on the structure of \mathbf{A} , \mathbf{v} or \mathbf{E}_k . In this section, we show a counter-intuitive result that, when \mathbf{E}_k arises from a truncated approximation of a special form of the matrix \mathbf{A} , the inexact scheme can be exact.

Consider a banded matrix, and assume that $\mathbf{v} = \mathbf{e}_1 = (1, 0, \dots, 0)^T$. The multiplication of \mathbf{A} with \mathbf{v} will therefore not involve the contribution of elements located in the trailing submatrices of \mathbf{A} . Generalizing this observation, we get the following result that arises frequently in the CME discussed in Section 2 and is therefore of wide interest.

Theorem 4.4. *Let $l \geq 0, k - l \geq 2$, assume that*

$$\mathbf{A} = \left(\begin{array}{c|c} \mathbf{A}_k & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{array} \right) \in \mathbb{R}^{n \times n},$$

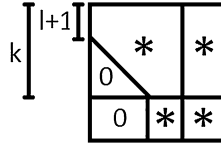
where

$$\mathbf{A}_k \in \mathbb{R}^{k \times k}, \quad (\mathbf{A}_k)_{ij} = 0 \text{ if } i > j + l,$$

and

$$\mathbf{C} \in \mathbb{R}^{(n-k) \times k}, \quad (\mathbf{C})_{ij} = 0 \text{ if } j \leq k - l - 1.$$

That is, \mathbf{A} visually has the form:



Also assume that $\mathbf{v} = \mathbf{e}_1$ and the relaxation matrices are identical,

$$\mathbf{E} = \left(\begin{array}{c|c} \mathbf{0} & -\mathbf{B} \\ \hline -\mathbf{C} & -\mathbf{D} \end{array} \right) \in \mathbb{R}^{n \times n}.$$

Then if $m \leq M = \max\{j : (j - 1)l + 1 \leq k - l - 1\}$, where m is the dimension of the basis built by the inexact Arnoldi algorithm based on matrix–vector products with $\mathbf{A} + \mathbf{E}$, we have

$$\mathcal{E}_m = \mathbf{0},$$

so that

$$\exp(\tau\mathbf{A}) = \exp(\tau\tilde{\mathbf{A}}),$$

and

$$\delta_m^{err} = \mathbf{0}.$$

Proof. Observe first that, because of the special forms of \mathbf{A} and \mathbf{v} , the \mathbf{v}_j produced by the Arnoldi process is such that

$$(\mathbf{v}_j)_i = 0, \quad i > (j - 1)l + 1, \quad 1 \leq j \leq M.$$

This means that, up to $j = M$, the exact Arnoldi process (where multiplications are performed with \mathbf{A}) and the inexact Arnoldi process (where they are instead performed with $\mathbf{A} + \mathbf{E}$) coincide, due to the fact that the multiplications only depend on the first $k - l - 1$ columns of \mathbf{A} and $\mathbf{A} + \mathbf{E}$, which are the same.

Because the first $k - l - 1$ columns of \mathbf{E} are zero, and only the first $(j - 1)l + 1 \leq k - l - 1$ elements of the vectors \mathbf{v}_j are nonzero for $1 \leq j \leq m \leq M$, we have $\mathbf{E}\mathbf{v}_j = \mathbf{0}$, and therefore

$$\mathcal{E}_m = \sum_{j=1}^m \mathbf{E}\mathbf{v}_j\mathbf{v}_j^T = \mathbf{0}.$$

Hence $\mathbf{A} = \tilde{\mathbf{A}}$ and naturally $\delta_m^{err} = \|\exp(\tau\mathbf{A})\mathbf{v} - \exp(\tau\tilde{\mathbf{A}})\mathbf{v}\| = 0. \quad \square$

Remark 4.5. The analysis in [Theorem 4.4](#) will explain results apparently intriguing in the experiments, as the reader will soon discover in the following section. Note however that it only works for $\mathbf{v} = \mathbf{e}_1$. Since the step-by-step integration scheme (11) is typically used, \mathbf{v} will not remain \mathbf{e}_1 past the first step, in which case the analysis of δ_m^{err} done in the previous section takes effect.

