

**ON THE ERROR IN COMPUTING LYAPUNOV EXPONENTS BY QR
METHODS ***

LUCA DIECI[†] AND ERIK S. VAN VLECK[‡]

Abstract. We consider the error introduced using QR methods to approximate Lyapunov exponents. We give a backward error statement for linear non-autonomous systems, and further discuss nonlinear autonomous problems. In particular, for linear systems we show that one approximates a “nearby” discontinuous problem where how nearby is measured in terms of local errors and a measure of non-normality. For nonlinear problems we use a type of shadowing result.

Key words. Lyapunov exponents, QR methods.

AMS subject classifications. 65L

1. The problem. Consider the non-autonomous linear system

$$(1.1) \quad \dot{x} = A(t)x, \quad t \geq 0,$$

where we will assume that the function $A : \mathbb{R}^+ \rightarrow \mathbb{R}^{n \times n}$ is bounded. To characterize the asymptotic growth behavior of (1.1), a commonly employed tool is that of *Lyapunov spectrum*, henceforth labeled Σ_L , which is based on upper and lower Lyapunov exponents.

Lyapunov exponents are also routinely used to study nonlinear dynamical systems via linearized analysis. Indeed, Lyapunov exponents are probably the most widely used quantities for detecting chaos, estimating dimensions of attractors, entropy, and so forth; e.g., see [5, 19, 2, 3].

There are some aspects about the popularity of Lyapunov exponents which are somewhat intriguing: (i) There are only a handful of non concocted problems for which the Lyapunov exponents are known analytically, and, as a consequence, (ii) Many studies using Lyapunov exponents are of numerical nature, but (iii) There is little error analysis of the techniques used to approximate Lyapunov exponents; see [7, 13, 15]. Our goal in this work is to rectify in part this situation, by providing a backward error result for QR techniques.

A plan of this paper is as follows. In the remainder of this Introduction, we recall the definition of Σ_L . In Section 2 we review the QR methods for approximating Lyapunov exponents. In Section 3 we give our main results on the error introduced by the QR methods when approximating Lyapunov exponents. These are largely results for linear problems, but we also discuss the nonlinear case with the help of a shadowing result. In Section 4 we illustrate our theoretical results on an example.

Lyapunov Spectrum: Σ_L . Let X be a fundamental matrix solution of (1.1) and consider

$$(1.2) \quad \lambda_i = \limsup_{t \rightarrow \infty} \frac{1}{t} \ln \|X(t)e_i\|, \quad i = 1, \dots, n,$$

where the e_i 's are the standard unit vectors. In (1.2) and everywhere in this paper the norm is the 2-norm for vectors. For a matrix, say $C \in \mathbb{R}^{n \times n}$, we will consider either the induced 2-norm, $\|C\|_2 = \max_{u \in \mathbb{R}^n: \|u\|_2=1} \|Cu\|_2$, or the Frobenius norm, $\|C\|_F = (\sum_{i,j=1}^n C_{ij}^2)^{1/2}$.

When $\sum_{i=1}^n \lambda_i$ is minimized with respect to all possible fundamental matrix solutions, then the λ_i 's are called (upper) Lyapunov exponents, and the corresponding fundamental matrix solution is called normal.

*This work was supported in part under NSF Grants DMS/FRG-0139895 and DMS/FRG-0139824

[†]School of Mathematics, Georgia Institute of Technology, Atlanta, Georgia 30332 (dieci@math.gatech.edu).

[‡]Department of Mathematics, University of Kansas, Lawrence, Kansas 66045 (evanvleck@math.ku.edu).

Consider also the adjoint equation

$$(1.3) \quad \dot{y}(t) = -A^T(t)y(t), \quad t \geq 0,$$

and let $\{-\mu_i\}_{i=1}^n$ be its Lyapunov exponents. We can assume that $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n$ and $\mu_1 \geq \mu_2 \geq \dots \geq \mu_n$. We define the Lyapunov spectrum Σ_L as

$$(1.4) \quad \Sigma_L := \bigcup_{j=1}^n [\mu_j, \lambda_j].$$

In general, for any $j = 1, \dots, n$, $\mu_j \leq \lambda_j$. If $\mu_j = \lambda_j$, for all $j = 1, \dots, n$, then the system is called regular.

Remark 1.1. Although in practice one may not need to find the entire spectrum of a system, but only a few of the most dominant spectral intervals, in this paper we consider the case in which we want to approximate the entire spectrum, and this is why we consider the entire fundamental matrix solution and not just a few of its columns.

The Lyapunov exponents (and hence Σ_L) are not necessarily stable; see [1]. In this context, the exponents are said to be stable if they are continuous with respect to perturbation in the coefficient matrix. More precisely, if “for any $\epsilon > 0$, there exists $\delta > 0$ such that $\sup_{t \in \mathbb{R}^+} \|E(t)\| < \delta$ implies

$$(1.5) \quad |\lambda_i - \hat{\lambda}_i| < \epsilon, \quad i = 1, \dots, n,$$

where the $\hat{\lambda}_i$'s are the (ordered) Lyapunov exponents of the perturbed system $\dot{x} = [A(t) + E(t)]x$.

If the Lyapunov exponents are distinct, then (again, see [1]) the exponents (and Σ_L) are stable if and only if X is an *integrally separated* fundamental matrix solution. Writing $X(t) = [X_1(t), \dots, X_n(t)]$, X is integrally separated if for $i = 1, \dots, n-1$, there exist $a > 0$ and $1 \geq d > 0$ such that

$$(1.6) \quad \frac{\|X_i(t)\|}{\|X_i(s)\|} \cdot \frac{\|X_{i+1}(s)\|}{\|X_{i+1}(t)\|} \geq de^{a(t-s)},$$

for all $t, s : t \geq s \geq 0$. Palmer [17, p. 21] considered the Banach space \mathcal{B} of continuous bounded matrix valued functions A , with norm $\|A\| = \sup_{t \geq 0} \|A(t)\|$, and employing results from

[14] showed that the systems with integral separation form an open and dense subset of \mathcal{B} . Thus, integral separation is a generic property in \mathcal{B} and it is thus reasonable to think (at least in first approximation) that one is trying to approximate stable and distinct Lyapunov exponents.

2. QR methods. In recent years, we have undertaken a systematic and focused effort in trying to numerically approximate Σ_L , and also the *Exponential Dichotomy* spectrum, by so-called *QR methods*; e.g., see [8, 9]. Ever since the work [4], these methods have been popular techniques to approximate Lyapunov exponents. QR methods come in two flavors: continuous or discrete. Although these are conceptually equivalent, their practical implementation and observed performance in finite precision often differ.

The bottom line of these techniques consist in extracting the required spectral information from the diagonal of the triangular factor in the QR factorization of the (normal) fundamental matrix solution X of (1.1). The function Q is orthogonal, while R is upper triangular with positive diagonal entries. Henceforth, whenever we refer to the QR factorization we always mean the unique one for which the diagonal of R is positive. It is well known that this factorization exists, is unique, and is as smooth as X .

Once R is available (but see below), the Lyapunov exponents are recovered from

$$(2.1) \quad \lambda_i = \limsup_{t \rightarrow \infty} \frac{1}{t} \log(R_{ii}(t)), \quad i = 1, \dots, n.$$

A similar formula can be used to obtain the μ_i 's. We remark (see [1]) that (2.1) gives the Lyapunov exponents for regular systems (and, then, the \limsup is really a limit). Moreover, in [10], we showed that (2.1) does indeed give the Lyapunov exponents as long as they are stable (not necessarily distinct). We hasten to stress that, from the numerical point of view, it is reasonable to approximate stable quantities.

2.1. Discrete QR method. Suppose we want the QR factorization at $X(t_k)$, for some sequence of points t_k , $k = 0, 1, 2, \dots$, with $t_0 = 0$. At any such point t_k , we can write

$$(2.2) \quad X(t_k) = \Phi(t_k, t_{k-1}) \dots \Phi(t_2, t_1) \Phi(t_1, 0) X_0,$$

where

$$(2.3) \quad \dot{\Phi}(t, t_{j-1}) = A(t) \Phi(t, t_{j-1}), \quad \Phi(t_{j-1}, t_{j-1}) = I, \quad t_{j-1} \leq t \leq t_j, \quad j = 1, 2, \dots, k.$$

Now, let $X_0 = Q(t_0)R(t_0)$, where $Q(t_0) \in \mathbb{R}^{n \times n}$ is orthogonal and $R(t_0) \in \mathbb{R}^{n \times n}$ is upper triangular with positive diagonal entries. Then, for $j = 1, 2, \dots, k$, recursively consider

$$(2.4) \quad \begin{aligned} \dot{\Psi}(t, t_{j-1}) &= A(t) \Psi(t, t_{j-1}), \quad \Psi(t_{j-1}, t_{j-1}) = Q(t_{j-1}) \\ \text{and factor } \Psi(t_j, t_{j-1}) &= Q(t_j)R(t_j, t_{j-1}), \end{aligned}$$

where $Q(t_j)$ are orthogonal and $R(t_j, t_{j-1})$ are upper triangular with positive diagonal. Then, we have the QR factorization of $X(t_k)$

$$(2.5) \quad X(t_k) = Q(t_k)R(t_k, t_{k-1}) \dots R(t_2, t_1)R(t_1, t_0)R(t_0).$$

In other words, we have

$$(2.6) \quad R(t_k) = \left(\prod_{j=k}^1 R(t_j, t_{j-1}) \right) R(t_0).$$

In order to access the diagonal of $R(t_k)$, we need to monitor only the diagonal entries of the factors $R(t_j, t_{j-1})$, so that (see (2.1))

$$(2.7) \quad \frac{1}{t_k} \log(R_{ii}(t_k)) = \frac{1}{t_k} \left(\sum_{j=k}^1 \log R_{ii}(t_j, t_{j-1}) + \log R_{ii}(t_0) \right).$$

Remark 2.1. Naturally, integral separation of X , see (1.6), can be phrased in terms of integral separation of R .

