

## Stability and convergence of numerical methods for ODEs

Consider the initial value problem for a system of ODEs

$$\begin{cases} \frac{d\mathbf{y}}{dt} = \mathbf{f}(\mathbf{y}, t) \\ \mathbf{y}(0) = \mathbf{y}_0 \end{cases} \quad (1)$$

and the perturbed problem

$$\begin{cases} \frac{dz}{dt} = \mathbf{f}(\mathbf{z}, t) + \boldsymbol{\delta}(t) \\ \mathbf{z}(0) = \mathbf{y}_0 + \boldsymbol{\delta}_0 \end{cases} \quad (2)$$

where  $\boldsymbol{\delta}(t)$  is an integrable function and  $\boldsymbol{\delta}_0 \in \mathbb{R}^n$ . Note that we replaced  $\mathbf{y}(t)$  with  $\mathbf{z}(t)$  in (2) to emphasize the fact that the solutions to (1) and (2) are (in general) different.

**Definition 1** (Stability of the Cauchy problem (see [1, 3])). The Cauchy problem (1) is said to be *stable* within the time interval  $[0, T]$  if for any perturbations  $\boldsymbol{\delta}_0$  and  $\boldsymbol{\delta}(t)$  such that<sup>1</sup>

$$\|\boldsymbol{\delta}_0\| \leq \epsilon, \quad \text{and} \quad \|\boldsymbol{\delta}(t)\| \leq \epsilon \quad \text{for all } t \in [0, T] \quad (3)$$

we have that

$$\|\mathbf{z}(t) - \mathbf{y}(t)\| \leq C\epsilon, \quad \text{for all } t \in [0, T], \quad (4)$$

where  $C$  is a finite constant that does not depend on  $\epsilon$ .

Based on this definition, the Cauchy problem (1) is “stable” if the difference between the solutions of (1) and (2) is bounded in  $[0, T]$  after we introduce a small perturbation  $\boldsymbol{\delta}_0$  in the initial condition  $\mathbf{y}_0$  and a perturbation  $\boldsymbol{\delta}(t)$  in  $\mathbf{f}(\mathbf{y}, t)$ . The definition of stability also implies that the difference between the solutions of (1) and (2) goes to zero as  $\epsilon \rightarrow 0$ . In fact, from (4) it follows that

$$\lim_{\epsilon \rightarrow 0} \|\mathbf{z}(t) - \mathbf{y}(t)\| \leq C \lim_{\epsilon \rightarrow 0} \epsilon = 0. \quad (5)$$

The constant  $C$  appearing in (4) may not be small. This is consistent with the fact that a small perturbations in the Cauchy problem (1) can introduce large perturbations in its solution.

Hereafter we show that any well-posed initial value problem (1) is stable, i.e., robust to perturbations in the limit of small perturbations.

**Theorem 1.** Let  $D \subseteq \mathbb{R}^n$  be an open set,  $\mathbf{y}_0 \in D$ . If  $\mathbf{f}(\mathbf{y}, t)$  is Lipschitz continuous in  $D$  and  $\boldsymbol{\delta}(t)$  is integrable then the initial value problem (1) is stable.

*Proof.* We need to show that for any  $\boldsymbol{\delta}_0$  and  $\boldsymbol{\delta}(t)$  the difference between the solutions of (1) and (2) is bounded in some time interval  $[0, T]$  and that the difference goes to zero as we send  $\epsilon$  to zero. First of all, we notice that if  $D$  is open and  $\epsilon$  is small enough then the initial condition  $(\mathbf{y}_0 + \boldsymbol{\delta}_0)$  is in  $D$ . If  $\mathbf{f}(\mathbf{y}, t)$  is Lipschitz continuous

$$\|\mathbf{f}(\mathbf{z}, t) - \mathbf{f}(\mathbf{y}, t)\| \leq L \|\mathbf{z} - \mathbf{y}\| \quad \forall \mathbf{z}, \mathbf{y} \in D, \quad (6)$$

<sup>1</sup>Note that in (3) we are bounding  $\boldsymbol{\delta}_0$  and  $\boldsymbol{\delta}(t)$  with the same constant  $\epsilon$ . Such a constant coincides with the radius of the largest ball centered at zero in  $\mathbb{R}^n$  that includes both  $\boldsymbol{\delta}_0$  and  $\boldsymbol{\delta}(t)$  (for all  $t \in [0, T]$ ).

and  $\delta(t)$  is integrable, then we have existence and uniqueness of the solution to both problems (1) and (2). Such problems can be equivalently written as

$$\mathbf{y}(t) = \mathbf{y}_0 + \int_0^t \mathbf{f}(\mathbf{y}(s), s) ds \quad (7)$$

$$\mathbf{z}(t) = \mathbf{y}_0 + \delta_0 + \int_0^t \mathbf{f}(\mathbf{z}(s), s) ds + \int_0^t \delta(s) ds \quad (8)$$

for all  $t \in [0, T]$ , where  $T$  is the smallest “exit time”, i.e., the time in which either  $\mathbf{y}(t)$  or  $\mathbf{z}(t)$  get out of  $D$ . Subtracting (7) from (8) and taking the norm yields

$$\begin{aligned} \|\mathbf{z}(t) - \mathbf{y}(t)\| &= \left\| \delta_0 + \int_0^t [\mathbf{f}(\mathbf{z}(s), s) - \mathbf{f}(\mathbf{y}(s), s)] ds + \int_0^t \delta(s) ds \right\| \\ &\leq \|\delta_0\| + \int_0^t \|\mathbf{f}(\mathbf{z}(s), s) - \mathbf{f}(\mathbf{y}(s), s)\| ds + \int_0^t \|\delta(s)\| ds \\ &\leq (1+t)\epsilon + L \int_0^t \|\mathbf{z}(s) - \mathbf{y}(s)\| ds \end{aligned} \quad (9)$$

where we used the triangle inequality, the inequalities (3), and the definition of Lipschitz continuity (6). At this point we use Grönwall’s inequality<sup>2</sup> to conclude that

$$\begin{aligned} \|\mathbf{z}(t) - \mathbf{y}(t)\| &\leq (1+t)e^{tL}\epsilon \\ &\leq \underbrace{(1+T)e^{TL}}_C \epsilon. \end{aligned} \quad (13)$$