Remark 4.6. Even though this theorem shows that we have $\mathcal{E}_m = \mathbf{0}$ and $\exp(\tau\mathbf{A}) = \exp(\tau\tilde{\mathbf{A}}), \forall \tau$, with a Krylov subspace of small dimension, this does not mean that the inexact Krylov approximation has no error. As shown in [Theorem 4.3](#), even if $\mathcal{E}_m = \mathbf{0}$, the expansion collapses to (20), which is generally not zero.

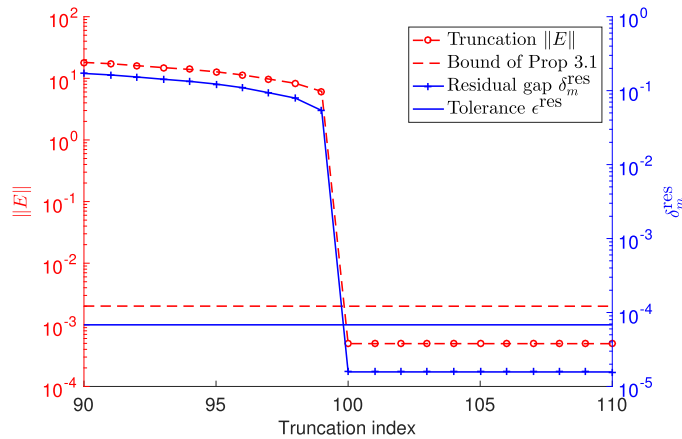


Fig. 1. Example 1: when the principal 100×100 submatrix becomes contained in the truncated matrix, the bound condition of Proposition 3.1 is satisfied, resulting in δ_m^{res} within the tolerance.

5. Numerical examples

We illustrate some of the theoretical results using three examples. The first two examples illustrate how Proposition 3.1 can be applied, by first showing how the relaxation scheme works according to the proposition, and secondly how the scheme fails when the condition of the proposition is not satisfied. The third example is a special one to demonstrate the peculiar behavior hinted in Remark 4.5: the scheme still works well even though the (sufficient) condition of Proposition 3.1 is not satisfied due to the reason given in Theorem 4.4. The examples were implemented in MATLAB.

We recall the following quantities that were given in the text to assess the inexact method:

$$\begin{aligned} \text{true residual} &= \|\mathbf{r}_m\| = \|\mathbf{V}_m \mathbf{H}_m \mathbf{y}_m - \mathbf{A} \mathbf{V}_m \mathbf{y}_m\| \\ \text{computed residual} &= \|\tilde{\mathbf{r}}_m\| = |h_{m+1,m} \mathbf{e}_m^T \mathbf{y}_m| \\ \text{residual gap} &= \delta_m^{\text{res}} = \|\mathbf{r}_m - \tilde{\mathbf{r}}_m\| = \|\mathcal{E}_m \mathbf{V}_m \mathbf{y}_m\| \\ \text{true error} &= \epsilon_m = \|\exp(\tau \mathbf{A}) \mathbf{v} - \mathbf{V}_m \mathbf{y}_m\|, \end{aligned}$$

where $\mathbf{y}_m = \exp(\tau \mathbf{H}_m) \beta \mathbf{e}_1$, with a given τ , and with \mathbf{H}_m and \mathbf{V}_m constructed by the inexact Krylov method using an initial given vector $\mathbf{v} = \mathbf{p}_0$ with $\beta = \|\mathbf{v}\|$.

5.1. Example 1 — illustration of the residual gap when Proposition 3.1 is satisfied

We take a 1000×1000 matrix, where the principal 100×100 submatrix is uniformly distributed in $[0, 1]$, and entries outside this submatrix are uniformly distributed in $[0, 10^{-6}]$. The initial vector is $\mathbf{v} = (10^{-3}, \dots, 10^{-3})^T$. We take the time point $\tau = 10^{-3}$, and tolerance on the residual gap $\epsilon^{\text{res}} = 10^{-3}$. The Krylov subspace is chosen to be of dimension $m = 15$.