2.2. Continuous QR method. The key difference between discrete and continuous QR methods is that in the continuous case one forms the triangular system satisfied by R , and in so doing bypasses forming explicitly the function R . From the QR factorization of X , we seek the function Q which performs the change of variables.

Now, it is well known (e.g., see [9]) that such Q is unique and satisfies

$$(2.8) \quad \dot{Q} = Q(t)H(Q, A), \quad Q(0) = Q_0,$$

where we have set $H := Q^T(t)\dot{Q}(t)$, with entries

$$(2.9) \quad H_{ij}(t) = \begin{cases} (Q^T(t)A(t)Q(t))_{ij}, & i > j, \\ 0, & i = j, \\ -(Q^T(t)A(t)Q(t))_{ji}, & i < j. \end{cases}$$

So, if Q is known, then R satisfies the transformed system

$$(2.10) \quad \dot{R} = B(t)R, \quad R(0) = R_0,$$

where we have set

$$(2.11) \quad B(t) := Q^T(t)A(t)Q(t) - Q^T(t)\dot{Q}(t),$$

and B is upper triangular by the way that H has been defined. The logarithm of the diagonal of $R(t_k)$ is now

$$\log\left(\text{diag}(R(t_k))\right) = \int_{t_0}^{t_k} \text{diag}(B(t))dt + \log(\text{diag}(R_0)).$$

and so (see (2.1))

$$(2.12) \quad \frac{1}{t_k} \log(R_{ii}(t_k)) = \frac{1}{t_k} \left(\int_0^{t_k} B_{ii}(t)dt + \log R_{ii}(t_0) \right).$$

Remark 2.2. Integral separation of the fundamental matrix solution X , see (1.6), can be equivalently phrased in terms of integral separation of the diagonal of B (see [10]):

$$(2.13) \quad \int_s^t (B_{ii}(\tau) - B_{i+1,i+1}(\tau))d\tau \geq d + a(t-s),$$

for all $t, s : t \geq s$, where $d \leq 0$ and $a > 0$.

In practice, any sensible implementation of the continuous QR method will form the (diagonal of the) function B in (2.11) and approximate the integral in (2.12) at the same time as the approximation of Q in (2.8) is done. We can think of this as follows.

From (2.11), we have that the diagonal entries of B are given by $B_{ii}(t) = [Q^T(t)A(t)Q(t)]_{ii}$, $i = 1, \dots, n$. Define the functions $\mu_i(t)$, $t \geq 0$, $i = 1, \dots, n$, as

$$(2.14) \quad \mu_i(t) = \int_0^t [Q^T(s)A(s)Q(s)]_{ii}ds, \quad i = 1, \dots, n.$$

Then, we can think of the continuous QR method as a technique where we approximate this differential system

$$(2.15) \quad \begin{cases} \dot{Q} &= QH, & Q(0) = Q_0, \\ \dot{\mu}_i &= [Q^T(t)A(t)Q(t)]_{ii}, & \mu_i(0) = 0, \quad i = 1, \dots, n. \end{cases}$$

The following simple observation will come in handy. Suppose we have values t_0, t_1, \dots , and let $t_{k-1} \leq t < t_k$. Clearly, the Q-factor in the QR factorization of $X(t)$ is the same as the unique Q-factor in the QR factorization of the function (recall (2.3))

$$(2.16) \quad \Phi(t, t_{k-1}) \dots \Phi(t_2, t_1)\Phi(t_1, 0)Q_0.$$

Therefore, the differential system (2.15) on which the continuous QR method is based can be formulated in a local way. That is, for $t_j \leq t \leq t_{j+1}$, $j = 0, 1, \dots, k-1$, we seek functions $Q(t, t_j)$ and $\mu_i(t, t_j)$ satisfying

$$(2.17) \quad \begin{cases} \dot{Q}(t, t_j) &= Q(t, t_j)H(Q(t, t_j), A(t)), & Q(t_j, t_j) = Q(t_j), \\ \dot{\mu}_i(t, t_j) &= [Q^T(t, t_j)A(t)Q(t, t_j)]_{ii}, & \mu_i(t_j, t_j) = 0, \quad i = 1, \dots, n, \end{cases}$$

and H is defined as in (2.9) but with $Q(t, t_j)$ replacing $Q(t)$ there. So doing, in (2.12) we will use

$$(2.18) \quad \frac{1}{t_k} \int_0^{t_k} B_{ii}(t) dt = \frac{1}{t_k} \sum_{j=1}^k \mu_i(t_j, t_{j-1}).$$

We conclude this section with two useful results. The following justifies the notation $R(t_j, t_{j-1})$ used in (2.4).

LEMMA 2.3. *For $j = 1, 2, \dots$, the matrices $R(t_j, t_{j-1})$ in (2.4) are the solution at t_j of*

$$\dot{R}(t, t_{j-1}) = B(t)R(t, t_{j-1}), \quad R(t_{j-1}, t_{j-1}) = I,$$

where $B(t)$ is given in (2.11).

Proof. This is a consequence of uniqueness of the QR factorization. The formal argument is by induction on j . For $j = 1$, from (2.5) we have that $R(t_1) = R(t_1, t_0)R(t_0)$, and from (2.10) we have this same formula as well, with $R(t_1, t_0)$ solution at t_1 of $\dot{R}(t, t_0) = B(t)R(t, t_0)$, $R(t_0, t_0) = I$. Uniqueness of the QR factorization of $X(t_1)$ gives the case of $j = 1$. Now, if the result is true for j , then it follows for $j + 1$, by the same argument we just used replacing $R(t_0)$ by $R(t_j)$, etc.. \square

We will also use the following result.

LEMMA 2.4. *Suppose that the sequence of points $\{t_j\}$, $j = 0, 1, \dots$, with $0 = t_0 < t_1 < t_2 < \dots$, is given, that $\lim_{j \rightarrow \infty} t_j = \infty$, and that all the matrices $R(t_{j+1}, t_j) \in \mathbb{R}^{n \times n}$ in (2.6) are upper triangular with positive diagonal. Then, at any t_k , $R(t_k)$, solution of (2.10), is the same matrix as the exact solution at t_k of the piecewise constant system*

$$(2.19) \quad \dot{\tilde{R}} = B_j \tilde{R}, \quad t_j \leq t < t_{j+1}, \quad \tilde{R}(0) = R(t_0), \quad j = 0, \dots, k-1,$$

where the matrices $B_j \in \mathbb{R}^{n \times n}$ are upper triangular and satisfy

$$(2.20) \quad R(t_{j+1}, t_j) = e^{h_j B_j}, \quad h_j := t_{j+1} - t_j, \quad j = 0, \dots, k-1.$$

Proof. The assumption on the matrices $R(t_{j+1}, t_j)$ implies that they are invertible with no negative eigenvalues. So, they have real logarithms, call them $h_j B_j$, which are upper triangular, and (2.20) is satisfied. Now the proof follows at once using Lemma 2.3, and the obvious equality $R(t_j, t_{j-1}) = \tilde{R}(t_j, t_{j-1})$, where we must interpret $\tilde{R}(t_j, t_{j-1})$ as the left limit: $t \rightarrow t_j^-$. \square

Remark 2.5. For us, the points t_j , $j = 0, 1, \dots$, in the Lemma will be the sequence of points found during numerical integration.

3. Error Analysis.

As far as we know, the first concrete effort to provide error analysis for QR methods was made in [7] where the authors, under the assumption of point spectrum¹, attempted an analysis of the errors introduced by both continuous and discrete QR methods and further specialized their results to the case of constant coefficient and periodic problems. The work [13] is a more recent effort for constant coefficients problems. For the case of algorithms based on the discrete Singular-Value-Decomposition (SVD) method, we refer to the work [15] (see also [20]). There, the authors consider the errors introduced by the finite precision computation of the SVDs of exact transition matrices. Interestingly, in [15], Oliveira and Stewart reach a perturbation result which rests on an assumption very close to one needed for stability of the continuous SVD method (cfr. [15, Formula (2.1)] with [9, Formulas (7.6), (7.7)]).

The question we want to address here is: when we approximate Lyapunov exponents by the QR methods, what kind of error result can we get?

¹That is, the Exponential Dichotomy spectrum reduces to n points, the Lyapunov exponents

There are always (at least) three sources of error: the one introduced by discretization of the relevant differential equations, the one introduced by the replacement of limits on the continuous time variable t by limits taken over a sequence of points $\{t_j\}$, and of course the one introduced by the need to truncate the limiting process itself. There are also the errors introduced by linear algebra computations in finite precision, but we may think that these are incorporated in the discretization errors.