This proves that the Cauchy problem (1) is stable. Note that the constant  $C$  appearing in (13) does not depend on  $\epsilon$  and it can be very big, depending on the Lipschitz constant  $L$  and the integration time  $T$  (integration time).

□

*Remark:* If we replace the initial value problem (1) by a numerical scheme we introduce errors that can accumulate in time. Such errors can be considered as perturbations in the ODE (1). Just think about representing the discrete numerical solution  $\mathbf{u}_k$  in terms of some local interpolant  $\mathbf{z}(t)$  (differentiable in time) and substituting it into (1). This yields an ODE in the form (2). If the initial value problem (1) is not stable, i.e., robust to small perturbations, then there is no hope for any numerical method to approximate the solution.

**Stability and zero-stability of numerical methods for ODEs.** The concept of stability we discussed in previous section for continuous-time dynamical systems can be extended to discrete-time dynamical

<sup>2</sup>The Grönwall’s inequality (see [3, Lemma 11.1]) can be stated as follows. Let  $u(t)$ ,  $g(t)$  and  $p(t)$  such that

$$u(t) \leq g(t) + \int_0^t p(s)u(s) ds. \quad (10)$$

If  $g(t)$  is non-decreasing and  $p(s)$  is strictly positive then

$$u(t) \leq g(t) \exp \left[ \int_0^t p(s)u(s) ds \right]. \quad (11)$$

In the case of equation (9) we have

$$u(t) = \|\mathbf{z}(t) - \mathbf{y}(t)\|, \quad g(t) = (1+t)\epsilon, \quad p(s) = L. \quad (12)$$

systems, i.e., to numerical schemes aiming at computing the approximate solution of the initial value problem (1). we have seen that such schemes can be written in the general form<sup>3</sup>

$$\begin{cases} \sum_{j=0}^q \alpha_j \mathbf{u}_{k+j} = \Delta t \Phi_{\mathbf{f}}(\mathbf{u}_{k+q}, \dots, \mathbf{u}_k, t_k, \Delta t), \\ \text{given } \{\mathbf{u}_0, \dots, \mathbf{u}_{q-1}\} \end{cases} \quad (14)$$

where  $\Phi_{\mathbf{f}}$  is some iteration function. By taking a perturbation of the initial condition  $\mathbf{u}_0$  and a “time-dependent” perturbation of  $\Phi_{\mathbf{f}}$  in (14) we obtain

$$\begin{cases} \sum_{j=0}^q \alpha_j \mathbf{z}_{k+j} = \Delta t [\Phi_{\mathbf{f}}(\mathbf{z}_{k+q}, \dots, \mathbf{z}_k, t_k, \Delta t) + \boldsymbol{\delta}_{k+q}], \\ \text{given } \{\mathbf{z}_0 = \mathbf{u}_0 + \boldsymbol{\delta}_0, \dots, \mathbf{z}_{q-1} = \mathbf{u}_{q-1} + \boldsymbol{\delta}_{q-1}\} \end{cases} \quad (15)$$

where  $\{\boldsymbol{\delta}_0, \dots, \boldsymbol{\delta}_{k+q}, \dots\}$  is a sequence of vectors in  $\mathbb{R}^n$  bounded by some constant  $\epsilon$ , i.e.,

$$\|\boldsymbol{\delta}_j\| \leq \epsilon \quad \text{for all } j = 0, 1, \dots \quad (16)$$

The perturbations  $\boldsymbol{\delta}_j$  can arise, e.g., because of round-off or truncation errors when performing floating point operations using double precision arithmetic. Clearly, the orbits generated by (14) and (15), i.e.,

$$\{\mathbf{u}_0, \dots, \mathbf{u}_N\} \quad \text{and} \quad \{\mathbf{z}_0, \dots, \mathbf{z}_N\} \quad (17)$$

are (in general) different. For any given iteration function  $\Phi_{\mathbf{f}}$  and any given  $\Delta t$  we can provide a definition of stability for the numerical scheme (14) that closely resembles Definition 1 for continuous time dynamical systems. To this end, let  $T$  be the period of integration, and  $N$  be the number of time steps, i.e.,

$$\Delta t^* = \frac{T}{N} \quad (18)$$

Of particular interest when performing convergence analysis, is the behavior of the scheme for small  $\Delta t$ , i.e., for all  $\Delta t$  smaller than  $\Delta t^*$ .

**Definition 2** (Zero-stability). We say that the numerical scheme (14) is *zero-stable* if there exists a  $\Delta t^* > 0$  such that for all  $\Delta t \leq \Delta t^*$  and for any perturbations  $\boldsymbol{\delta}_j$  ( $j = 0, \dots, N$ ) such that

$$\|\boldsymbol{\delta}_k\| \leq \epsilon, \quad \text{for all } j = 0, \dots, N \quad (19)$$

we have that

$$\|\mathbf{z}_k - \mathbf{u}_k\| \leq C\epsilon, \quad \text{for all } j = 0, \dots, N, \quad (20)$$

where  $\mathbf{u}_k$  and  $\mathbf{z}_k$  are defined by (14) and (15), and  $C$  is a finite constant that does not depend on  $\epsilon^4$ .

