# Stability of one-step and linear multistep methods – a matrix technique approach

## Miklós E. Mincsovics[✉ 1,2]

[1]MTA–ELTE Numerical Analysis and Large Networks Research Group,
Pázmány Péter sétány 1/C, Budapest H–1117, Hungary
[2]Budapest University of Technology and Economics, Department of Differential Equations,
Building H, Egry József utca 1, Budapest H–1111, Hungary

**Abstract.** We investigate the stability of one-step and linear multistep methods from a new direction. Our aim is to modify the long and technical proof which is consequently omitted in almost every textbook and make it user-friendly. In the literature the techniques of numerical solution of initial value problems and boundary value problems seem to have almost nothing in common which is quite surprising. Our new approach uses matrix techniques opposed to the usual recursion approach, thus applying the techniques of boundary value problems to initial value problems. Even though the proof remains long, it is easier to follow and connects two seemingly separated areas, consequently this approach might have educational profit.

**Keywords:** stability, linear multistep methods, M-matrix theory.

**2010 Mathematics Subject Classification:** 65L20, 65L06.

## 1  Introduction

Consider the initial value problem

$$\begin{cases} u(0) = u^0 \,, \\ u'(t) = f(u(t)) \,, \end{cases} \tag{1.1}$$

where $t \in [0,1]$, $u^0 \in \mathbb{R}$ is the initial value, $u : [0,1] \to \mathbb{R}$ is the unknown function and we assume that $f$ is Lipschitz continuous.

Since this problem is generally unsolvable, usually a numerical method is applied to approximate the solution. The most popular methods are the one-step and linear multistep methods. Both types use the grid $G_n = \{x_0 = 0, x_1, \ldots, x_{n+k-1} = 1\}$, where $h = x_{i+1} - x_i$ is the stepsize with $(n+k-1)h = 1$ (we investigate only the case when a uniform grid is used). The unknown function is approximated only at the gridpoints $u_i \approx u(x_i)$.

---

✉ Email: m.e.mincsovics@gmail.com

E.g. the *explicit Euler method* (EE) reads as

$$\begin{cases} u_0 = u^0, \\ n(u_i - u_{i-1}) = f(u_{i-1}), \quad i = 1, \ldots, n \end{cases} \tag{1.2}$$

and *linear multistep method*s (LMM) can be given in the following way

$$\begin{cases} u_i = c^i, & i = 0, \ldots, k-1 \\ \dfrac{1}{h} \sum_{j=0}^{k} \alpha_j u_{i-j} = \sum_{j=0}^{k} \beta_j f(u_{i-j}), & i = k, \ldots, n+k-1, \end{cases} \tag{1.3}$$

where $k$ denotes the number of steps. Note that EE can be viewed as a LMM with $k = 1$ step, $\alpha_0 = 1$, $\alpha_1 = -1$, $\beta_1 = 1$. We also note that there is a technical problem for $k > 1$, namely $c^i$, $i = 1, \ldots, k-1$ need to be determined, but this is beyond the scope of this present paper.

The usefulness of a method depends on whether it is convergent or not. The question of convergence can be split into two tasks, namely checking consistency and stability.

Consistency and its order can be determined by using the Taylor series theorem and the order conditions can be formalized by the help of the first and the second characteristic polynomial.

The *first characteristic polynomial* associated to (1.3) is defined as

$$\varrho(x) = \sum_{j=0}^{k} \alpha_j x^{k-j}, \tag{1.4}$$

while the *second characteristic polynomial* as

$$\sigma(x) = \sum_{j=0}^{k} \beta_j x^{k-j}. \tag{1.5}$$

For (at least first order) consistency the LMM must satisfy

$$\varrho(1) = 0 \quad \text{and} \quad \varrho'(1) = \sigma(1) = 1. \tag{1.6}$$

This latter usually appears in textbooks as $\varrho'(1) = \sigma(1)$ without being equal to 1, because LMMs can be scaled differently. However, we prefer this particular scaling since only in this case is true that (1.2) and (1.3) approximates (1.1) (and not a scalar times (1.1)).

Following the framework of [2, 6] we rewrite the methods (1.2) and (1.3) into the forms

$$(F_n(\mathbf{u}_n))_i = \begin{cases} u_0 - u^0, & i = 0 \\ n(u_i - u_{i-1}) - f(u_{i-1}), & i = 1, \ldots, n \end{cases} \tag{1.7}$$

and

$$(F_n(\mathbf{u}_n))_i = \begin{cases} u_i - c^i, & i = 0, \ldots, k-1 \\ \dfrac{1}{h} \sum_{j=0}^{k} \alpha_j u_{i-j} - \sum_{j=0}^{k} \beta_j f(u_{i-j}), & i = k, \ldots, n+k-1. \end{cases} \tag{1.8}$$

Exploiting the Lipschitz continuity stability simplifies to the following condition.