In this work, we analyze only the first source of error. To be justified in replacing the \limsup and \liminf by computations on a sequence of points rather than on a continuum, we may want to think that the system is regular, hence the Lyapunov exponents exist as limits. In this case, we can replace the limit along the positive real axis with the limit along any sequence of points $\{t_k\}$ converging to ∞ . To be conceptually justified in neglecting the truncation of time, we will think of the situation in which we use the QR methods taking them to the limit of $t \rightarrow \infty$.

With these simplifications, we will certainly expect to obtain good error statements for the Lyapunov exponents computed by QR methods if we could infer that the integration errors to approximate Q (and R) do not accumulate. But, in general, these errors do accumulate! However, in spite of this fact, one often obtains very accurate approximations to the Lyapunov exponents. Why? We believe that the reason lies in the following statement, which we will better qualify in this paper. In fact, our main results, Theorems 3.12 and 3.16, may be summarized as follows. “With the QR methods we are finding, at points t_k , the QR factorization of a matrix solution associated to a linear system close to the original one. The measure of closeness depends on several factors, the most relevant being: The magnitude of the local errors occurring during integration of the relevant differential equations, and how non diagonal is the function R (i.e., the *departure from normality* of R), the exact triangular factor in the QR factorization of the matrix solution X ”.

Backward error results have proven very important in the context of numerical linear algebra and geometric integration, see for example [21] and [12]. Our result is more in the flavor of the linear algebra type of results, and our is the first result of which we are aware that gives a backward error statement on the stability of the QR methods. It really says that, in theory, the QR methods can be used to approximate (in the infinite limit) the Lyapunov exponents of a system which we can make close to the original one, by controlling the integration stepsize, or the error tolerance. Once this is understood, then one may want to ask whether this type of backward error result implies that the computed Lyapunov exponents are close to the Lyapunov exponents of the original problem. We do not address this concern directly here. We simply remark that, for this to be true, a stability/continuity result for Lyapunov exponents is needed, which in the case of distinct Lyapunov exponents is equivalent to the existence of an integrally separated fundamental matrix solution.

3.1. Discrete QR error. In practice, we cannot solve for the transition matrices Φ in (2.2) exactly, and we will actually compute

$$(3.1) \quad X_k = X(t_k, t_{k-1}) \dots X(t_2, t_1) X(t_1, t_0) X_0,$$

where the matrices $X(t_j, t_{j-1})$ are approximations to $\Phi(t_j, t_{j-1})$, $j = 1, \dots, k$. Letting $Q(t_0) = Q_0$, and progressively setting

$$X(t_j, t_{j-1}) Q_{j-1} = Q_j R_j, \quad j = 1, \dots, k,$$

the numerical discrete QR method will read

$$(3.2) \quad X_k = Q_k R_k R_{k-1} \dots R_2 R_1 R(t_0).$$

The issue is: how does (3.2) relate to (2.5)? The next result is the key.

THEOREM 3.1. *For $k = 1, 2, \dots$, and $j = 1, \dots, k$, let $Q(t_j)$ and $R(t_j, t_{j-1})$ be the exact Q and R terms in (2.5), and let X_k be given in (3.1). We have*

$$(3.3) \quad X_k = Q(t_k) [R(t_k, t_{k-1}) + E_k] \dots [R(t_2, t_1) + E_2] [R(t_1, t_0) + E_1] R(t_0),$$

where

$$(3.4) \quad E_j = Q^T(t_j)N_jQ(t_{j-1}) , \quad j = 1, \dots, k ,$$

and N_j are the local errors obtained when approximating $\Phi(t_j, t_{j-1})$ by $X(t_j, t_{j-1})$: $N_j = X(t_j, t_{j-1}) - \Phi(t_j, t_{j-1})$, $j = 1, \dots, k$.

As a consequence of the above, the numerical realization of the discrete QR method as expressed by (3.2) finds the exact QR factorization of the sequence on the right-hand-side of (3.3).

Proof. Consider the first integration step. We have $X(t_1, t_0) = \Phi(t_1, t_0) + N_1$, and we form $X(t_1, t_0)Q_0 = Q_1R_1$. But, since we have $\Phi(t_1, t_0)Q_0 = Q(t_1)R(t_1, t_0)$, then we obtain $Q_1R_1 = Q(t_1)R(t_1, t_0) + N_1Q_0$, and so

$$X(t_1, t_0)Q_0 = Q(t_1)[R(t_1, t_0) + Q^T(t_1)N_1Q_0] .$$

Now consider the second integration step. We are looking at $X_2 = X(t_2, t_1)X(t_1, t_0)Q_0$. Arguing as we just did, we get

$$X(t_2, t_1)Q(t_1) = Q(t_2)[R(t_2, t_1) + Q^T(t_2)N_2Q(t_1)] ,$$

where N_2 is obtained from $X(t_2, t_1) = \Phi(t_2, t_1) + N_2$. Using this in the expression for X_2 gives

$$X_2 = Q(t_2)[R(t_2, t_1) + E_2][R(t_1, t_0) + E_1]R(t_0) .$$

Continuing this way, we obtain the sought result. Finally, the interpretation of what the numerical discrete QR method does is obvious. \square

Remark 3.2. To clarify, suppose that we are using a p -th order method for the numerical integration of the differential equations, enforcing a local error control, so that $\|N_j\| = \|\Phi(t_j, t_{j-1}) - X(t_j, t_{j-1})\| \leq \text{TOL}$, in the 2-norm or in the Frobenius norm. Then, we can assume that

$$\|E_j\| = \|N_j\| \leq \text{TOL} .$$

Of course, we could also say $\|E_j\| \leq c_j h_{j-1}^{p+1}$, where we have set $h_{j-1} = t_j - t_{j-1}$, $j = 1, 2, \dots$, and the c_j 's are constants (which depend on the function of coefficients A , and on the formula used). In all cases, we can control (either through the input tolerance, or through the choice of stepsize) the norm of E_j , in principle making it as small as we want, though in practice we can at best hope to make it as small as order `eps` (the machine precision)².

Let us now set

$$\widehat{R}(t_{j+1}, t_j) := R(t_{j+1}, t_j) + E_{j+1} ,$$

and more generally

$$\widehat{R}(t_k) = \left(\prod_{j=k}^1 \widehat{R}(t_j, t_{j-1}) \right) R(t_0) .$$

We must stress that these matrices $\widehat{R}(t_{j+1}, t_j)$ are not upper triangular, in general. Assume now that E_{j+1} is sufficiently small, so that $\widehat{R}(t_{j+1}, t_j)$ is nonsingular and has no eigenvalue on the negative real axis; see also Assumption 3.5 below. Then, in a similar way to what we did in Lemma 2.4, $\widehat{R}(t_k)$ is the exact solution at t_k of the piecewise constant problem

$$(3.5) \quad \dot{\widehat{R}} = \widehat{B}_j \widehat{R} , \quad t_j \leq t < t_{j+1} , \quad \widehat{R}(0) = R_0 , \quad j = 0, \dots, k-1 ,$$

²In double precision arithmetic, `eps` $\approx 2 \times 10^{-16}$

where the matrices $\widehat{B}_j \in \mathbb{R}^{n \times n}$ satisfy

$$(3.6) \quad \widehat{R}(t_{j+1}, t_j) = e^{h_j \widehat{B}_j}, \quad h_j := t_{j+1} - t_j, \quad j = 0, \dots, k-1.$$

By comparing the above result with Lemma 2.4, we now proceed to find bounds on the differences $B_j - \widehat{B}_j$. If we can make these differences small, then we would have found the sought backward error statement. Intuitively, since $\widehat{R}(t_{j+1}, t_j) = R(t_{j+1}, t_j) + E_{j+1}$, and $\|E_{j+1}\|$ is of the same order of magnitude as the local error, one may expect that $\|B_j - \widehat{B}_j\|$ is also of the same order. This is correct if logarithms of nearby matrices are near to one another. Precisely in which way this is the case is what we will investigate next. The following identity is useful for our analysis.

THEOREM 3.3. [6, Theorem 2.6] *Let $R \in \mathbb{R}^{n \times n}$ be an invertible matrix with no eigenvalues on the negative real axis, and let $E \in \mathbb{R}^{n \times n}$ be such that $R + E$ also has no eigenvalues on the negative real axis. Then, $R + E$ admits a principal logarithm $\log(R + E)$ and the following formula holds:*

$$(3.7) \quad \log(R + E) = \log(R) + \int_0^1 [(R - I)s + I]^{-1} E [(R + E - I)s + I]^{-1} ds.$$

Now, (3.7) indicates that $\log \widehat{R}(t_{j+1}, t_j)$ is close to $\log R(t_{j+1}, t_j)$ if the integral term

$$(3.8) \quad \int_0^1 [(R(t_{j+1}, t_j) - I)s + I]^{-1} E_{j+1} [(R(t_{j+1}, t_j) + E_{j+1} - I)s + I]^{-1} ds$$

will remain of the same order of magnitude as E_{j+1} . Naturally, we must have that this is true for all terms in the infinite sequence $\{t_j\}$. Our goal in the remainder of this section is to find an insightful expansion, and bounds, for the expression in (3.8), since as it stands (3.8) is not too revealing. We look for an expansion in E_{j+1} , and then first order bounds (that is, within terms of order $\|E_{j+1}\|^2$).

Let us introduce some notation. Define $D(t_{j+1}, t_j)$ (see Lemma 2.3), for $j = 0, 1, \dots$ as the solution at t_{j+1} of

$$(3.9) \quad \dot{D}(t, t_j) = \text{diag}(B(t))D(t, t_j), \quad D(t_j, t_j) = I,$$

where $B(t)$ is given in (2.11).

