



Universidad  
Carlos III de Madrid



This is a pre-copyedited, Frolilán Martínez Dopico produced PDF of an article accepted for publication in *IMA Journal of Numerical Analysis*. The version of record F. De Terán, F. M. Dopico, and J. Pérez (2015). Backward stability of polynomial root-finding using Fiedler companion matrices, in *IMA Journal of Numerical Analysis*, Numer Anal, pp. 1-41, is available online at: <http://dx.doi.org/10.1093/imanum/dru057>

© The author 2014. Oxford University Press

# Backward stability of polynomial root-finding using Fiedler companion matrices

FERNANDO DE TERÁN<sup>†</sup>, FROILÁN M. DOPICO<sup>‡</sup>, AND JAVIER PÉREZ<sup>§</sup>  
 DEPARTAMENTO DE MATEMÁTICAS, UNIVERSIDAD CARLOS III DE MADRID,  
 AVDA. UNIVERSIDAD 30, 28911 LEGANÉS, SPAIN

[Received on 20 August 2014]

Computing roots of scalar polynomials as the eigenvalues of Frobenius companion matrices using backward stable eigenvalue algorithms is a classical approach. The introduction of new families of companion matrices allows for the use of other matrices in the root-finding problem. In this paper, we analyze the backward stability of polynomial root-finding algorithms via Fiedler companion matrices. In other words, given a polynomial  $p(z)$ , the question is to determine whether the whole set of computed eigenvalues of the companion matrix, obtained with a backward stable algorithm for the standard eigenvalue problem, are the set of roots of a nearby polynomial or not. We show that, if the coefficients of  $p(z)$  are bounded in absolute value by a moderate number, then algorithms for polynomial root-finding using Fiedler matrices are backward stable, and Fiedler matrices are as good as the Frobenius companion matrices. This allows us to use Fiedler companion matrices with favorable structures in the polynomial root-finding problem. However, when some of the coefficients of the polynomial are large, Fiedler companion matrices may produce larger backward errors than Frobenius companion matrices, although in this case neither Frobenius nor Fiedler matrices lead to backward stable computations. To prove this we obtain explicit expressions for the change, to first order, of the characteristic polynomial coefficients of Fiedler matrices under small perturbations. We show that, for all Fiedler matrices except the Frobenius ones, this change involves quadratic terms in the coefficients of the characteristic polynomial of the original matrix, while for the Frobenius matrices it only involves linear terms. We present extensive numerical experiments that support these theoretical results. The effect of balancing these matrices is also investigated.

*Keywords:* roots of polynomials; eigenvalues; characteristic polynomial; Fiedler companion matrices; backward stability, conditioning

## 1. Introduction

Let  $p(z)$  be a monic polynomial of degree  $n$ ,

$$p(z) := z^n + \sum_{k=0}^{n-1} a_k z^k, \quad (1.1)$$

with  $a_k \in \mathbb{C}$ , for  $k = 0, \dots, n-1$ . The *first* and *second Frobenius companion matrices* of  $p(z)$  are defined as

$$C_1 := \begin{bmatrix} -a_{n-1} & -a_{n-2} & \cdots & -a_1 & -a_0 \\ 1 & 0 & \cdots & 0 & 0 \\ 0 & 1 & \ddots & 0 & 0 \\ \vdots & \ddots & \ddots & \vdots & \vdots \\ 0 & \cdots & 0 & 1 & 0 \end{bmatrix} \quad \text{and} \quad C_2 := \begin{bmatrix} -a_{n-1} & 1 & 0 & \cdots & 0 \\ -a_{n-2} & 0 & 1 & \ddots & \vdots \\ \vdots & \vdots & \ddots & \ddots & 0 \\ -a_1 & 0 & 0 & \cdots & 1 \\ -a_0 & 0 & 0 & \cdots & 0 \end{bmatrix}, \quad (1.2)$$

and they satisfy:  $\det(zI - C_1) = \det(zI - C_2) = p(z)$ . Hence, the eigenvalues of both  $C_1$  and  $C_2$  coincide with the roots of  $p(z)$ . Then, the root-finding problem for scalar monic polynomials (1.1) can be reformulated as an eigenvalue problem. However, these two problems present relevant differences from the numerical point of view regarding, in

<sup>†</sup>Corresponding author. Email: fteran@math.uc3m.es

<sup>‡</sup>Email: dopico@math.uc3m.es

<sup>§</sup>Email: jpalvaro@math.uc3m.es

particular, conditioning and backward errors. The difference relies on the fact that, due to perturbations, the companion matrix may become a dense matrix, which has not the structure of a companion matrix any more. In other words, small perturbations of the companion matrix might not correspond to small perturbations of the associated polynomial.

To be more precise, a standard way to compute the roots of  $p(z)$  is just by computing the eigenvalues of  $C_1$  (or  $C_2$ ). This is, for instance, the way followed by the MATLAB command `roots`, after balancing the Frobenius matrix. The MATLAB command `roots` then uses the QR-algorithm on the Frobenius matrix to get its eigenvalues. Though this may not be the best way to address the polynomial root-finding problem, from the point of view of efficiency and storage (see, for instance, Moler (1991)), it has been extensively used because of the advantages of the QR algorithm (robustness and backward stability). Nonetheless, to overcome the mentioned drawbacks on the efficiency (measured in number of operations) and storage, several fast variants of the QR method have been proposed, which take advantage of the structure of the companion matrix (see, for instance, Aurentz *et al.* (2013); Bini *et al.* (2004, 2005, 2010); Calvetti *et al.* (2002); Chandrasekaran *et al.* (2008); Gemignani (2007); Van Barel *et al.* (2010)), but none of them has been proved to be stable. In a different line of research, also variants of  $C_1, C_2$  have been proposed, devoted to improve the accuracy in the case of multiple roots, where the standard companion matrix gives less accurate results than for simple roots (see Brugnano & Trigianta (1995); Niu & Sakurai (2003)). In this paper, we are interested in the backward stability of the root-finding problem solved via an eigenvalue backward stable method, but for a wider class of companion matrices (namely, the Fiedler matrices, see Fiedler (2003)). Our work is motivated by Edelman & Murakami (1995) and Toh & Trefethen (1994), which address related issues for the Frobenius matrices.

Let us first focus on the root-finding problem for  $p(z)$  using the first Frobenius companion matrix  $C_1$ . Since the QR-algorithm is backward stable, the whole ensemble of computed eigenvalues is the whole ensemble of exact eigenvalues of a matrix  $C_1 + E$ , where  $E$  is a dense matrix such that

$$\|E\| = O(u)\|C_1\|, \quad (1.3)$$

for some matrix norm  $\|\cdot\|$ , and where  $u$  denotes the machine epsilon. However, this does not guarantee that these (computed) eigenvalues are the roots of a nearby polynomial of  $p(z)$  or, in other words, that the method is backward stable from the point of view of the polynomials. In this paper, we investigate this issue. In order for the method to be backward stable from the point of view of the polynomials in a normwise sense, the computed eigenvalues should be the exact roots of a polynomial  $\tilde{p}(z)$  such that

$$\frac{\|\tilde{p} - p\|}{\|p\|} = O(u),$$

for some polynomial norm  $\|\cdot\|$ . As we will see in Section 3, the backward stability of polynomial root-finding algorithms using companion matrices is closely related to the *conditioning* of the characteristic polynomial under perturbations of these matrices. This conditioning can be measured through the first order term of the Taylor expansion of the coefficients of the characteristic polynomial. In Edelman & Murakami (1995) it has been shown that, if

$$\tilde{p}(z) = \det(zI - C_1 - E) = z^n + \sum_{k=0}^{n-1} \tilde{a}_k z^k \quad (1.4)$$

then, to first order in (the entries of)  $E$ ,

$$\tilde{a}_k - a_k = \sum_{s=0}^k \sum_{j=1}^{n-k-1} a_s E_{j-s+k+1, j} - \sum_{s=k+1}^n \sum_{j=n-k}^n a_s E_{j-s+k+1, j}. \quad (1.5)$$

If the eigenvalues of  $C_1$  are computed with a backward stable algorithm, it may be proved from (1.5) that, to first order in  $E$ , the computed eigenvalues are the exact roots of a polynomial  $\tilde{p}(z)$  as in (1.4) such that

$$\frac{\|\tilde{p} - p\|}{\|p\|} = O(u)\|p\|, \quad (1.6)$$

with  $E$  satisfying (1.3). Note that (1.6) does not imply that computing the roots of  $p(z)$  using  $C_1$  (or  $C_2$ ) is a backward stable method from the point of view of the polynomials, since large values of  $\|p\|$  can give large backward errors. This had been already noticed, for instance, in Lemmonier & Van Dooren (2003), where the authors analyze diagonal scalings of the companion matrix to get small backward errors.

A key advantage in using Frobenius companion matrices in the root-finding problem is that they are easily constructible from the polynomial, without performing any arithmetic operation, by means of a uniform template valid for all polynomials. Any uniform template with these properties is what we mean by a *companion* matrix.

In Fiedler (2003), the author expanded the family of companion matrices associated with the monic polynomial  $p(z)$ . These matrices were named *Fiedler matrices* in De Terán *et al.* (2010). The family of Fiedler matrices includes  $C_1$  and  $C_2$  but, provided that  $n \geq 3$ , it contains some other different matrices and, in fact, many others when  $n$  is large. These matrices provide a new tool that could be used instead of  $C_1$  and  $C_2$  for computing the roots of  $p(z)$ . Some features of Fiedler matrices have been recently studied. For instance, in De Terán *et al.* (2013) the condition numbers for inversion of different Fiedler matrices have been compared, and it has been proved that, in many cases, some of the new Fiedler matrices have better conditioning than  $C_1$  and  $C_2$ . Also, in De Terán *et al.* (2014a), Fiedler matrices have been used to get new lower and upper bounds for the modulus of the roots of  $p(z)$ . We provide the formal definition of Fiedler matrices in Section 2. For the moment, the only relevant information is that, to construct them, we only need to know the polynomial  $p(z)$  and to fix a bijection  $\sigma : \{0, 1, \dots, n-1\} \rightarrow \{1, \dots, n\}$ , and that the Fiedler matrices contain, in different positions, exactly the same entries as  $C_1$  and  $C_2$ . We denote the Fiedler matrix associated with the polynomial  $p(z)$  and the bijection  $\sigma$  by  $M_\sigma(p)$ , or  $M_\sigma$  for brevity.