$\exists\, S \in \mathbb{R}$ and $\exists\, n_0 \in \mathbb{N}$ such that $\forall\, n \geq n_0$, $\forall\, \mathbf{u}_n^1, \mathbf{u}_n^2$ the estimate

$$\left\| \mathbf{u}_n^1 - \mathbf{u}_n^2 \right\|_{\mathcal{X}_n} \leq S \left\| F_n(\mathbf{u}_n^1) - F_n(\mathbf{u}_n^2) \right\|_{\mathcal{Y}_n} \tag{1.9}$$

holds.

If this condition is fulfilled for some method defined by $F_n$ and for the norms defined by $\mathcal{X}_n$ and $\mathcal{Y}_n$ we say that *the method is stable in the norm pair* $\|\cdot\|_{\mathcal{X}_n}$ *and* $\|\cdot\|_{\mathcal{Y}_n}$.

Naturally, the choice of $\mathcal{X}_n$ and $\mathcal{Y}_n$ is crucial in getting an admissible pair. One needs to take into consideration the original problem, usually some norm-consistency is required, see [6].

Ensuring stability is based on a technical result which states that stability is equivalent to the so-called root-condition. This is presented below.

The method is said to be *weakly stable* if for every root $\xi \in \mathbb{C}$ of the first characteristic polynomial $|\xi| \leq 1$ holds and if $|\xi| = 1$ then it is a simple root.

The method is said to be *strongly stable* if for every root $\xi \in \mathbb{C}$ of the first characteristic polynomial $|\xi| < 1$ holds except $\xi = 1$.

We note that for a consistent method $\varrho(1) = 0$ always holds. Weak stability corresponds to stability in the pair $\|\cdot\|_\infty$, $\|\cdot\|_\infty$ (we are interested in this case), while strong stability corresponds to stability in the pair $\|\cdot\|_\infty$, $\|\cdot\|_\$$, where the latter is the Spijker norm. If the connection of the root-condition and stability is proved then checking it is an easy task.

**Remarks and aims.**

- The "root-condition $\Rightarrow$ stability" part of the proof is technical and long. This is the reason why it is omitted in most of the textbooks and consequently in most of the courses.

- There are a few exceptions. The mostly referred books are [4] and [5] but these are hardly accessible nowadays. The book [6] contains much more but it is too detailed and really hard to read. More available is the book [3] which contains a proof based on the theory of linear difference equations, while [7] (in Hungarian) contains a proof based on using the transition matrix.

- Our aim is to give a new proof which is less technical and with this more followable.

- To do this we apply what we call matrix techniques. This is the technique which is usually applied in the case of boundary value problems. Using this direction we want to connect these areas and show the similarities between them which was concealed by the other proofs.

## 2  Stability of the explicit Euler method

To demonstrate the usefulness of the matrix technique, we first use it to prove the stability of the EE.

Throughout the paper we will use the following notations. If $\mathbf{u} = (u_0, u_1, \ldots, u_n)^T$ then $|\mathbf{u}| = (|u_0|, |u_1|, \ldots, |u_n|)^T$. $\mathbf{u} \geq \mathbf{v}$ is an elementwise relation i.e. it means $u_i \geq v_i$ for $i = 0, \ldots, n$. These notations will be used for matrices in the same sense.

Switching over to matrix form (1.7) can be written as

$$F_n(\mathbf{u}_n) = \mathbf{A}_n \mathbf{u}_n - \mathbf{B}_n \mathbf{f}(\mathbf{u}_n) - \mathbf{c}_n , \tag{2.1}$$

where $\mathbf{u}_n = (u_0, u_1, \ldots, u_n)^T$, $\mathbf{f}(\mathbf{u}_n) = (f(u_0), f(u_1), \ldots, f(u_n))^T$, $\mathbf{c}_n = (u^0, 0, \ldots, 0)^T$,

$$\mathbf{A}_n = \begin{pmatrix} 1 & 0 & \ldots & \ldots & 0 \\ -n & n & 0 & \ldots & 0 \\ 0 & -n & n & 0 & \ldots \\ \vdots & \ddots & \ddots & \ddots & \vdots \\ 0 & \ldots & 0 & -n & n \end{pmatrix} , \qquad \mathbf{B}_n = \begin{pmatrix} 0 & 0 & \ldots & \ldots & 0 \\ 1 & 0 & 0 & \ldots & 0 \\ 0 & 1 & 0 & 0 & \ldots \\ \vdots & \ddots & \ddots & \ddots & \vdots \\ 0 & \ldots & 0 & 1 & 0 \end{pmatrix} .$$