Notation 3.4. For all $j = 0, \dots, k-1$, and $k \geq 1$, we will set:

- For the diagonal and strictly upper triangular part of $R(t_{j+1}, t_j)$:

$$(3.10) \quad R(t_{j+1}, t_j) = D(t_{j+1}, t_j) + U_{j+1}.$$

- Also, for $s \in [0, 1]$, we set

$$(3.11) \quad R_j(s) := (R(t_{j+1}, t_j) - I)s + I,$$

and

$$(3.12) \quad D_j(s) := (D(t_{j+1}, t_j) - I)s + I.$$

We will also need Assumption 3.5 below.

ASSUMPTION 3.5. *Assume that*

$$(3.13) \quad \rho_j := \|E_{j+1}\| \cdot [\min_{1 \leq i \leq n} (1, D_{ii}(t_{j+1}, t_j))]^{-1} \leq \rho < 1, \quad \forall j = k, \dots, 0, \quad k = 1, 2, \dots.$$

Remark 3.6. Assumption 3.5 quantifies the interplay between the local errors (that is, $\|E_{j+1}\|$) and the solutions of the diagonal problems. In principle, since we can make $\|E_{j+1}\|$ as small as we like, then it is clear that we can enforce Assumption 3.5, for any given value of k . Alternatively, for given value of t_j , and desired bounds on $\|E_{j+1}\|$, we may choose the point t_{j+1} to enforce (3.13).

LEMMA 3.7. *Let (3.13) hold. For any given $j = k, \dots, 0$, and $k = 1, 2, \dots$, we have*

$$[R_j(s) + sE_{j+1}]^{-1} = R_j^{-1}(s) \sum_{l=0}^{\infty} (-1)^l (sE_{j+1}R_j^{-1}(s))^l, \quad s \in [0, 1].$$

Proof. Rewriting $[(R(t_{j+1}, t_j) + E_{j+1} - I)s + I]^{-1} = [R_j(s) + sE_{j+1}]^{-1} = R_j^{-1}(s)[I + sE_{j+1}R_j^{-1}(s)]^{-1}$, the result would follow upon expanding the latter inverse. Now, the given expansion is valid as long as the spectral radius of $sE_{j+1}R_j^{-1}(s)$ is less than 1. But, since the eigenvalues of $R_j(s)$ are $(D_{ii}(t_{j+1}, t_j) - 1)s + 1$, $i = 1, \dots, n$, and (3.13) holds, the claim follows. \square

Using the expansion of Lemma 3.7, and (3.7), we can thus write

$$(3.14) \quad \log(R(t_{j+1}, t_j) + E_{j+1}) = \log(R(t_{j+1}, t_j)) + \int_0^1 R_j^{-1}(s)E_{j+1}R_j^{-1}(s)ds + O(\|E_{j+1}\|^2),$$

where the constant hidden in the high order terms depends on bounds on $\max_{0 \leq s \leq 1} \|R_j^{-1}(s)\|$.

Our next tasks are to find a manageable expression (and bound) for $\int_0^1 R_j^{-1}(s)E_{j+1}R_j^{-1}(s)ds$ and to obtain uniform (in j and s) bounds on $\|R_j^{-1}(s)\|$ so that the quantity $O(\|E_{j+1}\|^2)$ is under control.

LEMMA 3.8. *We can rewrite $R_j^{-1}(s)$ as follows:*

$$(3.15) \quad R_j^{-1}(s) = D_j^{-1}(s) \sum_{k=0}^{n-1} (-1)^k V_j^k(s), \quad \text{where} \quad V_j(s) = sU_{j+1}D_j^{-1}(s),$$

and also as

$$(3.16) \quad R_j^{-1}(s) = \left[\sum_{k=0}^{n-1} (-1)^k W_j^k(s) \right] D_j^{-1}(s), \quad \text{where} \quad W_j(s) = sD_j^{-1}(s)U_{j+1},$$

where U_{j+1} is defined in (3.10). As a consequence, we can write

$$(3.17) \quad \int_0^1 R_j^{-1}(s)E_{j+1}R_j^{-1}(s)ds = \int_0^1 \left(\sum_{k=0}^{n-1} (-1)^k W_j^k(s) \right) D_j^{-1}(s)E_{j+1}D_j^{-1}(s) \left(\sum_{k=0}^{n-1} (-1)^k V_j^k(s) \right) ds.$$

Proof. The final rewriting (3.17) is obvious. We will only prove (3.15), the proof of (3.16) being similar. We have

$$R_j^{-1}(s) = (D_j(s) + sU_{j+1})^{-1} = D_j^{-1}(s)(I + sU_{j+1}D_j^{-1}(s))^{-1} = D_j^{-1}(s)(I + V_j(s))^{-1}.$$

Now, since U_{j+1} is strictly upper triangular, then $V_j(s)$ is strictly upper triangular and therefore we have $V_j^n(s) \equiv 0$ and

$$(I + V_j(s))^{-1} = \sum_{k=0}^{n-1} (-1)^k V_j^k(s).$$

\square

The expression in (3.17) elucidates the first order error expansion in (3.14). Now, we proceed to find bounds. To simplify notation, let us set

$$L_j := \int_0^1 R_j^{-1}(s) E_{j+1} R_j^{-1}(s) ds,$$

so that (3.14) reads

$$\log(R(t_{j+1}, t_j) + E_{j+1}) = \log(R(t_{j+1}, t_j)) + L_j + O(\|E_{j+1}\|^2),$$

and using (3.17) we will bound L_j in norm.

To begin with, we derive bounds for the matrix of absolute values. The triangular inequality gives

$$(3.18) \quad |L_j| \leq \max_{0 \leq s \leq 1} \sum_{k=0}^{n-1} \|W_j^k(s)\| \int_0^1 D_j^{-1}(s) |E_{j+1}| D_j^{-1}(s) ds \max_{0 \leq s \leq 1} \sum_{k=0}^{n-1} \|V_j^k(s)\|,$$

where the notation $|L_j|$, etc., refers to the matrix of absolute values (recall that $D_j^{-1}(s)$, see (3.12), are diagonal with positive diagonal entries). Next, we give bounds on the quantities in (3.18).

THEOREM 3.9. *Consider the matrix $\int_0^1 D_j^{-1}(s) |E_{j+1}| D_j^{-1}(s) ds$ of (3.18), call it F_{j+1} . Further, let $D(t_{j+1}, t_j) = \text{diag}(d_1, \dots, d_n)$. Then, for $p, q = 1, \dots, n$, F_{j+1} has entries*

$$(3.19) \quad (F_{j+1})_{pq} = |(E_{j+1})_{pq}| \times \begin{cases} \frac{\log(d_p) - \log(d_q)}{d_p - d_q}, & p \neq q \text{ and } d_p \neq d_q \\ \frac{1}{d_p}, & p = q \text{ or } p \neq q, \text{ but } d_p = d_q \end{cases}.$$

Moreover, for $p, q = 1, \dots, n$, we have

$$(3.20) \quad \frac{|(E_{j+1})_{pq}|}{\exp\left(\max\left(\int_{t_j}^{t_{j+1}} B_{pp}(t) dt, \int_{t_j}^{t_{j+1}} B_{qq}(t) dt\right)\right)} \leq (F_{j+1})_{pq} \leq \frac{|(E_{j+1})_{pq}|}{\exp\left(\min\left(\int_{t_j}^{t_{j+1}} B_{pp}(t) dt, \int_{t_j}^{t_{j+1}} B_{qq}(t) dt\right)\right)}.$$

Proof. Given the definition of F_{j+1} , we have $(F_{j+1})_{pq} = |(E_{j+1})_{pq}| \int_0^1 \frac{ds}{(D_j)_{pp}(s)(D_j)_{qq}(s)}$. Since $(D_j)_{pp}(s) = (d_p - 1)s + 1$, simple integrations give the stated form (3.19).

To get the bounds in (3.20), we proceed as follows. For $p = q$, or $p \neq q$ but $d_p = d_q$, then (3.20) is obviously sharp, since $d_p = \exp(\int_{t_j}^{t_{j+1}} B_{pp}(t) dt)$.

For $p \neq q$ and $d_p \neq d_q$, suppose $d_p > d_q$. Then, we have

$$(3.21) \quad \frac{\log(d_p) - \log(d_q)}{d_p - d_q} = \frac{1}{\exp(\int_{t_j}^{t_{j+1}} B_{qq}(t) dt)} \frac{\int_{t_j}^{t_{j+1}} (B_{pp}(t) - B_{qq}(t)) dt}{\exp(\int_{t_j}^{t_{j+1}} (B_{pp}(t) - B_{qq}(t)) dt) - 1}.$$

Now, consider the function $x/(e^x - 1)$ for $x > 0$. By expanding e^x , it is easy to see that

$$\frac{1}{e^x} \leq \frac{x}{e^x - 1} \leq 1, \quad x > 0.$$

Using this in (3.21) gives the claim in case $d_p > d_q$. Much the same argument for the case $d_q > d_p$ completes the proof. \square

Remark 3.10. As a consequence of Lemma 3.9, we see that the entries of F_{j+1} may differ significantly from those of $|E_{j+1}|$, only in case for some indices $q = 1, \dots, n$, $\int_{t_j}^{t_{j+1}} B_{qq}(t) dt \ll 0$. In general, this will betray that some Lyapunov exponents will be large, negative numbers.

