



# A convergence study for reduced rank extrapolation on nonlinear systems

Avram Sidi<sup>1</sup>

Received: 31 January 2019 / Accepted: 24 July 2019 / Published online: 20 August 2019  
© Springer Science+Business Media, LLC, part of Springer Nature 2019

## Abstract

Reduced Rank Extrapolation (RRE) is a polynomial type method used to accelerate the convergence of sequences of vectors  $\{\mathbf{x}_m\}$ . It is applied successfully in different disciplines of science and engineering in the solution of large and sparse systems of linear and nonlinear equations of very large dimension. If  $s$  is the solution to the system of equations  $\mathbf{x} = \mathbf{f}(\mathbf{x})$ , first, a vector sequence  $\{\mathbf{x}_m\}$  is generated via the fixed-point iterative scheme  $\mathbf{x}_{m+1} = \mathbf{f}(\mathbf{x}_m)$ ,  $m = 0, 1, \dots$ , and next, RRE is applied to this sequence to accelerate its convergence. RRE produces approximations  $s_{n,k}$  to  $s$  that are of the form  $s_{n,k} = \sum_{i=0}^k \gamma_i \mathbf{x}_{n+i}$  for some scalars  $\gamma_i$  depending (nonlinearly) on  $\mathbf{x}_n, \mathbf{x}_{n+1}, \dots, \mathbf{x}_{n+k+1}$  and satisfying  $\sum_{i=0}^k \gamma_i = 1$ . The convergence properties of RRE when applied in conjunction with linear  $\mathbf{f}(\mathbf{x})$  have been analyzed in different publications. In this work, we discuss the convergence of the  $s_{n,k}$  obtained from RRE with nonlinear  $\mathbf{f}(\mathbf{x})$  (i) when  $n \rightarrow \infty$  with fixed  $k$ , and (ii) in two so-called *cycling* modes.

**Keywords** Vector extrapolation methods · Minimal polynomial extrapolation (MPE) · Reduced rank extrapolation (RRE) · Krylov subspace methods · Nonlinear equations · Cycling mode

**Mathematics Subject Classification (2010)** Primary 65B05 · 65H10; Secondary 65F10

## 1 Introduction

Consider a system of nonlinear algebraic equations of dimension  $N$ , which we choose to write as

$$\mathbf{x} = \mathbf{f}(\mathbf{x}), \quad \mathbf{f} : \mathbb{C}^N \rightarrow \mathbb{C}^N; \quad s \text{ solution}, \quad (1.1)$$

---

✉ Avram Sidi  
asidi@cs.technion.ac.il  
<http://www.cs.technion.ac.il/~asidi>

<sup>1</sup> Computer Science Department, Technion - Israel Institute of Technology, Haifa 32000, Israel

where

$$\mathbf{x} = [x^{(1)}, \dots, x^{(N)}]^T, \quad \mathbf{s} = [s^{(1)}, \dots, s^{(N)}]^T; \quad x^{(i)}, s^{(i)} \text{ scalars}, \quad (1.2)$$

and

$$\mathbf{f}(\mathbf{x}) = [f_1(\mathbf{x}), \dots, f_N(\mathbf{x})]^T; \quad f_i(\mathbf{x}) = f_i(x^{(1)}, \dots, x^{(N)}) \text{ scalar functions}. \quad (1.3)$$

One immediate way of solving this system is via the fixed-point iterative scheme

$$\mathbf{x}_{m+1} = \mathbf{f}(\mathbf{x}_m), \quad m = 0, 1, \dots; \quad \text{for some } \mathbf{x}_0, \quad (1.4)$$

provided the sequence  $\{\mathbf{x}_m\}$  converges. Let  $\mathbf{f}(\mathbf{x})$  be twice continuously differentiable in a neighborhood of  $\mathbf{s}$ , and let  $\mathbf{F}(\mathbf{x})$  be the Jacobian matrix of  $\mathbf{f}$  evaluated at  $\mathbf{x}$ , that is,

$$\mathbf{F}(\mathbf{x}) = \begin{bmatrix} f_{1,1}(\mathbf{x}) & f_{1,2}(\mathbf{x}) & \cdots & f_{1,N}(\mathbf{x}) \\ f_{2,1}(\mathbf{x}) & f_{2,2}(\mathbf{x}) & \cdots & f_{2,N}(\mathbf{x}) \\ \vdots & \vdots & & \vdots \\ f_{N,1}(\mathbf{x}) & f_{N,2}(\mathbf{x}) & \cdots & f_{N,N}(\mathbf{x}) \end{bmatrix}; \quad f_{i,j}(\mathbf{x}) = \frac{\partial f_i}{\partial x^{(j)}}(\mathbf{x}). \quad (1.5)$$

It is known that (see Ortega and Rheinboldt [22], for example) if  $\rho(\mathbf{F}(\mathbf{x}))$ , the spectral radius of  $\mathbf{F}(\mathbf{x})$ , is such that  $\rho(\mathbf{F}(\mathbf{s})) < 1$  and if  $\mathbf{x}_0$  is sufficiently close to  $\mathbf{s}$ , then the sequence  $\{\mathbf{x}_m\}$  converges to  $\mathbf{s}$ . The closer  $\rho(\mathbf{F}(\mathbf{s}))$  is to one, the slower is the convergence of  $\{\mathbf{x}_m\}$  to  $\mathbf{s}$ ; this is the case in most practical engineering applications.

The convergence of  $\{\mathbf{x}_m\}$  to  $\mathbf{s}$  can be accelerated substantially by applying to it a vector extrapolation method. When applied to  $\{\mathbf{x}_m\}$ , an extrapolation method produces approximations  $s_{n,k}$  to  $\mathbf{s}$  that are, either directly or indirectly, of the form

$$s_{n,k} = \sum_{i=0}^k \gamma_i \mathbf{x}_{n+i}; \quad \gamma_i \text{ some scalars,} \quad \sum_{i=0}^k \gamma_i = 1, \quad (1.6)$$

the  $\gamma_i$  depending nonlinearly on the  $\mathbf{x}_m$  used in constructing  $s_{n,k}$ . Let  $M$  be the number of the  $\mathbf{x}_m$  needed to construct  $s_{n,k}$ . (Of course,  $M$  is not necessarily the same for all vector extrapolation methods.)<sup>1</sup>

For the sake of completeness, here we mention briefly those vector extrapolation methods that have been shown to be useful in applications.

1. **Polynomial type methods:** These are *minimal polynomial extrapolation (MPE)*, *reduced rank extrapolation (RRE)*, *modified minimal polynomial extrapolation*

---

<sup>1</sup>It is clear that the integers  $n$  and  $k$  are chosen by the user and that  $M$  is determined by  $n$ ,  $k$ , and the extrapolation method being used.

(MMPE), and the most recent *singular value decomposition-based minimal polynomial extrapolation (SVD-MPE)*. MPE was introduced by Cabay and Jackson [9], RRE was introduced independently by Kaniel and Stein [18], Eddy [11], and Mešina [21].<sup>2</sup> MMPE was introduced independently by Brezinski [6], Pugachev [23], and Sidi, Ford, and Smith [35]. SVD-MPE is a new method by Sidi [31].

2. **Epsilon algorithms:** These are the *scalar epsilon algorithm (SEA)*, the *vector epsilon algorithm (VEA)*, and the *topological epsilon algorithm (TEA)*. SEA is a method that is based entirely on the famous *epsilon algorithm* of Wynn [45] that implements the transformation of Shanks [24] for scalar sequences. VEA was introduced by Wynn [46]. TEA was introduced by Brezinski [6].

For an earlier account of the epsilon algorithms, see the book by Brezinski [7]. For a comprehensive survey covering the developments that took place until the 1980s, see the survey paper by Smith, Ford, and Sidi [39] and the book by Brezinski and Redivo Zaglia [8]. For a geometric approach to the treatment of vector extrapolation methods as these are being applied to linear systems, see Jbilou and Sadok [16]. For a more recent review of MPE and RRE, see Sidi [30]. For a detailed and up-to-date treatment, including development, analysis, numerical implementation, and various applications, of all these methods, see the recent book of Sidi [33].

Numerically stable and efficient algorithms for implementing polynomial methods have been proposed by Sidi [27], [31] for MPE, RRE, and SVD-MPE and by Jbilou and Sadok [17] for MMPE. The epsilon algorithms are normally implemented via their definitions, which involve recursion relations. When applied to sequences  $\{\mathbf{x}_m\}$  generated via fixed-point iterative schemes from systems of linear equations, MPE, RRE, and TEA turn out to be equivalent to known Krylov subspace methods for linear systems. This is explored in Sidi [26]. Yet another recent paper by Sidi [32] shows that MPE and RRE are very closely related in more than one way.

Now, all the methods mentioned above have interesting convergence and convergence acceleration properties that concern the precise asymptotic behavior of the sequences  $\{s_{n,k}\}_{n=0}^\infty$ , with fixed  $k$ , when the sequences  $\{\mathbf{x}_m\}$  are generated via fixed-point iterative schemes from systems of linear equations; see Sidi [25], [28], Sidi, Ford, and Smith [35], and Sidi and Bridger [34], and also Sidi [33, Chapter 6] for the methods MPE, RRE, MMPE, and TEA, Wynn [47] and Sidi [29] for SEA, and Graves-Morris and Saff [14] for VEA. We shall call this mode of usage of vector extrapolation methods the *n-Mode*.

Unfortunately, the *n-Mode* convergence theories that apply to the case in which  $f(\mathbf{x})$  is linear do not apply to the case in which  $f(\mathbf{x})$  is nonlinear. This is one of the topics we would like to study here, RRE being the extrapolation method used. That is, we would like to investigate the convergence properties of the sequences  $\{s_{n,k}\}_{n=0}^\infty$ , with fixed  $k$ , obtained by applying RRE to  $\{\mathbf{x}_m\}$  generated as in (1.4), where  $f(\mathbf{x})$  is nonlinear.

---

<sup>2</sup>The approaches of [18] and [21] to RRE are almost identical, in the sense that  $s_{n,k} = \sum_{i=0}^k \gamma_i \mathbf{x}_{n+i}$  in [21], while  $s_{n,k} = \sum_{i=0}^k \gamma_i \mathbf{x}_{n+i+1}$  in [18], the  $\gamma_i$  being the same for both. The approaches of [11] and [21] are completely different, however; their equivalence was proved in the review paper of Smith, Ford, and Sidi [39].

The numerical implementations of polynomial extrapolation methods and of epsilon algorithms, when generating the vectors  $s_{n,k}$ , necessitate the keeping of resp.  $k + 2$  and  $2k + 1$  vectors in core memory simultaneously. In case we would like to increase  $k$  to improve the quality of the  $s_{n,k}$ , this may pose a serious problem when we are dealing with very high dimensional vectors, which is the case in most large scale applications. Within the context described via (1.1)–(1.4) in the first paragraph of this section, it is best to apply vector extrapolation methods in the so-called *cycling* mode, and this has been the usual practice. This mode of usage of vector extrapolation methods, which we shall call the *C-Mode*, can be described via the following steps:

**C-Mode**

- C0. Choose integers  $n \geq 0$  and  $k \geq 1$  and an initial vector  $x_0$ .
- C1. Compute the vectors  $x_1, x_2, \dots, x_M$  [via  $x_{m+1} = f(x_m)$ ].<sup>3</sup>
- C2. Apply the extrapolation method to the vectors  $x_n, x_{n+1}, \dots, x_M$ , and compute  $s_{n,k}$ .
- C3. If  $s_{n,k}$  satisfies the accuracy test, stop.  
Otherwise, set  $x_0 = s_{n,k}$  and go to step C1.

We call each application of steps C1–C3 a *cycle* and denote by  $s^{(r)}$  the  $s_{n,k}$  computed in the  $r$ th cycle. We will also denote the initial vector  $x_0$  in step C0 by  $s^{(0)}$ . Under suitable conditions, it has been shown rigorously for MPE and RRE that the sequence  $\{s^{(r)}\}_{r=0}^\infty$  has very good convergence properties when  $f(x)$  is linear. See [36], [37]. See also [33, Chapter 7]. The case in which  $f(x)$  is nonlinear has proved to be complicated and has not been resolved till the present.

In some cases, RRE stalls if applied in the C-Mode with  $n = 0$ , in the sense that it takes too many iterations until one sees meaningful convergence; in such cases, even a moderate  $n > 0$  can be very helpful to accelerate convergence effectively. See the numerical examples in [36], [37].