The definition “zero-stability” follows from the fact that we require  $\|\mathbf{z}_k - \mathbf{u}_k\| \leq C\epsilon$  for all  $\Delta t \leq \Delta t^*$ , and in particular for  $\Delta t \rightarrow 0$ . Hence the “zero” part in “zero-stability” refers to the stability of the scheme in the limit  $\Delta t \rightarrow 0$ .

- Zero stability is a property of the numerical scheme, not of the ODE system (1). We have seen, in fact, that a well-posed Cauchy problem is always stable.

<sup>3</sup>In (14) we set  $\alpha_q = 1$  to remove the non-uniqueness of  $\alpha_j$  and  $\beta_j$  due to possible rescaling by a constant.

<sup>4</sup>The constant  $C$  in (20) may depend also on  $T$ ,  $\Delta t$  or other constants, but it cannot depend on  $k$  or  $\epsilon$ .

- Numerical methods that are not zero-stable have no hope to reliably approximate the solution of (1). In fact, even if the method is consistent, i.e., if the truncation error goes to zero as  $\Delta t \rightarrow 0$ , we have that perturbations due to finite-arithmetic may rapidly propagate in schemes that are not zero-stable, and therefore generate instabilities. In other words, consistent schemes that are not zero-stable may not converge as  $\Delta t \rightarrow 0$ . For example, the numerical scheme in equation (32) hereafter is consistent but not zero stable. Another example of a consistent scheme that is not zero stable is discussed in [2, p. 32].

**The root condition and zero-stability.** The numerical method (14) is said to satisfy the *root condition* if all roots of the first characteristic polynomial

$$\rho(z) = \sum_{j=0}^q \alpha_j z^j \quad (21)$$

are within the unit circle, and those of modulus one (i.e., the ones on the unit circle) are simple. The following fundamental theorem relates zero stability of the numerical method (14) to the root condition.

**Theorem 2.** The numerical method (14) is zero-stable if and only if it satisfies the root condition.

A detailed proof of this theorem is provided at the end of this note<sup>5</sup>. Recall that a necessary condition for consistency is that  $\rho(1) = 0$ , i.e.,  $z = 1$  is a root of (21). Such a root must be simple in order for the method to satisfy the root condition. Let us now study zero-stability of all schemes we have considered so far.

- **One-step methods:** The most general form of a one-step method is

$$\mathbf{u}_{k+1} = \mathbf{u}_k + \Delta t \Phi_{\mathbf{f}}(\mathbf{u}_{k+1}, \mathbf{u}_k, t_k, \Delta t). \quad (22)$$

The characteristic polynomial for this class of methods is

$$\rho(z) = z - 1 \quad (23)$$

Clearly,  $\rho(z)$  has a simple root at  $z = 1$  and therefore (22) satisfies the root condition. This implies that *all one-step methods are zero-stable*. Recall that all Runge-Kutta methods are one-step methods.

- **Adams-Bashforth and Adams-Moulton methods:** A  $q$ -step Adams method can be written in the general form

$$\mathbf{u}_{k+q} = \mathbf{u}_{k+q-1} + \Delta t \sum_{j=0}^q \beta_j \mathbf{f}(\mathbf{u}_{k+j}, t_{k+j}). \quad (24)$$

For Adams-Bashforth methods (explicit) we have  $\beta_q = 0$ ; for Adams-Moulton (implicit)  $\beta_q \neq 0$ . The characteristic polynomial associated with (24) is

$$\rho(z) = z^q - z^{q-1} = z^{q-1}(z - 1). \quad (25)$$

This polynomial has as a simple root at  $z = 1$  and a root with algebraic multiplicity  $q - 1$  at  $z = 0$ . Therefore it satisfies the root condition. By Theorem (2) we have that *all Adams-Bashforth and all Adams-Moulton methods are zero-stable*.

---

<sup>5</sup>For whatever reason, none of the books I came across in my career provides concise and direct proof of Theorem 2 in the general case we are considering here, i.e., for vector-valued ODEs and numerical methods of the form (14). Hence, I decided to provide my own version of the proof.

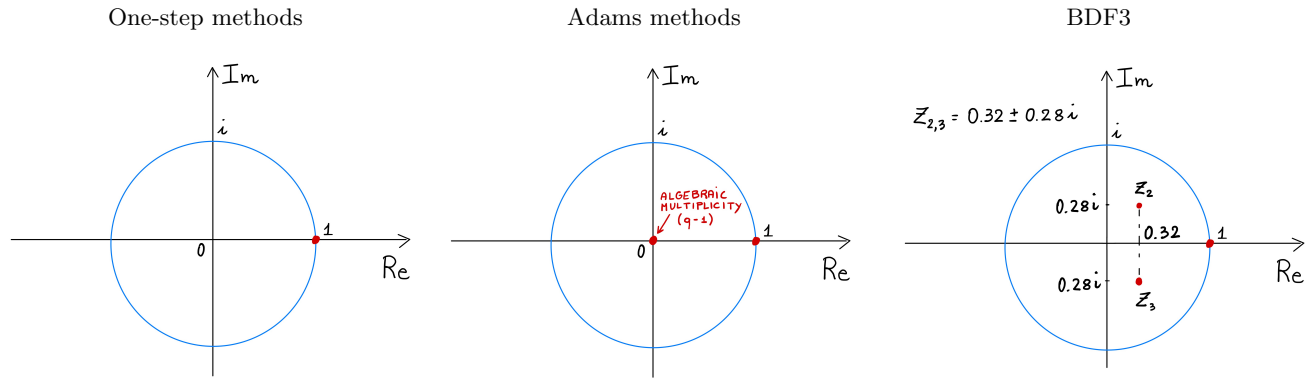


Figure 1: Roots of the characteristic polynomial (21). If all roots are within the unit circle and those modulus one (i.e., the ones on the unit circle) are simple (i.e., they have algebraic multiplicity one) then the method is zero-stable. All methods sketched in this figure are zero-stable.