Let us continue to bound the terms in (3.18).

THEOREM 3.11. *Recall the notation of (3.10) and (3.11). For all $j = 0, 1, \dots$, let*

$$(3.22) \quad \delta_j := \min_{1 \leq p \leq n} \frac{1}{\min(1, \exp(\int_{t_j}^{t_{j+1}} B_{pp}(t) dt))},$$

and let $\nu_j = \|U_{j+1}\|$ be the defect from normality³ of $R(t_{j+1}, t_j)$. Then, the following bounds hold

$$(3.23) \quad \max_{0 \leq s \leq 1} \sum_{k=0}^{n-1} \|W_j^k(s)\| \leq \frac{1 - (\delta_j \nu_j)^n}{1 - \delta_j \nu_j}$$

and also

$$(3.24) \quad \max_{0 \leq s \leq 1} \sum_{k=0}^{n-1} \|V_j^k(s)\| \leq \frac{1 - (\delta_j \nu_j)^n}{1 - \delta_j \nu_j}.$$

Finally, the following bound on $\|R_j^{-1}(s)\|$ holds as well

$$(3.25) \quad \max_{0 \leq s \leq 1} \|R_j^{-1}(s)\| \leq \delta_j \frac{1 - (\delta_j \nu_j)^n}{1 - \delta_j \nu_j}.$$

Proof. With $R_j^{-1}(s)$ given by (3.15) or (3.16), (3.25) follows immediately from (3.24) or (3.23). To witness, using (3.15), we have $\|R_j^{-1}(s)\| \leq \|D_j^{-1}(s)\| \sum_{k=0}^{n-1} \|V_j^k(s)\|$. Now, since

$$D_j^{-1}(s) = \text{diag}\left(\frac{1}{1 + s(\exp(\int_{t_j}^{t_{j+1}} B_{ii} d\tau) - 1)}, i = 1, \dots, n\right),$$

given the definition of δ_j and (3.24), then (3.25) follows.

The argument to show (3.23) and (3.24) is the same, so we only show the latter. Recall that $V_j(s) = sU_{j+1}D_j^{-1}(s)$. So, norm bounds and the triangular inequality, given the definitions of ν_j and δ_j , give

$$\sum_{k=0}^{n-1} \|V_j^k(s)\| \leq \sum_{k=0}^{n-1} (\delta_j \nu_j)^k$$

and (3.24) follows. \square

By putting together the bounds in Theorems 3.9 and 3.11, we can finally summarize our backward error statement.

THEOREM 3.12. *Consider the system (1.1). Let $\{t_j\}$, $j = 0, 1, \dots$, $t_0 = 0 < t_1 < t_2 < \dots$, be the sequence of points (converging to ∞) generated by the numerical realization of the discrete QR method.*

At each t_j , $j = 1, 2, \dots$, the exact discrete QR method delivers the factorization (2.5), where $R(t_{j+1}, t_j) = e^{h_j B_j}$, $h_j = t_{j+1} - t_j$, $j = 0, 1, \dots$, is the solution of the upper triangular system (2.19). The numerical discrete QR method, instead, gives the QR factorization of the matrix in (3.1), that is in (3.3). Assume that $\|E_{j+1}\| = \|\Phi(t_{j+1}, t_j) - X(t_{j+1}, t_j)\| \leq \text{TOL}$, for all $j = 0, 1, \dots$. Finally, let (3.13) hold, and recall the notation in (3.10-3.11-3.12) and the notation from Theorems 3.9 and 3.11.

Then, the numerical discrete QR method finds the (exact) QR factorization of the system (3.5), where

$$(3.26) \quad h_j \widehat{B}_j = h_j B_j + L_j + O(\|E_{j+1}\|^2),$$

³The meaning of *normal* here is the one in common usage in the linear algebra community, it has little to do with that of a normal fundamental matrix solution

and we have the bounds

$$(3.27) \quad |L_j| \leq \left[\frac{1 - (\delta_j \nu_j)^n}{1 - \delta_j \nu_j} \right]^2 |F_{j+1}|$$

where the general (p,q) -entry of $|F_{j+1}|$ is given in (3.19) and further bounded as in (3.20):

$$(|F_{j+1}|)_{(p,q)} \leq \text{TOL} / \min_{i=p,q} \left(\exp \left(\int_{t_j}^{t_{j+1}} B_{ii}(t) dt \right) \right),$$

δ_j is given in (3.22) and ν_j is the departure from normality of the exact triangular transition matrix $R(t_{j+1}, t_j)$. \square

Remarks 3.13.

- In a way, the above result is good news: We have obtained the exact realization of the method we wanted on a problem which we can make close to the original one by decreasing the error tolerance (or the stepsize), and this is true regardless of how long we compute. So, if we have stable exponents, we will get answers close to the true values. However, this is correct only in part. The result we found has two factors in the estimates. The first was predictable: $1/\exp(\int_{t_j}^{t_{j+1}} B_{ii}(t) dt)$ tells us that it will be harder to approximate large and negative exponents. The second factor hints at an inherent difficulty in approximating Lyapunov exponents by the QR methods when the (exact) factor R is far from diagonal; this is reflected in the quantity ν_j in (3.27). This is bothersome, since it has little to do with having stable exponents.
- Increasing the accuracy (equivalently, decreasing the stepsize) can alleviate these difficulties, specifically the one due to lack of normality. In principle, in fact, we can make the factors ν_j small, by decreasing the stepsize. This is because, see Lemma 2.4, we can write

$$\|U_{j+1}\| \leq \|R(t_{j+1}, t_j)\| = \|e^{h_j B_j}\|$$

and thus $\|U_{j+1}\| \leq \|I\| + h_j \|B_j\| + O(h_j^2)$. Thus, for h_j sufficiently small, within terms of $O(h_j^2)$, in (3.27) we can claim that $|L_j| \leq (\|I\| + 2\delta_j h_j \|B_j\| + O(h_j^2)) |F_{j+1}|$. However, it may be impossible in practice to push the stepsize to be very small while at the same time try to compute on long intervals of time. This impasse is real, and tells us that it may be hard to accurately compute (by QR methods) Lyapunov exponents for highly non-normal problems, those for which severe lack of normality is not localized.

- If there exists an integrally separated fundamental matrix solution for (1.1), then by Theorem 5.1 of [9] there exists $R(t_0)$ (see (2.5)) such that $R(t_k)$ in (2.6) approaches a diagonal matrix as $k \rightarrow \infty$. In addition, if the diagonal of B is integrally separated, see (2.13), then by Lemma 7.4 of [9] $R(t_k)$ approaches $\text{diag}(R)\bar{Z}$ where \bar{Z} is unit upper triangular, and hence \bar{Z} determines the asymptotic lack of normality.
- What is obtained in Theorem 3.12 for the non-autonomous linear system is that the computed solution is the exact solution to a nearby piecewise continuous system. This is in contrast (e.g., see [12]) to the case in backward error analysis for nonlinear autonomous problems, in which the computed solution is close to the solution of a smooth modified equation. We remark that, since the nearby problem we obtain is piecewise continuous, as opposed to continuous, the stability/continuity results for Lyapunov exponents summarized in [1, pp. 172-173] cannot be applied directly.

3.2. Continuous QR error. The situation for the continuous QR method is similar to the discrete QR case, though in some way the situation is better, since one does not have to compute the triangular factor directly. The following observations help to understand the error behavior of the continuous QR method.

Observations 3.14.

(a) Recall that the Q-factor of $X(t)$ is the same as that of (2.16). But, the Q-factor of the expression in (2.16) is the same as the Q-factor in the QR factorization of

$$(3.28) \quad \Phi(t, t_{k-1}) Q_{k-1} Q_{k-1}^T \Phi(t_{k-1}, t_{k-2}) \dots Q_1^T \Phi(t_1, 0) Q_0,$$

where Q_j , $j = 1, \dots, k-1$, are (any) orthogonal matrices. Consider (3.28). For $j = 0, 1, \dots$, define $U(t, t_j) := \Phi(t, t_j) Q_j$, and let $\widehat{Q}(t, t_j)$ be the Q-factor in the QR factorization of $U(t, t_j)$. Then, by direct differentiation, it is easy to see that $\widehat{Q}(t, t_j)$ satisfies the differential equation (2.8) with H defined as in (2.9), but with $\widehat{Q}(t, t_j)$ replacing $Q(t)$ there. In other words, $\partial_t \widehat{Q}(t, t_j) = \widehat{Q}(t, t_j) H(\widehat{Q}(t, t_j), A(t))$, and $\widehat{Q}(t_j, t_j) = Q_j$. Notice that this fact is true for any choice of matrices Q_j 's, not only if $Q_j = Q(t_j)$, $j = 1, \dots, k-1$, in (3.28), in which case we already knew it: The first equation of (2.17). In particular, this means that at any value t : $t_{k-1} \leq t \leq t_k$, $k = 1, 2, \dots$, the Q-factor in the QR factorization of $X(t)$ or –which is the same– of (3.28) can be written as

$$(3.29) \quad Q(t) = \widehat{Q}(t, t_{k-1}) [Q_{k-1}^T \widehat{Q}(t_{k-1}, t_{k-2})] \dots [Q_2^T \widehat{Q}(t_2, t_1)] [Q_1^T \widehat{Q}(t_1, t_0)].$$