A different cycling procedure involving the minimal polynomial of the (constant) Jacobian matrix  $F(s)$  with respect to a nonzero vector<sup>4</sup> has been considered in various publications. The description of this procedure, which we shall call the *MC-Mode*, is as follows:

**MC-Mode**

- MC0. Choose an integer  $n \geq 0$  and an initial vector  $x_0$ .

---

<sup>3</sup>Note that  $M = n + k + 1$  for MPE, RRE, MMPE, and SVD-MPE, while  $M = n + 2k$  for SEA, VEA, and TEA.

<sup>4</sup>Given a nonzero vector  $u \in \mathbb{C}^N$ , the monic polynomial  $P(\lambda)$  is said to be a *minimal polynomial of the matrix  $T \in \mathbb{C}^{N \times N}$  with respect to  $u$*  if  $P(T)u = 0$  and if  $P(\lambda)$  has smallest degree.

The polynomial  $P(\lambda)$  exists and is unique. Moreover, if  $P_1(T)u = 0$  for some polynomial  $P_1(\lambda)$  with  $\deg P_1 > \deg P$ , then  $P(\lambda)$  divides  $P_1(\lambda)$ . In particular,  $P(\lambda)$  divides the minimal polynomial of  $T$ , which in turn divides the characteristic polynomial of  $T$ . [Thus, the degree of  $P(\lambda)$  is at most  $N$  and its zeros are some or all of the eigenvalues of  $T$ .]

- MC1. Compute the vectors  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_M$  [via  $\mathbf{x}_{m+1} = \mathbf{f}(\mathbf{x}_m)$ ],  $M$  being as explained in footnote<sup>3</sup>, with  $k$  there being the degree of the minimal polynomial of  $\mathbf{F}(\mathbf{s})$  with respect to  $\boldsymbol{\epsilon}_n = \mathbf{x}_n - \mathbf{s}$ .<sup>5</sup>
- MC2. Apply the extrapolation method to the vectors  $\mathbf{x}_n, \mathbf{x}_{n+1}, \dots, \mathbf{x}_M$ , and compute  $\mathbf{s}_{n,k}$ .
- MC3. If  $\mathbf{s}_{n,k}$  satisfies the accuracy test, stop.  
 Otherwise, set  $\mathbf{x}_0 = \mathbf{s}_{n,k}$  and go to step MC1.

As before, we call each application of steps MC1–MC3 a *cycle* and denote by  $\mathbf{s}^{(r)}$  the  $\mathbf{s}_{n,k}$  computed in the  $r$ th cycle.<sup>6</sup> We will also denote the initial vector  $\mathbf{x}_0$  in step MC0 by  $\mathbf{s}^{(0)}$ . It is observed in many numerical examples that the sequence  $\{\mathbf{s}^{(r)}\}_{r=0}^\infty$  converges quadratically to the solution  $\mathbf{s}$  of the system  $\mathbf{x} = \mathbf{f}(\mathbf{x})$  when  $\mathbf{f}(\mathbf{x})$  is nonlinear.<sup>7</sup> The first papers dealing with this topic (that is the MC-Mode with  $\mathbf{s}_{0,k}$  only) are those by Brezinski [4], [5], Gekeler [12], and Skelboe [38]. Of these, [4], [5], and [12] consider the application of the epsilon algorithms, while [38] also considers the application of MPE and RRE. The quadratic convergence proofs in all of these papers have a gap in that they all end up with the relation

$$\|\mathbf{s}^{(r+1)} - \mathbf{s}\|_2 \leq K_r \|\mathbf{s}^{(r)} - \mathbf{s}\|_2^2,$$

from which they conclude that  $\{\mathbf{s}^{(r)}\}_{r=0}^\infty$  converges quadratically. However,  $K_r$  is a scalar that depends on  $r$  through  $\mathbf{s}^{(r)}$ , and the proofs do not show how it depends on  $r$ . In particular, they do not show whether  $K_r$  is bounded in  $r$  or how it grows with  $r$  if it is not bounded. This gap was disclosed in the review paper of Smith, Ford, and Sidi [39].

A more recent paper by Jbilou and Sadok [15] deals with the same MC-Mode cycling via MPE and RRE. Yet another paper by Le Ferrand [20] treats TEA. Both these works provide proofs of quadratic convergence by imposing some global conditions on the whole sequence  $\{\mathbf{s}^{(r)}\}_{r=0}^\infty$  as well as on  $\mathbf{f}(\mathbf{x})$ . (See also Laurens and Le Ferrand [19].)

In this work, we present a new convergence study of RRE when it is being applied to nonlinear systems. Specifically, we treat the convergence of RRE (i) in the  $n$ -Mode, and (ii) in the two cycling modes mentioned above. By making a *global* assumption, we are able to prove convergence in all cases. We can justify heuristically the plausibility of this assumption; we do not have a rigorous justification for it, however. This difficulty is inherent to all studies. We explore the source of this difficulty here. It must be mentioned that the difficulties that exist in the previous papers mentioned above are similar to ours, although they take different forms. Whether and how we can circumvent these difficulties is not clear at this time.

<sup>5</sup>It is clear that to apply any of the extrapolation methods in this mode, one needs to know the matrix  $\mathbf{F}(\mathbf{s})$ , for which one also needs to know the solution  $\mathbf{s}$ .

<sup>6</sup>Note that  $k$  is not necessarily fixed in this mode of cycling; it may vary from one cycle to the next. It always satisfies  $k \leq N$ , however.

<sup>7</sup>Quadratic convergence is relevant only when  $\mathbf{f}(\mathbf{x})$  is nonlinear. When  $\mathbf{f}(\mathbf{x})$  is linear, that is,  $\mathbf{f}(\mathbf{x}) = \mathbf{T}\mathbf{x} + \mathbf{d}$ , where  $\mathbf{T}$  is a fixed  $N \times N$  matrix and  $\mathbf{d}$  is a fixed vector, hence  $\mathbf{F}(\mathbf{s}) = \mathbf{T}$ , the solution  $\mathbf{s}$  is obtained already at the end of step MC2 of the first cycle, that is, we have  $\mathbf{s}^{(1)} = \mathbf{s}$ . Therefore, there is nothing to analyze when  $\mathbf{f}(\mathbf{x})$  is linear.

The plan of this paper is as follows: In Section 2, we give a brief description of RRE, which is needed throughout. In Section 3, we derive a formula for the error vector  $s_{n,k} - s$  when the vectors  $x_m$  are generated via (1.4) with a nonlinear  $f(x)$ . In Section 4, we use this error formula to derive an upper bound on  $\|s_{n,k} - s\|$ . In Section 5, we complete the convergence studies of RRE in the different modes mentioned above. We mention that our results concerning the convergence of RRE in the  $n$ -Mode and the C-Mode are the first ones in the literature of extrapolation methods. In our study, we make much use of the results presented in Sidi [26] throughout these studies. In Section 6, we discuss the nature of the problem/difficulty mentioned above and compare our global assumption with that of [15]. In the appendix, we review some known theorems concerning Moore–Penrose generalized inverses of perturbed matrices, which we use in Section 4. (For generalized inverses, see Ben-Israel and Greville [3] and Campbell and Meyer [10], for example.)

Throughout this work, we will use lowercase boldface italic letters to denote vectors and we will use uppercase boldface italic letters to denote matrices.

Finally, we mention that in our study of RRE, we employ two different vector norms:

- (i) The standard  $l_2$  vector norm defined via  $\|z\|_2 = \sqrt{z^*z}$  and the matrix norm induced by it, namely,  $\|A\|_2 = \sigma_{\max}(A)$ , where  $\sigma_{\max}(A)$  is the largest singular value of the matrix  $A$ .
- (ii) The  $G$  norm defined via  $\|z\| = \|Gz\|_2$  and the matrix norm  $\|A\|$  induced by it, where  $G = I - F(s)$ . Note that  $G$  is nonsingular since  $F(s)$  does not have unity as an eigenvalue; therefore, the  $G$  norm is a true vector norm.

Of course, the two vector norms are equivalent and we have

$$\frac{1}{\|G^{-1}\|_2} \|z\|_2 \leq \|z\| \leq \|G\|_2 \|z\|_2. \tag{1.7}$$

When  $A$  is an  $N \times N$  (square) matrix, we have  $\|A\| = \|GAG^{-1}\|_2$ . We make extensive use of these connections between the two norms,  $\|\cdot\|$  and  $\|\cdot\|_2$ , in the sequel.

## 2 Description of RRE

Consider the system of equations given in (1.1)–(1.3), and let the sequence  $\{x_m\}$  be generated via the fixed-point iterative scheme in (1.4).

Define the first and second order differences of the  $x_m$  as in

$$u_m = x_{m+1} - x_m, \quad w_m = u_{m+1} - u_m = x_{m+2} - 2x_{m+1} + x_m, \quad m = 0, 1, \dots, \tag{2.1}$$

and, for some fixed  $n \geq 0$ , form the  $N \times (j + 1)$  matrices

$$U_j = [u_n | u_{n+1} | \dots | u_{n+j}], \quad W_j = [w_n | w_{n+1} | \dots | w_{n+j}], \quad j = 0, 1, \dots \tag{2.2}$$

Then the  $\gamma_i$  in (1.6) for RRE are the solution to the constrained standard  $l_2$  minimization problem

$$\min_{\gamma_0, \gamma_1, \dots, \gamma_k} \left\| \sum_{i=0}^k \gamma_i \mathbf{u}_{n+i} \right\|_2 \quad \text{subject to} \quad \sum_{i=0}^k \gamma_i = 1, \tag{2.3}$$

which can also be expressed in matrix terms as

$$\min_{\boldsymbol{\gamma}} \|\mathbf{U}_k \boldsymbol{\gamma}\|_2 \quad \text{subject to} \quad \sum_{i=0}^k \gamma_i = 1; \quad \boldsymbol{\gamma} = [\gamma_0, \gamma_1, \dots, \gamma_k]^T \in \mathbb{C}^{k+1}. \tag{2.4}$$

Then, with the solution  $\boldsymbol{\gamma}$  of this problem, the RRE approximation  $s_{n,k}$  is given as in

$$s_{n,k} = \sum_{i=0}^k \gamma_i \mathbf{x}_{n+i}. \tag{2.5}$$

Noting that

$$\mathbf{x}_{n+m} = \mathbf{x}_n + \sum_{j=0}^{m-1} \mathbf{u}_{n+j}, \quad \mathbf{u}_{n+m} = \mathbf{u}_n + \sum_{j=0}^{m-1} \mathbf{w}_{n+j}, \quad m = 0, 1, \dots,$$

we can reexpress  $s_{n,k}$  and  $\mathbf{U}_k \boldsymbol{\gamma}$  as

$$s_{n,k} = \mathbf{x}_n + \sum_{j=0}^{k-1} \xi_j \mathbf{u}_{n+j} = \mathbf{x}_n + \mathbf{U}_{k-1} \boldsymbol{\xi}, \quad \mathbf{U}_k \boldsymbol{\gamma} = \mathbf{u}_n + \sum_{j=0}^{k-1} \xi_j \mathbf{w}_{n+j} = \mathbf{u}_n + \mathbf{W}_{k-1} \boldsymbol{\xi}, \tag{2.6}$$

where

$$\boldsymbol{\xi} = [\xi_0, \xi_1, \dots, \xi_{k-1}]^T \in \mathbb{C}^k; \quad \xi_j = \sum_{i=j+1}^k \gamma_i, \quad j = 0, 1, \dots, k-1. \tag{2.7}$$

The (constrained) minimization problem for the vector  $\boldsymbol{\gamma}$  in (2.4) can now be replaced by the following (unconstrained) minimization problem for the vector  $\boldsymbol{\xi}$  in (2.6):

$$\min_{\boldsymbol{\xi}} \|\mathbf{u}_n + \mathbf{W}_{k-1} \boldsymbol{\xi}\|_2, \quad \boldsymbol{\xi} = [\xi_0, \xi_1, \dots, \xi_{k-1}]^T \in \mathbb{C}^k. \tag{2.8}$$

Now, the solution to this problem (for  $\boldsymbol{\xi}$ ) is simply  $-\mathbf{W}_{k-1}^+ \mathbf{u}_n$ , where  $\mathbf{K}^+$  stands for the Moore–Penrose generalized inverse of the matrix  $\mathbf{K}$ . Upon substituting this into (2.6), we obtain

$$\boxed{s_{n,k} = \mathbf{x}_n - \mathbf{U}_{k-1} \mathbf{W}_{k-1}^+ \mathbf{u}_n}. \tag{2.9}$$

We will be making use of this representation of  $s_{n,k}$  in the sequel. For the above developments, see Sidi [25].