- **BDF methods:** We know that a  $q$ -step BDF method can be written in the form

$$\sum_{j=0}^q \alpha_j \mathbf{u}_{k+j} = c\Delta t \mathbf{f}(\mathbf{u}_{k+q}, t_{k+q}). \quad (26)$$

The characteristic polynomial associated with (26) is

$$\rho(z) = z^q + \alpha_{q-1}z^{q-1} + \dots + \alpha_0. \quad (27)$$

It can be shown that a  $q$ -step BDF method satisfies the root condition and therefore it is zero-stable if and only if  $q \leq 6$ .

- **2-step midpoint method:** The 2-step midpoint method

$$\mathbf{u}_{k+2} = \mathbf{u}_k + 2\Delta t \mathbf{f}(\mathbf{u}_{k+1}, t_{k+1}) \quad (28)$$

satisfies the root condition and therefore it is zero-stable. In fact, the characteristic polynomial associated with (28) is

$$\rho(z) = z^2 - 1. \quad (29)$$

The roots  $z = \pm 1$  are both simple and sitting at the boundary of the unit circle in the complex plane. As we will see, a scheme that satisfies the root condition with simple eigenvalues at boundary of the unit circle is theoretically zero-stable, but in practical applications it can generate instabilities.

- **2-step LMM method:** The following two-step explicit linear multi-step method<sup>6</sup>

$$\mathbf{u}_{k+2} - 4\mathbf{u}_{k+1} + 3\mathbf{u}_k = -2\Delta t \mathbf{f}(\mathbf{u}_k, t_k) \quad (32)$$

is consistent but *not zero-stable*. The characteristic polynomial is

$$\rho(z) = z^2 - 4z + 3. \quad (33)$$

<sup>6</sup>The method (32) is *not* a BDF method, and is obtained by approximating  $d\mathbf{y}(t_k)/dt$  with a second-order forward finite difference formula:

$$\frac{d\mathbf{y}(t_k)}{dt} \simeq \frac{-3\mathbf{y}(t_k) + 4\mathbf{y}(t_{k+1}) - \mathbf{y}(t_{k+2})}{2\Delta t}. \quad (30)$$

and setting the equality

$$\frac{-3\mathbf{u}_k + 4\mathbf{u}_{k+1} - \mathbf{u}_{k+2}}{2\Delta t} = \mathbf{f}(\mathbf{u}_k, t_k). \quad (31)$$

Consistency can be checked immediately, since (see course note number 3)

$$\rho(1) = 0 \quad \frac{\Phi(\mathbf{u}_k, t_k, 0)}{\rho'(z)} = f(\mathbf{u}_k, t_k). \quad (34)$$

The polynomial (33) has roots  $z = 1$  and  $z = 3$ . Therefore the method (32) is not zero-stable.

- **General LMM methods:** We have seen in the course note 3 that the maximal order of consistency of a linear  $q$ -step method of the form

$$\sum_{j=0}^q \alpha_j \mathbf{u}_{k+j} = \Delta t \sum_{j=0}^q \beta_j \mathbf{f}(\mathbf{u}_{k+j}, t_{k+1}), \quad (35)$$

is  $2q$  (implicit methods) or  $2q - 1$  (explicit methods). At this point we notice that such maximal order LMM methods are, in general, *zero-unstable*, i.e., they do not satisfy the root-condition (see [2, §3.4]). In fact the following theorem holds true.

**Theorem 3** (First Dahlquist barrier - 1956). There is no zero-stable linear  $q$ -step method with consistency order exceeding  $q + 1$  ( $q$  odd) or  $q + 2$  ( $q$  even).

Zero-stable linear  $q$ -step implicit methods with order  $q + 2$  are called *optimal*. These methods have all roots with algebraic multiplicity one sitting on the boundary of the unit circle. This can yield stability issues.

**Convergence.** Let  $T = N\Delta t$  be period of integration. We say that the scheme (14) is convergent if the error (in any norm)

$$\max_{k \in \{0, \dots, N\}} \|\mathbf{u}_k - \mathbf{y}_k\| \quad (36)$$

goes to zero as  $\Delta t \rightarrow 0$ . Here  $\mathbf{y}_k = \mathbf{y}(t_k)$  represents the analytical solution of the ODE system (1) evaluated at  $t = t_k$ , while  $\mathbf{u}_k$  is the numerical solution produced by the scheme (14). If the error decreases as  $\Delta t^p$  then we say that the scheme converges with order  $p$ .

If a numerical scheme is convergent then the order of convergence is the same as the order of consistency (see the proof of theorem 4 at the end of this note). Indeed the error (36) can be bounded by the norm of the global truncation error, which goes to zero to some order in  $\Delta t$  (if the scheme is consistent) The following fundamental theorem provides necessary and sufficient conditions for convergence of numerical method for a system of ODEs.

**Theorem 4** (Convergence). The numerical method (14) is convergent if and only if it is consistent and zero-stable. In other words,

$$\text{convergence} \Leftrightarrow \text{consistency} + \text{zero stability}. \quad (37)$$

Moreover, the convergence order coincides with the consistency order.