Thus

$$F_n(\mathbf{u}_n^1) - F_n(\mathbf{u}_n^2) = \mathbf{A}_n(\mathbf{u}_n^1 - \mathbf{u}_n^2) - \mathbf{B}_n(f(\mathbf{u}_n^1) - f(\mathbf{u}_n^2)) ,$$

multiplying by $\mathbf{A}_n^{-1}$, which exists, we have

$$\mathbf{A}_n^{-1} \left( F_n(\mathbf{u}_n^1) - F_n(\mathbf{u}_n^2) \right) = (\mathbf{u}_n^1 - \mathbf{u}_n^2) - \mathbf{A}_n^{-1} \mathbf{B}_n(f(\mathbf{u}_n^1) - f(\mathbf{u}_n^2)) ,$$

Taking absolute value and exploiting the Lipschitz continuity of $f$ we can estimate the right side

$$\begin{aligned} |\mathbf{A}_n^{-1} \left( F_n(\mathbf{u}_n^1) - F_n(\mathbf{u}_n^2) \right)| &= |(\mathbf{u}_n^1 - \mathbf{u}_n^2) - \mathbf{A}_n^{-1} \mathbf{B}_n(f(\mathbf{u}_n^1) - f(\mathbf{u}_n^2))| \\ &\geq |\mathbf{u}_n^1 - \mathbf{u}_n^2| - |\mathbf{A}_n^{-1} \mathbf{B}_n(f(\mathbf{u}_n^1) - f(\mathbf{u}_n^2))| \\ &\geq |\mathbf{u}_n^1 - \mathbf{u}_n^2| - |\mathbf{A}_n^{-1}||\mathbf{B}_n||f(\mathbf{u}_n^1) - f(\mathbf{u}_n^2)| \\ &\geq |\mathbf{u}_n^1 - \mathbf{u}_n^2| - |\mathbf{A}_n^{-1}||\mathbf{B}_n|L|\mathbf{u}_n^1 - \mathbf{u}_n^2| \\ &= (\mathbf{I} - L|\mathbf{A}_n^{-1}||\mathbf{B}_n|)|\mathbf{u}_n^1 - \mathbf{u}_n^2| . \end{aligned}$$

If $\mathbf{X}_n = \mathbf{I} - L|\mathbf{A}_n^{-1}||\mathbf{B}_n|$ is inverse nonnegative then

$$\mathbf{X}_n^{-1}|\mathbf{A}_n^{-1} \left( F_n(\mathbf{u}_n^1) - F_n(\mathbf{u}_n^2) \right)| \geq |\mathbf{u}_n^1 - \mathbf{u}_n^2| ,$$

thus

$$\left\|\mathbf{X}_n^{-1}\right\|_\infty \left\|\mathbf{A}_n^{-1}\right\|_\infty \left\|F_n(\mathbf{u}_n^1) - F_n(\mathbf{u}_n^2)\right\|_\infty \geq \left\|\mathbf{u}_n^1 - \mathbf{u}_n^2\right\|_\infty .$$

If both of $\left\|\mathbf{X}_n^{-1}\right\|_\infty$ and $\left\|\mathbf{A}_n^{-1}\right\|_\infty$ are bounded independently of $n$ then we got stability in the $\|\cdot\|_\infty, \|\cdot\|_\infty$ pair. So we have two tasks.

**1.**

$$\mathbf{A}_n^{-1} = \begin{pmatrix} 1 & 0 & \ldots & \ldots & 0 \\ 1 & h & 0 & \ldots & 0 \\ 1 & h & h & 0 & \ldots \\ \vdots & \ddots & \ddots & \ddots & \vdots \\ 1 & h & \ldots & h & h \end{pmatrix} , \tag{2.2}$$

thus

$$\|\mathbf{A}_n^{-1}\|_\infty = 2 .$$

Alternatively we can use M-matrix theory. For the Reader's convenience we collected the necessary information on M-matrices in the Appendix.

We choose $d(t) = e^t$ and so the dominant-vector is $\mathbf{d}_n$ with $(\mathbf{d}_n)_i = e^{t_i} > 0$ and $\|\mathbf{d}_n\|_\infty = e$.
Then

$$(\mathbf{A}_n\mathbf{d}_n)_i = \begin{cases} 1, & i = 0, \\ n\left(-e^{t_{i-1}} + e^{t_i}\right), & i = 1,\ldots,n, \end{cases}$$

$n\left(-e^{t_{i-1}} + e^{t_i}\right) = \frac{e^h - 1}{h}e^{t_{i-1}} \geq e^{t_{i-1}} \geq 1$, thus $\|\mathbf{A}_n^{-1}\|_\infty \leq e$.