### 3 An error formula for RRE

#### 3.1 RRE on the linear system $x = s + F(s)(x - s)$

Let us now consider the linear system

$$x = \tilde{f}(x), \quad \tilde{f}(x) = s + F(s)(x - s), \tag{3.1}$$

where  $F(s)$  is the Jacobian matrix of  $f$  evaluated at  $s$ , as given in (1.5). Note that  $\tilde{f}(x)$  is simply the linear part of the Taylor series of  $f(x)$  in (1.1) about  $s$ . Clearly,  $s$  is the solution to (3.1) since  $\tilde{f}(s) = s$ .

With the vectors  $x_0, x_1, \dots, x_n$  generated nonlinearly as in (1.4) of the preceding section, let

$$\tilde{x}_n = x_n \quad \text{and} \quad \tilde{x}_{m+1} = \tilde{f}(\tilde{x}_m), \quad m = n, n + 1, \dots \tag{3.2}$$

Following this, define

$$\tilde{\epsilon}_m = \tilde{x}_m - s, \quad \tilde{u}_m = \tilde{x}_{m+1} - \tilde{x}_m, \quad \tilde{w}_m = \tilde{u}_{m+1} - \tilde{u}_m, \quad m = n, n + 1, \dots, \tag{3.3}$$

$$\tilde{U}_j = [\tilde{u}_n \mid \tilde{u}_{n+1} \mid \dots \mid \tilde{u}_{n+j}], \quad \tilde{W}_j = [\tilde{w}_n \mid \tilde{w}_{n+1} \mid \dots \mid \tilde{w}_{n+j}], \quad j = 0, 1, \dots \tag{3.4}$$

Then, by (2.9), the vector  $\tilde{s}_{n,k}$  produced by applying RRE to the sequence  $\{\tilde{x}_m\}$  is

$$\tilde{s}_{n,k} = \tilde{x}_n - \tilde{U}_{k-1} \tilde{W}_{k-1}^+ \tilde{u}_n. \tag{3.5}$$

Upon subtracting  $s$  from both sides of this equality and invoking  $\tilde{\epsilon}_n = \tilde{x}_n - s$ , we obtain the error formula

$$\tilde{s}_{n,k} - s = \tilde{\epsilon}_n - \tilde{U}_{k-1} \tilde{W}_{k-1}^+ \tilde{u}_n. \tag{3.6}$$

The error  $\tilde{s}_{n,k} - s$  has been studied in detail in [25], [28], [34], [36]; for a summary, see [33, Chapters 6,7].<sup>8</sup>

The following result from [26, Theorem 4.2] concerning the application of RRE to vector sequences from fixed-point iteration of *linear* systems will be crucial in our analysis of RRE concerning *nonlinear* systems in Section 5.

**Theorem 3.1** *Denote  $\tilde{F} = F(s)$  for short; thus  $G = I - \tilde{F}$ . Then the vector  $\tilde{s}_{n,k}$  is the solution to the optimization problem*

$$\begin{aligned} \|\tilde{s}_{n,k} - s\| &= \|G(\tilde{s}_{n,k} - s)\|_2 = \min_{g \in \tilde{P}_k} \|g(\tilde{F})G(\tilde{x}_n - s)\|_2, \\ \tilde{P}_k &= \left\{ g(z) = \sum_{j=0}^k \alpha_j z^j : g(1) = 1 \right\}. \end{aligned} \tag{3.7}$$

<sup>8</sup>See also Sidi and Shapira [37] concerning a modified version of restarted GMRES with prior Richardson iterations, that is very closely related to RRE.



*Remarks*

1. If  $k$  is the degree of the minimal polynomial of  $F(s)$  with respect to the vector  $\tilde{\epsilon}_n$ , then  $\tilde{s}_{n,k} = s$ , the solution to (3.1). See footnote<sup>7</sup>.
2. Concerning Theorem 3.1, note that the vector  $G(y - s)$  is simply the *residual* of the vector  $y$  for the linear system  $x = \tilde{f}(x)$  because

$$G(y - s) = y - \tilde{f}(y),$$

and also

$$y - \tilde{f}(y) = 0 \iff y = s, \quad \text{since } G \text{ is nonsingular.}$$

Thus, what Theorem 3.1 means is that  $\|G(\tilde{s}_{n,k} - s)\|_2$ , the  $l_2$  norm of the residual vector of  $\tilde{s}_{n,k} = \sum_{i=0}^k \tilde{\gamma}_i \tilde{x}_{n+i}$  subject to  $\sum_{i=0}^k \tilde{\gamma}_i = 1$ , is the smallest of all the  $l_2$  norms of the residuals of the vectors  $\sum_{i=0}^k \alpha_i \tilde{x}_{n+i}$  subject to  $\sum_{i=0}^k \alpha_i = 1$ . Here we also recall that  $\tilde{x}_n = x_n$  by (3.2).

**3.2 RRE on the nonlinear system  $x = f(x)$**

In the Introduction, we assumed that  $f(x)$  is twice continuously differentiable in a neighborhood of the solution  $s$ . We also assumed that  $\rho(\tilde{F}) < 1$ , where we recall  $\tilde{F} = F(s)$ , thus ensuring the convergence of the sequence  $\{x_m\}$  to  $s$ . We now assume, in addition, that  $\|\tilde{F}\|_2 < 1$  too and define the ball  $B(s, \delta)$  containing  $s$  in its interior via

$$B(s, \delta) = \{x : \|x - s\| \equiv \|G(x - s)\|_2 \leq \delta\}. \tag{3.8}$$

Clearly,  $B(s, \delta)$  is a convex set. In addition, we assume  $f(x)$  is twice continuously differentiable in  $B(s, \delta)$ .

**Lemma 3.2** *For all  $\delta$  sufficiently small, there exists a positive constant  $L < 1$  independent of  $\delta$ , such that*

$$\|x_{m+1} - s\| \leq L \|x_m - s\|, \quad m = 0, 1, \dots, \quad \text{provided } x_0 \in B(s, \delta). \tag{3.9}$$

*Consequently, the whole sequence  $\{x_m\}$  is in  $B(s, \delta)$  and converges to  $s$ .*

*Proof* We begin with the following result that follows from Ortega and Rheinboldt [22, p. 69]:

$$f(x) - f(s) = \int_0^1 F(s + t(x - s))(x - s) dt \quad \text{provided } x \in B(s, \delta).$$

It is important to note that  $s + t(x - s)$ , with  $t \in [0, 1]$ , is a convex combination of  $x$  and  $s$  hence is also in  $B(s, \delta)$ . Multiplying both sides of this equality on the left by  $G$ , we obtain

$$G[f(x) - f(s)] = \int_0^1 [GF(s + t(x - s))G^{-1}][G(x - s)] dt,$$

which, upon taking  $l_2$  norms on both sides and invoking the known fact that

$$\left\| \int_a^b \mathbf{u}(\xi) d\xi \right\|_2 \leq \int_a^b \|\mathbf{u}(\xi)\|_2 d\xi, \quad \mathbf{u}(\xi) \in \mathbb{C}^N,$$

gives

$$\|\mathbf{G}[\mathbf{f}(\mathbf{x}) - \mathbf{f}(\mathbf{s})]\|_2 \leq \int_0^1 \|\mathbf{G}\mathbf{F}(\mathbf{s} + t(\mathbf{x} - \mathbf{s}))\mathbf{G}^{-1}\|_2 \|\mathbf{G}(\mathbf{x} - \mathbf{s})\|_2 dt.$$

Finally, invoking in this last inequality  $\|\mathbf{G}\mathbf{z}\|_2 = \|\mathbf{z}\|$  and the fact that  $\|\mathbf{G}\mathbf{A}\mathbf{G}^{-1}\|_2 = \|\mathbf{A}\|$ , we obtain

$$\begin{aligned} \|\mathbf{f}(\mathbf{x}) - \mathbf{f}(\mathbf{s})\| &\leq \left( \int_0^1 \|\mathbf{F}(\mathbf{s} + t(\mathbf{x} - \mathbf{s}))\| dt \right) \|\mathbf{x} - \mathbf{s}\| \\ &\leq \left[ \max_{0 \leq t \leq 1} \|\mathbf{F}(\mathbf{s} + t(\mathbf{x} - \mathbf{s}))\| \right] \|\mathbf{x} - \mathbf{s}\| \\ &\leq \left[ \max_{\mathbf{z} \in B(\mathbf{s}, \delta)} \|\mathbf{F}(\mathbf{z})\| \right] \|\mathbf{x} - \mathbf{s}\|. \end{aligned} \tag{3.10}$$

Now, by the fact that  $\mathbf{f}(\mathbf{x})$  is twice differentiable in  $B(\mathbf{s}, \delta)$ , it follows that

$$\mathbf{F}(\mathbf{x}) = \mathbf{F}(\mathbf{s} + (\mathbf{x} - \mathbf{s})) = \mathbf{F}(\mathbf{s}) + \mathbf{\Delta}(\mathbf{x} - \mathbf{s}), \tag{3.11}$$

where the matrix  $\mathbf{\Delta}(\mathbf{x} - \mathbf{s})$  satisfies

$$\|\mathbf{\Delta}(\mathbf{x} - \mathbf{s})\| \leq \alpha \|\mathbf{x} - \mathbf{s}\| \quad \text{for some } \alpha > 0 \text{ independent of } \mathbf{x} \in B(\mathbf{s}, \delta). \tag{3.12}$$

Taking norms on both sides of (3.11), realizing that  $\|\mathbf{F}(\mathbf{s})\| = \|\mathbf{F}(\mathbf{s})\|_2$  because  $\mathbf{F}(\mathbf{s})$  and  $\mathbf{G} = \mathbf{I} - \mathbf{F}(\mathbf{s})$  commute, and invoking  $\mathbf{x} \in B(\mathbf{s}, \delta)$ , we have

$$\|\mathbf{F}(\mathbf{x})\| \leq \|\mathbf{F}(\mathbf{s})\|_2 + \alpha \|\mathbf{x} - \mathbf{s}\| \leq \|\mathbf{F}(\mathbf{s})\|_2 + \alpha\delta \quad \forall \mathbf{x} \in B(\mathbf{s}, \delta). \tag{3.13}$$

Since we have assumed that  $\|\mathbf{F}(\mathbf{s})\|_2 < 1$ , we can choose  $\delta$  sufficiently small to cause

$$\max_{\mathbf{x} \in B(\mathbf{s}, \delta)} \|\mathbf{F}(\mathbf{x})\| = L < 1. \tag{3.14}$$

With this, (3.10) becomes

$$\|\mathbf{f}(\mathbf{x}) - \mathbf{f}(\mathbf{s})\| \leq L \|\mathbf{x} - \mathbf{s}\| \quad \forall \mathbf{x} \in B(\mathbf{s}, \delta). \tag{3.15}$$

The proof of (3.9) for the sequence  $\{\mathbf{x}_m\}$  can now be carried out by letting  $\mathbf{x} = \mathbf{x}_m$  in (3.15), recalling that  $\mathbf{f}(\mathbf{x}_m) = \mathbf{x}_{m+1}$  and  $\mathbf{f}(\mathbf{s}) = \mathbf{s}$ , and then proceeding by induction on  $m$ . □

In the sequel, we adopt the shorthand notation

$$\boldsymbol{\epsilon}_m = \mathbf{x}_m - \mathbf{s}, \quad m = 0, 1, \dots; \quad \tilde{\mathbf{F}} = \mathbf{F}(\mathbf{s}). \tag{3.16}$$

We also make use of the fact that  $\|\tilde{\mathbf{F}}\| \leq L < 1$ , which follows from (3.14), and, along with (3.9), guarantees that the sequence  $\{\|\boldsymbol{\epsilon}_m\|\}$  decreases monotonically and converges to zero.