The proof of this theorem follows exactly the same steps as the proof of theorem 2, and it is briefly discussed at the end of this note. This theorem has several corollaries. For instance, we have just seen that all one-step methods are zero-stable and therefore we have that:

**Corollary 1.** A one-step method is convergent if and only if it is consistent.

This means that in order to prove convergence of a one-step method it is necessary and sufficient to prove consistency. Hence, in the case of RK methods a necessary and sufficient condition for convergence is

$$\sum_{i=1}^s b_i = 1. \quad (38)$$

**Corollary 2.** Adams methods are convergent if and only if they are consistent.

In fact, we have seen that Adams methods are always zero-stable and therefore consistency implies convergence. Recall that Adams-Bashforth and Adams-Moulton methods are consistent if and only if

$$\sum_{j=0}^q \beta_j = 1. \quad (39)$$

Hence, if (39) is satisfied then Adams-Bashforth ( $\beta_q = 0$ ) or Adams-Moulton ( $\beta_q \neq 0$ ) methods are convergent.

*Example:* The numerical scheme (32) is not convergent. In fact, it is consistent, but not zero-stable.

**Estimating the convergence order of a numerical method.** To estimate the convergence order of the scheme (14) numerically it is sufficient to compute the error  $\|\mathbf{y}(t_k) - \mathbf{u}_k\|$  (in any norm) relative to an analytical solution  $\mathbf{y}(t)$  for various (sufficiently small)  $\Delta t$ , and then plot

$$\max_{k=1, \dots, N} \|\mathbf{y}(t_k) - \mathbf{u}_k\|$$

versus  $\Delta t$  in a logarithmic scale. The slope of the line obtained in this way represents the order of the method. In fact, suppose that for sufficiently small  $\Delta t$  we have

$$\max_{k=1, \dots, N} \|\mathbf{y}(t_k) - \mathbf{u}_k\| \simeq C \Delta t^p. \quad (40)$$

Taking the logarithm yields

$$\log \left( \max_{k=1, \dots, N} \|\mathbf{y}(t_k) - \mathbf{u}_k\| \right) \simeq \log(C) + p \log(\Delta t) \quad (41)$$

which represents a line with slope  $p$  in a log-log plot. To compute the error, we need of course the analytical solution to the initial value problem (1), which is not always available. However, it is very easy to manufacture an ODE with a time-dependent right hand side that has any desired solution  $\mathbf{y}(t)$ . To this end, choose any continuously differentiable vector  $\mathbf{y}(t)$  and any Lipschitz continuous function  $\mathbf{f}(\mathbf{y})$ . Compute the time forcing term

$$\mathbf{h}(t) = \frac{d\mathbf{y}(t)}{dt} - \mathbf{f}(\mathbf{y}(t)). \quad (42)$$

Then the chosen  $\mathbf{y}(t)$  is the analytical solution to the initial value problem

$$\begin{cases} \frac{d\mathbf{y}}{dt} = \mathbf{f}(\mathbf{y}) + \mathbf{h}(t) \\ \mathbf{y}(0) = \mathbf{y}_0 \end{cases} \quad (43)$$

In this way, for each given  $\Delta t$  we can solve (43) using the numerical method (14) and compute the error  $\|\mathbf{u}_k - \mathbf{y}(t_k)\|$ .

**Proof of Theorem 2.** Let us consider the  $m$ -th component of the perturbed scheme (14)

$$\sum_{j=0}^q \alpha_j z_{k+j}^m = \Delta t [\Phi_{\mathbf{f}}^m(\mathbf{z}_{k+q}, \dots, \mathbf{z}_k, t_k, \Delta t) + \delta_{q+k}^m]. \quad (44)$$

and the unperturbed one

$$\sum_{j=0}^q \alpha_j u_{k+j}^m = \Delta t \Phi_{\mathbf{f}}^m(\mathbf{u}_{k+q}, \dots, \mathbf{u}_k, t_k, \Delta t). \quad (45)$$

Subtracting (44) from (45) yields

$$\sum_{j=0}^q \alpha_j e_{k+j}^m = \Delta t [\Phi_{\mathbf{f}}^m(\mathbf{z}_{k+q}, \dots, \mathbf{z}_k, t_k, \Delta t) - \Phi_{\mathbf{f}}^m(\mathbf{u}_{k+q}, \dots, \mathbf{u}_k, t_k, \Delta t) + \delta_{q+k}^m], \quad (46)$$

where

$$e_{k+j}^m = z_{k+j}^m - u_{k+j}^m. \quad (47)$$

Upon definition of

$$\mathbf{e}_k^m = \begin{bmatrix} e_k^m \\ e_{k+1}^m \\ \vdots \\ e_{k+q-1}^m \end{bmatrix}, \quad \mathbf{b}_k^m = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ \Phi_{\mathbf{f}}^m(\mathbf{z}_{k+q}, \dots, \mathbf{z}_k, t_k, \Delta t) - \Phi_{\mathbf{f}}^m(\mathbf{u}_{k+q}, \dots, \mathbf{u}_k, t_k, \Delta t) \end{bmatrix}, \quad \mathbf{d}_k^m = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ \delta_{q+k}^m \end{bmatrix} \quad (48)$$

we see that we can write (46) in a compact form as<sup>7</sup>

$$\mathbf{e}_{k+1}^m = \mathbf{A} \mathbf{e}_k^m + \Delta t (\mathbf{b}_k^m + \mathbf{d}_k^m), \quad (49)$$

where

$$\mathbf{A} = \begin{bmatrix} 0 & 1 & 0 & \cdots & 0 \\ 0 & 0 & 1 & \cdots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \cdots & 1 \\ -\alpha_0 & -\alpha_1 & -\alpha_2 & \cdots & -\alpha_{q-1} \end{bmatrix}. \quad (50)$$