**2.**

$$|\mathbf{A}_n^{-1}||\mathbf{B}_n| = \mathbf{A}_n^{-1}\mathbf{B}_n = \begin{pmatrix} 0 & 0 & \ldots & \ldots & 0 \\ h & 0 & 0 & \ldots & 0 \\ h & h & 0 & 0 & \ldots \\ \vdots & \ddots & \ddots & \ddots & \vdots \\ h & h & \ldots & h & 0 \end{pmatrix}, \tag{2.3}$$

which is "small" enough to ensure for

$$\mathbf{X}_n = \begin{pmatrix} 1 & 0 & \ldots & \ldots & 0 \\ -Lh & 1 & 0 & \ldots & 0 \\ -Lh & -Lh & 1 & 0 & \ldots \\ \vdots & \ddots & \ddots & \ddots & \vdots \\ -Lh & -Lh & \ldots & -Lh & 1 \end{pmatrix} \tag{2.4}$$

to be an M-matrix. It is clearly a Z-matrix. We choose $d(t) = e^{Lt}$ and so the dominant-vector
is $\mathbf{d}_n$ with $(\mathbf{d}_n)_i = e^{Lt_i} > 0$ and $\|\mathbf{d}_n\|_\infty = e^L$. Then

$$(\mathbf{X}_n\mathbf{d}_n)_i = \begin{cases} 1, & i = 0, \\ e^{Lt_i} - Lh\sum\limits_{j=0}^{i-1} e^{Lt_j}, & i = 1,\ldots,n, \end{cases}$$

$$e^{Lt_i} - Lh\sum_{j=0}^{i-1} e^{Lt_j} = e^{Lt_i} - Lh\frac{e^{Lt_i} - 1}{e^{Lh} - 1} = e^{Lt_i} - L\frac{1}{\frac{e^{Lh}-1}{h}}(e^{Lt_i} - 1) \geq 1,$$

thus $\|\mathbf{X}_n^{-1}\|_\infty \leq e^L$.

With this we proved the stability of the EE in the $\|\cdot\|_\infty, \|\cdot\|_\infty$ pair.

## 3 Stability of linear multistep methods

We proceed similarly to the EE and we use the matrix form corresponding to (1.3)

$$F_n(\mathbf{u}_n) = \mathbf{A}_n\mathbf{u}_n - \mathbf{B}_n\mathbf{f}(\mathbf{u}_n) - \mathbf{c}_n, \tag{3.1}$$

where

$$\mathbf{u}_n = (u_0, u_1, \ldots, u_{n+k-1})^T,$$
$$\mathbf{f}(\mathbf{u}_n) = (f(u_0), f(u_1), \ldots, f(u_{n+k-1}))^T,$$
$$\mathbf{c}_n = (c^0, c^1, \ldots, c^{k-1}, 0, \ldots, 0)^T,$$
$$\mathbf{A}_n = \begin{pmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{A}_{n,\partial} & \mathbf{A}_{n,0} \end{pmatrix}, \qquad \mathbf{B}_n = \begin{pmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{B}_{n,\partial} & \mathbf{B}_{n,0} \end{pmatrix},$$

where $\mathbf{I} \in \mathbb{R}^{k \times k}$ is the identity matrix, $\mathbf{A}_{n,0}, \mathbf{B}_{n,0} \in \mathbb{R}^{n \times n}$ and

$$
\mathbf{A}_{n,\partial} = \frac{1}{h}
\begin{pmatrix}
\alpha_k & \cdots & \alpha_2 & \alpha_1 \\
0 & \alpha_k & \cdots & \alpha_2 \\
\vdots & \ddots & \ddots & \vdots \\
0 & \cdots & \cdots & \alpha_k \\
0 & \cdots & \cdots & 0 \\
\vdots & \ddots & \ddots & \vdots \\
0 & \cdots & \cdots & 0
\end{pmatrix}
\qquad
\mathbf{A}_{n,0} = \frac{1}{h}
\begin{pmatrix}
\alpha_0 & 0 & \cdots & \cdots & \cdots & 0 \\
\alpha_1 & \alpha_0 & 0 & \cdots & \cdots & 0 \\
\alpha_2 & \alpha_1 & \alpha_0 & 0 & \cdots & 0 \\
\vdots & \ddots & \ddots & \ddots & \ddots & \vdots \\
\vdots & & \ddots & \ddots & \ddots & \vdots \\
\vdots & & & \ddots & \ddots & \vdots \\
0 & \cdots & 0 & \alpha_k & \cdots & \alpha_0
\end{pmatrix},
$$