(b) Now, when using the continuous QR method numerically, we find approximate values Q_j to the exact values $Q(t_j)$, $j = 1, 2, \dots$. We can assume that these Q_j 's are orthogonal matrices, since this is the case when we use any of a host of schemes which maintain orthogonality at the discrete level, see [11] and references there. It is these approximations that we must think we are using in (3.28). So, by virtue of the above point (a), on a step $t_j \leq t \leq t_{j+1}$, $j = 0, 1, 2, \dots$, instead of (2.17) we will be approximating the solution of

$$(3.30) \quad \begin{cases} \partial_t \widehat{Q}(t, t_j) &= \widehat{Q}(t, t_j) H(\widehat{Q}(t, t_j), A(t)) , & \widehat{Q}(t_j, t_j) = Q_j , \\ \partial_t \widehat{\mu}_i(t, t_j) &= [\widehat{Q}^T(t, t_j) A(t) \widehat{Q}(t, t_j)]_{ii} , & \widehat{\mu}_i(t_j, t_j) = 0 , \quad i = 1, \dots, n , \end{cases}$$

where H is defined as in (2.9) using $\widehat{Q}(t, t_j)$. Using (3.29) at $t = t_k$, and comparing with Q_k , we get that

$$Q(t_k) = Q_k [Q_k^T \widehat{Q}(t_k, t_{k-1})] \dots [Q_2^T \widehat{Q}(t_2, t_1)] [Q_1^T \widehat{Q}(t_1, t_0)]$$

where all terms in brackets are defects from the identity of the size of the local errors incurred in approximating Q : That is, if we let $\widehat{N}_j = Q_j - \widehat{Q}(t_j, t_{j-1})$ for the local errors, then $\widehat{Q}^T(t_j, t_{j-1}) Q_j = I + \widehat{Q}^T(t_j, t_{j-1}) \widehat{N}_j$. Notice that the above expression for $Q(t_k)$ shows that we cannot generally expect small global errors $Q(t_k) - Q_k$.

Now, from (3.30), we will obtain approximations Q_{j+1} and $\mu_i^c(t_{j+1}, t_j)$ instead of the exact values $\widehat{Q}(t_{j+1}, t_j)$ and $\widehat{\mu}_i(t_{j+1}, t_j)$. Then, at t_k , we will form

$$\frac{1}{t_k} \sum_{j=1}^k \mu_i^c(t_j, t_{j-1}),$$

instead of the analogous formula (2.18). The question is: How does this formula compare to (2.18)? We give an answer to this question in the same flavor of what we did for the discrete QR method, that is we argue that essentially Theorem 3.1 holds.

We need to get around the fact that with the continuous QR method we do not find the triangular factor, only the orthogonal one, and the diagonal of the triangular factor. So, there is some extra freedom in choosing a triangular factor amongst all those with given diagonal. The next Lemma gives two possibilities.

LEMMA 3.15. *Let $\eta > 0$ be given, $\eta \ll 1$. For $j = 1, 2, \dots$, let $\widehat{Q}(t_j, t_{j-1}) \widehat{R}(t_j, t_{j-1})$ be the exact QR factorization of $\Phi(t_j, t_{j-1}) Q_{j-1}$, and let Q_j be the numerical approximation to $\widehat{Q}(t_j, t_{j-1})$; see*

(3.30). Assume that $Q_j - \widehat{Q}(t_j, t_{j-1}) = O(\eta)$. Let R_j 's be upper triangular matrices whose strictly upper triangular part is not determined, but whose diagonal is assigned and satisfies $\text{diag}(R_j - \widehat{R}(t_j, t_{j-1})) = O(\eta)$.

By choosing either one of (a) or (b) below for the strictly upper part of the R_j 's, we obtain that

$$Q_j R_j - \widehat{Q}(t_j, t_{j-1}) \widehat{R}(t_j, t_{j-1}) = O(\eta).$$

(a) “Frobenius-optimal”. Choose R_j to minimize the Frobenius norm $\|\widehat{Q}(t_j, t_{j-1}) \widehat{R}(t_j, t_{j-1}) - Q_j R_j\|_F$.

(b) “Upper-exact”. For $p = 1, \dots, n-1$, $q = p+1, \dots, n$, take $(R_j)_{pq} = (\widehat{R}(t_j, t_{j-1}))_{pq}$.

Proof. The fact that choosing (b) we obtain $Q_j R_j - \widehat{Q}(t_j, t_{j-1}) \widehat{R}(t_j, t_{j-1}) = O(\eta)$ is clear, since we would have $R_j = \text{diag}(R_j - \widehat{R}(t_j, t_{j-1})) + \widehat{R}(t_j, t_{j-1})$.

Since the choice (a) would be optimal, the $O(\eta)$ estimate will also hold for (a), as long as (a) is solvable. We show this next. We have

$$\|\widehat{Q}(t_j, t_{j-1}) \widehat{R}(t_j, t_{j-1}) - Q_j R_j\|_F^2 = \|\widehat{R}(t_j, t_{j-1}) - (E + I)R_j\|_F^2$$

where we have set $E = \widehat{Q}^T(t_j, t_{j-1})Q_j - I$, and hence $E = O(\eta)$. Now,

$$\|\widehat{R}(t_j, t_{j-1}) - (E + I)R_j\|_F^2 = \text{trace}[(\widehat{R}(t_j, t_{j-1}) - (E + I)R_j)^T(\widehat{R}(t_j, t_{j-1}) - (E + I)R_j)]$$

which is the same as

$$\sum_{p=1}^n e_p^T [R_j^T R_j + \widehat{R}^T(t_j, t_{j-1}) \widehat{R}(t_j, t_{j-1}) - 2\widehat{R}^T(t_j, t_{j-1})(E + I)R_j] e_p,$$

where e_p are the standard unit vectors.

Differentiating each term in the sum and setting the derivatives to zero lead to the optimal solution

$$(R_j)_{1:p-1,p} = [(I + E^T)(\widehat{R}(t_j, t_{j-1}))_{:,p}]_{1:p-1}.$$

□

There may be other useful possibilities for “completing” the matrices R_j 's of Lemma 3.15. However, the basic fact remains: When we use the continuous QR method we do not find a triangular factor of an associated transition matrix, but only its diagonal.

THEOREM 3.16. For $j = 1, 2, \dots$, let Q_j be the numerical approximations generated by the continuous QR method to the exact Q -factors in the QR factorizations of $X(t_j)$, let $\mu_i^c(t_j, t_{j-1})$, $i = 1, \dots, n$, be the approximations to the values $\widehat{\mu}_i(t_j, t_{j-1})$ in (3.30), and let $\widehat{Q}(t_j, t_{j-1}) \widehat{R}(t_j, t_{j-1})$ be the exact QR factorization of $\Phi(t_j, t_{j-1})Q_{j-1}$. Let \widehat{N}_j be the local error in Q : $\widehat{N}_j = Q_j - \widehat{Q}(t_j, t_{j-1})$. Finally, let R_j 's be upper triangular matrices with diagonal given by the exponential of

$$(3.31) \quad \log(\text{diag}(R_j)) = (\mu_1^c(t_j, t_{j-1}), \dots, \mu_n^c(t_j, t_{j-1})),$$

and otherwise the R_j 's are chosen so that the entries of R_j differ from those of $\widehat{R}(t_j, t_{j-1})$ by the same order of magnitude as the differences on the diagonal (e.g., this is the case for (a) and (b) of Lemma 3.15). Let $\Delta_j := R_j - \widehat{R}(t_j, t_{j-1})$.

Then, the continuous QR method is giving the QR factorization of $X(t_j, t_{j-1})Q_{j-1}$, where $X(t_j, t_{j-1})$ approximates $\Phi(t_j, t_{j-1})$ with error of the same norm as that of the local errors in Q_j and R_j :

$$X(t_j, t_{j-1}) - \Phi(t_j, t_{j-1}) = N_j,$$

where we have set

$$N_j = \widehat{N}_j \widehat{R}(t_j, t_{j-1}) Q_{j-1}^T + \widehat{Q}(t_j, t_{j-1}) \Delta_j Q_{j-1}^T + \widehat{N}_j \Delta_j.$$

Proof. We define $X(t_j, t_{j-1})$ indirectly from the relation

$$X(t_j, t_{j-1}) Q_{j-1} = Q_j R_j.$$

Then, since

$$Q_j R_j = (\widehat{Q}(t_j, t_{j-1}) + \widehat{N}_j) (\widehat{R}(t_j, t_{j-1}) + \Delta_j),$$

the result follows. \square

Using Theorem 3.16, we can now invoke the same results as we had for the discrete QR method. In particular, Theorem 3.1 all the way to Theorem 3.12 hold, and the continuous QR method is therefore delivering the QR factorization (or, better, the Q-factor and the logarithm of the diagonal of R) of the fundamental matrix solution of a problem close to the original one, in the sense of the backward error statement of Theorem 3.12.

Remark 3.17. We do not want to give the impression that discrete and continuous QR methods are one and the same. In practice, they are not, but at the theoretical level they allow for a unified treatment insofar as error analysis. The chief reason for their practical difference is that the differential equations one ends up solving have different stability types, and when using these methods with variable stepsizes one controls different local errors! With the discrete QR method one controls directly the error on the transition matrices, and indirectly obtains a control on the R-factor (and Q-factor, though it is not used). In the continuous QR method, instead, one controls directly the local error in the Q-factor of the transition matrices, and essentially the local error in the R-factor, and only indirectly the error on the transition matrices. Often, it is easier to control the error of the Q factor, and of $\text{diag}(\log(R))$, than it is to control the error on the local transition matrices. As a consequence, when choosing the stepsize by enforcing a local error control, one may end up taking fewer steps with the continuous QR method than with the discrete QR method.