By using the discrete variation of constant formula (in which we treat  $\Delta t (\mathbf{b}_k^m + \mathbf{d}_k^m)$  as a “forcing term”) we write the formal solution of (49) as

$$\mathbf{e}_{k+1}^m = \mathbf{A}^{k+1} \mathbf{e}_0^m + \Delta t \sum_{p=0}^k \mathbf{A}^{k-p} (\mathbf{b}_p^m + \mathbf{d}_p^m), \quad (51)$$

As we shall see hereafter, the zero-stability of the numerical scheme is essentially determined by the properties of the matrix  $\mathbf{A}$ , in particular by the behavior of the matrix powers  $\mathbf{A}^k$  as  $k$  is increased. If the norm of the matrix powers can be bounded by a constant that is independent of  $k$  then zero-stability follows rather straightforwardly. The properties of the matrix powers  $\mathbf{A}^k$  are fully determined by the roots of the first characteristic polynomial (21).

---

<sup>7</sup>Recall that  $\alpha_q = 1$ .



**Lemma 1.** Let  $\|\cdot\|$  be any matrix norm compatible with a vector norm. Then  $\|\mathbf{A}^k\|$  can be bounded by a quantity  $M$  that does not depend on  $k$ , i.e.,

$$\|\mathbf{A}^k\| \leq M \quad \text{for all } k \in \mathbb{N} \quad (52)$$

if and only if the root condition is satisfied.

*Proof.* The matrix  $\mathbf{A}$  in (50) is the transpose of the companion matrix associated with the characteristic polynomial (21). This means that the eigenvalues of  $\mathbf{A}$  coincide with the roots of the polynomial (21). Moreover, companion matrices are *non-derogatory*, which means that there exists only one eigenvector corresponding to each eigenvalue  $\lambda$ . Such eigenvector is explicitly obtained as

$$\mathbf{h} = \begin{bmatrix} 1 \\ \lambda \\ \lambda^2 \\ \vdots \\ \lambda^{q-1} \end{bmatrix}. \quad (53)$$

The non-derogatory property of  $\mathbf{A}$  implies that if there exists any eigenvalue with algebraic multiplicity  $r_j > 1$ , then the corresponding eigenspace has dimension  $r_j - 1$ . This means that the matrix  $\mathbf{A}$  is diagonalizable (i.e., similar to a diagonal matrix), if and only if all the eigenvalues are simple. If there exist any eigenvalue with multiplicity larger than one then the matrix  $\mathbf{A}$  is similar to a (block-diagonal) Jordan matrix  $\mathbf{J}$

$$\mathbf{A} = \mathbf{P}\mathbf{J}\mathbf{P}^{-1} \quad (54)$$

where  $\mathbf{P}$  is the matrix that has the generalized eigenvectors of  $\mathbf{A}$  columnwise and

$$\mathbf{J} = \begin{bmatrix} \mathbf{J}_1 & \mathbf{0} & \cdots & \mathbf{0} \\ \mathbf{0} & \mathbf{J}_2 & \cdots & \mathbf{0} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{0} & \mathbf{0} & \cdots & \mathbf{J}_l \end{bmatrix}, \quad \mathbf{J}_i = \lambda_i \mathbf{I}_{r_i} + \mathbf{B}_{r_i}. \quad (55)$$

In this equation,  $\mathbf{J}_i$  denotes the Jordan block corresponding to the eigenvalue  $\lambda_i$  (which has algebraic multiplicity  $r_i$ ),  $\mathbf{I}_{r_i}$  is a  $r_i \times r_i$  identity matrix and  $\mathbf{B}_{r_i}$  is a  $r_i \times r_i$  matrix with ones above the main diagonal. For instance, if  $\lambda_i$  has algebraic multiplicity  $r_i = 3$  then the geometric multiplicity is 2 and we have

$$\mathbf{J}_i = \begin{bmatrix} \lambda_i & 1 & 0 \\ 0 & \lambda_i & 1 \\ 0 & 0 & \lambda_i \end{bmatrix}, \quad \mathbf{I}_3 = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}, \quad \mathbf{B}_3 = \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 0 & 0 & 0 \end{bmatrix}. \quad (56)$$

The matrix power  $\mathbf{A}^k$  can be written as

$$\mathbf{A}^k = \mathbf{P}\mathbf{J}^k\mathbf{P}^{-1}, \quad (57)$$

where

$$\mathbf{J}^k = \begin{bmatrix} \mathbf{J}_1^k & \mathbf{0} & \cdots & \mathbf{0} \\ \mathbf{0} & \mathbf{J}_2^k & \cdots & \mathbf{0} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{0} & \mathbf{0} & \cdots & \mathbf{J}_l^k \end{bmatrix}, \quad \mathbf{J}_i^k = (\lambda_i \mathbf{I}_{r_i} + \mathbf{B}_{r_i})^k \quad (58)$$

Let us compute  $\mathbf{J}_i^k$  for  $r_i = 1$  (simple eigenvalue)

$$\mathbf{J}_i^k = \lambda_i^k. \quad (59)$$

On the other hand, for  $r_i = 2$ , (eigenvalue with algebraic multiplicity 2 and geometric multiplicity 2) we have have

$$J_i = \lambda_i \mathbf{I}_2 + \mathbf{B}_2, \quad (60)$$

$$J_i^2 = (\lambda_i \mathbf{I}_2 + \mathbf{B}_2)^2 = \lambda_i^2 \mathbf{I}_2 + 2\lambda_i \mathbf{B}_2, \quad (61)$$