A natural question is whether or not computing the roots of  $p(z)$  using a Fiedler matrix  $M_\sigma$  and a backward stable eigenvalue algorithm is backward stable from the point of view of the polynomials, that is, whether or not the computed roots are the exact roots of a polynomial  $\tilde{p}(z)$  such that  $\|\tilde{p} - p\| = O(u)\|p\|$ . As it happens with Frobenius matrices, if we compute the roots of  $p(z)$  as the eigenvalues of  $M_\sigma$  with a backward stable algorithm (like the QR algorithm), then the computed roots are the exact eigenvalues of  $M_\sigma + E$ , where  $\|E\| = O(u)\|M_\sigma\|$ . However, again, this does not guarantee backward stability from the point of view of the polynomials. The goal of this paper is to analyze this issue.

To accomplish this task we need to know how the coefficients of the characteristic polynomial of  $M_\sigma$  change when the matrix is perturbed as  $M_\sigma + E$ , with  $E$  an arbitrary perturbation with no special structure. This change can be estimated, up to first order in  $E$ , through the gradients  $\nabla a_k(M_\sigma)$ , where  $a_k(X) : \mathbb{C}^{n^2} \rightarrow \mathbb{C}$  is the  $k$ th coefficient of the characteristic polynomial of a matrix  $X \in \mathbb{C}^{n \times n}$ , considered as a function of its entries. In particular, we find explicitly  $\nabla a_k(M_\sigma)$  in terms of the coefficients of  $p(z)$ . This allows us to get, up to first order, a formula for the variation of the characteristic polynomial of  $M_\sigma$  under small perturbations of  $M_\sigma$ . From this formula, we analyze the backward stability of the polynomial root-finding problem solved by applying backward stable eigensolvers to Fiedler matrices. In the recent reference Lawrence & Corless (2014), the authors address the same problem as in the present paper, namely, to know whether or not solving the polynomial root-finding problem as an eigenvalue problem is backward stable, but they use a suitable companion matrix for the polynomial expressed in barycentric form. In that reference the polynomials are not necessarily monic, but the authors follow a similar approach to ours.

To get an expression for  $\nabla a_k(M_\sigma)$ , we first prove that its coordinates are the entries of the  $(k+1)$ th coefficient of the adjoint  $\text{adj}(zI - M_\sigma)$ . Then, we get an explicit formula of  $\text{adj}(zI - M_\sigma)$ . This is a general theoretical result on Fiedler matrices that may be useful in the future to analyze other features of this family of matrices.

For a precedent on the perturbation analysis of the characteristic polynomial, we refer the reader to Ipsen & Rehman (2008). In that paper, several bounds are derived for the variation of the characteristic polynomial of an arbitrary matrix  $A$  under perturbations, in terms of symmetric functions of the singular values of  $A$ . The bounds there are very pessimistic for general matrices. However, here we take advantage of the sparsity and the structure of the Fiedler matrices to get more specific bounds depending on the coefficients of  $p(z)$ .

Throughout this paper, if  $A \in \mathbb{C}^{n \times n}$  is a matrix, then  $\|A\|_\infty$  denotes the usual matrix  $\infty$ -norm (see (Higham, 2002, p. 108)). In particular, for a vector  $v = [v_1 \ \dots \ v_n]^T \in \mathbb{C}^n$ , we have  $\|v\|_\infty = \max\{|v_1|, \dots, |v_n|\}$ . Similarly, for a polynomial  $p(z) = \sum_{k=0}^n a_k z^k$  (not necessarily monic),  $\|p\|_\infty$  is the norm on the vector space of scalar polynomials of degree less than or equal to  $n$  defined as

$$\|p\|_\infty := \max\{|a_n|, |a_{n-1}|, \dots, |a_1|, |a_0|\}.$$

Notice that, since we deal in this paper with monic polynomials,  $a_n = 1$  and we always have  $\|p\|_\infty \geq 1$ .

The main results of this work are Theorem 3.3 and Corollary 3.2. Theorem 3.3 gives, to first order in  $E$ , the coefficients of the characteristic polynomial of  $M_\sigma + E$ , and Corollary 3.2 tells us that if we compute the roots of a monic polynomial  $p(z)$  as the eigenvalues of a Fiedler matrix  $M_\sigma$  other than the Frobenius companion matrices using a backward stable eigenvalue algorithm, then the computed roots are the exact roots of a monic polynomial  $\tilde{p}(z)$  with

$$\frac{\|\tilde{p} - p\|_\infty}{\|p\|_\infty} = O(u)\|p\|_\infty^2, \quad (1.7)$$

which implies that computing the roots of  $p(z)$  using any of the Fiedler matrices of  $p(z)$  is not backward stable if  $\|p\|_\infty$  is large. For the Frobenius companion matrices, Corollary 3.2 recovers (1.6). In Section 4 we provide numerical experiments that support this theoretical result.

Our results are even more general because our formulation allows us to translate the backward errors of any algorithm for computing the eigenvalues of  $M_\sigma$  to the backward error of the polynomial root-finding problem, even when the algorithm is not backward stable. This is particularly interesting for any fast algorithm that has been or might be developed in the future for computing eigenvalues of special Fiedler matrices. To be more precise, if instead of an expression like (1.3) the eigensolver computes the eigenvalues of a matrix  $M_\sigma + E$ , with

$$\|E\| = c(p)O(u)\|M_\sigma\|,$$

where  $c(p)$  is some quantity depending on  $p(z)$ , then (1.7) is replaced by

$$\frac{\|\tilde{p} - p\|_\infty}{\|p\|_\infty} = c(p)O(u)\|p\|_\infty^2.$$

As a consequence of (1.6) and (1.7) we get the following conclusions:

- (C1) From the point of view of the normwise backward errors in the (monic) polynomial  $p(z)$ , any Fiedler matrix can be used for solving the root-finding problem with the same reliability as Frobenius companion matrices when  $\|p\|_\infty = O(1)$ . In this case, the root-finding problem solved by applying a backward stable eigenvalue algorithm on any Fiedler companion matrix is a backward stable method.
- (C2) However, when  $\|p\|_\infty$  is large none of the Fiedler matrices leads to a backward stable algorithm for the root-finding problem and, moreover, any Fiedler matrix other than Frobenius companion matrices may produce much larger backward errors than the ones produced when using Frobenius matrices.

Note, in particular, that since  $\|p\|_\infty \geq 1$ , no Fiedler matrix can improve the behavior of Frobenius matrices in the root-finding problem from the point of view of backward errors. Anyway, the particular structure of some Fiedler matrices can make their use more efficient than the use of classical Frobenius companion matrices. For instance, we could take advantage of the pentadiagonal structure of some Fiedler matrices (which exist for any value of  $n$ , see De Terán *et al.* (2010)) to devise structured versions of the *LR* algorithm to get its eigenvalues in  $O(n^2)$  flops (see, for instance, Zhlobich (2012)). However, as for all structured methods for the root-finding problem, stability can not yet be guaranteed.

We have also considered the effect of balancing (see Parlett & Reinsch (1969)) Fiedler companion matrices on the backward errors of the root-finding problem for  $p(z)$  using a Fiedler matrix  $M_\sigma$ . The numerical experiments carried out in Section 4 indicate that balancing very often improves the backward errors for general polynomials, including some polynomials for which the backward error without balancing is quite large. However, we prove that, when  $|a_{n-1}|$  is much larger than  $|a_{n-2}|$ , the condition number of  $p(z)$  using any balanced Fiedler matrix is large, and so is the backward error. Some experiments on polynomials with  $|a_{n-1}|$  much larger than  $|a_{n-2}|$  show that, indeed, balancing the Fiedler matrices does not guarantee backward stability for the root-finding polynomial problem.

The paper is organized as follows. In Section 2 we introduce Fiedler matrices and their basic properties. In Section 3 we analyze, to first order, the change of the coefficients of the characteristic polynomial of Fiedler matrices under matrix perturbations, and we connect it with the backward error of the polynomial root-finding problem solved via an eigenvalue algorithm. This section contains the main results of the paper. Due to the length and technical nature of this section, some proofs have been omitted or reduced. For more detailed proofs, we refer the reader to De Terán *et al.* (2014b), which is an extended version of this paper. Section 4 is devoted to numerical experiments that illustrate the theoretical results obtained in Section 3. In Section 5 we provide a geometric interpretation of the change, to first order, of the characteristic polynomial of Fiedler matrices in terms of the orbit space under similarity of these matrices. This is motivated by the one in Edelman & Murakami (1995) for Frobenius companion matrices, and gives a decomposition of  $\mathbb{C}^{n \times n}$  as the sum of the tangent space to the similarity orbit of a Fiedler matrix and the Sylvester space of matrices associated to it. Section 6 presents a summary of the main contributions of the paper.

## 2. Fiedler matrices. Definition and basic properties

For a given polynomial  $p(z)$  as in (1.1), we define the  $n \times n$  matrices

$$M_0 := \begin{bmatrix} I_{n-1} & 0 \\ 0 & -a_0 \end{bmatrix} \quad \text{and} \quad M_k := \begin{bmatrix} I_{n-k-1} & & & \\ & -a_k & 1 & \\ & 1 & 0 & \\ & & & I_{k-1} \end{bmatrix}, \quad k = 1, \dots, n-1, \quad (2.1)$$

which are the basic factors used to build all Fiedler matrices. Here and in the rest of the paper  $I_j$  denotes the  $j \times j$  identity matrix. In Fiedler (2003) Fiedler matrices are constructed as the product  $M_{i_1}M_{i_2}\cdots M_{i_n}$ , where  $(i_1, i_2, \dots, i_n)$  is any possible permutation of the  $n$ -tuple  $(0, 1, \dots, n-1)$ . In order to better express certain key properties of this permutation and the resulting Fiedler matrix, in De Terán *et al.* (2010) the authors index the product of the  $M_i$  factors in a slightly different way, as it is described in the following definition.

**DEFINITION 2.1** Let  $p(z) = z^n + \sum_{k=0}^{n-1} a_k z^k$ , with  $n \geq 2$ , and let  $M_i$ , for  $i = 0, 1, \dots, n-1$ , be the matrices in (2.1). Given any bijection  $\sigma : \{0, 1, \dots, n-1\} \rightarrow \{1, \dots, n\}$ , the Fiedler matrix of  $p(z)$  associated with  $\sigma$  is the  $n \times n$  matrix

$$M_\sigma(p) := M_{\sigma^{-1}(1)} \cdots M_{\sigma^{-1}(n)}. \quad (2.2)$$

We want to notice that  $\sigma(i)$  in (2.2) describes the position of the factor  $M_i$  in the product  $M_{\sigma^{-1}(1)} \cdots M_{\sigma^{-1}(n)}$ , i.e.,  $\sigma(i) = j$  means that  $M_i$  is the  $j$ th factor in the product. We want to note also that the building factors (2.1) of (2.2) depend also on  $p(z)$  (to be precise, they depend on its coefficients). However, in this case we do not write explicitly this dependence for the sake of simplicity. For the same reason, we will also drop the dependence on  $p$  in  $M_\sigma$  when there is no risk of confusion (namely, until Section 5).