For the inexact scheme, we define \mathbf{A}_k to be a 1000×1000 matrix containing the principal $k \times k$ submatrix of \mathbf{A} and 0 outside. Inexact multiplications with \mathbf{A} are then performed with \mathbf{A}_k instead. Observations are shown in Fig. 1. Since the key entries of \mathbf{A} are in the principal 100×100 submatrix, it is clear that if $k < 100$, then \mathbf{A}_k will leave out vital entries and inflate the error matrix $\mathbf{E} = \mathbf{A} - \mathbf{A}_k$, which in turn will cause the residual gap to be large. However, if $k \geq 100$, then $\|\mathbf{E}\| \leq \frac{\epsilon^{\text{res}}}{m\beta \|\exp(\tau \mathbf{H}_m)\|}$, which is the bound condition in Proposition 3.1, and therefore $\delta_m^{\text{res}} \leq \epsilon^{\text{res}}$ as guaranteed there. Fig. 1 illustrates this with the truncation index $90 \leq k \leq 110$.

5.2. Example 2 — illustration of the residual gap when Proposition 3.1 is not satisfied

We now consider another 1000×1000 matrix, where the main diagonal is uniformly distributed in $[0, 1]$, and the off-diagonal entries are uniformly distributed in $[0, 10^{-6}]$. As in the first example, we keep $\mathbf{v} = (10^{-3}, \dots,$

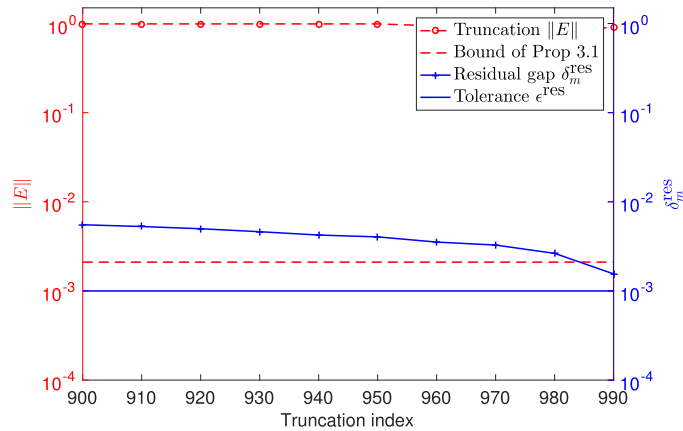


Fig. 2. Example 2: The bound condition of Proposition 3.1 is never satisfied, and δ_m^{res} is greater than the tolerance.

Table 1
Michaelis–Menten reactions and propensities.

	Reaction	Propensity	Rate constant (s^{-1})
1.	$E + S \xrightarrow{\kappa_1} ES$	$\alpha_1 = \kappa_1 [E] [S]$	$\kappa_1 = 1.0$
2.	$ES \xrightarrow{\kappa_2} E + S$	$\alpha_2 = \kappa_2 [ES]$	$\kappa_2 = 1.0$
3.	$ES \xrightarrow{\kappa_3} E + P$	$\alpha_3 = \kappa_3 [ES]$	$\kappa_3 = 0.1$

$10^{-3})^T$, $\tau = 10^{-3}$, $\epsilon^{\text{res}} = 10^{-3}$, and $m = 15$. The inexact setup is also the same, with inexact matrix–vector products against A performed using A_k , where A_k contains the $k \times k$ principal submatrix of A and zeros elsewhere.

As Fig. 2 shows, in this example, the residual gap δ_m^{res} is consistently above the tolerance on the residual gap ϵ^{res} even when k is almost the size of A . The reason for this is that since the diagonal elements are significant to be omitted, truncating even only one or two of them would make $\|E\|$ large and not satisfy the bound condition of Proposition 3.1.

The conclusion to draw from these two examples is that when inexactness is achieved in a reasonable way, the inexact Krylov method works well. Depending on the particular problem at hand, one can decide how to relax the matrix–vector multiplications to satisfy Proposition 3.1, in which case the residual gap can be controlled by ϵ^{res} , ensuring that the computed residual $\|\tilde{r}_m\|$ serves as a reliable approximation of the true residual $\|r_m\|$.