Expanding  $f(x)$  in a Taylor series about the solution  $s$  and using the fact that  $f(s) = s$  and  $f \in C^2(B(s, \delta))$ , we have

$$f(x) = s + \tilde{F} \cdot (x - s) + \mu(x - s), \tag{3.17}$$

where

$$\|\mu(x - s)\| \leq a \|x - s\|^2 \quad \forall x \in B(s, \delta), \quad \text{for some } a > 0. \tag{3.18}$$

Consequently,

$$x_{m+1} = f(x_m) = s + \tilde{F}\epsilon_m + \mu(\epsilon_m) \Rightarrow \epsilon_{m+1} = \tilde{F}\epsilon_m + \mu(\epsilon_m). \tag{3.19}$$

Then, by induction,

$$\epsilon_{n+i} = \tilde{F}^i \epsilon_n + \sum_{j=0}^{i-1} \tilde{F}^{i-j-1} \mu(\epsilon_{n+j}), \quad i = 0, 1, 2, \dots \tag{3.20}$$

**Lemma 3.3** *The vectors  $\epsilon_m, u_m$ , and  $w_m$  satisfy*

$$\epsilon_{n+i} = \tilde{F}^i \epsilon_n + \check{\epsilon}_{n+i}; \quad \|\check{\epsilon}_{n+i}\| \leq C_i \|\epsilon_n\|^2, \quad C_i = a \frac{1 - L^i}{1 - L} L^{i-1}, \tag{3.21}$$

$$u_{n+i} = (\tilde{F} - I)\tilde{F}^i \epsilon_n + \check{u}_{n+i}; \quad \|\check{u}_{n+i}\| \leq D_i \|\epsilon_n\|^2, \quad D_i = C_i + C_{i+1}, \tag{3.22}$$

$$w_{n+i} = (\tilde{F} - I)^2 \tilde{F}^i \epsilon_n + \check{w}_{n+i}; \quad \|\check{w}_{n+i}\| \leq E_i \|\epsilon_n\|^2, \quad E_i = C_i + 2C_{i+1} + C_{i+2}. \tag{3.23}$$

*Remark* Note that  $C_0 = 0$  and  $C_1 = a$  by (3.21). Therefore,  $D_0 = a$  by (3.22).

*Proof* We start by noting that, by (3.20),

$$\check{\epsilon}_{n+i} = \sum_{j=0}^{i-1} \tilde{F}^{i-j-1} \mu(\epsilon_{n+j}),$$

which, upon taking norms and invoking  $\|\tilde{F}\| \leq L$  and (3.18), gives

$$\begin{aligned} \|\check{\epsilon}_{n+i}\| &\leq \sum_{j=0}^{i-1} \|\tilde{F}^{i-j-1}\| \|\mu(\epsilon_{n+j})\| \leq \sum_{j=0}^{i-1} L^{i-j-1} a (L^j \|\epsilon_n\|)^2 \\ &= a \left( \sum_{j=0}^{i-1} L^{i+j-1} \right) \|\epsilon_n\|^2, \end{aligned}$$

from which (3.21) follows.

The proofs of (3.22)–(3.23) follow from (3.21) and the observation that

$$\check{u}_m = \check{\epsilon}_{m+1} - \check{\epsilon}_m \quad \text{and} \quad \check{w}_m = \check{\epsilon}_{m+2} - 2\check{\epsilon}_{m+1} + \check{\epsilon}_m.$$

We leave the details to the reader. □

Let us now go back to the linear system  $x = \tilde{f}(x)$  in (3.1), recalling that  $F(s) = \tilde{F}$ . As already explained,  $\tilde{f}(x)$  is simply the linear part of the Taylor series of  $f(x)$

about  $s$ , obtained from (3.17) by letting  $\mu(y) \equiv \mathbf{0}$  there. In addition,  $\tilde{f}(s) = s$ , that is,  $s$  is the solution to  $x = \tilde{f}(x)$ , as well as  $x = f(x)$ . Let us now note that  $\mu(y) \equiv \mathbf{0}$  also implies that  $\check{\epsilon}_m = \mathbf{0}$ ,  $\check{u}_m = \mathbf{0}$ , and  $\check{w}_m = \mathbf{0}$  in (3.21), (3.22), and (3.23), respectively. Recalling also that  $\check{\epsilon}_n = \epsilon_n$ , we finally realize that, for  $i = 0, 1, \dots$ ,

$$\begin{aligned} \tilde{\epsilon}_{n+i} &= \tilde{F}^i \tilde{\epsilon}_n = \tilde{F}^i \epsilon_n, & \tilde{u}_{n+i} &= \tilde{F}^i \tilde{u}_n = (\tilde{F} - I) \tilde{F}^i \epsilon_n, \\ \tilde{w}_{n+i} &= \tilde{F}^i \tilde{w}_n = (\tilde{F} - I)^2 \tilde{F}^i \epsilon_n; \end{aligned} \tag{3.24}$$

consequently,

$$u_{n+i} = \tilde{u}_{n+i} + \check{u}_{n+i}, \quad w_{n+i} = \tilde{w}_{n+i} + \check{w}_{n+i}. \tag{3.25}$$

As a result of all this, we have

$$U_{k-1} = \tilde{U}_{k-1} + \check{U}_{k-1}, \quad \check{U}_{k-1} = [\check{u}_n \mid \check{u}_{n+1} \mid \dots \mid \check{u}_{n+k-1}] \tag{3.26}$$

and

$$W_{k-1} = \tilde{W}_{k-1} + \check{W}_{k-1}, \quad \check{W}_{k-1} = [\check{w}_n \mid \check{w}_{n+1} \mid \dots \mid \check{w}_{n+k-1}], \tag{3.27}$$

with  $U_j$  and  $W_j$  as in (2.2). For simplicity of notation, in what follows, we drop the subscript  $k - 1$  from the matrices  $U_{k-1}$ ,  $W_{k-1}$ ,  $\check{U}_{k-1}$ ,  $\check{W}_{k-1}$ , etc. With these, (2.9) becomes

$$\begin{aligned} s_{n,k} &= x_n - U W^+ u_n \\ &= x_n - (\tilde{U} + \check{U})(\tilde{W} + \check{W})^+ (\tilde{u}_n + \check{u}_n). \end{aligned} \tag{3.28}$$

Letting also

$$H = W^+ - \check{W}^+ = (\tilde{W} + \check{W})^+ - \check{W}^+, \tag{3.29}$$

we rewrite (3.28) in the form

$$s_{n,k} = x_n - (\tilde{U} + \check{U})(\tilde{W}^+ + H)(\tilde{u}_n + \check{u}_n). \tag{3.30}$$

Next, opening the parentheses in (3.30), we obtain the equality

$$s_{n,k} = x_n - \tilde{U} \tilde{W}^+ \tilde{u}_n - \tilde{U} \check{W}^+ \check{u}_n - (\tilde{U} H + \check{U} \tilde{W}^+ + \check{U} H)(\tilde{u}_n + \check{u}_n). \tag{3.31}$$

Now, by the fact that  $x_n = \tilde{x}_n$  and by (3.5), we have that  $x_n - \tilde{U} \tilde{W}^+ \tilde{u}_n = \tilde{s}_{n,k}$  in this equality. Next, we invoke  $u_n = \tilde{u}_n + \check{u}_n$  and  $U = \tilde{U} + \check{U}$  again, and obtain a convenient representation of  $s_{n,k}$  and the error in it. We summarize all this in the following lemma.

**Lemma 3.4** *Let*

$$\check{s}_{n,k} = -\tilde{U} \check{W}^+ \check{u}_n - (U H + \check{U} \tilde{W}^+) u_n. \tag{3.32}$$

*Then,  $s_{n,k}$  is given by the equality*

$$s_{n,k} = \tilde{s}_{n,k} + \check{s}_{n,k}. \tag{3.33}$$

*Subtracting  $s$  from both sides of this equality, we also obtain the error formula*

$$s_{n,k} - s = (\tilde{s}_{n,k} - s) + \check{s}_{n,k}. \tag{3.34}$$

## 4 Derivation of upper bounds for $\|s_{n,k} - s\|$

### 4.1 Preliminaries

We now turn to the study of  $s_{n,k} - s$ . Multiplying both sides of (3.34) on the left by  $\mathbf{G}$  and taking  $l_2$  norms, and also invoking  $\|z\|_2 \leq \|\mathbf{G}^{-1}\|_2 \|z\|$ , we obtain

$$\begin{aligned} \|s_{n,k} - s\| &\leq \|\check{s}_{n,k} - s\| + \|\check{s}_{n,k}\|, \\ \frac{\|\check{s}_{n,k}\|}{\|\mathbf{G}^{-1}\|_2} &\leq \|\mathbf{G}\check{\mathbf{U}}\|_2 \|\check{\mathbf{W}}^+\|_2 \|\check{\mathbf{u}}_n\| + \|\mathbf{G}\mathbf{U}\|_2 \|\mathbf{H}\|_2 \|u_n\| + \|\mathbf{G}\check{\mathbf{U}}\|_2 \|\check{\mathbf{W}}^+\|_2 \|u_n\|. \end{aligned} \tag{4.1}$$

Thus, we need to study the behavior of each one of the terms in this bound. We begin with the following lemma.

**Lemma 4.1** *The following are true:*

$$\|\mathbf{G}\mathbf{U}\|_2 \leq K_1 \|\epsilon_n\|, \quad \|\mathbf{G}\check{\mathbf{U}}\|_2 \leq K_2 \|\epsilon_n\|, \quad \|\mathbf{G}\check{\mathbf{U}}\|_2 \leq K_3 \|\epsilon_n\|^2, \tag{4.2}$$

$$\|\mathbf{W}\|_2 \leq K'_1 \|\epsilon_n\|, \quad \|\check{\mathbf{W}}\|_2 \leq K'_2 \|\epsilon_n\|, \quad \|\check{\mathbf{W}}\|_2 \leq K'_3 \|\epsilon_n\|^2, \tag{4.3}$$

with  $K_i, K'_i, i = 1, 2, 3$ , positive constants independent of  $k$  and  $n$ .

*Proof* To achieve the proof, we make use of (3.21)–(3.24) and

$$\begin{aligned} \|u_m\| &\leq (1 + L)\|\epsilon_m\| \quad \text{and} \quad \|w_m\| \leq (1 + L)^2\|\epsilon_m\|, \\ \|\tilde{u}_m\| &\leq (1 + L)\|\epsilon_m\| \quad \text{and} \quad \|\tilde{w}_m\| \leq (1 + L)^2\|\epsilon_m\|. \end{aligned} \tag{4.4}$$

We prove the validity of the bound on  $\|\mathbf{G}\mathbf{U}\|_2$  only; the others can be proved in exactly the same way.

We start by analyzing  $\|\mathbf{G}\mathbf{U}\|_F$ , the Frobenius norm of  $\mathbf{G}\mathbf{U}$ . Noting that

$$\mathbf{G}\mathbf{U} = [\mathbf{G}u_n \mid \mathbf{G}u_{n+1} \mid \cdots \mid \mathbf{G}u_{n+k-1}],$$

we have

$$\begin{aligned} \|\mathbf{G}\mathbf{U}\|_F^2 &= \sum_{j=0}^{k-1} \|\mathbf{G}u_{n+j}\|_2^2 = \sum_{j=0}^{k-1} \|u_{n+j}\|^2 \leq \sum_{j=0}^{k-1} [(1 + L)\|\epsilon_{n+j}\|]^2 \quad \text{by (4.4)} \\ &\leq (1 + L)^2 \sum_{j=0}^{k-1} (L^j \|\epsilon_n\|)^2 \quad \text{by (3.9)} \\ &= \frac{1 + L}{1 - L} (1 - L^{2k}) \|\epsilon_n\|^2 \\ &< \frac{1 + L}{1 - L} \|\epsilon_n\|^2. \end{aligned}$$

The result  $\|\mathbf{G}\mathbf{U}\|_2 \leq K_1 \|\epsilon_n\|$ , with  $K_1 = \sqrt{(1 + L)/(1 - L)}$ , now follows by invoking  $\|\mathbf{G}\mathbf{U}\|_2 \leq \|\mathbf{G}\mathbf{U}\|_F$ .<sup>9</sup> □

<sup>9</sup>Recall that, for any matrix  $\mathbf{K}$  with  $\text{rank}(\mathbf{K}) = r$ , we have  $\|\mathbf{K}\|_2 \leq \|\mathbf{K}\|_F \leq r\|\mathbf{K}\|_2$ . See Golub and Van Loan [13].

### 4.2 Upper bounds for $\|\tilde{W}^+\|_2$ and $\|H\|_2$

Next, by Theorem A.3 in the [Appendix](#), we can bound  $\|H\|_2$  as in

$$\|H\|_2 \leq \sqrt{2} \frac{\Delta}{1 - \Delta} \|\tilde{W}^+\|_2 \quad \text{provided } \Delta = \|\tilde{W}^+\|_2 \|\check{W}\|_2 < 1. \tag{4.5}$$

We realize that all we need is a suitable upper bound on  $\|\tilde{W}^+\|_2$  since we already have an upper bound on  $\|\check{W}\|_2$  from (4.13). We turn to this issue next.