The family of matrices  $\{M_k\}_{k=0}^{n-1}$  satisfies the following commutativity relations

$$M_i M_j = M_j M_i \quad \text{for } |i - j| \neq 1. \quad (2.3)$$

It is proved in Fiedler (2003) that all Fiedler matrices of  $p(z)$  are similar, so they have  $p(z)$  as characteristic polynomial. Frobenius companion matrices of  $p(z)$  are particular cases of Fiedler matrices, namely,  $C_1 = M_{n-1}M_{n-2}\cdots M_1M_0$  and  $C_2 = M_0M_1\cdots M_{n-2}M_{n-1}$ . Observe that the matrices  $M_i$  are symmetric, and therefore the transpose of any Fiedler matrix is another Fiedler matrix, obtained by reversing the order of the  $M_i$  factors in (2.2).

The relations (2.3) imply that some Fiedler matrices associated with different bijections  $\sigma$  are equal. For example, for  $n = 3$ , the Fiedler matrices  $M_0M_2M_1$  and  $M_2M_0M_1$  are equal. These relations suggest that the relative positions of the matrices  $M_i$  and  $M_{i+1}$  in the product  $M_\sigma$  are of fundamental interest in studying Fiedler matrices. This motivates Definition 2.2, partially introduced in De Terán *et al.* (2010).

**DEFINITION 2.2** Let  $\sigma : \{0, 1, \dots, n-1\} \rightarrow \{1, \dots, n\}$  be a bijection.

- (a) For  $i = 0, \dots, n-2$ , we say that  $\sigma$  has a consecution at  $i$  if  $\sigma(i) < \sigma(i+1)$  and that  $\sigma$  has an inversion at  $i$  if  $\sigma(i) > \sigma(i+1)$ .
- (b) The positional consecution-inversion sequence of  $\sigma$ , denoted by  $\text{PCIS}(\sigma)$ , is the  $(n-1)$ -tuple  $(v_0, \dots, v_{n-2})$  such that  $v_j = 1$  if  $\sigma$  has a consecution at  $j$  and  $v_j = 0$  otherwise.

**REMARK 2.1** We note that  $\sigma$  has a consecution at  $i$ , that is  $v_i = 1$ , if and only if  $M_i$  is to the left of  $M_{i+1}$  in the product defining the Fiedler matrix  $M_\sigma$ , while  $\sigma$  has an inversion at  $i$ , that is  $v_i = 0$ , if and only if  $M_i$  is to the right of  $M_{i+1}$  in  $M_\sigma$ . This simple observation on Definition 2.2 will be used freely.

In order to keep the notation in future sections reasonably simple we introduce the following definitions.

**DEFINITION 2.3** Let  $\sigma : \{0, 1, \dots, n-1\} \rightarrow \{1, \dots, n\}$  be a bijection with  $\text{PCIS}(\sigma) = (v_0, v_1, \dots, v_{n-2})$ , then:

- (a) The extended positional consecution-inversion sequence of  $\sigma$ , denoted by  $\text{EPCIS}(\sigma)$ , is the  $n$ -tuple  $(v_0, v_1, \dots, v_{n-1})$ , where  $v_{n-1} = v_{n-2}$ .
- (b) For  $0 \leq i \leq j \leq n-2$ , we set

$$i_\sigma(i : j) := \sum_{k=i}^j (1 - v_k) \quad \text{and} \quad c_\sigma(i : j) := \sum_{k=i}^j v_k$$

for, respectively, the number of inversions and consecutions of  $\sigma$  from  $i$  to  $j$ . We also set  $i_\sigma(i : j) := c_\sigma(i : j) := 0$  for  $i > j$ .

The following immediate identities will be used several times along the paper:

$$i_\sigma(i : j) + c_\sigma(i : j) = j - i + 1, \quad \text{for } 0 \leq i \leq j \leq n - 2, \quad (2.4)$$

$$i_\sigma(0 : i) + c_\sigma(0 : j) \leq n - 1, \quad \text{for } 0 \leq i, j \leq n - 2. \quad (2.5)$$

We close this section with the following notion, that will be used along the paper.

**DEFINITION 2.4** Let  $p(z) = z^n + \sum_{k=0}^{n-1} a_k z^k$  be a monic polynomial of degree  $n$ . For  $d = 0, 1, \dots, n$ , the degree  $d$  Horner shift of  $p(z)$  is the polynomial  $p_d(z) = z^d + a_{n-1}z^{d-1} + \dots + a_{n-d+1}z + a_{n-d}$ .

Notice that the Horner shifts of  $p(z)$  satisfy the following recurrence relation

$$\begin{cases} p_0(z) = 1, & \text{and} \\ p_d(z) = zp_{d-1}(z) + a_{n-d}, & \text{for } d = 1, 2, \dots, n. \end{cases} \quad (2.6)$$

### 3. Backward error, conditioning, and first order perturbation terms of the characteristic polynomial

A natural definition of the normwise *backward error* of the computed roots,  $\tilde{\lambda}_1, \dots, \tilde{\lambda}_n$ , of the monic polynomial (1.1) via a certain algorithm is

$$\eta_\infty(\tilde{\lambda}_1, \dots, \tilde{\lambda}_n) := \frac{\|\tilde{p} - p\|_\infty}{\|p\|_\infty},$$

where  $\tilde{p}(z) = \prod_{i=1}^n (z - \tilde{\lambda}_i)$ . This notion of backward error coincides with the relative distance, in the  $\infty$ -norm, between the original polynomial  $p(z)$  and the monic polynomial  $\tilde{p}(z)$  whose roots are  $\tilde{\lambda}_1, \dots, \tilde{\lambda}_n$ . The key in our approach is that the roots are computed as the eigenvalues of a (companion) matrix,  $A$ , so that the computed roots are the exact eigenvalues of some perturbation of  $A$ , say  $A + E$ . In other words,  $p(z) = \det(zI - A)$  and, following (1.4) for a general companion matrix  $A$ , we also have  $\tilde{p}(z) = \det(zI - (A + E))$ . Hence, the difference between  $p(z)$  and  $\tilde{p}(z)$  can be measured from the variation of the coefficients of the characteristic polynomial of  $A$  under small perturbations of  $A$ .

Hence, we consider the  $k$ th coefficient of the characteristic polynomial of a matrix  $X = [x_{ij}] \in \mathbb{C}^{n \times n}$  as a function of the entries of  $X$ ,  $a_k(X) : \mathbb{C}^{n^2} \rightarrow \mathbb{C}$ , for  $k = 0, 1, \dots, n - 1$ . Equivalently:

$$\det(zI - X) = z^n + \sum_{k=0}^{n-1} a_k(X) z^k.$$

The function  $a_k(X)$  is a multivariable polynomial function of the entries of  $X$ . Therefore, the first order term in  $E$  of its Taylor polynomial centered at  $A$  is (see, for instance (Grauert & Fritzsche, 1976, Th. 3.8)) for functions of several complex variables)

$$a_k(A + E) = a_k(A) + \sum_{i,j=1}^n \frac{\partial a_k(X)}{\partial x_{ij}} \Big|_{X=A} E_{ij} = a_k(A) + \nabla a_k(A) \cdot \text{vec}(E), \quad \text{for } k = 0, 1, \dots, n - 1, \quad (3.1)$$

where, for a given  $m \times n$  matrix  $M = [m_{ij}]$ ,  $\text{vec}(M)$  is the *vectorization* of  $M$ , namely, the column vector

$$\text{vec}(M) := [m_{11} \dots m_{m1} m_{12} \dots m_{m2} \dots m_{1n} \dots m_{mn}]^T$$

(see (Horn & Johnson, 1985, Def. 4.2.9), for instance), and

$$\nabla a_k(A) = \left[ \frac{\partial a_k(X)}{\partial x_{11}} \Big|_{X=A} \dots \frac{\partial a_k(X)}{\partial x_{n1}} \Big|_{X=A} \quad \frac{\partial a_k(X)}{\partial x_{12}} \Big|_{X=A} \dots \frac{\partial a_k(X)}{\partial x_{n2}} \Big|_{X=A} \quad \dots \quad \frac{\partial a_k(X)}{\partial x_{1n}} \Big|_{X=A} \dots \frac{\partial a_k(X)}{\partial x_{nn}} \Big|_{X=A} \right].$$

Therefore, to first order in  $E$ , we have

$$|a_k(A + E) - a_k(A)| = |\nabla a_k(A) \cdot \text{vec}(E)|.$$

For any Fiedler matrix  $M_\sigma$ , we get an explicit expression of  $\nabla a_k(M_\sigma)$  in terms of the entries of  $M_\sigma$  or, equivalently, in terms of the coefficients of its characteristic polynomial  $p(z)$ . The corresponding expression was given in Edelman & Murakami (1995) for Frobenius companion matrices, which are particular cases of Fiedler matrices. The general expression we provide here, valid for any Fiedler matrix, requires different techniques to the ones in that paper.

The following well-know result (known as Jacobi's formula, see Bhatia & Jain (2009)) provides us a description of the gradient of the determinant. We include a short proof here for completeness.

LEMMA 3.1 Let  $A \in \mathbb{C}^{n \times n}$ , and consider a small perturbation  $A + E$ , with  $E \in \mathbb{C}^{n \times n}$ . Then, the function

$$\begin{aligned} \det : \mathbb{C}^{n \times n} &\longrightarrow \mathbb{C} \\ X &\longmapsto \det(X), \end{aligned}$$

is analytic in a neighborhood of  $A$ , and

$$\det(A + E) = \det(A) + \text{tr}(\text{adj}(A)E) + O(\|E\|^2),$$

where  $\|\cdot\|$  is any norm in  $\mathbb{C}^{n \times n}$ ,  $\text{adj}(A)$  is the adjugate matrix of  $A$  (see Bernstein (2009)), and  $\text{tr}(B)$  is the trace of  $B$ .

*Proof.* The function  $\det : \mathbb{C}^{n \times n} \rightarrow \mathbb{C}$  is clearly analytic in a neighborhood of  $A$ , since it is a polynomial function on the entries of  $X \in \mathbb{C}^{n \times n}$ . Moreover, analogously to (3.1), with the function  $\det$  instead of  $a_k$ , we get

$$\det(A + E) = \det(A) + \nabla \det(A) \cdot \text{vec}(E) + O(\|E\|^2).$$

Now, it is straightforward to check that

$$\left. \frac{\partial \det(X)}{\partial x_{ij}} \right|_{X=A} = (\text{adj}(A))_{ji}$$

(see also (Bernstein, 2009, Fact 10.11.21)). The result now follows from the identity  $\text{tr}(AB) = \text{vec}(A^T)^T \cdot \text{vec}(B)$ , which is valid for every  $A, B \in \mathbb{C}^{n \times n}$ .  $\square$

As an immediate consequence of Lemma 3.1, applied to  $p(z) = \det(zI - A)$ , we get Proposition 3.1, which gives a description of the gradient of the coefficients of the characteristic polynomial of  $A$  and, as a consequence, an expression for the variation of the characteristic polynomial under small perturbations, up to first order.