3.3. Other Errors. As we said already, there are other obvious sources of errors if one tries to approximate Lyapunov exponents.

To begin with, it is impossible in general to give sharp results on the error one commits when truncating time, unless some extra assumptions are placed on the function A , such as some form of recurrence. This is simply because A may change completely its character past the time where one

stops computing. For example, just take a scalar problem with $A(t) = \begin{cases} 0, & t \leq T \\ t - T, & T \leq t \leq T + 1 \\ 1, & T + 1 < t \end{cases}$.

Any computation which stops at, before, or shortly after T , is bound to give wrong results insofar as the Lyapunov exponents.

Furthermore, it is important to appreciate that the asymptotic behavior can be approached in many different ways, also within the same system. For example, it may take a short time to approximate some exponents, and a long time to approximate some other ones, regardless of the integral separation in the system. To witness, we can consider the following example. Take a diagonal system, where (for t large) A is of the form $A_{ii}(t) = i - 1 + \frac{1}{t^{1/p_i}}$, $p_i > 0$, $i = 1, 2, \dots, n$. Clearly, the exponents are $\{0, 1, 2, \dots, n - 1\}$, and the system is regular and integrally separated. However, if $0 < p_i \ll 1$ the exponent is reached very quickly, whereas if $p_i \gg 1$ it is approached very slowly.

Also, we have not considered the errors induced by the finite precision computation. This is (at least in part) justified by the fact that linear algebra errors can be incorporated with the integration errors. For example, in (3.4), the terms N_j can be made to comprise both the local errors obtained when approximating $\Phi(t_j, t_{j-1})$ by $X(t_j, t_{j-1})$ and the finite precision errors arising from the QR factorization of $X(t_j, t_{j-1})$.

3.4. Nonlinear case. The setup now is the following. We have the nonlinear problem

$$(3.32) \quad \dot{x} = f(x), \quad x(0) = x_0,$$

with flow (solution) $\phi^t(x_0)$. So, we would like to obtain the Lyapunov exponents for the linear variational problem

$$(3.33) \quad \dot{X} = Df(\phi^t(x_0))X, \quad \text{or} \quad \dot{X} = A(t)X,$$

subject to some initial conditions X_0 . By the exact QR techniques on this problem (recall (2.5)), we would thus find a sequence of t -values $\{t_j\}$ such that

$$(3.34) \quad X(t_k) = Q(t_k)R(t_k, t_{k-1}) \cdots R(t_1, t_0)R_0.$$

In practice, we will have a numerical approximation to the flow $\phi^t(x_0)$, call it $\psi_h^\tau(x_0)$. [We can think of $\psi_h^\tau(x_0)$ as being defined at grid points by the numerical scheme, and everywhere else by some interpolation process]. Here we highlight the fact that the exact flow and the approximate flow may not be comparable on the same time scale as in the case of approximation of a hyperbolic periodic orbit or more general hyperbolic attractor.

Thus, we will end up attempting to approximate the spectra of

$$(3.35) \quad Y' = Df(\psi_h^\tau(x_0))Y, \quad \text{or} \quad Y' = C(\tau)Y, \quad " \equiv \frac{d}{d\tau}$$

and by the exact QR techniques on this problem, we would find a sequence of values $\{\tau_j\}$ such that

$$(3.36) \quad Y(\tau_k) = Z(\tau_k)U(\tau_k, \tau_{k-1}) \cdots U(\tau_1, \tau_0)R_0.$$

Now, assume that there exist smooth monotone functions $w_j(t)$ such that

$$w_j(t_j) = \tau_j, \quad w_j(t_{j+1}) = \tau_{j+1}$$

and define $w(t) = w_j(t)$ for $t \in [t_j, t_{j+1})$. Assume in addition that for all $t \geq 0$, there exists $\epsilon(t) > 0$ and $\delta > 0$ such that

$$(3.37) \quad \begin{aligned} (a) \quad & |(w(t_{j+1}) - w(t_j)) - (t_{j+1} - t_j)| \leq \delta, \\ (b) \quad & \|\phi^t(\tilde{x}_0) - \psi_h^{w(t)}(x_0)\| \leq \epsilon(t), \end{aligned}$$

where \tilde{x}_0 is some initial condition (IC), which is the same for all $t \geq 0$. We remark that the assumptions in (3.37) hold in the context of continuous shadowing with rescaling of time (e.g., see [18]) in which case δ and $\epsilon(t)$ are proportional to the absolute local error in approximating the nonlinear differential equation.

Now, relative to the ICs \tilde{x}_0 , we will also have the linearized problem

$$(3.38) \quad \dot{\tilde{X}} = Df(\phi^t(\tilde{x}_0))\tilde{X}, \quad \text{or} \quad \dot{\tilde{X}} = \tilde{A}(t)\tilde{X},$$

and by the exact QR techniques on this problem at the t -values $\{t_j\}$ we have

$$(3.39) \quad \tilde{X}(t_k) = \tilde{Q}(t_k)\tilde{R}(t_k, t_{k-1}) \cdots \tilde{R}(t_1, t_0)\tilde{R}_0.$$

So, quite clearly, there are two aspects to consider.

- (1) Firstly, there is the error caused by the difference between the two linear problems (3.35) and (3.38). This is in essence an issue of comparing the functions \tilde{A} and C .

(2) Then, there is the issue of the difference between \tilde{A} and A . In general, \tilde{A} and A will not be close (at least, not in a pointwise sense). However, under some important circumstances, A and \tilde{A} will lead to spectra which are close to one another. For this to happen, one needs some measure of ergodicity.

We have the following result, which addresses (1) above.

THEOREM 3.18. *Suppose (3.37) hold, and let f be \mathcal{C}^k , $k \geq 2$. Then we have*

$$(3.40) \quad \|\tilde{A}(t) - C(w(t))\| \leq M \epsilon(t).$$

where M is a bound on the second derivative f_{xx} evaluated along the path $\phi^t(\tilde{x}_0) - s(\phi^t(\tilde{x}_0) - \psi_h^{w(t)}(x_0))$ for $0 \leq s \leq 1$ and $t \geq 0$.

Proof. We only need to notice that (3.37-b) implies a similar inequality for the Jacobians. In fact, from the mean value theorem in integral form we have that

$$\begin{aligned} & f_x(\phi^t(\tilde{x}_0)) - f_x(\psi_h^{w(t)}(x_0)) \\ &= \int_0^1 f_{xx}(\phi^t(\tilde{x}_0) - s(\phi^t(\tilde{x}_0) - \psi_h^{w(t)}(x_0))) ds (\phi^t(\tilde{x}_0) - \psi_h^{w(t)}(x_0)). \end{aligned}$$

so that (3.40) follows.

□

Remarks 3.19.

- (1) The bound (3.40) together with a stability/continuity result for Lyapunov exponents (see e.g. [1]) implies that for $M\epsilon(t)$ small enough, uniformly in t , the Lyapunov exponents of (3.35) and (3.38) are close.
- (2) With an assumption of ergodicity [16] we have that with probability one the Lyapunov exponents of (3.33) and (3.38) are the same (if x_0 and \tilde{x}_0 both lie on the compact invariant set with an ergodic probability measure) and thus a stability result for Lyapunov exponents allows (with probability one) a comparison of the Lyapunov exponents of (3.33) and (3.35) using the previous remark (1).

4. An Example: Numerical Results.

We build an example where we vary the departure from normality of the exact triangular factor. Take the following upper triangular function $B(t) = D(t) + U(t)$, with

$$(4.1) \quad D(t) = \text{diag}(D_{11}(t), D_{22}(t), D_{33}(t), D_{44}(t)),$$

where we take $D_{11}(t) = 1$, $D_{22}(t) = \cos(t)$, $D_{33}(t) = -\frac{1}{\sqrt{t+1}}$, $D_{44}(t) = -10$, and

$$(4.2) \quad U(t) = \alpha \begin{pmatrix} 0 & \cos(t) & \sin(t) & \cos(t) \\ 0 & 0 & \cos(t) & \sin(t) \\ 0 & 0 & 0 & \cos(t) \\ 0 & 0 & 0 & 0 \end{pmatrix},$$

and we will be interested in the two cases of $\alpha = 0$ and $\alpha = 10^4$.