Now, by (3.4) and (3.24), we have

$$\tilde{W} = [(\tilde{F} - I)^2 \epsilon_n \mid (\tilde{F} - I)^2 \tilde{F} \epsilon_n \mid \dots \mid (\tilde{F} - I)^2 \tilde{F}^{k-1} \epsilon_n], \tag{4.6}$$

which can be written in the form

$$\tilde{W} = \|\epsilon_n\|_2 \overset{\circ}{W}, \quad \overset{\circ}{W} = RS(e_n), \tag{4.7}$$

where

$$R = (\tilde{F} - I)^2 = G^2, \quad S(y) = [y \mid \tilde{F}y \mid \dots \mid \tilde{F}^{k-1}y], \quad e_n = \frac{\epsilon_n}{\|\epsilon_n\|_2}. \tag{4.8}$$

[Note that the columns of  $S(y)$  span the Krylov subspace  $\mathcal{K}_k(\tilde{F}; y) = \text{span}\{y, \tilde{F}y, \dots, \tilde{F}^{k-1}y\}$ .] First,  $R$  is  $N \times N$ , constant, and nonsingular since  $G$  is. Next, we recall that  $k$  is at most the degree of the minimal polynomial of  $\tilde{F}$  with respect to the vector  $\epsilon_n$ , which implies that the vectors  $\tilde{F}^i \epsilon_n, i = 0, 1, \dots, k - 1$ , are linearly independent and, therefore,  $\text{rank}(S(e_n)) = k$ . As a result,  $\text{rank}(\overset{\circ}{W}) = k = \text{rank}(\tilde{W})$  since  $R$  is nonsingular. By the fact that  $(aK)^+ = a^{-1}K^+$  for every nonzero scalar  $a \in \mathbb{C}$ , and by Theorem A.1 in the [Appendix](#), we thus have

$$\tilde{W}^+ = \frac{1}{\|\epsilon_n\|_2} \overset{\circ}{W}^+ \Rightarrow \|\tilde{W}^+\|_2 = \frac{1}{\|\epsilon_n\|_2} \|\overset{\circ}{W}^+\|_2 \tag{4.9}$$

and

$$\|\overset{\circ}{W}^+\|_2 \leq \|R^{-1}\|_2 \|S(e_n)^+\|_2. \tag{4.10}$$

We need to bound only  $\|S(e_n)^+\|_2$  *uniformly* (i) for all  $n = 1, 2, \dots$  in the  $n$ -Mode, and (ii) for all unit vectors  $e_n^{(r)} = \epsilon_n^{(r)} / \|\epsilon_n^{(r)}\|_2$  arising in the different cycles of the C-Mode and the MC-Mode. Unfortunately, we are not able to prove the existence of such uniform bounds. In what follows, concerning the application of RRE in the  $n$ -Mode and in two cycling modes, we *assume* that, at each step of the different modes of usage of RRE,  $\|S(e_n)^+\|_2$  is bounded uniformly throughout, that is, we assume that, for some constant  $\tilde{\eta} > 0$ ,

$$\|S(e_n)^+\|_2 \leq \tilde{\eta}. \tag{4.11}$$

Combining (4.9)–(4.11) and invoking also  $\|\epsilon_n\| \leq \|G\|_2 \|\epsilon_n\|_2$ , we obtain

$$\boxed{\|\tilde{W}^+\|_2 \leq \frac{\eta}{\|\epsilon_n\|}, \quad \eta = \tilde{\eta} \|G\|_2 \|R^{-1}\|_2.} \tag{4.12}$$

We shall comment on this assumption concerning the uniform upper bound for  $\|S(e_n)^+\|_2$  in Section 6.

The first thing to do now is to guarantee that  $\Delta = \|\tilde{\mathbf{W}}^+\|_2 \|\check{\mathbf{W}}\|_2 < 1$  in (4.5) is satisfied under the assumption in (4.12) concerning  $\|\tilde{\mathbf{W}}^+\|_2$ . By (4.12) and (4.3) and the fact that  $\|\epsilon_0\| \leq \delta$  since  $\mathbf{x}_0 \in B(\mathbf{s}, \delta)$ , we have

$$\Delta \leq K'_3 \eta \|\epsilon_n\| \leq K'_3 \eta L^n \|\epsilon_0\| \leq K'_3 \eta L^n \delta. \tag{4.13}$$

Clearly, by making  $\delta$  sufficiently small, we can make the upper bound on  $\Delta$  smaller than one. The closer  $\delta$  is to zero, the closer  $\mathbf{x}_0$  is to  $\mathbf{s}$ . This is precisely what is needed in order to develop a *local* convergence theory for any extrapolation method.

Next, by (4.5), (4.12), and (4.13),

$$\|\mathbf{H}\|_2 \leq \lambda_n, \quad \lambda_n = \sqrt{2} \frac{K'_3 \eta^2}{1 - K'_3 \eta \|\epsilon_n\|}. \tag{4.14}$$

As we will show later,  $\epsilon_n$  is bounded in all three modes (*n*-Mode, C-Mode, and MC-Mode) we study here, which implies that  $\lambda_n$  is bounded too.

*Remark* Before proceeding further, we would like to discuss an interesting consequence of the global assumption we have made concerning  $\tilde{\mathbf{W}}^+$ . By (4.12) and (4.14) and also by (3.29), namely, that  $\mathbf{W}^+ = \tilde{\mathbf{W}}^+ + \mathbf{H}$ , we have

$$\|\mathbf{W}^+\|_2 \leq \|\tilde{\mathbf{W}}^+\|_2 + \|\mathbf{H}\|_2 \leq \frac{\eta}{\|\epsilon_n\|} + \lambda_n.$$

As a result, the vector  $\xi = -\mathbf{W}^+ \mathbf{u}_n$  defined via (2.8), satisfies

$$\|\xi\|_2 \leq \|\mathbf{W}^+\|_2 \|\mathbf{u}_n\|_2 \leq (1 + L)(\eta + \lambda_n \|\mathbf{G}^{-1}\|_2 \|\epsilon_n\|).$$

Here we have made use of (4.4) too. Since  $\|\epsilon_n\|$  and  $\lambda_n$  are bounded, so is  $\lambda_n \|\epsilon_n\|$ , in all three modes. This implies that  $\xi$  is bounded, which causes  $\boldsymbol{\gamma}$  in (2.3)–(2.5) to be bounded as well. This can be seen by expressing the  $\gamma_i$  in terms of the  $\xi_i$  by employing (2.7) as in

$$\gamma_0 = 1 - \xi_0; \quad \gamma_i = \xi_{i-1} - \xi_i, \quad i = 1, \dots, k - 1; \quad \gamma_k = \xi_{k-1}.$$

Thus, we have globally

$$\sum_{i=0}^k |\gamma_i| \leq \Gamma \text{ for some } \Gamma > 0 \text{ throughout all three modes.}$$

Interestingly, this is analogous to the global assumption made by Toth and Kelly [41] in the convergence analysis of the acceleration method of Anderson [1]. Note that, when applied to linear systems, Anderson acceleration is equivalent to GMRES (see Walker and Ni [43]), which is equivalent to RRE applied to linear systems (see Sidi [26]).

### 4.3 Upper bound for $\|s_{n,k} - s\|$

With the different matrices in (4.1) bounded as above, we turn to  $s_{n,k} - s$ . By (4.2), (4.3), (4.4), and (4.12), we have

$$\|G\tilde{U}\|_2 \|\tilde{W}^+\|_2 \|\tilde{u}_n\| \leq K_2\eta D_0 \|\epsilon_n\|^2, \tag{4.15}$$

$$\|GU\|_2 \|H\|_2 \|u_n\| \leq K_1\lambda_n(1 + L) \|\epsilon_n\|^2, \tag{4.16}$$

$$\|G\check{U}\|_2 \|\check{W}^+\|_2 \|u_n\| \leq K_3\eta(1 + L) \|\epsilon_n\|^2. \tag{4.17}$$

Substituting these into (4.1), we obtain

$$\|\check{s}_{n,k}\| \leq \tau_n \|\epsilon_n\|^2, \quad \tau_n = [K_2\eta D_0 + (K_1\lambda_n + K_3\eta)(1 + L)] \|G^{-1}\|_2, \tag{4.18}$$

and this leads to the bound on  $\|s_{n,k} - s\|$  in the next lemma:

**Lemma 4.2** *The norm of the error vector  $s_{n,k} - s$  can be bounded as in*

$$\|s_{n,k} - s\| \leq \|\check{s}_{n,k} - s\| + \tau_n \|\epsilon_n\|^2, \quad \tau_n = [K_2\eta D_0 + (K_1\lambda_n + K_3\eta)(1 + L)] \|G^{-1}\|_2. \tag{4.19}$$

*Remark* Note that, by (4.14) and (4.19),  $\lim_{n \rightarrow \infty} \tau_n$  is finite since  $\lim_{n \rightarrow \infty} \lambda_n$  is finite. Therefore,  $\|s_{n,k} - s\|$  cannot be smaller than  $\|\check{s}_{n,k}\| \leq \tau_n \|\epsilon_n\|^2$ , even though  $\|\check{s}_{n,k} - s\|$  may be smaller. In other words, the term  $\|\check{s}_{n,k}\|$  limits the accuracy of  $s_{n,k}$  as an approximation to  $s$ .

## 5 Convergence analysis

### 5.1 Preliminaries

We start by studying the term  $\|\check{s}_{n,k} - s\|$ . We recall that  $\check{s}_{n,k}$  is the vector obtained by applying RRE to the vectors  $\check{x}_m, m = n, n + 1, \dots, n + k$ , with  $\check{x}_n = x_n$ , as described in Section 3.1. Our study will be based on the developments of [26, 36], and [33, Chapters 6,7].

We first have

$$\begin{aligned} \|\check{s}_{n,k} - s\| &= \|G(\check{s}_{n,k} - s)\|_2 = \min_{g \in \tilde{\mathcal{P}}_k} \|g(\tilde{F})G(\check{x}_n - s)\|_2 \quad \text{by Theorem 3.1} \\ &= \min_{g \in \tilde{\mathcal{P}}_k} \|g(\tilde{F})G(x_n - s)\|_2 \quad \text{because } \check{x}_n = x_n \\ &\leq \left[ \min_{g \in \tilde{\mathcal{P}}_k} \|g(\tilde{F})\|_2 \right] \|G(x_n - s)\|_2 \\ &= \theta_k \|\epsilon_n\|, \end{aligned} \tag{5.1}$$

recalling that  $x_n - s = \epsilon_n$  and defining

$$\theta_k = \min_{g \in \tilde{\mathcal{P}}_k} \|g(\tilde{F})\|_2. \tag{5.2}$$



(Note that  $\theta_k$  depends only on  $\tilde{F}$  and  $k$ .) Of course, we also have

$$\theta_k \leq \|g(\tilde{F})\|_2 \quad \forall g \in \tilde{\mathcal{P}}_k. \tag{5.3}$$

We now would like to bound  $\theta_k$  appropriately. Choosing  $g(z) = z^k$  in (5.3), we obtain,

$$\theta_k = \min_{g \in \tilde{\mathcal{P}}_k} \|g(\tilde{F})\|_2 \leq \|\tilde{F}^k\|_2 \leq \|\tilde{F}\|_2^k \leq L^k, \quad \text{at worst.} \tag{5.4}^{10}$$

With all these developments, (5.1) and (4.18) together give the result in the next lemma:

**Lemma 5.1** *The error vector  $s_{n,k} - s$  satisfies*

$$\|s_{n,k} - s\| \leq \theta_k \|\epsilon_n\| + \tau_n \|\epsilon_n\|^2. \tag{5.5}$$

*Remark* By choosing  $g(z) \in \tilde{\mathcal{P}}_k$  suitably, upper bounds on  $\theta_k$  that are smaller than  $L^k$  can be given for different cases. We give such bounds for two such cases here. For additional cases involving orthogonal polynomials, such as Jacobi polynomials, we refer the reader to Sidi and Shapira [36].

- If the hermitian part of  $G = I - \tilde{F}$ , namely, the matrix  $G_H = \frac{1}{2}(G + G^*)$ , is positive definite, then

$$\theta_k \leq (1 - \nu^2/\sigma^2)^{k/2},$$

where  $\sigma$  is the largest singular value of  $G$  and  $\nu$  is the smallest eigenvalue of  $G_H$ . Of course,  $0 < \nu < \sigma$ . See [26].