PROPOSITION 3.1 Let  $A \in \mathbb{C}^{n \times n}$  and  $z \in \mathbb{C}$ . Let us write the adjugate matrix of  $zI - A$  as

$$\text{adj}(zI - A) = \sum_{k=0}^{n-1} z^k P_{k+1}, \quad (3.2)$$

with  $P_{k+1} \in \mathbb{C}^{n \times n}$ , for  $k = 0, 1, \dots, n-1$ . Let  $a_k(X) : \mathbb{C} \rightarrow \mathbb{C}$  be the  $k$ th coefficient of the characteristic polynomial of  $X = [x_{ij}] \in \mathbb{C}^{n \times n}$ , and let  $\nabla a_k(A)$  be the gradient of the function  $a_k(X)$  evaluated at  $A$ . Then, for  $k = 0, 1, \dots, n-1$ ,

$$\nabla a_k(A) = - [\text{vec}(P_{k+1}^T)]^T.$$

As a consequence, if  $A + E$  is a small perturbation of  $A$ , with  $E \in \mathbb{C}^{n \times n}$ , then

$$\det(zI - (A + E)) - \det(zI - A) = - \sum_{k=0}^{n-1} z^k [\text{vec}(P_{k+1}^T)]^T \cdot \text{vec}(E) + O(\|E\|^2) = - \sum_{k=0}^{n-1} z^k \text{tr}(P_{k+1}E) + O(\|E\|^2),$$

where  $\|\cdot\|$  is any norm in  $\mathbb{C}^{n \times n}$ .

*Proof.* From Lemma 3.1 and (3.2), we have

$$\begin{aligned} \det(zI - (A + E)) &= \det(zI - A) - \text{tr}(\text{adj}(zI - A)E) + O(\|E\|^2) \\ &= \det(zI - A) - \sum_{k=0}^{n-1} z^k \text{tr}(P_{k+1}E) + O(\|E\|^2) \\ &= \det(zI - A) - \sum_{k=0}^{n-1} z^k [\text{vec}(P_{k+1}^T)]^T \cdot \text{vec}(E) + O(\|E\|^2), \end{aligned}$$

and the expression for  $\nabla a_k(A)$  follows immediately from this. Note that in the last identity we have used that  $\text{tr}(AB) = \text{vec}(A^T)^T \cdot \text{vec}(B)$ , as in the proof of Lemma 3.1.  $\square$

Proposition 3.1 tells us that the variation of the characteristic polynomial of  $A \in \mathbb{C}^{n \times n}$  is given, to first order, by the trace of  $\text{adj}(zI - A)$ . This adjugate matrix is an  $n \times n$  matrix whose entries are polynomials of degree at most  $n-1$  or, equivalently, a matrix polynomial of size  $n \times n$  with degree at most  $n-1$ . Actually, its degree is exactly  $n-1$ , because of the identity:  $(zI - A) \cdot \text{adj}(zI - A) = \det(zI - A)I_n$ . In Section 3.1 we give an explicit expression for the entries of  $\text{adj}(zI - A)$ , for  $A$  being an arbitrary Fiedler matrix  $M_\sigma$ . Then, in Section 3.2, we use this information, following Proposition 3.1, to present an explicit expression for the variation, up to first order, of the coefficients of the characteristic polynomial of  $M_\sigma$  or, in other words, an explicit expression for  $\nabla a_k(M_\sigma)$ .



### 3.1 Adjugate matrix of $zI - M_\sigma$

The main result of this section is Theorem 3.2, which gives an explicit expression for  $\text{adj}(zI - M_\sigma)$ . As we have seen in (3.2), this is a matrix polynomial in the variable  $z$ . We use the notation  $\mathbb{C}^{n \times n}[z]$  for the set of  $n \times n$  matrix polynomials.

An explicit expression for the adjugate in the case of first and second Frobenius companion matrices was already known (see (Gantmacher, 1959, Ch. IV §4) or (Edelman & Murakami, 1995, p. 768)):

$$\text{adj}(zI - C_2) = \begin{bmatrix} p_0(z) \\ p_1(z) \\ \vdots \\ p_{n-1}(z) \end{bmatrix} \begin{bmatrix} z^{n-1} & \cdots & z & 1 \end{bmatrix} - p(z) \begin{bmatrix} 0 & & & & & \\ 1 & 0 & & & & \\ & z & 1 & & & \\ \vdots & & z & \ddots & & \\ \vdots & & & \ddots & \ddots & \\ z^{n-2} & z^{n-3} & \cdots & z & 1 & 0 \end{bmatrix}, \quad (3.3)$$

and  $\text{adj}(zI - C_1) = (\text{adj}(zI - C_2))^T$ . Here  $p_0(z), \dots, p_{n-1}(z)$  are the Horner shifts introduced in Definition 2.4. Equation (3.3) has a very particular structure: it is a sum of a rank-1 matrix plus a matrix whose  $(i, j)$  entry is of the form  $p(z)p_{ij}(z)$ , where  $p_{ij}(z)$  is a polynomial of degree at most  $n - 2$ . We will prove that this structure is shared also by  $\text{adj}(zI - M_\sigma)$ , for any Fiedler matrix  $M_\sigma$ . For example, if we consider the Fiedler matrix  $M_\sigma$  of a degree-6 monic polynomial  $p(z) = z^6 + \sum_{k=0}^5 a_k z^k$ , with  $\text{PCIS}(\sigma) = (1, 0, 1, 0, 1)$ , we will show that

$$\text{adj}(zI - M_\sigma) = \begin{bmatrix} z^2 \\ z^2 p_1(z) \\ z \\ z p_3(z) \\ 1 \\ p_5(z) \end{bmatrix} \begin{bmatrix} z^3 p_0(z) & z^2 & z^2 p_2(z) & z & z p_4(z) & 1 \end{bmatrix} - p(z) \begin{bmatrix} 0 & 0 & 1 & 0 & z & 0 \\ 1 & 0 & p_1(z) & 0 & z p_1(z) & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 \\ z & 1 & p_2(z) & 0 & p_3(z) & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ z^2 & z & z p_2(z) & 1 & p_4(z) & 0 \end{bmatrix}.$$

**THEOREM 3.2** Let  $p(z) = z^n + \sum_{k=0}^{n-1} a_k z^k$  be a polynomial and  $p_d(z)$ , for  $d = 0, 1, \dots, n-1$ , the degree  $d$  Horner shift of  $p(z)$ . Let  $\sigma : \{0, 1, \dots, n-1\} \rightarrow \{1, \dots, n\}$  be a bijection with  $\text{EPCIS}(\sigma) = (v_0, v_1, \dots, v_{n-1})$  and let  $M_\sigma$  be the Fiedler matrix of  $p(z)$  associated with  $\sigma$ . Let  $x_\sigma, y_\sigma \in \mathbb{C}^n[z]$  be the vector polynomials whose  $k$ th entry is

$$x_\sigma(k) = \begin{cases} z^{i_\sigma(0:n-k-1)} p_{k-1}(z) & \text{if } v_{n-k} = 1, \\ z^{i_\sigma(0:n-k-1)} & \text{if } v_{n-k} = 0, \end{cases} \quad \text{and} \quad y_\sigma(k) = \begin{cases} z^{c_\sigma(0:n-k-1)} p_{k-1}(z) & \text{if } v_{n-k} = 0, \\ z^{c_\sigma(0:n-k-1)} & \text{if } v_{n-k} = 1, \end{cases} \quad (3.4)$$

for  $k = 1, 2, \dots, n$ , and let  $A_\sigma \in \mathbb{C}^{n \times n}[z]$  be the matrix polynomial whose  $(i, j)$  entry is

$$A_\sigma(i, j) = \begin{cases} 0 & \text{if } v_{n-i} = v_{n-j} = 0 \text{ and } i \geq j, \\ z^{i_\sigma(n-j+1:n-i-1)} & \text{if } v_{n-i} = v_{n-j} = 0 \text{ and } i < j, \\ z^{c_\sigma(n-i+1:n-j-1)} & \text{if } v_{n-i} = v_{n-j} = 1 \text{ and } i > j, \\ 0 & \text{if } v_{n-i} = v_{n-j} = 1 \text{ and } i \leq j, \\ 0 & \text{if } v_{n-i} = 0 \text{ and } v_{n-j} = 1, \\ z^{c_\sigma(n-i+1:n-j-1)} p_{j-1}(z) & \text{if } v_{n-i} = 1, v_{n-j} = 0 \text{ and } i > j, \\ z^{i_\sigma(n-j+1:n-i-1)} p_{i-1}(z) & \text{if } v_{n-i} = 1, v_{n-j} = 0 \text{ and } i < j, \end{cases} \quad (3.5)$$

for  $i, j = 1, 2, \dots, n$ . Then,

$$\text{adj}(zI - M_\sigma) = x_\sigma y_\sigma^T - p(z) A_\sigma.$$

Note that  $x_\sigma, y_\sigma$  and  $A_\sigma$  depend on the variable  $z$ , though we drop it for the ease of notation.

Before proving Theorem 3.2 we state and prove some technical lemmas.

**LEMMA 3.2** Let  $x_\sigma$  and  $y_\sigma$  be the vectors defined in (3.4), and  $A_\sigma$  be the matrix defined in (3.5). Then,  $A_\sigma$  is the unique  $n \times n$  matrix satisfying the following two properties:

- (i) The entries of  $A_\sigma$  are polynomials in  $z$ , and
- (ii) all entries of  $x_\sigma y_\sigma^T - p(z) A_\sigma$  are polynomials of degree less than or equal to  $n - 1$ .

*Proof.* To prove that the entries of  $A_\sigma$  are polynomials, it suffices to see that the exponents of the powers of  $z$  appearing in the entries of (3.5) are nonnegative. This is immediate by Definition 2.3. To prove that the  $(i, j)$  entry of  $x_\sigma y_\sigma^T - p(z)A_\sigma$  is a polynomial of degree less than or equal to  $n-1$  is straightforward using (2.4) and (2.5). We show here the proof of just one case in (3.5) and refer the reader to De Terán *et al.* (2014b) for more details on the remaining cases. In particular, we assume that  $v_{n-i} = v_{n-j} = 0$  and  $i < j$ . In this case, using (2.4), the  $(i, j)$  entry of  $x_\sigma y_\sigma^T - p(z)A_\sigma$  is equal to

$$\begin{aligned} x_\sigma(i)y_\sigma(j) - p(z)A_\sigma(i, j) &= z^{i_\sigma(0:n-i-1)+c_\sigma(0:n-j-1)} p_{j-1}(z) - p(z)z^{i_\sigma(n-j+1:n-i-1)} \\ &= z^{i_\sigma(n-j+1:n-i-1)} (z^{n-j+1} p_{j-1}(z) - p(z)) \\ &= z^{i_\sigma(n-j+1:n-i-1)} (-a_{n-j} z^{n-j} - a_{n-j-1} z^{n-j-1} - \dots - a_1 z - a_0), \end{aligned}$$

which is a polynomial of degree less than  $n-1$ , because  $i_\sigma(n-j+1:n-i-1) + n-j \leq n-i-1 < n-1$ .