Then, we rotate B , and consider the linear system (1.1) with

$$A(t) = Q(t)B(t)Q^T(t) + \dot{Q}(t)Q^T(t),$$

and

$$Q(t) = \text{diag}(1, Q_\beta(t), 1) \cdot \text{diag}(Q_\eta(t), Q_\eta(t)).$$

We set

$$Q_\gamma(t) = \begin{pmatrix} \cos(\gamma t) & \sin(\gamma t) \\ -\sin(\gamma t) & \cos(\gamma t) \end{pmatrix}, \quad \eta = 1, \quad \beta = \sqrt{2}.$$

$$\alpha = 0 \text{ in (4.2), } T = 10^3.$$

TOL	Method	Steps	e_1	e_2	e_3	e_4
1.E-3	Disc	12750	2.66E-7	2.53E-7	2.42E-8	2.20E-3
1.E-6	Disc	47248	3.28E-10	3.17E-10	3.10E-11	1.42E-6
1.E-9	Disc	185519	1.48E-12	3.35E-13	3.28E-14	1.19E-9
1.E-3	Cont	5962	4.00E-4	5.72E-4	1.31E-4	1.10E-3
1.E-6	Cont	21328	1.26E-7	1.77E-7	4.36E-8	3.46E-7
1.E-9	Cont	82592	3.16E-11	4.62E-11	1.19E-11	8.98E-11

TABLE 4.1
Error in QR methods when the triangular factor is diagonal.

$$\alpha = 10^4 \text{ in (4.2), } T = 10^3.$$

TOL	Method	Steps	e_1	e_2	e_3	e_4
1.E-3	Disc	537164	1.06E+0	6.34E-1	3.11E+0	2.93E+0
1.E-6	Disc	1606256	2.45E-3	2.95E-2	4.42E-3	2.30E-2
1.E-9	Disc	5703057	2.83E-4	6.10E-4	8.89E-4	4.50E-6
$h=1.E-3$	Disc	1000000	1.50E-1	3.71E-1	9.10E-1	1.14E+0
$h=1.E-4$	Disc	10000000	2.58E-01	2.30E-1	4.90E-1	2.01E-3
1.E-3	Cont	334986	8.26E+0	9.78E+0	2.50E+0	1.55E+1
1.E-6	Cont	1032455	4.25E-2	4.58E-2	1.69E-3	1.58E-3
1.E-9	Cont	4087009	5.95E-4	3.52E-4	2.46E-4	4.00E-6
$h=1.E-3$	Cont	1000000	7.28E+0	1.01E+1	8.64E+0	8.70E+0
$h=1.E-4$	Cont	10000000	2.62E-1	1.03E+0	8.05E-1	4.85E-1

TABLE 4.2
Error in QR methods when the triangular factor is highly non-diagonal.

Regardless of the value of α in (4.2), this is a regular system with stable Lyapunov exponents given by the limits of $\lambda_i(t) := \frac{1}{t} \int_0^t D_{ii}(s) ds$, $i = 1, 2, 3, 4$, that is: $\{1, 0, 0, -10\}$.

All results on this problem have been obtained using the code **leslis**, which we wrote and is public domain: <http://www.math.gatech.edu/~dieci>

In particular, we use the continuous QR method using the projected 5th order scheme (**IPAR(8)=0** in **LESLIS**), with error control on the Q-factor and the μ_i 's in (2.15), and the 5th order discrete QR method (**IPAR(8)=4**) with error control on the μ_i 's. In all cases, we call **TOL** the required tolerance values.

In Tables 4.1 and 4.2 we report on selected experiments in the case of $\alpha = 0$, respectively $\alpha = 10^4$, in (4.2). The computations of Tables 4.1 and 4.2 have been carried out up to $T = 10^3$, and we show the error between the finite-time computed and finite-time exact Lyapunov exponents, $e_i(T) := |\lambda_i(T) - \lambda_i^c(T)|$ where $\lambda_i^c(T)$ are the computed values at T of $\lambda_i(T)$, $i = 1, 2, 3, 4$. The heading **Method** refers to either the Discrete or Continuous QR method. Scientific notation is used throughout. In Table 4.2, we also report on some experiments made with constant stepsize, under the heading of **TOL**.

From Table 4.1, it is apparent that both discrete and continuous QR methods deliver approximations with the same error as the error tolerances. This is in agreement with the theory, of course, since the triangular factor is actually diagonal. Incidentally, we also observe that the continuous QR method takes fewer steps and is actually less expensive (and at least as accurate) than the discrete QR method.

The results in Table 4.2 show unarguably the impact of lack of normality for the upper triangular matrix solution. There is a deterioration in accuracy with respect to the required error

$$\alpha = 10^4 \text{ in (4.2), TOL=1.E-6.}$$

T	Method	Steps	e ₁	e ₂	e ₃	e ₄
10 ⁴	Disc	16059751	1.07E-2	6.3E-3	1.83E-2	2.31E-2
10 ⁵	Disc	160823338	7.63E-2	1.51E-1	2.50E-1	2.31E-2
10 ⁶	Disc	1618405426	6.27E-1	2.66E-1	9.34E-1	4.18E-2
10 ⁴	Cont	10346418	5.26E-2	4.91E-2	4.98E-3	1.46E-3
10 ⁵	Cont	104059732	9.30E-3	2.21E-1	2.33E-1	2.31E-3

TABLE 4.3
Error when the triangular factor is non-diagonal, in function of T.

tolerances which is proportional to the departure from normality of the triangular factor. This is in agreement with our backward error results, in particular with (3.27) and Remark 3.13. We notice that –with respect to Table 4.1– we are taking far more steps for the same error tolerance. This is an indication that the methods try to compensate for the departure of normality by restricting the stepsize, see Remark 3.13. For the record, the continuous QR method again takes fewer steps than the discrete QR method, but in the end (for this case) it takes a longer time since each step is more expensive. We also notice that the methods with constant stepsize perform much worse for comparable cost; e.g., compare the results obtained with $h = 1.E-4$ with those obtained with $\text{TOL} = 1.E-9$.

Finally, in Table 4.3 we show that there is really no appreciable difference resulting from increasing the length of the interval, i.e., T . In fact, the errors remain of the same order as the departure from normality times TOL . These results in Table 4.3 were all obtained with $\text{TOL}=1.E-6$.

REFERENCES

- [1] L. Ya. Adrianova, *Introduction to Linear Systems of Differential Equations*, Translations of Mathematical Monographs Vol. 146, AMS, Providence, R.I. (1995).
- [2] L. Arnold and V. Wihstutz Eds. *Lyapunov Exponents. Proceedings, Bremen 1984*. Springer-Verlag, Berlin, 1986. Lecture Notes in Mathematics 1186.
- [3] L. Arnold, H. Crauel, and J.P. Eckmann Eds. *Lyapunov Exponents. Proceedings, Oberwolfach 1990*. Springer-Verlag, Berlin, 1991. Lecture Notes in Mathematics 1486.
- [4] G. Benettin, L. Galgani, A. Giorgilli and J.-M. Strelcyn, “Lyapunov Exponents for Smooth Dynamical Systems and for Hamiltonian Systems; A Method for Computing All of Them. Part 1: Theory”, and “...Part 2: Numerical Applications”, *Meccanica* **15** (1980), pp. 9-20, 21-30.
- [5] P. Constantin and C. Foias, “Global Lyapunov exponents, Kaplan-Yorke formulas and the dimension of the attractors for 2D Navier-Stokes equations,” *Comm. Pure Applied Mathematics* **38** (1985), pp. 1–27.
- [6] L. Dieci, B. Morini, A. Papini, and A. Pasquali. “On real logarithms of nearby matrices and structured matrix interpolation”, *Appl. Numer. Math.*, 19:145–165, 1999.
- [7] L. Dieci, R. D. Russell, and E. S. Van Vleck. “On the computation of Lyapunov exponents for continuous dynamical systems”, *SIAM J. Numer. Anal.*, 34:402–423, 1997.
- [8] L. Dieci and E.S. Van Vleck, “Lyapunov and other spectra: a survey,” *Preservation of Stability under Discretization*, D. Estep and S. Tavener Ed.s, SIAM Publications (2002).
- [9] L. Dieci and E.S. Van Vleck, “Lyapunov spectral intervals: theory and computation,” *SIAM J. Numer. Anal.*, 40:516–542, 2003.
- [10] L. Dieci and E.S. Van Vleck, “Lyapunov and Sacker-Sell spectral intervals,” (2003) submitted.
- [11] L. Dieci and E.S. Van Vleck, “LESLIS and LESLIL: Codes for approximating Lyapunov exponents of linear systems”, Technical Report (2004): www.math.gatech.edu/~dieci.
- [12] E. Hairer, C. Lubich, “The life-span of backward error analysis for numerical integrators,” *Numer. Math.* **76** (1997), pp. 441–462.
- [13] E. McDonald and D. Higham, “Error analysis of QR algorithms for computing Lyapunov exponents”, *ETNA*, 12:234–251, 2001.
- [14] V.M. Millionshchikov, “Systems with integral division are everywhere dense in the set of all linear systems of differential equations,” *Differentsial'nye Uravneniya* **5** (1969), pp. 1167–1170.
- [15] S. Oliveira and D.E. Stewart. “Exponential splitting of products of matrices and accurately computing singular values of long products”, *LAA*, 309:175–190, 2000.

- [16] V. I. Oseledec, “A multiplicative ergodic theorem. Lyapunov characteristic numbers for dynamical systems”, *Trans. Moscow Mathem. Society*, 19:197, 1968.
- [17] K.J. Palmer, “The structurally stable systems on the half-line are those with exponential dichotomy,” *J. Diff. Eqn.* **33** (1979), pp. 16–25.
- [18] K.J. Palmer, *Shadowing in Dynamical Systems*, Mathematics and Its Applications, Vol. 501, Kluwer, 2000.
- [19] D. Ruelle. *Chaotic evolution and strange attractors*. Cambridge University Press, Cambridge, 1989.
- [20] D. E. Stewart. “A new algorithm for the SVD of a long product of matrices and the stability of products”, *ETNA*, 5:29–47, 1997.
- [21] L. N. Trefethen and D. Bau, *Numerical Linear algebra*. Society for Industrial and Applied Mathematics (SIAM), Philadelphia, PA, 1997.