- If  $\tilde{F}$  is hermitian with eigenvalues in the (real) interval  $[\alpha, \beta]$ ,  $-1 < \alpha < \beta < 1$ , then

$$\theta_k \leq \frac{1}{T_k\left(\frac{2-\alpha-\beta}{\beta-\alpha}\right)} < 2 \left(\frac{\sqrt{k}-1}{\sqrt{k}+1}\right)^k, \quad \kappa = \frac{1-\alpha}{1-\beta} > 1.$$

Here  $T_k(z)$  is the Chebyshev polynomial of the first kind of degree  $k$ . (See Varga [42, Chapter 5], for example.) Note that, in this case,  $\theta_k < L^k$ , with  $L = \max(|\alpha|, |\beta|) < 1$ .

### 5.1.1 Main assumptions

Before delving into the local convergence analyses of the different modes of usage of RRE, we would like to summarize the assumptions we have made so far. We will be referring to them in the statements of our (local) convergence theorems below.

- A1.  $f \in C^2(B(s, \delta))$  for some  $\delta > 0$ . (We can assume  $\delta$  to be as small as needed in our proofs.)
- A2.  $\|\tilde{F}\| \leq \max_{x \in B(s, \delta)} \|F(x)\| = L < 1$ , which also implies that  $\rho(\tilde{F}) \leq L$ .

---

<sup>10</sup>Clearly,  $g(z) = z^k$  is in  $\tilde{\mathcal{P}}_k$  and  $\theta_k < 1$  since  $L < 1$ . Next, in general, the polynomial  $g(z)$  that gives the optimum in (5.4) is different from  $z^k$ . Thus, generally speaking,  $\theta_k < L^k$

- A3. The very first vector  $\mathbf{x}_0$ , with which we start any of the modes, is in  $B(\mathbf{s}, \delta)$ . Thus,  $\|\mathbf{x}_0 - \mathbf{s}\| < \delta$ .
- A4.  $\|\tilde{\mathbf{W}}^+\| \leq \eta/\|\epsilon_n\|$  for every  $n$  in the  $n$ -Mode and for every cycle in the C-Mode and the MC-Mode. ( $\eta > 0$  is fixed.)

### 5.2 Convergence in $n$ -Mode

We recall that, in the  $n$ -Mode, we are applying RRE, with  $k \geq 1$  fixed throughout, to the infinite sequence  $\{\mathbf{x}_m\}$  that is generated as in (1.4). (That is, no cycling is involved.)

**Theorem 5.2** *Under the assumptions A1–A4, RRE converges in the  $n$ -Mode. Actually, we have*

$$\limsup_{n \rightarrow \infty} \frac{\|\mathbf{s}_{n,k} - \mathbf{s}\|}{\|\epsilon_n\|} \leq \theta_k < 1. \tag{5.6}$$

*Proof* Since  $\lim_{n \rightarrow \infty} \|\epsilon_n\| = 0$  and  $\lim_{n \rightarrow \infty} \tau_n < \infty$ , it is clear from Lemma 5.1 that  $\lim_{n \rightarrow \infty} \|\mathbf{s}_{n,k} - \mathbf{s}\| = 0$ , hence  $\lim_{n \rightarrow \infty} \mathbf{s}_{n,k} = \mathbf{s}$ .

Next, again by Lemma 5.1, we have

$$\frac{\|\mathbf{s}_{n,k} - \mathbf{s}\|}{\|\epsilon_n\|} \leq \theta_k + \tau_n \|\epsilon_n\|.$$

Taking the limsup as  $n \rightarrow \infty$  of both sides and recalling again that  $\lim_{n \rightarrow \infty} \tau_n < \infty$  and  $\lim_{n \rightarrow \infty} \|\epsilon_n\| = 0$ , the result in (5.6) follows. □

*Remark* Let us also rewrite (4.19) as

$$\|\mathbf{s}_{n,k} - \mathbf{s}\| = O(\psi_n) \quad \text{as } n \rightarrow \infty; \quad \psi_n = \max\{\|\tilde{\mathbf{s}}_{n,k} - \mathbf{s}\|, \|\epsilon_n\|^2\}. \tag{5.7}$$

This is possible since  $\lim_{n \rightarrow \infty} \tau_n$  is finite. It is thus clear that  $\|\mathbf{s}_{n,k} - \mathbf{s}\|$  cannot be less than  $O(\|\epsilon_n\|^2) = O(L^{2n}) \approx O(\rho(\tilde{\mathbf{F}})^{2n})$  as  $n \rightarrow \infty$ , no matter what  $\|\tilde{\mathbf{s}}_{n,k} - \mathbf{s}\|$  is. [See the remark following (4.19).]

### 5.3 Convergence in C-Mode cycling

In C-Mode cycling, we keep  $n \geq 0$  and  $k \geq 1$  fixed throughout,  $k$  always being assumed to be less than the degree of the minimal polynomial of  $\tilde{\mathbf{F}}$  with respect to the vector  $\epsilon_n$  in every cycle.

**Theorem 5.3** *Under the assumptions A1–A4, RRE converges linearly in the C-Mode. Actually, we have*

$$\limsup_{r \rightarrow \infty} \frac{\|\mathbf{s}_{n,k}^{(r+1)} - \mathbf{s}\|}{\|\mathbf{s}_{n,k}^{(r)} - \mathbf{s}\|} \leq \theta_k L^n < 1. \tag{5.8}$$

*Proof* We start by observing that, by Lemma 3.2, there holds  $\|\epsilon_n\| \leq L^n \|\epsilon_0\|$ . With this, (5.5) becomes

$$\|s_{n,k} - s\| \leq \left( \theta_k L^n + \tau_n L^{2n} \|\epsilon_0\| \right) \|\epsilon_0\|. \tag{5.9}$$

Let us now denote the vectors  $x_m$ ,  $\epsilon_m = x_m - s$ , and  $s_{n,k}$  used/computed in cycle  $i$  by  $x_m^{(i)}$ ,  $\epsilon_m^{(i)}$ , and  $s_{n,k}^{(i)}$ , respectively, and rewrite (5.9) that is relevant to the cycle  $(r + 1)$  as

$$\|s_{n,k}^{(r+1)} - s\| \leq \left( \theta_k L^n + \tau_n L^{2n} \|\epsilon_0^{(r+1)}\| \right) \|\epsilon_0^{(r+1)}\|. \tag{5.10}$$

Let us also recall that, in the C-Mode,  $x_0^{(r+1)} = s_{n,k}^{(r)}$ , and hence  $\epsilon_0^{(r+1)} = s_{n,k}^{(r)} - s$ . With these, (5.10) becomes

$$\|s_{n,k}^{(r+1)} - s\| \leq \mu_r \|s_{n,k}^{(r)} - s\|, \quad \mu_r = \theta_k L^n + \tau_n L^{2n} \|s_{n,k}^{(r)} - s\|. \tag{5.11}$$

We now show by induction that for each  $r$ ,  $s_{n,k}^{(r)}$  is in the ball  $B(s, \delta)$  and tends to  $s$  as  $r \rightarrow \infty$ , provided  $x_0$  in step C0 of C-Mode cycling is sufficiently close to  $s$ .

For  $r = 0$ , we have  $x_0 = s_{n,k}^{(0)} \in B(s, \delta)$  by choice; therefore,

$$\mu_0 \leq \theta_k L^n + \tau_n L^{2n} \delta.$$

Since  $\theta_k L^n < 1$ , we can force  $\mu_0 < 1$  by choosing  $\delta$  sufficiently small or by choosing  $s_{n,k}^{(0)}$  sufficiently close to  $s$ . This, in turn, forces

$$\|s_{n,k}^{(1)} - s\| \leq \mu_0 \|s_{n,k}^{(0)} - s\| \leq \mu_0 \delta < \delta \quad \Rightarrow \quad s_{n,k}^{(1)} \in B(s, \delta).$$

Continuing by induction on  $r$ , we see that  $\mu_{r+1} < \mu_r$ ,  $\|s_{n,k}^{(r+1)} - s\| < \|s_{n,k}^{(r)} - s\|$ , hence  $s_{n,k}^{(r+1)} \in B(s, \delta)$  since  $s_{n,k}^{(r)} \in B(s, \delta)$ . We also have  $\lim_{r \rightarrow \infty} \|s_{n,k}^{(r)} - s\| = 0$ , hence  $\lim_{r \rightarrow \infty} s_{n,k}^{(r)} = s$ . With the convergence of  $\{s_{n,k}^{(r)}\}_{r=0}^\infty$  to  $s$  established, let us now rewrite (5.11) as

$$\frac{\|s_{n,k}^{(r+1)} - s\|}{\|s_{n,k}^{(r)} - s\|} \leq \mu_r. \tag{5.12}$$

Taking the limsup as  $r \rightarrow \infty$  on both sides of this inequality, we obtain (5.8). □

### 5.4 Convergence in MC-Mode cycling

We recall that in MC-Mode cycling, we keep  $n$  fixed while  $k = k_r$  is the degree of the minimal polynomial of  $\tilde{F}$  with respect to  $\epsilon_n$  in the  $r$ th cycle.

**Theorem 5.4** *Under the assumptions A1–A4, RRE converges quadratically in the MC-Mode. Actually, we have*

$$\limsup_{r \rightarrow \infty} \frac{\|s_{n,k_{r+1}}^{(r+1)} - s\|}{\|s_{n,k_r}^{(r)} - s\|^2} \leq \tau_n L^{2n}. \tag{5.13}$$

*Proof* We start by noting that  $\tilde{s}_{n,k_r} = s$  in each cycle, as mentioned at the end of Section 3.1. Thus, (4.19) becomes

$$\|s_{n,k} - s\| \leq \tau_n \|\epsilon_n\|^2. \tag{5.14}$$

Proceeding precisely as in the proof of Theorem 5.3 concerning the C-Mode cycling, we next obtain

$$\|s_{n,k} - s\| \leq \tau_n L^{2n} \|\epsilon_0\|^2. \tag{5.15}$$

As in the case of the C-Mode, noting that  $\epsilon_0^{(r+1)} = s_{n,k_r}^{(r)} - s$ , we write (5.15) in the MC-Mode as

$$\|s_{n,k_{r+1}}^{(r+1)} - s\| \leq \tau_n L^{2n} \|s_{n,k_r}^{(r)} - s\|^2 = \phi_r \|s_{n,k_r}^{(r)} - s\|, \quad \phi_r = \tau_n L^{2n} \|s_{n,k_r}^{(r)} - s\|. \tag{5.16}$$

We now show, by induction on  $r$ , that  $s_{n,k_r}^{(r)}$  is in the ball  $B(s, \delta)$  and tends to  $s$  as  $r \rightarrow \infty$ , provided  $x_0$  in step MC0 of MC-Mode is sufficiently close to  $s$ .

For  $r = 0$ , we have  $x_0 = s_{n,k_0}^{(0)} \in B(s, \delta)$  by choice; therefore,

$$\phi_0 \leq \tau_n L^{2n} \delta \quad \Rightarrow \quad \phi_0 < 1 \text{ provided } \delta \text{ sufficiently small.}$$

This implies that  $\|s_{n,k_1}^{(1)} - s\| < \|s_{n,k_0}^{(0)} - s\|$ ; therefore,  $s_{n,k_1}^{(1)} \in B(s, \delta)$ . In addition, we also have  $\phi_1 < \phi_0$ . Continuing by induction on  $r$ , we see that  $\phi_r < \phi_{r-1} < 1$  hence  $\|s_{n,k_{r+1}}^{(r+1)} - s\| < \|s_{n,k_r}^{(r)} - s\|$ , which implies that  $s_{n,k_{r+1}}^{(r+1)} \in B(s, \delta)$  since  $s_{n,k_r}^{(r)} \in B(s, \delta)$ , and that  $\lim_{r \rightarrow \infty} \|s_{n,k_r}^{(r)} - s\| = 0$ , meaning that  $\lim_{r \rightarrow \infty} s_{n,k_r}^{(r)} = s$ . With the convergence of  $\{s_{n,k}^{(r)}\}_{r=0}^\infty$  to  $s$  established, let us now rewrite (5.16) as

$$\frac{\|s_{n,k_{r+1}}^{(r+1)} - s\|}{\|s_{n,k_r}^{(r)} - s\|^2} \leq \tau_n L^{2n}. \tag{5.17}$$

Taking the limsup as  $r \rightarrow \infty$  on both sides of this inequality, we obtain (5.13). Thus, the convergence of the sequence  $\{s_{n,k}^{(r)}\}_{r=0}^\infty$  is quadratic. □