Now, suppose that there is another matrix  $B$ , whose entries are polynomials in  $z$ , and such that the entries of the matrix  $x_\sigma y_\sigma^T - p(z)B$  are polynomials in  $z$  of degree at most  $n-1$ . Let  $W_1 = x_\sigma y_\sigma^T - p(z)A_\sigma$  and let  $W_2 = x_\sigma y_\sigma^T - p(z)B$ , then,  $W_1 - W_2 = p(z)(B - A_\sigma)$  is a matrix whose entries are polynomials of degree at most  $n-1$ , but if  $A_\sigma \neq B$ , then  $p(z)(B - A_\sigma)$  has, at least, one entry which is a polynomial of degree at least  $n$ , hence  $A_\sigma = B$ .  $\square$

Lemma 3.3 is key to prove Theorem 3.2. It allows us to relate  $\text{adj}(zI - M_\sigma)$  with the adjugate of an  $(n-1) \times (n-1)$  matrix obtained by deflating  $zI - M_\sigma$  in a certain way. In the following, a matrix polynomial  $P(z) \in \mathbb{C}^{n \times n}[z]$  is said to be *unimodular* if  $\det P(z)$  is a nonzero constant. In other words,  $P(z)$  has a polynomial inverse.

LEMMA 3.3 Let  $p(z) = z^n + \sum_{k=0}^{n-1} a_k z^k$ , let  $\sigma : \{0, 1, \dots, n-1\} \rightarrow \{1, \dots, n\}$  be a bijection with  $\text{PCIS}(\sigma) = (v_0, v_1, \dots, v_{n-2})$ , let  $M_\sigma$  be the Fiedler matrix of  $p(z)$  associated with  $\sigma$ , and define unimodular  $Q(z), R(z) \in \mathbb{C}^{n \times n}[z]$  as

$$Q(z) := \begin{bmatrix} 1 & 0 & & \\ z & 1 & & \\ & & I_{n-2} & \end{bmatrix} \quad \text{and} \quad R(z) := \begin{bmatrix} 0 & 1 & & \\ -1 & p_1(z) & & \\ & & I_{n-2} & \end{bmatrix}.$$

Then,

- (a) if  $\sigma$  has a consecution at  $n-2$ ,

$$Q(z)(zI_n - M_\sigma)R(z) = \begin{bmatrix} 1 & & \\ & zI_{n-1} - \tilde{M}_\rho & \end{bmatrix},$$

- (b) if  $\sigma$  has an inversion at  $n-2$ ,

$$R(z)^T(zI_n - M_\sigma)Q(z)^T = \begin{bmatrix} 1 & & \\ & zI_{n-1} - \tilde{M}_\rho & \end{bmatrix},$$

where  $\rho : \{0, 1, \dots, n-2\} \rightarrow \{1, \dots, n-1\}$  is a bijection such that  $\text{PCIS}(\rho) = (v_0, v_1, \dots, v_{n-3})$ , and  $\tilde{M}_\rho = \tilde{M}_{\rho^{-1}(1)} \tilde{M}_{\rho^{-1}(2)} \cdots \tilde{M}_{\rho^{-1}(n-1)}$ , with  $\tilde{M}_0 = \text{diag}(I_{n-2}, -a_0)$ , and

$$\tilde{M}_k = \begin{bmatrix} I_{n-k-2} & & & \\ & -a_k & 1 & \\ & 1 & 0 & \\ & & & I_{k-1} \end{bmatrix}, \quad \text{for } k = 1, 2, \dots, n-3, \quad \tilde{M}_{n-2} = \begin{bmatrix} -p_2(z) + z & 1 & & \\ 1 & 0 & & \\ & & & I_{n-3} \end{bmatrix}.$$

*Proof.* We only prove part (a) because part (b) is similar. So, let us assume that  $\sigma$  has a consecution at  $n-2$ . Then, using (2.3), the factors of  $M_\sigma$  can be rearranged so that  $M_\sigma = XM_{n-2}M_{n-1}Y$ , where  $X, Y$  are products of  $M_i$  matrices, with  $i < n-2$ . Now, since  $Q(z)$  and  $R(z)$  commute with  $M_i$ , for  $i < n-2$ , we have

$$\begin{aligned}
Q(z)(zI_n - M_\sigma)R(z) &= zQ(z)R(z) - XQ(z)M_{n-2}M_{n-1}R(z)Y \\
&= \begin{bmatrix} 0 & z & 0 \\ -z & z^2 + zp_1(z) & 0 \\ 0 & 0 & z \end{bmatrix} zI_{n-3} - X \begin{bmatrix} -1 & z & 0 \\ -z & z^2 - a_{n-2} & 1 \\ 0 & 1 & 0 \end{bmatrix} I_{n-3} Y \\
&= \begin{bmatrix} 0 & z & 0 \\ -z & z^2 & 0 \\ 0 & 0 & z \end{bmatrix} zI_{n-3} - X \left( \begin{bmatrix} -1 & z & 0 \\ -z & z^2 - z & 0 \\ 0 & 0 & 0 \end{bmatrix} 0_{n-3} + \begin{bmatrix} 0 & 0 & 0 \\ 0 & -p_2(z) + z & 1 \\ 0 & 1 & 0 \end{bmatrix} I_{n-3} \right) Y \\
&= \begin{bmatrix} 1 & & & \\ & z & & \\ & & z & \\ & & & zI_{n-3} \end{bmatrix} - X \begin{bmatrix} 0 & 0 & 0 \\ 0 & -p_2(z) + z & 1 \\ 0 & 1 & 0 \end{bmatrix} I_{n-3} Y \\
&= \begin{bmatrix} 1 & & & \\ & zI_{n-1} & & \\ & & 0 & \\ & & \tilde{M}_{\rho^{-1}(1)} \tilde{M}_{\rho^{-1}(2)} \cdots \tilde{M}_{\rho^{-1}(n-1)} & \end{bmatrix} = \begin{bmatrix} 1 & & & \\ & zI_{n-1} - \tilde{M}_\rho & & \end{bmatrix},
\end{aligned}$$

where we have used that  $p_2(z) = zp_1(z) + a_{n-2}$  and the fact that multiplying any matrix of the form  $\text{diag}(A, 0_{n-2})$ , with  $A \in \mathbb{C}^{2 \times 2}$ , by  $M_k$ , for  $k = 0, 1, \dots, n-3$ , keeps that matrix unchanged. Finally, note that the relative positions of  $\tilde{M}_0, \tilde{M}_1, \dots, \tilde{M}_{n-2}$  in  $\tilde{M}_\rho$  coincide with the ones of  $M_0, M_1, \dots, M_{n-2}$  in  $M_\sigma$ , so  $\text{PCIS}(\rho) = (v_0, v_1, \dots, v_{n-3})$ .  $\square$

REMARK 3.1 Some important observations about the matrix  $\tilde{M}_\rho$  in Lemma 3.3 are in order:

- (a) The matrix  $\tilde{M}_i$ , for  $i = 0, \dots, n-3$  is obtained from  $M_i$  by removing the first row and column.
- (b) The matrix  $\tilde{M}_\rho$  can be seen formally as a Fiedler matrix of the polynomial  $r(z) := z^{n-1} + \sum_{k=0}^{n-2} b_k z^k$ , where  $b_{n-2} = p_2(z) - z$  and  $b_k = a_k$  for  $k = 0, 1, \dots, n-3$ . Notice that  $r(z) = p(z)$  for all  $z \in \mathbb{C}$ . We also want to emphasize that the formal  $(n-2)$ th coefficient of  $r(z)$  is not an scalar, but a polynomial in  $z$ .
- (c) The formal Horner shifts of  $r(z)$  satisfy:  $r_0(z) = p_0(z) = 1$  and  $r_k(z) = p_{k+1}(z)$  for  $k = 1, 2, \dots, n-2$ .

Now, armed with Lemmas 3.2 and 3.3, we are in the position to prove Theorem 3.2.

*Proof.* (of **Theorem 3.2**) The proof proceeds by induction in  $n$ . For  $n = 2$  there are only two Fiedler matrices, namely the first and second Frobenius companion matrices. For these two matrices we have

$$\text{adj}(zI - C_2) = \text{adj} \left( \begin{bmatrix} a_1 + z & -1 \\ a_0 & z \end{bmatrix} \right) = \begin{bmatrix} z & 1 \\ -a_0 & a_1 + z \end{bmatrix} = \begin{bmatrix} 1 \\ p_1(z) \end{bmatrix} \begin{bmatrix} z & 1 \end{bmatrix} - p(z) \begin{bmatrix} 0 & 0 \\ 1 & 0 \end{bmatrix}$$

and  $\text{adj}(zI - C_1)$  is the transpose of  $\text{adj}(zI - C_2)$ . These are the matrices in the statement with  $\text{PCIS}(\sigma) = (1)$  and  $\text{PCIS}(\sigma) = (0)$ , respectively. Assume that the result is true for Fiedler matrices of size  $(n-1) \times (n-1)$ . To prove it for size  $n \times n$ , we assume that  $\sigma$  has a consecution at  $n-2$  (the proof when  $\sigma$  has an inversion at  $n-2$  is similar and we omit it). Then, from Lemma 3.3 (a), we have that

$$zI_n - M_\sigma = Q(z)^{-1} \begin{bmatrix} 1 & \\ & zI_{n-1} - \tilde{M}_\rho \end{bmatrix} R(z)^{-1},$$

therefore

$$\text{adj}(zI_n - M_\sigma) = \text{adj}(R(z)^{-1}) \text{adj} \left( \begin{bmatrix} 1 & \\ & zI_{n-1} - \tilde{M}_\rho \end{bmatrix} \right) \text{adj}(Q(z)^{-1}) = R(z) \begin{bmatrix} p(z) & \\ & \text{adj}(zI_{n-1} - \tilde{M}_\rho) \end{bmatrix} Q(z),$$

where we have used the identities  $\text{adj}(AB) = \text{adj}(B)\text{adj}(A)$ ,  $\det R(z) = \det Q(z) = 1$ , and  $\det(zI_{n-1} - \tilde{M}_\rho) = p(z)$ . By the induction hypothesis

$$\text{adj}(zI_n - M_\sigma) = R(z) \begin{bmatrix} p(z) & \\ & x_\rho y_\rho^T - p(z)A_\rho \end{bmatrix} Q(z) = R(z) \begin{bmatrix} 0 \\ x_\rho \end{bmatrix} \begin{bmatrix} 0 & y_\rho^T \end{bmatrix} Q(z) - p(z)R(z) \begin{bmatrix} -1 & \\ & A_\rho \end{bmatrix} Q(z).$$