$$
\mathbf{B}_{n,\partial} =
\begin{pmatrix}
\beta_k & \cdots & \beta_2 & \beta_1 \\
0 & \beta_k & \cdots & \beta_2 \\
\vdots & \ddots & \ddots & \vdots \\
0 & \cdots & \cdots & \beta_k \\
0 & \cdots & \cdots & 0 \\
\vdots & \ddots & \ddots & \vdots \\
0 & \cdots & \cdots & 0
\end{pmatrix}
\qquad
\mathbf{B}_{n,0} =
\begin{pmatrix}
\beta_0 & 0 & \cdots & \cdots & \cdots & 0 \\
\beta_1 & \beta_0 & 0 & \cdots & \cdots & 0 \\
\beta_2 & \beta_1 & \beta_0 & 0 & \cdots & 0 \\
\vdots & \ddots & \ddots & \ddots & \ddots & \vdots \\
\vdots & & \ddots & \ddots & \ddots & \vdots \\
\vdots & & & \ddots & \ddots & \vdots \\
0 & \cdots & 0 & \beta_k & \cdots & \beta_0
\end{pmatrix}.
$$

Following the way we calculated the stability of the EE, we have

$$
\left| \mathbf{A}_n^{-1} \big( F_n(\mathbf{u}_n^1) - F_n(\mathbf{u}_n^2) \big) \right| \geq \big( \mathbf{I} - L|\mathbf{A}_n^{-1}||\mathbf{B}_n| \big) |\mathbf{u}_n^1 - \mathbf{u}_n^2|,
$$

where the problem is that $|\mathbf{A}_n^{-1}||\mathbf{B}_n|$ is difficult to calculate – even determining $\mathbf{A}_n^{-1}$ is more difficult than previously – so we use an estimate

$$
\left| \mathbf{A}_n^{-1} \big( F_n(\mathbf{u}_n^1) - F_n(\mathbf{u}_n^2) \big) \right| \geq \big( \mathbf{I} - \mathbf{W}_n \big) |\mathbf{u}_n^1 - \mathbf{u}_n^2|,
$$

where $L|\mathbf{A}_n^{-1}||\mathbf{B}_n| \leq \mathbf{W}_n$ for some $\mathbf{W}_n$ which is still small enough for $\bar{\mathbf{X}}_n = \mathbf{I} - \mathbf{W}_n$ to be an M-matrix. The finishing is the same as in the case of the EE. Thus

$$
(\bar{\mathbf{X}}_n)^{-1} \left| \mathbf{A}_n^{-1} \big( F_n(\mathbf{u}_n^1) - F_n(\mathbf{u}_n^2) \big) \right| \geq |\mathbf{u}_n^1 - \mathbf{u}_n^2|,
$$

and taking norms we get

$$
\left\| (\bar{\mathbf{X}}_n)^{-1} \right\|_\infty \left\| \mathbf{A}_n^{-1} \right\|_\infty \left\| F_n(\mathbf{u}_n^1) - F_n(\mathbf{u}_n^2) \right\|_\infty \geq \left\| \mathbf{u}_n^1 - \mathbf{u}_n^2 \right\|_\infty.
$$

So we have two tasks.

1. Giving an upper bound for $\left\| \mathbf{A}_n^{-1} \right\|_\infty$.

2. Giving an upper bound for $\left\| (\bar{\mathbf{X}}_n)^{-1} \right\|_\infty$, which includes the following subtasks.

   (a) Finding an appropriate $\mathbf{W}_n$ which is an upper estimate for $L|\mathbf{A}_n^{-1}||\mathbf{B}_n|$ and

   (b) proving that $\bar{\mathbf{X}}_n = \mathbf{I} - \mathbf{W}_n$ is still an M-matrix using a dominant vector so that we get an upper bound for $\left\| (\bar{\mathbf{X}}_n)^{-1} \right\|_\infty$ independent of $n$.

**1.** Note that

$$
\mathbf{A}_n^{-1} =
\begin{pmatrix}
\mathbf{I} & \mathbf{0} \\
-\mathbf{A}_{n,0}^{-1}\mathbf{A}_{n,\partial} & \mathbf{A}_{n,0}^{-1}
\end{pmatrix}.
$$

We split the task here as well by first calculating $\mathbf{A}_{n,0}^{-1}$ then estimating the term $-\mathbf{A}_{n,0}^{-1}\mathbf{A}_{n,\partial}$.