### 6 Remarks on $\|S(e_n)^+\|_2$

Let us observe that  $S(y)$  can be written as the product of two matrices as

$$S(y) = P Q(y), \tag{6.1}$$

where  $P \in \mathbb{C}^{N \times kN}$  and  $Q(y) \in \mathbb{C}^{kN \times k}$  are given as

$$P = [I \mid \tilde{F} \mid \dots \mid \tilde{F}^{k-1}]; \quad Q(y) = \begin{bmatrix} y & \mathbf{0} & \dots & \mathbf{0} \\ \mathbf{0} & y & \dots & \mathbf{0} \\ \vdots & & \ddots & \vdots \\ \mathbf{0} & \mathbf{0} & \dots & y \end{bmatrix}, \quad y \in \mathbb{C}^k. \tag{6.2}$$

Clearly,  $P$  is a constant matrix and has full row rank, while  $Q(y)$  has full column rank for all nonzero  $y$ , that is,

$$\text{rank}(P) = N, \quad \text{rank}(Q(y)) = k \quad \forall y \neq \mathbf{0}. \tag{6.3}$$

Before going on, we recall that if  $\mathbf{K} \in \mathbb{C}^{m \times k}$  with  $\text{rank}(\mathbf{K}) = k$ , then it has  $k$  nonzero singular values, which we order such that

$$\sigma_1(\mathbf{K}) \geq \sigma_2(\mathbf{K}) \geq \dots \geq \sigma_k(\mathbf{K}) > 0,$$

and

$$\sigma_k(\mathbf{K}) = \min_{\mathbf{z} \in \mathbb{C}^k, \|\mathbf{z}\|_2=1} \|\mathbf{K}\mathbf{z}\|_2 \quad \text{and} \quad \|\mathbf{K}^+\|_2 = 1/\sigma_k(\mathbf{K}).$$

Now,  $\mathbf{P}$  has  $N$  positive singular values, and therefore

$$\|\mathbf{P}^+\|_2 = 1/\sigma_N(\mathbf{P}).$$

Next,  $\mathbf{Q}(\mathbf{y})$  is unitary when  $\|\mathbf{y}\| = 1$ , in the sense that

$$\mathbf{Q}(\mathbf{y})^* \mathbf{Q}(\mathbf{y}) = \mathbf{I}_{k \times k} \quad \forall \mathbf{y} \in \mathbb{C}^k, \|\mathbf{y}\|_2 = 1, \tag{6.4}$$

hence so is  $\mathbf{Q}(\mathbf{e}_n)$  since  $\|\mathbf{e}_n\|_2 = 1$ . As a result  $\mathbf{Q}(\mathbf{y})^+ = \mathbf{Q}(\mathbf{y})^*$  and  $\mathbf{Q}(\mathbf{y})$  has  $k$  singular values, all equal to one, for all  $\mathbf{y}, \|\mathbf{y}\|_2 = 1$ . Consequently,

$$\|\mathbf{Q}(\mathbf{y})^+\|_2 = 1 \quad \forall \mathbf{y} \in \mathbb{C}^k, \|\mathbf{y}\|_2 = 1. \tag{6.5}$$

Despite these interesting facts—that  $\|\mathbf{P}^+\|_2$  is fixed and that  $\|\mathbf{Q}(\mathbf{e}_n^{(r)})^+\|_2 = 1$  throughout the cycling process—we are not able to prove that  $\|\mathbf{S}(\mathbf{e}_n^{(r)})^+\|_2 = \|[\mathbf{P} \mathbf{Q}(\mathbf{e}_n^{(r)})]^+\|_2 \leq \alpha$  for some fixed  $\alpha > 0$ , for all  $r = 0, 1, \dots$ , where  $\mathbf{e}_n^{(r)} = \boldsymbol{\epsilon}_n^{(r)} / \|\boldsymbol{\epsilon}_n^{(r)}\|_2$  in the  $r$ th cycle.

For example, (A.3) in the appendix, which would be extremely useful if applicable, does not apply to  $\mathbf{S}(\mathbf{y})$ . If it did, then we would have  $\mathbf{S}(\mathbf{y})^+ = \mathbf{Q}(\mathbf{y})^* \mathbf{P}^+$  hence  $\|\mathbf{S}(\mathbf{y})^+\|_2 \leq \|\mathbf{P}^+\|_2$ , very conveniently.

We might think that Theorem A.4 in the Appendix would apply to the  $n$ -Mode and C-Mode (it does not necessarily apply to the MC-Mode since the  $\text{rank}(\mathbf{S}(\mathbf{e}_n^{(r)})) = k_r$  may vary with  $r$ ), but this too is problematic. Theorem A.4 requires the following:

- In the  $n$ -Mode, the sequence  $\{\mathbf{e}_n\}_{n=0}^\infty$ , where  $\mathbf{e}_n = \boldsymbol{\epsilon}_n / \|\boldsymbol{\epsilon}_n\|_2$ , must have a limit  $\mathbf{e}_\infty$  such that  $\text{rank}(\mathbf{S}(\mathbf{e}_\infty)) = k$ . It is obvious from (3.20)–(3.21) that it is very difficult to determine whether such a vector  $\mathbf{e}_\infty$  exists when  $\mathbf{f}(\mathbf{x})$  is nonlinear.<sup>11</sup>
- In the C-Mode, the sequence  $\{\mathbf{e}_n^{(r)}\}_{r=0}^\infty$ , where  $\mathbf{e}_n^{(r)} = \boldsymbol{\epsilon}_n^{(r)} / \|\boldsymbol{\epsilon}_n^{(r)}\|_2$ , must have a limit  $\mathbf{e}_n^{(\infty)}$  such that  $\text{rank}(\mathbf{S}(\mathbf{e}_n^{(\infty)})) = k$ . It is obvious again from (3.20)–(3.21) that it is very difficult to ascertain whether such a limit exists when  $\mathbf{f}(\mathbf{x})$  is nonlinear.

A different approach to the issue, for the C-Mode, would be as follows: Since  $\mathbf{S}(\mathbf{e}_n^{(r)})$  has full column rank,  $\|\mathbf{S}(\mathbf{e}_n^{(r)})^+\|_2 = 1/\sigma_k(\mathbf{S}(\mathbf{e}_n^{(r)})) > 0$  for every  $r = 1, 2, \dots$ . Defining the vector  $\boldsymbol{\zeta}(\mathbf{y}) \in \mathbb{C}^k, \|\boldsymbol{\zeta}(\mathbf{y})\|_2 = 1$ , via

$$\min_{\mathbf{z} \in \mathbb{C}^k, \|\mathbf{z}\|_2=1} \|\mathbf{S}(\mathbf{y})\mathbf{z}\|_2 = \|\mathbf{S}(\mathbf{y})\boldsymbol{\zeta}(\mathbf{y})\|_2, \tag{6.6}$$

<sup>11</sup>For the linear system  $\mathbf{x} = \tilde{\mathbf{f}}(\mathbf{x})$ , we have  $\boldsymbol{\epsilon}_{n+1} = \tilde{\mathbf{F}}\boldsymbol{\epsilon}_n, n = 0, 1, \dots$ , as power iterations. Thus, in some cases,  $\mathbf{e}_\infty = \lim_{n \rightarrow \infty} \mathbf{e}_n$  exists and is an eigenvector of  $\tilde{\mathbf{F}}$ , hence causes  $\text{rank}(\mathbf{S}(\mathbf{e}_\infty)) = 1$  at most. Clearly, this is a problem when  $\text{rank}(\mathbf{S}(\mathbf{e}_n)) = k > 1$ , for  $n = 0, 1, \dots$

we thus have

$$\sigma_k(\mathbf{S}(\mathbf{e}_n^{(r)})) = \min_{z \in \mathbb{C}^k, \|z\|_2=1} \|\mathbf{S}(\mathbf{e}_n^{(r)})z\|_2 = \|\mathbf{S}(\mathbf{e}_n^{(r)})\xi(\mathbf{e}_n^{(r)})\|_2 > 0 \quad \forall r = 1, 2, \dots, \tag{6.7}$$

from which, we obtain

$$\sigma_k(\mathbf{S}(\mathbf{e}_n^{(r)})) \geq \liminf_{r \rightarrow \infty} \|\mathbf{S}(\mathbf{e}_n^{(r)})\xi(\mathbf{e}_n^{(r)})\|_2 = \alpha \geq 0. \tag{6.8}$$

Clearly,  $\alpha$  is independent of  $r$ . Now, if we can show that  $\alpha > 0$ , we will have shown that  $\|\mathbf{S}(\mathbf{e}_n)^+\|_2 \leq 1/\alpha$ , hence that  $\|\mathbf{S}(\mathbf{e}_n)^+\|_2$  is bounded uniformly throughout the cycling process. Unfortunately, this does not seem to be the case in general; the best we can say is that  $\alpha \geq 0$ .

Thus, even though  $\sigma_k(\mathbf{S}(\mathbf{e}_n^{(r)})) > 0$  for  $r = 0, 1, \dots$ , it seems we cannot guarantee the existence of a fixed positive constant  $\tilde{\alpha}$  such that, when applying RRE in the cycling mode,  $\sigma_k(\mathbf{S}(\mathbf{e}_n^{(r)})) \geq \tilde{\alpha}$  uniformly in every cycle. Therefore, we can only *assume* that such a constant exists for the C-Mode cycling process being studied, for which  $k$  is fixed throughout, namely,

$$\boxed{\|\mathbf{S}(\mathbf{e}_n^{(r)})^+\|_2 \leq 1/\tilde{\alpha} < \infty \quad \forall r, \text{rank}(\mathbf{S}(\mathbf{e}_n^{(r)})) = k \leq k_r, \quad r = 0, 1, \dots,} \tag{6.9}$$

where  $k_r$  is the degree of the minimal polynomial of  $\tilde{\mathbf{F}}$  with respect to  $\mathbf{e}_n^{(r)}$ .

As for the MC-Mode cycling process, we can, similarly, only *assume* that

$$\boxed{\|\mathbf{S}(\mathbf{e}_n^{(r)})^+\|_2 \leq 1/\tilde{\alpha} < \infty \quad \forall r, \text{rank}(\mathbf{S}(\mathbf{e}_n^{(r)})) = k_r, \quad r = 0, 1, \dots} \tag{6.10}$$

(This is reasonable because there are only finitely many  $k_r$  as  $1 \leq k_r \leq N$ .) Precisely (6.9) and (6.10) are what we have assumed in (4.11).

Finally, we note that the global condition in (4.12) we have imposed on the three modes for RRE discussed in this work is formulated in terms of  $\tilde{\mathbf{F}}$ , the Jacobian matrix of  $\mathbf{f}(x)$  at the solution  $\mathbf{s}$  only, and it concerns  $s_{n,k}$  with arbitrary  $n$ . This should be contrasted with the global condition introduced in [15] for the MC-Mode only that is formulated in terms of  $\mathbf{f}(x)$ , and concerns  $s_{0,k}$ . Denoting the  $\mathbf{x}_i$  and the  $\mathbf{u}_i = \mathbf{x}_{i+1} - \mathbf{x}_i$  generated at the  $r$ th cycle by  $\mathbf{x}_i^{(r)}$  and  $\mathbf{u}_i^{(r)}$ , respectively, with  $\mathbf{x}_0^{(r)} = \mathbf{s}_{0,k_{r-1}}^{(r-1)}$ , the condition of [15] reads as follows:

$$\sqrt{\det(\mathbf{Y}_r^* \mathbf{Y}_r)} \geq \alpha > 0 \quad \forall r; \quad \mathbf{Y}_r = [\hat{\mathbf{u}}_0^{(r)} \mid \hat{\mathbf{u}}_1^{(r)} \mid \dots \mid \hat{\mathbf{u}}_{k_r-1}^{(r)}], \quad \hat{\mathbf{u}}_i^{(r)} = \mathbf{u}_i^{(r)} / \|\mathbf{u}_i^{(r)}\|_2.$$

**Acknowledgements** The author would like to thank one of the anonymous referees for his/her remarks that helped to improve the presentation and results of this work substantially.

### Appendix : Some properties of Moore–Penrose inverses

First , we recall the well-known facts

$$A \in \mathbb{C}^{m \times n}, \quad \text{rank}(A) = n \quad \Rightarrow \quad A^+ = (A^*A)^{-1}A^* \quad \Rightarrow \quad A^+A = I_{n \times n}, \tag{A.1}$$

$$A \in \mathbb{C}^{m \times n}, \quad \text{rank}(A) = m \quad \Rightarrow \quad A^+ = A^*(AA^*)^{-1} \quad \Rightarrow \quad AA^+ = I_{m \times m}, \tag{A.2}$$

and

$$A \in \mathbb{C}^{m \times n}, \quad B \in \mathbb{C}^{n \times p}, \quad \text{rank}(A) = \text{rank}(B) = n \quad \Rightarrow \quad (AB)^+ = B^+A^+. \tag{A.3}$$

The following theorems on Moore–Penrose inverses of perturbed matrices can be found in Ben-Israel and Greville [2], Wedin [44], and Stewart [40]. Here we give independent proofs of two of them.