$\vdots$

$$J_i^k = (\lambda_i \mathbf{I}_2 + \mathbf{B}_2)^k = \lambda_i^k \mathbf{I}_2 + k\lambda_i^{k-1} \mathbf{B}_2, \quad (62)$$

where we used the fact that  $\mathbf{B}_2^k = 0$  for all  $k \geq 2$ . Similarly, for  $r_i = 3$  it can be shown that  $\mathbf{B}_3^k = 0$  for all  $k \geq 3$ , which yields

$$J_i^k = \lambda_i^k \mathbf{I}_3 + k\lambda_i^{k-1} \mathbf{B}_3 + k\lambda_i^{k-2} \mathbf{B}_3^2 \quad k \geq 3. \quad (63)$$

By taking the norm of (57) we obtain

$$\|\mathbf{A}^k\| \leq K \|\mathbf{J}^k\|, \quad K = \|\mathbf{P}\| \|\mathbf{P}^{-1}\|. \quad (64)$$

At this point we recall that for any matrix norm compatible with a vector norm and for any block-diagonal matrix such as  $\mathbf{J}$  or  $\mathbf{J}^k$  we have

$$\|\mathbf{J}^k\| = \max \left\{ \|\mathbf{J}_1^k\|, \dots, \|\mathbf{J}_l^k\| \right\}. \quad (65)$$

If the eigenvalues of  $\mathbf{A}$  sitting at the boundary of the unit circle are *simple* then, by equation (59) we have

$$\|\mathbf{J}_i^k\| = 1. \quad (66)$$

On the other hand, if  $|\lambda_i| < 1$  (eigenvalue within the unit circle or arbitrary multiplicity) then by equation (62) or (63) we have that

$$\|\mathbf{J}_i^k\| \rightarrow 0 \quad \text{for } k \rightarrow \infty. \quad (67)$$

Since  $\|\mathbf{J}_i^k\|$  is finite for all  $k$ , there exists a finite  $M$  such that  $\|\mathbf{J}_i^k\| \leq M$  for all  $k$ .

Finally, if there exists a non-simple eigenvalue  $\lambda_i$  (eigenvalue with algebraic multiplicity larger than one) at the boundary of the unit circle then we can no longer guarantee that  $\|\mathbf{J}_i^k\|$  is bounded independently of  $k$ . In fact, suppose that the algebraic multiplicity of the eigenvalue  $\lambda_i$  at the boundary of the unit circle (i.e.,  $|\lambda_i| = 1$ ) is  $r_i = 2$ . Then by using (62) we see that

$$\|\mathbf{J}_i^k\|_1 = |\lambda_i|^k + k|\lambda_i|^{k-1} = 1 + k \quad \text{for all } k \geq 2. \quad (68)$$

In summary, if the root condition is satisfied, i.e., if all the eigenvalues of  $\mathbf{A}$  are within the unit circle with the exception of a finite number of *simple* eigenvalues sitting at the boundary of the unit circle then

$$\|\mathbf{A}^k\| \leq K \|\mathbf{J}^k\| \leq M \quad \text{for all } k \in \mathbb{N}, \quad (69)$$

where  $M > 0$  is independent of  $k$ . This completes the proof of Lemma 1. □

We now have all elements to show that if a scheme satisfies the root condition then it is zero-stable. To this end, let us take the infinity norm of (51), and use (52) (or (69)) to obtain

$$|e_{k+q}^m| \leq \|e_{k+1}^m\|_\infty \leq M \left( \|e_0^m\|_\infty + \Delta t \sum_{p=0}^k \|b_p^m\|_\infty + \Delta t \sum_{p=0}^k \|d_p^m\|_\infty \right). \quad (70)$$

By using definition (48) and (47) we see that

$$\begin{aligned}
\sum_{m=1}^n \|\mathbf{b}_p^m\|_\infty &= \sum_{m=1}^n |\Phi_{\mathbf{f}}^m(\mathbf{z}_{p+q}, \dots, \mathbf{z}_p, t_p, \Delta t) - \Phi_{\mathbf{f}}^m(\mathbf{u}_{p+q}, \dots, \mathbf{u}_p, t_p, \Delta t)| \\
&= \|\Phi_{\mathbf{f}}(\mathbf{z}_{p+q}, \dots, \mathbf{z}_p, t_p, \Delta t) - \Phi_{\mathbf{f}}(\mathbf{u}_{p+q}, \dots, \mathbf{u}_p, t_p, \Delta t)\|_1 \\
&\leq L \sum_{s=0}^q \|\mathbf{z}_{p+s} - \mathbf{u}_{p+s}\|_1 \\
&= L \sum_{s=0}^q \sum_{m=1}^n |z_{p+s}^m - u_{p+s}^m| \\
&= L \sum_{s=0}^q \sum_{m=1}^n |e_{p+s}^m|, \tag{71}
\end{aligned}$$

where we assumed that  $\Phi_{\mathbf{f}}$  is Lipschitz continuous. Next, define

$$g_{k+q} = \sum_{m=1}^n |e_{k+q}^m|. \tag{72}$$

Note that  $g_{k+q}$  is the 1-norm of the vector  $\mathbf{z}_{k+q} - \mathbf{u}_{k+q}$  (see Eq. (47)). Substituting (71) into the inequality (70) (summed-up in  $m$ ) yields

$$\begin{aligned}
g_{k+q} &\leq M \left( \sum_{m=1}^n \|e_0^m\|_\infty + L\Delta t \sum_{s=0}^{k+q} g_s + \Delta t \sum_{p=0}^k \sum_{m=1}^n \|\mathbf{d}_p^m\|_\infty \right) \\
&\leq Mn\epsilon(1 + k\Delta t) + ML\Delta t \sum_{s=0}^{k+q} g_s. \tag{73}
\end{aligned}$$