Note that in the induction step we may see  $\tilde{M}_\rho$  as a Fiedler matrix associated with  $r(z) = z^{n-1} + \sum_{k=0}^{n-2} b_k z^k$ , with  $b_i$ , for  $i = 0, \dots, n-2$ , as in Remark 3.1, part (b). To finish the proof it suffices to prove the following three identities:

$$(i) \quad x_\sigma = R(z) \begin{bmatrix} 0 \\ x_\rho \end{bmatrix}, \quad (ii) \quad y_\sigma = Q^T(z) \begin{bmatrix} 0 \\ y_\rho \end{bmatrix}, \quad \text{and} \quad (iii) \quad A_\sigma = R(z) \begin{bmatrix} -1 & \\ & A_\rho \end{bmatrix} Q(z).$$

- (i) From the expressions of PCIS( $\sigma$ ) and PCIS( $\rho$ ) we have  $i_\rho(0 : k-1) = i_\sigma(0 : k-1)$ , for  $k = 1, 2, \dots, n-2$ . Also, the Horner shifts corresponding to  $\tilde{M}_\rho$  are  $p_0(z), p_2(z), \dots, p_{n-1}(z)$ . These observations imply that  $x_\rho(k) = x_\sigma(k+1)$ , for  $k = 2, 3, \dots, n-1$  (note that, for the permutation  $\rho$ ,  $n$  must be replaced by  $n-1$  in (3.4)). Therefore

$$R(z) \begin{bmatrix} 0 \\ x_\rho \end{bmatrix} = \begin{bmatrix} 0 & 1 \\ -1 & p_1(z) \\ & & I_{n-2} \end{bmatrix} \begin{bmatrix} 0 \\ z^{i_\rho(0:n-3)} \\ x_\rho(2:n-1) \end{bmatrix} = \begin{bmatrix} z^{i_\rho(0:n-3)} \\ z^{i_\rho(0:n-3)} p_1(z) \\ x_\rho(2:n-1) \end{bmatrix} = \begin{bmatrix} z^{i_\sigma(0:n-2)} p_0(z) \\ z^{i_\sigma(0:n-3)} p_1(z) \\ x_\sigma(3:n) \end{bmatrix} = x_\sigma,$$

where we have used, since  $v_{n-2} = 1$ , that  $i_\sigma(0 : n-3) = i_\sigma(0 : n-2)$  and  $p_0(z) = 1$ .

- (ii) This can be proved in a similar way as (i). We refer the reader to De Terán *et al.* (2014b) for more details.  
 (iii) We prove this using Lemma 3.2. From (i) and (ii) we know that

$$\text{adj}(zI - M_\sigma) = x_\sigma y_\sigma^T - p(z) R(z) \begin{bmatrix} -1 & \\ & A_\rho \end{bmatrix} Q(z).$$

But the entries of  $R(z) \text{diag}(-1, A_\rho) Q(z)$  are polynomials in  $z$  and, moreover, the entries of  $\text{adj}(zI - M_\sigma)$  are polynomials of degree less than or equal to  $n-1$ . By the uniqueness proved in Lemma 3.2, we get (iii).  $\square$

### 3.2 First-order perturbation of the coefficients of the polynomial $\det(zI - M_\sigma)$

In this section, we derive formulas, to first order in  $E$ , for the coefficients of the characteristic polynomial of  $M_\sigma + E$ , where  $E$  is an arbitrary dense matrix.

**THEOREM 3.3** Let  $p(z) = z^n + \sum_{k=0}^{n-1} a_k z^k$  be a monic polynomial, let  $\sigma : \{0, 1, \dots, n-1\} \rightarrow \{1, \dots, n\}$  be a bijection with EPCIS( $\sigma$ ) =  $(v_0, v_1, \dots, v_{n-1})$ , let  $M_\sigma$  be the Fiedler companion matrix of  $p(z)$  associated with  $\sigma$ , and let  $E \in \mathbb{C}^{n \times n}$  be an arbitrary matrix. If the characteristic polynomial of  $M_\sigma + E$  is denoted by  $\tilde{p}(z) = z^n + \sum_{k=0}^{n-1} \tilde{a}_k z^k$ , then, to first order in  $E$ ,

$$\tilde{a}_k - a_k = - \sum_{i,j=1}^n p_{ij}^{(\sigma,k)}(a_0, a_1, \dots, a_{n-1}) E_{ij}, \quad k = 0, 1, \dots, n-1, \quad (3.6)$$

where, for  $i, j = 1, 2, \dots, n$ , the function  $p_{ij}^{(\sigma,k)}(a_0, a_1, \dots, a_{n-1})$  is a multivariable polynomial in the coefficients of  $p(z)$ . More precisely,  $p_{ij}^{(\sigma,k)}(a_0, a_1, \dots, a_{n-1})$  is equal to:

- (a) if  $v_{n-i} = v_{n-j} = 0$ :

- $a_{k+i_\sigma(n-j:n-i)}$ , if  $j \geq i$  and  $n-k-i+1 \leq i_\sigma(n-j:n-i) \leq n-k$ ;
- $-a_{k+1-i_\sigma(n-i:n-j-1)}$ , if  $j < i$  and  $k+1+i-n \leq i_\sigma(n-i:n-j-1) \leq k+1$ ;
- 0, otherwise;

- (b) if  $v_{n-i} = v_{n-j} = 1$ :

- $a_{k+c_\sigma(n-i:n-j)}$ , if  $j \leq i$  and  $n-k-j+1 \leq c_\sigma(n-i:n-j) \leq n-k$ ;
- $-a_{k+1-c_\sigma(n-j:n-i-1)}$ , if  $j > i$  and  $k+1+j-n \leq c_\sigma(n-j:n-i-1) \leq k+1$ ;
- 0, otherwise;

- (c) if  $v_{n-i} = 1$  and  $v_{n-j} = 0$ :

- 1, if  $i_\sigma(0:n-j-1) + c_\sigma(0:n-i-1) = k$ ,
- 0, otherwise;

(d) if  $v_{n-i} = 0$  and  $v_{n-j} = 1$ :

- $$\sum_{l=\max\{0, k+1+j-c_\sigma(n-j:n-i-1)-n\}}^{l=\min\{k+1-c_\sigma(n-j:n-i-1), i-1\}} -(a_{n+1-i+l} a_{k+1-c_\sigma(n-j:n-i-1)-l}),$$

if  $j > i$  and  $k+2+j-i-n \leq c_\sigma(n-j:n-i-1) \leq k+1$ ;
- $$\sum_{l=\max\{0, k+1+i-i_\sigma(n-i:n-j-1)-n\}}^{l=\min\{k+1-i_\sigma(n-i:n-j-1), j-1\}} -(a_{n+1-j+l} a_{k+1-i_\sigma(n-i:n-j-1)-l}),$$

if  $j < i$  and  $k+2+i-j-n \leq i_\sigma(n-i:n-j-1) \leq k+1$ ;
- 0, otherwise;

where we set  $a_n := 1$ .

*Proof.* From Proposition 3.1, the coefficients of the characteristic polynomial of  $M_\sigma + E$  satisfy, to first order in  $E$ ,

$$\tilde{a}_k - a_k = - \sum_{i,j=1}^n P_{k+1}(j,i) E_{ij},$$

where  $P_{k+1}(j,i)$  is the  $(j,i)$  entry of  $P_{k+1}$  which, according to (3.2) is the  $k$ th matrix coefficient of the matrix polynomial  $\text{adj}(zI - M_\sigma)$ . Hence  $p_{ij}^{(\sigma,k)}(a_0, a_1, \dots, a_{n-1})$  is the  $k$ th coefficient of the  $(j,i)$  entry of  $\text{adj}(zI - M_\sigma)$ . From Theorem 3.2, we know an explicit expression for the  $(j,i)$  entry of  $\text{adj}(zI - M_\sigma)$ . By analyzing separately each case in the statement, it is straightforward to check that the  $k$ th coefficient of this entry coincides with the expression given in this theorem (see De Terán *et al.* (2014b) for more details).  $\square$

REMARK 3.2 According to the notation in (3.1), we have

$$\nabla a_k(M_\sigma) = - \left[ p_{11}^{(\sigma,k)} \cdots p_{n1}^{(\sigma,k)} p_{12}^{(\sigma,k)} \cdots p_{n2}^{(\sigma,k)} \cdots p_{1n}^{(\sigma,k)} \cdots p_{nn}^{(\sigma,k)} \right],$$

where we have dropped the dependence on  $a_0, \dots, a_{n-1}$  for brevity.

REMARK 3.3 For  $k = n-1$ , and  $\sigma$  an arbitrary bijection, a direct verification in Theorem 3.3 gives

$$p_{ij}^{(\sigma, n-1)}(a_0, \dots, a_{n-1}) = \begin{cases} 1 & \text{if } i = j \\ 0 & \text{otherwise} \end{cases}.$$

Then, for any Fiedler matrix  $M_\sigma$ , it follows from (3.6) that

$$a_{n-1}(M_\sigma + E) - a_{n-1}(M_\sigma) = - \sum_{i=1}^n E_{ii}.$$

But, since the  $(n-1)$ th coefficient of the characteristic polynomial of  $A$  is equal to  $-\text{tr}(A)$ , this is a restatement of the well-know identity:  $\text{tr}(M_\sigma + E) = \text{tr}(M_\sigma) + \text{tr}(E)$ .

We emphasize that  $p_{ij}^{(\sigma,k)}(a_0, a_1, \dots, a_{n-1})$  are always linear or quadratic polynomials in the coefficients  $a_0, \dots, a_{n-1}$ . They depend, at a first stage, on whether the bijection  $\sigma$  has a consecution or an inversion at  $n-i$  and  $n-j$ . In particular,  $p_{ij}^{(\sigma,k)}(a_0, a_1, \dots, a_{n-1})$  can only be quadratic when there is a consecution at  $n-j$  and an inversion at  $n-i$ .