5.3. Example 3 — illustration of the error and residuals when Theorem 4.4 is satisfied

We consider the CME arising from the Michaelis–Menten enzyme kinetics, which is a well known system of biochemical reactions in cell biology. There are four species: substrates (S), enzymes (E), enzyme–substrate complexes (ES) and products (P), interacting according to the three chemical reactions listed in Table 1, and we took reaction rates as in [7]. The state vector is $x = ([P], [E], [S], [ES])^T$, where $[X]$ is the current number of copies of species X. If we start with $x_1 = (0, 50, 50, 0)^T$, i.e., a maximum of 50 substrates, the resulting matrix A of the underlying CME is of dimension $n = 1,326$. We use MATLAB’s expm command to check for correctness. Fig. 3 shows the sparsity pattern of A .

Since we know that the system starts in state x_1 , the initial probability vector is $p_0 = e_1$, and in the spirit of (6), we compute the approximation

$$\exp(\tau A)p_0 \approx \exp(\tau A_J)p_0,$$

where A_J is a padded truncation of A . We simply take a principal submatrix instead of the more general scheme where J can be an arbitrary subset of $\{1, \dots, n\}$. There is no loss of generality because it can be assumed that there has been a reordering $P^T A P$ where P is an appropriate permutation matrix. Our aim is to test the theory developed here and not really to craft efficient implementation details with elaborate sparse data structures that are best communicated

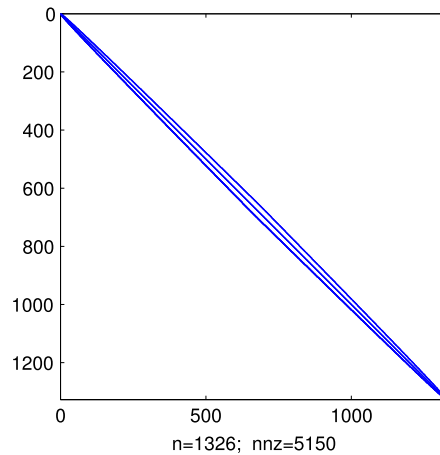


Fig. 3. Sparsity pattern of the matrix A from the CME of the Michaelis–Menten enzyme kinetics in Example 3.

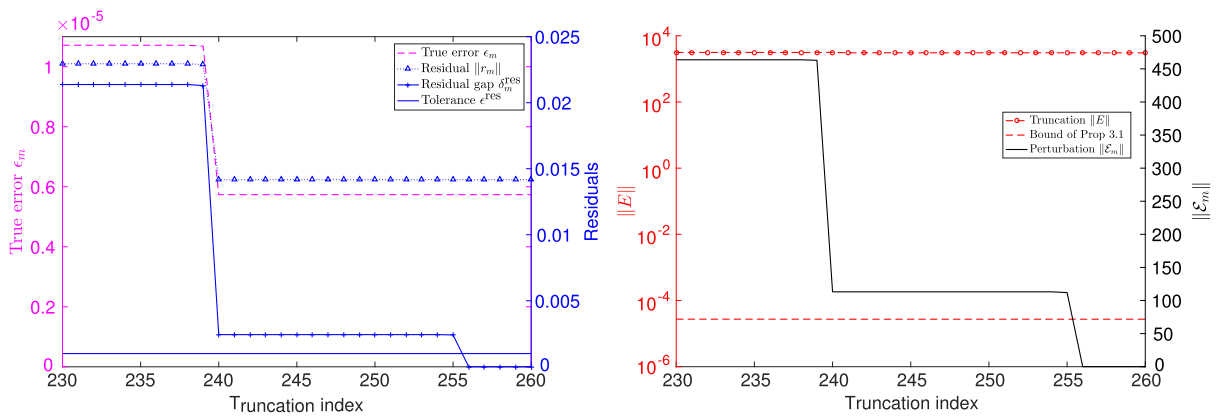


Fig. 4. Example 3. (Left) the true error (on the left y-axis), and the true residual and residual gap (on the right y-axis); (Right) $\|E\|$ and $\|E_m\|$; computed as $|J|$ increases from 230 to 280; $v = e_1$, $\tau = 10^{-2}$, and $m = 30$.

elsewhere. We vary $|J|$, the cardinality of J that determines the size of the truncation from 230 to 280, so that $|J|$ ranges from about 17% to 21% of n in this example. We take $v = e_1$, $\tau = 10^{-2}$, and $m = 30$.