**(a)** $\mathbf{A}_{n,0}$ is a lower triangular Toeplitz matrix with inverse of the same type.

**Lemma 3.1.**

$$\mathbf{A}_{n,0}^{-1} = h \begin{pmatrix} a_1 & 0 & \ldots & \ldots & 0 \\ a_2 & a_1 & 0 & \ldots & 0 \\ \vdots & \ddots & \ddots & \ddots & \vdots \\ \vdots & & \ddots & \ddots & \vdots \\ a_n & a_{n-1} & \ldots & \ldots & a_1 \end{pmatrix},$$

*where*

$$a_l = \sum_{i=1}^{\hat{k}} \frac{(l+k-2)!}{(l+k-k_i-1)!} \frac{\xi_i^{l+k-k_i-1}}{\alpha_0 \prod_{j \neq i}(\xi_i - \xi_j)}, \qquad l = 1, \ldots, n \tag{3.2}$$

*where $k_i$ denotes the multiplicity of $\xi_i$, the roots of the first characteristic polynomial $\varrho$; $\sum k_i = k$, and the number of the different roots is $\hat{k}$.*

*If all of the roots of $\varrho$ are simple then (3.2) simplifies to*

$$a_l = \sum_{i=1}^{k} \frac{\xi_i^{l+k-2}}{\alpha_0 \prod_{j \neq i}(\xi_i - \xi_j)} = \sum_{i=1}^{k} \frac{\xi_i^{l+k-2}}{\varrho'(\xi_i)}, \qquad l = 1, \ldots, n. \tag{3.3}$$

*Proof.* Introducing $\mathbf{H} \in \mathbb{R}^{n \times n}$

$$\mathbf{H} = \begin{pmatrix} 0 & 0 & \ldots & \ldots & 0 \\ 1 & 0 & 0 & \ldots & 0 \\ 0 & 1 & 0 & 0 & \ldots \\ \vdots & \ddots & \ddots & \ddots & \vdots \\ 0 & \ldots & 0 & 1 & 0 \end{pmatrix},$$

and using the identity $(\mathbf{I} - x\mathbf{H})(\mathbf{I} + x\mathbf{H} + \ldots + (x\mathbf{H})^{n-1}) = \mathbf{I} + (x\mathbf{H})^n = \mathbf{I}$, we get

$$(\mathbf{I} - x\mathbf{H})^{-1} = \mathbf{I} + x\mathbf{H} + \ldots + (x\mathbf{H})^{n-1}. \tag{3.4}$$

$$h\mathbf{A}_{n,0} = \alpha_0 \mathbf{I} + \alpha_1 \mathbf{H} + \alpha_2 \mathbf{H}^2 + \ldots + \alpha_k \mathbf{H}^k = \alpha_k \prod_{i=1}^{k}(\mathbf{H} - x_i \mathbf{I})$$

$$= \alpha_k(-1)^k \left(\prod_{i=1}^{k} x_i\right) \prod_{i=1}^{k}\left(\mathbf{I} - \frac{1}{x_i}\mathbf{H}\right) = \alpha_0 \prod_{i=1}^{k}\left(\mathbf{I} - \frac{1}{x_i}\mathbf{H}\right)$$

$$= \alpha_0 \prod_{i=1}^{k}(\mathbf{I} - \xi_i \mathbf{H}),$$

where $\xi_i$ are the roots of the first characteristic polynomial $\varrho$, since $\alpha_0 + \alpha_1 x + \alpha_2 x^2 + \cdots + \alpha_k x^k$ is the reciprocal polynomial of $\varrho$. Note that the $(\mathbf{I} - \xi_i \mathbf{H})$-s commute.

Using (3.4), we get

$$\mathbf{A}_{n,0}^{-1} = \frac{h}{\alpha_0} \prod_{i=1}^{k} \begin{pmatrix} 1 & 0 & \ldots & \ldots & 0 \\ \xi_i & 1 & 0 & \ldots & 0 \\ \vdots & \ddots & \ddots & \ddots & \vdots \\ \vdots & & \ddots & \ddots & \vdots \\ \xi_i^{n-1} & \xi_i^{n-2} & \ldots & \ldots & 1 \end{pmatrix}. \tag{3.5}$$

Finally we use induction to get the formula (3.2). $\qquad\square$

We remark that Lemma 3.1 corresponds to the solution formula for homogeneous linear difference equations which is used in other proofs.

Formula (3.2) has some immediate profit. One can see that the weak stability is necessary and sufficient to have a constant $K_1$ for which $|a_l| < K_1$ holds for all $n$ and $l = 1, \ldots, n$. As a consequence we have that the weak stability is necessary and sufficient to have a constant $K_2$ for which $\|\mathbf{A}_{n,0}^{-1}\|_\infty < K_2$ holds for all $n$. In this case $K_2$ can be chosen as $K_2 = K_1$.