COROLLARY 3.1 Let  $M_\sigma$  be  $C_1$  or  $C_2$  in the statement of Theorem 3.3. Then  $p_{ij}^{(\sigma,k)}(a_0, a_1, \dots, a_{n-1})$  in (3.6) is a polynomial of degree at most 1 in  $a_0, \dots, a_{n-1}$ , for all  $k = 0, 1, \dots, n-1$ , and all  $1 \leq i, j \leq n$ . For the remaining Fiedler matrices  $M_\sigma$ , there is always some  $k$  and some  $i, j$  such that  $p_{ij}^{(\sigma,k)}(a_0, a_1, \dots, a_{n-1})$  is a quadratic polynomial in  $a_0, a_1, \dots, a_{n-1}$ .

*Proof.* Let us first recall that the bijection associated with  $C_1$  is  $\sigma_1 = (\sigma_1(0), \sigma_1(1), \dots, \sigma_1(n-1)) = (n, n-1, \dots, 1)$ , whereas the bijection associated with  $C_2$  is  $\sigma_2 = (\sigma_2(0), \sigma_2(1), \dots, \sigma_2(n-1)) = (1, 2, \dots, n)$ . Hence,  $\sigma_1$  has no consecutions, whereas  $\sigma_2$  has no inversions.

Then, it remains to show that, if  $\sigma : \{0, 1, \dots, n-1\} \rightarrow \{1, \dots, n\}$  is a bijection having a consecution at  $n-j$  and an inversion at  $n-i$ , for some  $2 \leq i, j \leq n$ , then there is some  $0 \leq k \leq n-1$  such that  $p_{ij}^{(\sigma, k)}(a_0, a_1, \dots, a_{n-1})$  has degree 2. Note, first, that it must be  $i \neq j$ . Without loss of generality, let us assume that  $j > i$ . The proof for the case  $j < i$  is analogous. We need to prove that, in the sum defining  $p_{ij}^{(\sigma, k)}(a_0, a_1, \dots, a_{n-1})$  in the first bullet of case (d) in Theorem 3.3 there is at least one monomial  $a_r a_s$  such that  $0 \leq r, s \leq n-1$ . More precisely, we need to prove:

- (i) There is some  $0 \leq k \leq n-1$  such that  $k+2+j-i-n \leq c_\sigma(n-j : n-i-1) \leq k+1$ .
- (ii) There is some  $l$ , with  $\max\{0, k+1+j-c_\sigma(n-j : n-i-1)-n\} \leq l \leq \min\{k+1-c_\sigma(n-j : n-i-1), i-1\}$ , such that  $0 \leq n+1-i+l \leq n-1$  and  $0 \leq k+1-c_\sigma(n-j : n-i-1)-l \leq n-1$ .

For this, it suffices to take  $k = c_\sigma(n-j : n-i-1) - 1 = c_\sigma(n-j+1 : n-i-1)$  and  $l = 0$ . Note that (ii) is fulfilled for these values of  $k$  and  $l$ , because  $i \geq 2$ .  $\square$

The expressions given in Theorem 3.3 for the variation of the coefficients of the characteristic polynomial of  $M_\sigma$  are involved in general (that is, for arbitrary Fiedler matrices). We will show them explicitly in Section 3.2.2 for some particularly relevant Fiedler matrices, including the Frobenius companion matrices.

The following result, which is a direct consequence of Theorem 3.3 (see De Terán *et al.* (2014b)), describes one property of the polynomials  $p_{ij}^{(\sigma, k)}(a_0, \dots, a_{n-1})$  that will be used later.

LEMMA 3.4 Let  $p_{ij}^{(\sigma, k)}(a_0, a_1, \dots, a_{n-1})$  be the polynomial defined in (3.6), and set  $a_n = 1$ . Then:

- (a) For  $k = 0, 1, \dots, n-1$ ,

$$p_{ii}^{(\sigma, k)}(a_0, a_1, \dots, a_{n-1}) = \begin{cases} a_{k+1} & \text{if } i \geq n-k, \\ 0 & \text{if } i < n-k. \end{cases}$$

- (b) If  $\sigma$  has a consecution at  $n-2$ , then  $p_{12}^{(\sigma, 0)}(a_0, a_1, \dots, a_{n-1}) = -a_0$ , and if  $\sigma$  has an inversion at  $n-2$ , then  $p_{21}^{(\sigma, 0)}(a_0, a_1, \dots, a_{n-1}) = -a_0$ .

To identify those indices  $k$  for which  $\nabla a_k(M_\sigma)$  contains quadratic terms in  $a_0, \dots, a_{n-1}$  may be interesting in practice. The presence of such terms implies that the sensitivity of the coefficient  $a_k(M_\sigma)$  to perturbations of  $M_\sigma$  is quadratic in  $a_0, \dots, a_{n-1}$ , instead of linear. This implies in turn that, for large values of  $a_0, \dots, a_{n-1}$ , we can expect much larger changes after small perturbations in these coefficients than in the ones where  $\nabla a_k(M_\sigma)$  contains only linear terms. We have seen in Corollary 3.1 that, for all Fiedler matrices but the Frobenius ones, there is always at least one  $k$  such that  $\nabla a_k(M_\sigma)$  contains quadratic terms. Moreover, the proof of Corollary 3.1 tells us that if  $i, j$  are such that  $\sigma$  has a consecution at  $n-j$  and an inversion at  $n-i$ , and  $j > i$  (respectively,  $j < i$ ), then for  $k = c_\sigma(n-j+1 : n-i-1)$  (resp.,  $k = i_\sigma(n-i+1 : n-j-1)$ ) the gradient  $\nabla a_k(M_\sigma)$  contains quadratic terms. In particular, Lemma 3.5 states that, for all Fiedler matrices but the Frobenius ones,  $\nabla a_0(M_\sigma)$  contains always quadratic polynomials in  $a_0, \dots, a_{n-1}$ . Its proof is a direct consequence of Theorem 3.3 (see De Terán *et al.* (2014b) for more details).

LEMMA 3.5 Let  $p_{ij}^{(\sigma, k)}(a_0, a_1, \dots, a_{n-1})$  be the polynomial defined in (3.6), and let  $t \in \{0, 1, \dots, n-3\}$ .

- (a) If  $\text{PCIS}(\sigma) = (v_0, v_1, \dots, v_t = 1, v_{t+1} = 0, v_{t+2} = 0, \dots, v_{n-2} = 0)$  then

$$p_{2, n-t}^{(\sigma, 0)}(a_0, a_1, \dots, a_{n-1}) = -a_{n-1}a_0.$$

- (b) If  $\text{PCIS}(\sigma) = (v_0, v_1, \dots, v_t = 0, v_{t+1} = 1, v_{t+2} = 1, \dots, v_{n-2} = 1)$  then

$$p_{n-t, 2}^{(\sigma, 0)}(a_0, a_1, \dots, a_{n-1}) = -a_{n-1}a_0.$$

The main result, from the theoretical point of view, in this section is a direct consequence of Theorem 3.3.

COROLLARY 3.2 Let  $p(z) = z^n + \sum_{k=0}^{n-1} a_k z^k$  be a monic polynomial, and  $M_\sigma$  be a Fiedler companion matrix of  $p(z)$ . Assume that the roots of  $p(z)$  are computed as the eigenvalues of  $M_\sigma$  with a backward stable algorithm i. e., an algorithm that computes the exact eigenvalues of some matrix  $M_\sigma + E$ , with  $\|E\|_\infty = O(u)\|M_\sigma\|_\infty$ . Then the computed roots are the exact roots of a polynomial  $\tilde{p}(z)$  such that:

(a) If  $M_\sigma = C_1, C_2$ ,

$$\frac{\|\tilde{p} - p\|_\infty}{\|p\|_\infty} = O(u)\|p\|_\infty, \quad (3.7)$$

(b) if  $M_\sigma \neq C_1, C_2$ ,

$$\frac{\|\tilde{p} - p\|_\infty}{\|p\|_\infty} = O(u)\|p\|_\infty^2, \quad (3.8)$$

where  $u$  is the machine epsilon. In other words, the backward error of the computed roots  $\tilde{\lambda}_1, \dots, \tilde{\lambda}_n$  is

$$\eta_\infty(\tilde{\lambda}_1, \dots, \tilde{\lambda}_n) = \begin{cases} O(u)\|p\|_\infty, & \text{if } M_\sigma = C_1, C_2, \\ O(u)\|p\|_\infty^2, & \text{if } M_\sigma \neq C_1, C_2. \end{cases}$$

*Proof.* If the eigenvalues of  $M_\sigma$  are computed with a backward stable algorithm, the computed eigenvalues are the exact eigenvalues of a matrix  $M_\sigma + E$ , for some  $E \in \mathbb{C}^{n \times n}$  with  $\|E\|_\infty = O(u)\|M_\sigma\|_\infty$ . Thus, the computed eigenvalues are the exact roots of  $\tilde{p}(z) = z^n + \sum_{k=0}^{n-1} \tilde{a}_k z^k = \det(zI - M_\sigma - E)$ . From Theorem 3.3, to first order in  $E$ ,

$$\begin{aligned} |\tilde{a}_k - a_k| &= \left| \sum_{i,j=1}^n p_{ij}^{(\sigma,k)}(a_0, a_1, \dots, a_{n-1}) E_{ij} \right| \leq \sum_{i,j=1}^n \left| p_{ij}^{(\sigma,k)}(a_0, a_1, \dots, a_{n-1}) \right| \cdot |E_{ij}| \\ &\leq (\max_{1 \leq i, j \leq n} |E_{ij}|) \cdot \left( \sum_{i,j=1}^n |p_{ij}^{(\sigma,k)}(a_0, a_1, \dots, a_{n-1})| \right). \end{aligned} \quad (3.9)$$

Notice, also from Theorem 3.3, that the absolute value of every polynomial  $p_{ij}^{(\sigma,k)}(a_0, a_1, \dots, a_{n-1})$  is bounded by  $n\|p\|_\infty^2$  and that, by Corollary 3.1, the square in the norm of  $p$  is necessary in all Fiedler matrices except the Frobenius companion matrices, where it can be replaced by 1. Therefore,

$$\max_{k=0,1,\dots,n-1} |\tilde{a}_k - a_k| = \|\tilde{p} - p\|_\infty = O(u)\|M_\sigma\|_\infty\|p\|_\infty^2 = O(u)\|p\|_\infty^3,$$

using that  $\max_{i,j=1,2,\dots,n} |E_{ij}| = O(u)\|M_\sigma\|_\infty$  and  $\|M_\sigma\|_\infty = O(1)\|p\|_\infty$  (see (De Terán *et al.*, 2014a, Th. 3.3)).  $\square$

It is worth to remark that if the matrix  $E$  in the statement of Corollary 3.2 satisfies  $\|E\|_\infty = c(p)O(u)\|M_\sigma\|_\infty$ , with  $c(p)$  being some positive quantity depending on  $p(z)$  then, with the appropriate changes in (3.9), we could replace (3.7) and (3.8) by, respectively:

$$\frac{\|\tilde{p} - p\|_\infty}{\|p\|_\infty} = c(p)O(u)\|p\|_\infty \quad \text{and} \quad \frac{\|\tilde{p} - p\|_\infty}{\|p\|_\infty} = c(p)O(u)\|p\|_\infty^2.$$

Hence, even for eigensolvers whose backward stability can not be guaranteed (like the fast QR-like algorithms mentioned in the Introduction for the Frobenius companion matrix or those that can be applied to other Fiedler matrices) our developments allow us to provide backward error estimates for the polynomial root-finding problem using Fiedler companion matrices.