Fig. 4 shows the results. We see on the right plot a decrease of $\|E_m\|$ as the truncation size $|J|$ increases. This example has the particularity of illustrating the phenomenon described in Section 4.4, where E_m becomes a zero matrix for a big enough truncation size $|J|$, even though $|J|$ is still only about 20% of n . Because of this phenomenon, the residual gap becomes 0, as the left plot shows. These results agree with the theory presented here. Note also on the left plot how the stair functions change in unison, illustrating the connection between the residuals and the error.

6. Conclusion

This work has analyzed inexact Krylov methods for approximating the action of the matrix exponential and provided insights into why they can be successful. We obtained results that in hindsight connect well with previous results, but it is worth recalling that it was unclear at the beginning exactly how such inexact methods related to previous works. The rigorous treatment presented here made the connection clear and established the details. This therefore fills a gap in the literature. We also brought into focus a particularly attractive aspect of inexact methods: they set a framework that encompasses model reduction methods in a generic way. We gave the important application of solving the chemical master equation (CME) as an example to motivate this viewpoint, with truncation methods such as the finite state projection (FSP) method that fit naturally in the framework. The study included numerical experiments to illustrate the theory.

Acknowledgments

This work was supported in part by NSF grant DMS-1320849. We thank the anonymous reviewers for their comments that helped improve the paper.

References

- [1] A.H. Al-Mohy, N.J. Higham, Computing the action of the matrix exponential, with an application to exponential integrators, *SIAM J. Sci. Comput.* 33 (2) (2011) 488–511.
- [2] M.A. Botchev, V. Grimm, M. Hochbruck, Residual, restarting and Richardson iteration for the matrix exponential, *SIAM J. Sci. Comput.* 35 (3) (2013) A1376–A1397.
- [3] A. Bouras, V. Frayssé, Inexact matrix–vector products in Krylov methods for solving linear systems: A relaxation strategy, *SIAM J. Matrix Anal. Appl.* 26 (3) (2005) 660–678.
- [4] J.V.D. Eshof, G.L.G. Sleijpen, Inexact Krylov subspace methods for linear systems, *SIAM J. Matrix Anal. Appl.* 26 (1) (2004) 125–153.
- [5] E. Gallopoulos, Y. Saad, Efficient solution of parabolic equations by Krylov approximation methods, *SIAM J. Sci. Stat. Comput.* 13 (5) (1992) 1236–1264.
- [6] L. Giraud, S. Gratton, J. Langou, Convergence in backward error of relaxed GMRES, *SIAM J. Sci. Comput.* 29 (2) (2007) 710–728.
- [7] J. Goutsias, Quasiequilibrium approximation of fast reaction kinetics in stochastic biochemical systems, *J. Chem. Phys.* 122 (18) (2005) 184102.
- [8] N.J. Higham, *Functions of Matrices: Theory and Computation*, SIAM, Philadelphia, PA, USA, 2008.
- [9] M. Hochbruck, C. Lubich, On Krylov subspace approximations to the matrix exponential operator, *SIAM J. Numer. Anal.* 34 (1997) 1911–1925.
- [10] B. Munsky, M. Khammash, The finite state projection algorithm for the solution of the chemical master equation, *J. Chem. Phys.* 124 (4) (2006) 044104.
- [11] Y. Saad, Analysis of some Krylov subspace approximations to the matrix exponential operator, *SIAM J. Numer. Anal.* 29 (1) (1992) 209–228.
- [12] R.B. Sidje, EXPOKIT: A software package for computing matrix exponentials, *ACM Trans. Math. Software* 24 (1) (1998) 130–156.
- [13] R.B. Sidje, Inexact uniformization and GMRES methods for large Markov chains, *Numer. Linear Algebra Appl.* 18 (2011) 947–960.
- [14] R.B. Sidje, N. Winkles, Evaluation of the performance of inexact GMRES, *J. Comput. Appl. Math.* 235 (2011) 1956–1975.
- [15] V. Simoncini, D.B. Szyld, Theory of inexact Krylov subspace methods and applications to scientific computing, *SIAM J. Sci. Comput.* 25 (2) (2003) 454–477.