**(b)** If we assume the weak stability one can see that

$$|(-\mathbf{A}_{n,0}^{-1}\mathbf{A}_{n,\partial})_{ij}| < K_1 \alpha k$$

holds, where $\alpha = \max |\alpha_i|$. This means that $\| - \mathbf{A}_{n,0}^{-1}\mathbf{A}_{n,\partial}\|_\infty < K_1 \alpha k^2$.

Consequently the weak stability is necessary and sufficient to have a constant $\hat{K}$ for which $\|\mathbf{A}_n^{-1}\|_\infty < \hat{K}$ holds for all $n$. In this case $\hat{K}$ can be chosen as $\hat{K} = \max\{1, K_1(1 + \alpha k^2)\}$.

Choosing $f \equiv 0$ we get that the weak stability is necessary to the stability in the $\|\cdot\|_\infty, \|\cdot\|_\infty$ pair.

**2.**

**(a)** Note that

$$|\mathbf{A}_n^{-1}||\mathbf{B}_n| = \begin{pmatrix} \mathbf{0} & \mathbf{0} \\ |\mathbf{A}_{n,0}^{-1}||\mathbf{B}_{n,\partial}| & |\mathbf{A}_{n,0}^{-1}||\mathbf{B}_{n,0}| \end{pmatrix}$$

is a lower triangular matrix. If the weak root-condition holds one can see that its entries can be estimated similarly as in the last paragraph.

$$(|\mathbf{A}_n^{-1}||\mathbf{B}_n|)_{ij} < h\,K_1 \beta k\,,$$

where $\beta = \max |\beta_i|$.

Thus we can choose

$$\mathbf{W}_n = \begin{pmatrix} \bar{L}h & 0 & \ldots & 0 \\ \bar{L}h & \bar{L}h & 0 & \ldots \\ \vdots & \ddots & \ddots & \vdots \\ \bar{L}h & \ldots & \bar{L}h & \bar{L}h \end{pmatrix}\,,$$

with $\bar{L} = K_1 \beta k$.

**(b)**

$$\bar{\mathbf{X}}_n = \mathbf{I} - \mathbf{W}_n = \begin{pmatrix} 1 - \bar{L}h & 0 & \ldots & 0 \\ -\bar{L}h & 1 - \bar{L}h & 0 & \ldots \\ \vdots & \ddots & \ddots & \vdots \\ -\bar{L}h & \ldots & -\bar{L}h & 1 - \bar{L}h \end{pmatrix}$$

is inverse nonnegative for large enough $n$-s. To prove that we choose $d(t) = e^{\bar{L}t}$ and so the dominant-vector is $\mathbf{d}_n$ with $(\mathbf{d}_n)_i = e^{\bar{L}t_i} > 0$ and $\|\mathbf{d}_n\|_\infty = e^{\bar{L}}$. Then

$$(\bar{\mathbf{X}}_n \mathbf{d}_n)_i = \begin{cases} 1 - \bar{L}h, & i = 0, \\ e^{\bar{L}t_i} - \bar{L}h \sum_{j=0}^{i} e^{\bar{L}t_j}, & i = 1, \ldots, n + k - 1, \end{cases}$$

$$e^{\bar{L}t_i} - \bar{L}h \sum_{j=0}^{i} e^{\bar{L}t_j} = e^{\bar{L}t_i} - \bar{L}h \frac{e^{\bar{L}(t_i+h)} - 1}{e^{\bar{L}h} - 1}$$

$$= e^{\bar{L}t_i} - \bar{L}\frac{1}{\frac{e^{\bar{L}h}-1}{h}}(e^{\bar{L}(t_i+h)} - 1) \geq e^{\bar{L}t_i} - (e^{\bar{L}(t_i+h)} - 1)$$

$$= 1 - he^{\bar{L}t_i}\frac{e^{\bar{L}h} - 1}{h} \geq 1 - he^{\bar{L}}2\bar{L},$$

if $h$ is small enough. Thus $(\bar{\mathbf{X}}_n\mathbf{d}_n)_i \geq \frac{1}{2}$ if $h$ is small enough, thus for these $h$-s $\bar{\mathbf{X}}_n^{-1} \geq \mathbf{0}$ and $\|\bar{\mathbf{X}}_n^{-1}\|_\infty \leq 2e^{\bar{L}}$.

Summarizing the results we can state the following.

**Theorem 3.2.** *The weak stability is necessary and sufficient to the stability in the $\|\cdot\|_\infty$, $\|\cdot\|_\infty$ pair.*

## 4  Remarks

Beyond that we have proved we wanted, some part of the proof can be exploited to get interesting additional results.

- Using the formulas (1.6), (3.2) and (3.3) and assuming strong stability we have

$$\lim_{n,l\to\infty} a_l = 1. \tag{4.1}$$

  Which is not surprising since the matrix $\mathbf{A}_n^{-1}$ is expected to represent some numerical quadrature formula.