*Remark* For convenience of notation, throughout this appendix only, we will use  $\|\cdot\|$  to denote the  $l_2$  norm. (Thus,  $\|\cdot\|$  here does *not* stand for the  $G$  norm we have used in Sections 1–6.)

**Theorem A.1** *Let  $A \in \mathbb{C}^{m \times n}$ ,  $\text{rank}(A) = n$ , and let  $G \in \mathbb{C}^{m \times m}$  be nonsingular and define  $B = GA$ . Then  $\text{rank}(B) = n$  too, and*

$$\|B^+\| \leq \|G^{-1}\| \|A^+\|.$$

*Proof* That  $\text{rank}(B) = n$  is clear since  $G$  is nonsingular. Starting now with  $A = G^{-1}B$ , we first have

$$Ax = G^{-1}(Bx) \quad \Rightarrow \quad \|Ax\| \leq \|G^{-1}\| \|Bx\| \quad \forall x \in \mathbb{C}^n, \quad \|x\| = 1.$$

Let  $x'$  and  $x''$ , with  $\|x'\| = 1$  and  $\|x''\| = 1$ , be such that

$$\sigma_{\min}(A) = \min_{\|x\|=1} \|Ax\| = \|Ax'\| \quad \text{and} \quad \sigma_{\min}(B) = \min_{\|x\|=1} \|Bx\| = \|Bx''\|,$$

where  $\sigma_{\min}(K)$  denotes the smallest singular value of a matrix  $K$ . Then

$$\sigma_{\min}(A) = \|Ax'\| \leq \|Ax''\| \leq \|G^{-1}\| \|Bx''\| = \|G^{-1}\| \sigma_{\min}(B).$$

The result follows by recalling that  $\|K^+\| = 1/\sigma_{\min}(K)$  when  $K$  has full column rank, which implies that  $\sigma_{\min}(K) > 0$ . □

**Theorem A.2** *Let  $A \in \mathbb{C}^{m \times n}$  and  $(A + E) \in \mathbb{C}^{m \times n}$ ,  $m \geq n$ , such that  $\text{rank}(A) = n$  and  $\|EA^+\| < 1$ . Then*

$$\|(A + E)^+\| \leq \frac{\|A^+\|}{1 - \|EA^+\|}.$$

If  $\Delta = \|E\| \|A^+\| < 1$  in addition, then this result can be expressed as

$$\|(A + E)^+\| \leq \frac{1}{1 - \Delta} \|A^+\|.$$

*Proof* First, because  $A$  is of full column rank, we have that  $A^+A = I_{n \times n}$ . Consequently,

$$A + E = (I + EA^+)A.$$

Since  $\|EA^+\| < 1$  by assumption, the matrix  $G = I + EA^+$  is nonsingular. The first result now follows from Theorem A.1 and by the fact that  $\|G^{-1}\| \leq 1/(1 - \|EA^+\|)$ . The second result follows by invoking  $\|EA^+\| \leq \|E\| \|A^+\| = \Delta$  and the additional assumption that  $\Delta < 1$ .  $\square$

**Theorem A.3** Let  $A$  and  $E$  be as in Theorem A.2,  $\Delta = \|E\| \|A^+\| < 1$ , and let  $H = (A + E)^+ - A^+$ . Then

$$\|H\| \leq \sqrt{2} \frac{\Delta}{1 - \Delta} \|A^+\|.$$

*Proof* By Wedin [44, Theorem 4.1], there holds

$$\|H\| \leq \sqrt{2} \|(A + E)^+\| \|A^+\| \|E\|.$$

Invoking now Theorem A.2, the result follows.  $\square$

The following theorem is due to Stewart [40].

**Theorem A.4** Let  $A_1, A_2, \dots$ , and  $A$  be such that  $\lim_{n \rightarrow \infty} A_n = A$ . Then  $\lim_{n \rightarrow \infty} A_n^+ = A^+$  if and only if  $\text{rank}(A_n) = \text{rank}(A)$ ,  $n \geq n_0$ , for some integer  $n_0$ .

## References

1. Anderson, D.G.: Iterative procedures for nonlinear integral equations. *J. ACM* **12**, 547–560 (1965)
2. Ben-Israel, A.: On error bounds for generalized inverses. *SIAM J. Numer. Anal.* **3**, 585–592 (1966)
3. Ben-Israel, A., Greville, T.N.E. *Generalized Inverses: Theory and Applications*. CMS Books in Mathematics, 2nd edn. Springer, New York (2003)
4. Brezinski, C.: Application de l' $\epsilon$ -algorithme à la résolution des systèmes non linéaires. *C. R. Acad. Sci. Paris* **271 A**, 1174–1177 (1970)
5. Brezinski, C.: Sur un algorithme de résolution des systèmes non linéaires. *C. R. Acad. Sci. Paris* **272 A**, 145–148 (1971)
6. Brezinski, C.: Généralisations de la transformation de Shanks, de la table de Padé, et de l' $\epsilon$ -algorithme. *Calcolo* **12**, 317–360 (1975)
7. Brezinski, C.: Accélération de la Convergence en Analyse Numérique. Number 584 in *Lecture Notes in Mathematics*, Springer, Berlin (1977)
8. Brezinski, C., Redivo Zaglia, M.: *Extrapolation Methods: Theory and Practice*. North-Holland, Amsterdam (1991)
9. Cabay, S., Jackson, L.W.: A polynomial extrapolation method for finding limits and antilimits of vector sequences. *SIAM J. Numer. Anal.* **13**, 734–752 (1976)
10. Campbell, S.L., Meyer, C.D. Jr.: *Generalized Inverses of Linear Transformations*. Dover, New York (1991)



11. Eddy, R.P.: Extrapolating to the limit of a vector sequence. In: Wang, P.C.C. (ed.) *Information Linkage Between Applied Mathematics and Industry*, pp. 387–396. Academic Press, New York (1979)
12. Gekeler, E.: On the solution of systems of equations by the epsilon algorithm of Wynn. *Math. Comp.* **26**, 427–436 (1972)
13. Golub, G.H., Van Loan, C.F. *Matrix Computations*, 4th edn. The Johns Hopkins University Press, Baltimore (2013)
14. Graves-Morris, P.R., Saff, E.B.: Row convergence theorems for generalised inverse vector-valued Padé, approximants. *J. Comp. Appl. Math.* **23**, 63–85 (1988)
15. Jbilou, K., Sadok, H.: Some results about vector extrapolation methods and related fixed-point iterations. *J. Comp. Appl. Math.* **36**, 385–398 (1991)
16. Jbilou, K., Sadok, H.: Analysis of some vector extrapolation methods for linear systems. *Numer. Math.* **70**, 73–89 (1995)
17. Jbilou, K., Sadok, H.: LU-implementation of the modified minimal polynomial LU-extrapolation method. *IMA J. Numer. Anal.* **19**, 549–561 (1999)
18. Kaniel, S., Stein, J.: Least-square acceleration of iterative methods for linear equations. *J. Optimization Theory Appl.* **14**, 431–437 (1974)
19. Laurens, J., Le Ferrand, H.: Fonctions d’itérations vectorielles, itérations rationnelles. *C. R. Acad. Sci. Paris* **321 I**, 631–636 (1995)
20. Le Ferrand, H.: Convergence of the topological  $\epsilon$ -algorithm for solving systems of nonlinear equations. *Numer. Algorithms* **3**, 273–283 (1992)
21. Mešina, M.: Convergence acceleration for the iterative solution of the equations  $x = AX + f$ . *Comput. Methods Appl. Mech. Engrg.* **10**, 165–173 (1977)
22. Ortega, J., Rheinboldt, W.: *Iterative Solution of Nonlinear Equations in Several Variables*. Academic Press, New York (1970)
23. Pugachev, B.P.: Acceleration of the convergence of iterative processes and a method of solving systems of nonlinear equations. *U.S.S.R. Comput. Math. Math. Phys.* **17**, 199–207 (1978)
24. Shanks, D.: Nonlinear transformations of divergent and slowly convergent sequences. *J. Math. and Phys.* **34**, 1–42 (1955)
25. Sidi, A.: Convergence and stability properties of minimal polynomial and reduced rank extrapolation algorithms. *SIAM J. Numer. Anal.* **23**, 197–209 (1986). Originally appeared as NASA TM-83443 (1983)
26. Sidi, A.: Extrapolation vs. projection methods for linear systems of equations. *J. Comp. Appl. Math.* **22**, 71–88 (1988)
27. Sidi, A.: Efficient implementation of minimal polynomial and reduced rank extrapolation methods. *J. Comp. Appl. Math.* **36**, 305–337 (1991). Originally appeared as NASA TM-103240 ICOMP-90-20 (1990)
28. Sidi, A.: Convergence of intermediate rows of minimal polynomial and reduced rank extrapolation tables. *Numer. Algorithms* **6**, 229–244 (1994)
29. Sidi, A.: Extension and completion of Wynn’s theory on convergence of columns of the epsilon table. *J. Approx. Theory* **86**, 21–40 (1996)
30. Sidi, A.: Review of two vector extrapolation methods of polynomial type with applications to large-scale problems. *J. Comput. Sci.* **3**, 92–101 (2012)
31. Sidi, A.: SVD-MPE: An SVD-based vector extrapolation method of polynomial type. *Appl. Math.* **7**, 1260–1278 (2016). Special issue on Applied Iterative Methods
32. Sidi, A.: Minimal polynomial and reduced rank extrapolation methods are related. *Adv. Comput. Math.* **43**, 151–170 (2017)
33. Sidi, A.: *Vector Extrapolation Methods with Applications*. Number 17 in SIAM Series on Computational Science and Engineering. SIAM, Philadelphia (2017)
34. Sidi, A., Bridger, J.: Convergence and stability analyses for some vector extrapolation methods in the presence of defective iteration matrices. *J. Comp. Appl. Math.* **22**, 35–61 (1988)
35. Sidi, A., Ford, W.F., Smith, D.A.: Acceleration of convergence of vector sequences. *SIAM J. Numer. Anal.* **23**, 178–196 (1986). Originally appeared as NASA TP-2193 (1983)
36. Sidi, A., Shapira, Y.: Upper bounds for convergence rates of vector extrapolation methods on linear systems with initial iterations. Technical Report 701, Computer Science Dept., Technion–Israel Institute of Technology, 1991. Appeared also as NASA TM-105608 ICOMP-92-09 (1992)
37. Sidi, A., Shapira, Y.: Upper bounds for convergence rates of acceleration methods with initial iterations. *Numer. Algorithms* **18**, 113–132 (1998)

38. Skelboe, S.: Computation of the periodic steady-state response of nonlinear networks by extrapolation methods. *IEEE Trans. Circuits Syst.* **27**, 161–175 (1980)
39. Smith, D.A., Ford, W.F., Sidi, A.: Extrapolation methods for vector sequences. *SIAM Rev.* **29**, 199–233 (1987). Erratum: *SIAM Rev.* **30**, 623–634 (1988)
40. Stewart, G.W.: On the continuity of the generalized inverse. **17**, 33–45 (1969)
41. Toth, A., Kelly, C.T.: Convergence analysis for Anderson acceleration. *SIAM J. Numer. Anal.* **53**, 805–819 (2015)
42. Varga, R.S. *Matrix Iterative Analysis*. Number 27 in Springer Series in Computational Mathematics, 2nd edn. Springer, New York (2000)
43. Walker, H.F., Ni, P.: Anderson acceleration for fixed-point iterations. *SIAM J. Numer. Anal.* **49**, 1715–1735 (2011)
44. Wedin, P.Å.: Perturbation theory for pseudo-inverses. *BIT* **13**, 217–232 (1973)
45. Wynn, P.: On a device for computing the  $e_m(S_n)$  transformation. *Math. Tables Aids to Comput.* **10**, 91–96 (1956)
46. Wynn, P.: Acceleration techniques for iterated vector and matrix problems. *Math. Comp.* **16**, 301–322 (1962)
47. Wynn, P.: On the convergence and stability of the epsilon algorithm. *SIAM J. Numer. Anal.* **3**, 91–122 (1966)

**Publisher's note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.