Now we can use the discrete Grönwall lemma (see, e.g., [3, Lemma 11.2]) to conclude that

$$g_{k+q} \leq \epsilon \left( nM(1 + (k+1)\Delta t) e^{ML(k+q+1)\Delta t} \right) \leq \underbrace{\epsilon nM(1+T)e^{MLT}}_C, \tag{74}$$

where  $T \geq (q+k+1)\Delta t$  is some integration period. Recalling that  $g_{k+p}$  is the 1-norm of the vector  $\mathbf{z}_{k+p} - \mathbf{u}_{k+p}$  (see Eq. (72)) we see that (74) can be written as

$$\|\mathbf{z}_{k+q} - \mathbf{u}_{k+q}\|_1 \leq C\epsilon, \tag{75}$$

for all  $k$  such that  $(q+k+1)\Delta t \leq T$ . Alternatively, if we set a maximum number of time steps  $N \geq k$  and an integration period  $T$  then (75) holds for all  $k \leq N$  and for all  $\Delta t \leq T/(N+q) = \Delta t^*$ . This is were the definition of zero-stability kicks in, i.e., conditions (74) and (75) are satisfied for all  $\Delta t \leq T/(N+q) = \Delta t^*$ . Based on definition (19) we conclude that the root condition implies zero-stability. The converse statement, i.e., zero-stability implies root condition, is straightforward. Indeed, if the scheme is zero stable then (20) is satisfied for all  $\epsilon$ . This implies that (see Equation 51)

$$\left\| \mathbf{A}^{k+1} \mathbf{e}_0^m + \Delta t \sum_{p=0}^k \mathbf{A}^{k-p} (\mathbf{b}_p^m + \mathbf{d}_p^m) \right\|_\infty \leq C\epsilon \tag{76}$$

Recalling that  $C$  must be independent of  $k$ , this condition can be satisfied for all  $\epsilon$  if and only if  $\|\mathbf{A}^k\| \leq M$ .

**Proof of theorem 4.** Let  $\mathbf{y}_k = \mathbf{y}(t_k)$  be the solution of the ODE (1) evaluated at  $t = t_k$ . A substitution of such solution into the scheme (14) yields the truncation error  $\boldsymbol{\tau}_{k+q}$ , hereafter written in a componentwise form ( $m = 1, \dots, n$ )

$$\sum_{j=0}^q \alpha_j y_{k+j}^m = \Delta t (\Phi_{\mathbf{f}}^m(\mathbf{y}_{k+q}, \dots, \mathbf{y}_k, t_k, \Delta t) + \tau_{q+k}^m). \quad (77)$$

Similarly, the numerical solution  $\mathbf{u}_k$  satisfies

$$\sum_{j=0}^q \alpha_j u_{k+j}^m = \Delta t \Phi_{\mathbf{f}}^m(\mathbf{u}_{k+q}, \dots, \mathbf{u}_k, t_k, \Delta t). \quad (78)$$

Subtracting (78) from (77) yields

$$\sum_{j=0}^q \alpha_j e_{k+j}^m = \Delta t [\Phi_{\mathbf{f}}^m(\mathbf{y}_{k+q}, \dots, \mathbf{y}_k, t_k, \Delta t) - \Phi_{\mathbf{f}}^m(\mathbf{u}_{k+q}, \dots, \mathbf{u}_k, t_k, \Delta t) + \tau_{q+k}^m], \quad (79)$$

where

$$e_{k+j}^m = y_{k+j}^m - u_{k+j}^m. \quad (80)$$

By following exactly same steps that took as from equation (46) to (75) in the proof of theorem 2 we obtain the error bound

$$\|\mathbf{y}(t_k) - \mathbf{u}_k\|_1 \leq MT e^{MLT} \|\boldsymbol{\tau}(\Delta t)\|_1, \quad (81)$$

where the global truncation error  $\|\boldsymbol{\tau}\|$  is a function of  $\Delta t$ . To obtain (81) we replaced  $(q+k)\Delta t$  with  $T$ , which implies that (81) holds for all  $\Delta t \leq T/(N+q)$  (this is where zero stability comes in) where  $N$  any fixed number larger or equal than  $(k+q)$ .

Moreover, for  $\Delta t$  small enough, we have seen that  $\|\boldsymbol{\tau}\|_1$  goes to zero as some power of  $\Delta t$  (otherwise the method is not consistent). Equation (81) says that the convergence order of the method is the same as the order of consistency.

To obtain the bound (81) we assumed that the initial condition has no error, and that the numerical computation of  $\Phi_{\mathbf{f}}$  and all arithmetic operations in the schemes are exact. Clearly this is not the case in practice. It is possible to repeat the proof above, by assuming that all these numerical inaccuracies are bounded, e.g., as a function of the machine precision  $\epsilon$ , and develop a more detailed bound that depends on  $\epsilon$ .

## References

- [1] W. Hahn. *Stability of motion*. Springer, 1967.
- [2] J. D. Lambert. *Numerical methods for ordinary differential systems: the initial value problem*. Wiley, 1991.
- [3] A. Quarteroni, R. Sacco, and F. Saleri. *Numerical mathematics*. Springer, 2007.