  Weak stability is not enough to ensure (4.1) as the following example shows. Consider the Milne method

$$(F_n(\mathbf{u}_n))_i = \begin{cases} u_i - c^i, & i = 0,1 \\ \frac{1}{h}\left(\frac{1}{2}u_i - \frac{1}{2}u_{i-2}\right) - \left(\frac{1}{6}f(u_i) + \frac{4}{6}f(u_{i-1}) + \frac{1}{6}f(u_{i-2})\right), & i = 2,\ldots,n+1. \end{cases}$$

  For this method $a_l = 2$, if $l$ is odd and $a_l = 0$, if $l$ is even.

- Based on (4.1) and assuming strong stability we can conclude that

$$\lim_{n\to\infty} \|\mathbf{A}_{n,0}^{-1}\|_\infty = 1.$$

- Assuming strong stability we have another consequence. If $\mathbf{A}_{n,0}$ is inverse nonnegative for small $n$-s then it is inverse nonnegative for all $n$.

  It is trivial that for the Adams methods $\mathbf{A}_{n,0}$ is inverse nonnegative since their matrix is identical to the matrix of EE. It can be checked that $\mathbf{A}_{n,0}$ is inverse nonnegative for the BDF methods ($k = 1,\ldots,6$) as well, in spite of not being a Z-matrix for $k > 1$.

- For $k > 1$ $\mathbf{A}_n$ is not a Z-matrix. We note that the norm estimate of Lemma 5.3 holds not only for M-matrices, but it is true for inverse nonnegative matrices as well. Knowing this, it is tempting to try to get an upper bound for $\|\mathbf{A}_n^{-1}\|_\infty$ similarly we did in the case of EE. We might use the same function $d(t) = e^t$ to construct the dominant vector for which it is easy to prove that $\mathbf{A}_n\mathbf{d}_n > \mathbf{0}$ holds. But the problem is that $\mathbf{A}_n$ is not inverse nonnegative any more.

## 5 Appendix

We collected here the necessary information on Z- and M-matrices. The Reader can find more details in [1,8].

**Definition 5.1.** A matrix **M** is said to be a *Z-matrix* if its offdiagonal entries are nonpositive. A matrix **M** is said to be a regular *M-matrix* if it is a regular Z-matrix, moreover, $\mathbf{M}^{-1} \geq \mathbf{0}$ holds.

**Theorem 5.2.** *The matrix* **M** *is assumed to be a Z-matrix. Then the following are equivalent.*

1. **M** *is a regular M-matrix.*

2. $\exists \, \mathbf{d} > \mathbf{0}$: $\mathbf{Md} > \mathbf{0}$.

**Lemma 5.3.** *The matrix* **M** *is assumed to be an M-matrix and* **d** *is a corresponding dominant vector (i.e.* $\mathbf{d} > \mathbf{0} : \mathbf{Md} > \mathbf{0}$*). Then the following estimate holds*

$$\|\mathbf{M}^{-1}\|_\infty \leq \frac{\|\mathbf{d}\|_\infty}{\min_i (\mathbf{Md})_i} \, . \tag{5.1}$$

## References

[1] A. Berman, R. J. Plemmons, *Nonnegative matrices in the mathematical sciences*, Classics in Applied Mathematics, Vol. 9, Society for Industrial and Applied Mathematics (SIAM), Philadelphia, PA, Revised reprint of the 1979 original, 1994. MR1298430

[2] I. Faragó, M. E. Mincsovics, I. Fekete, Notes on the basic notions in nonlinear numerical analysis *Electron. J. Qual. Theory Differ. Equ., Proc. 9'th Coll. QTDE* **2012**, No. 6, 1–22. MR3338525; url

[3] W. Gautschi, *Numerical analysis*, Birkhäuser, 2011.

[4] P. Henrici, *Discrete variable methods in ordinary differential equations*, John Wiley & Sons, Inc., New York–London, 1962. MR0135724

[5] K. W. Morton, *Numerical solution of ordinary differential equations*, Oxford University Computing Laboratory, 1987.

[6] H. J. Stetter, *Analysis of discretization methods for ordinary differential equations*, Springer Tracts in Natural Philosophy, Springer-Verlag, New-York–Heidelberg, 1973. MR0426438

[7] G. Stoyan, G. Takó, *Numerikus módszerek 2.* (in Hungarian) [Numerical methods 2.], ELTE-Typotex, Budapest, 1995

[8] R. S. Varga, *Matrix iterative analysis*, Springer Series in Computational Mathematics, Vol. 27, Springer-Verlag, Berlin, 2000. MR1753713; url