**3.2.1 Recursive formula for the derivatives of the coefficient of the characteristic polynomial.** In Section 3.1 we have given an explicit formula for the entries of  $\text{adj}(zI - M_\sigma)$ . The aim of this subsection is to provide, in Proposition 3.4, a recursive formula for the coefficients of  $\text{adj}(zI - A)$  when viewed as a matrix polynomial in  $z$ , for arbitrary  $A \in \mathbb{C}^{n \times n}$ . This is an interesting theoretical result that gives an alternative description of the coefficients of  $\text{adj}(zI - A)$  and, as a consequence of Lemma 3.1, of the gradient of the characteristic polynomial of  $A$ . It may also have a practical interest, as it provides a recursive construction of these coefficients. This construction is related to the Faddev-Leverrier method to compute the coefficients of the characteristic polynomial (Gantmacher, 1959, Ch. 4, §5).

**PROPOSITION 3.4** (Gantmacher, 1959, Ch. 4, §4) Let  $A \in \mathbb{C}^{n \times n}$  and let  $p(z) = z^n + \sum_{k=0}^{n-1} a_k z^k$  be the characteristic polynomial of  $A$ . Let the matrices  $A_1, A_2, \dots, A_n \in \mathbb{C}^{n \times n}$  be defined by the following recurrence relation

$$\begin{cases} A_n = I, & \text{and} \\ A_k = A \cdot A_{k+1} + a_k I, & \text{for } k = n-1, n-2, \dots, 1. \end{cases} \quad (3.10)$$

Then,

$$\text{adj}(zI - A) = \sum_{k=0}^{n-1} z^k A_{k+1}.$$

We note that, as a consequence of the recursive relations of the Horner shifts (2.6), the matrices  $A_k$  are the Horner shifts of  $p(z) = \det(zI - A)$  evaluated at  $A$ . More precisely:

$$A_k = p_{n-k}(A) = A^{n-k} + a_{n-1}A^{n-k-1} + \cdots + a_{k+1}A + a_k I.$$

With this in mind, Proposition 3.1 gives the following expression for the gradient of the  $k$ th coefficient of the characteristic polynomial of  $A$ :

$$\nabla a_k(A) = -[\text{vec}(p_{n-k-1}(A^T))]^T, \quad \text{for } k = 0, 1, \dots, n-1. \quad (3.11)$$

Proposition 3.4 has been used in Edelman & Murakami (1995) to get an explicit formula for the derivatives of the coefficients of  $\det(zI - C)$ , with  $C$  being a Frobenius companion matrix. For this, the authors take advantage of the explicit expression of the matrices  $A_k$  defined in (3.10) with  $A = C$ , which are very simple in this case (see (Edelman & Murakami, 1995, p. 768)). However, for  $A$  being an arbitrary Fiedler matrix, the matrices  $A_k$  become much more involved, and it is not easy to get an explicit expression of these matrices just using (3.10). For this reason, we have obtained the expression of the entries of  $\text{adj}(zI - A)$  by other means. However, Proposition 3.4 gives us an alternative way to get  $\text{adj}(zI - A)$  using the Horner shifts of  $A$ .

As a consequence of the previous remarks, the polynomial  $p_{ij}^{(\sigma,k)}(a_0, a_1, \dots, a_{n-1})$  in Theorem 3.3 corresponds to the  $(j, i)$  entry of  $p_{n-k-1}(M_\sigma)$ . In the following section, we display these matrices for some particular cases, including the Frobenius matrices. Corollary 3.1 implies that the first and second Frobenius matrices are the only Fiedler matrices  $M_\sigma$  for which all Horner shifts  $p_k(M_\sigma)$  have entries which are linear multivariable polynomials in the coefficients of  $p(z)$ . For all other Fiedler matrices  $M_\sigma$ , there is at least one  $k$  such that  $p_k(M_\sigma)$  contains some quadratic entries.

**3.2.2 Some particular cases.** We obtain in this section the explicit expression (3.6) for some particular Fiedler matrices. We start with the classical Frobenius companion matrices in Theorem 3.5, where we get analogous formulas to the ones obtained in Edelman & Murakami (1995). The proof is a direct consequence of Theorem 3.3, so we omit it. We refer the reader to De Terán *et al.* (2014b) for more details.

**THEOREM 3.5** Let  $p(z) = z^n + \sum_{k=0}^{n-1} a_k z^k$  be a monic polynomial of degree  $n$ , let  $C = C_1$  or  $C_2$  be the first or second Frobenius companion matrix of  $p(z)$ , and let  $E \in \mathbb{C}^{n \times n}$ . If  $\tilde{p}(z) = z^n + \sum_{k=0}^{n-1} \tilde{a}_k z^k$  is the characteristic polynomial of  $C + E$ . Then, to first order in  $E$ , for  $k = 0, 1, \dots, n-1$ :

(i) If  $C = C_1$ :

$$\tilde{a}_k - a_k = \sum_{s=0}^k \sum_{j=1}^{n-k-1} a_s E_{j-s+k+1, j} - \sum_{s=k+1}^n \sum_{j=n-k}^n a_s E_{j-s+k+1, j}. \quad (3.12)$$

(ii) If  $C = C_2$ :

$$\tilde{a}_k - a_k = \sum_{s=0}^k \sum_{i=1}^{n-k-1} a_s E_{i, i-s+k+1} - \sum_{s=k+1}^n \sum_{i=n-k}^n a_s E_{i, i-s+k+1}. \quad (3.13)$$

According to (3.11), the matrix  $p_{n-k-1}(A^T)$  encodes the information about  $\nabla a_k(A)$ . In particular, the  $(i, j)$  entry of  $p_{n-k-1}(A^T)$  is the coefficient of  $E_{ij}$  in (3.1). In the case of Frobenius companion matrices, these Horner shifts can be computed without too much effort, since they are equal to:

$$p_{n-k-1}(C_1^T) = p_{n-k-1}(C_2) = \left[ \begin{array}{ccc|ccc} 0 & \cdots & 0 & 1 & & 0 \\ -a_k & & & a_{n-1} & 1 & \\ \vdots & \ddots & & \vdots & a_{n-1} & \ddots \\ -a_1 & \ddots & -a_k & a_{k+1} & \vdots & \ddots & 1 \\ -a_0 & \ddots & \vdots & & a_{k+1} & \ddots & a_{n-1} \\ & \ddots & -a_1 & & & \ddots & \vdots \\ 0 & & -a_0 & 0 & & & a_{k+1} \end{array} \right], \quad \text{for } k = 0, 1, \dots, n-1, \quad (3.14)$$



where the first block-column contains  $n - k - 1$  columns, and the second block-column contains  $k + 1$  columns. The reader may check that the  $(i, j)$  entry of (3.14) is the coefficient of  $E_{ij}$  in (3.12). The same happens with the transpose of (3.14) and formula (3.13).

Excluding the Frobenius companion matrices, the simplest Fiedler matrices are those corresponding to bijections with just one inversion (resp., consecution) at 0, and consecutions (resp., inversions) elsewhere. These particular Fiedler matrices present several numerical advantages that may be of interest in new enhancements of the current codes for the Polynomial Eigenvalue Problem (like MATLAB's `polyeig`). To be precise, one of these matrices is

$$F = M_0(M_{n-1}M_{n-2}\cdots M_1) = \begin{bmatrix} -a_{n-1} & 1 & & & \\ -a_{n-2} & 0 & 1 & & \\ \vdots & & \ddots & \ddots & \\ -a_1 & & & 0 & -a_0 \\ 1 & 0 & \cdots & 0 & 0 \end{bmatrix},$$

and the other one is  $F^T$ . Theorem 3.6 is again a direct consequence of Theorem 3.3 (see De Terán *et al.* (2014b)).

**THEOREM 3.6** Let  $p(z) = z^n + \sum_{k=0}^{n-1} a_k z^k$  be a monic polynomial of degree  $n$ , let  $M_\sigma = F$  be the Fiedler companion matrix of  $p(z)$  with PCIS( $\sigma$ ) =  $(0, 1, 1, \dots, 1)$  and let  $E \in \mathbb{C}^{n \times n}$ . If  $\tilde{p}(z) = z^n + \sum_{k=0}^{n-1} \tilde{a}_k z^k$  is the characteristic polynomial of  $F + E$ , then, to first order in  $E$ ,

$$\begin{aligned} \tilde{a}_k - a_k = & \sum_{j=k+1}^{n-1} a_0 a_{n+k+1-j} E_{nj} + \sum_{s=0}^k \sum_{i=1}^{n-k-2} a_s E_{i,i+k+1-s} + \sum_{s=1}^k a_s E_{n-k-1,n-s} - E_{n-k-1,n} \\ & - \sum_{s=k+1}^n \sum_{i=n-k}^{n-1} a_s E_{i,i+k+1-s} - E_{n-k-1,n} - a_{k+1} E_{nn}. \end{aligned} \quad (3.15)$$

Theorem 3.6 illustrates how a single change in the PCIS of the Frobenius companion matrix (just the position of the factor  $M_0$  in the product defining  $C_1$  and  $C_2$ ) implies the appearance of quadratic terms in the coefficients of  $p(z)$  in the formula for the gradient of the coefficients of the characteristic polynomial (see the first summand in the right-hand-side of (3.15)). As before, this can also be seen by explicitly displaying the Horner shifts evaluated at  $F$ :

$$p_{n-k-1}(F) = \left[ \begin{array}{cccc|ccc} 0 & & & & 1 & & 0 \\ -a_k & & & & a_{n-1} & \ddots & \vdots \\ \vdots & \ddots & & & \vdots & \ddots & 1 \\ -a_1 & & -a_k & & a_{k+2} & a_{n-1} & -a_0 \\ -a_0 & \ddots & \vdots & -a_k & a_{k+1} & \ddots & \vdots \\ & \ddots & -a_1 & \vdots & & \ddots & a_{k+2} \\ & & -a_0 & -a_1 & & a_{k+1} & -a_0 a_{k+2} \\ & & & 1 & & & a_{k+1} \end{array} \right], \quad \text{for } k = 0, 1, \dots, n-3,$$

$$p_1(F) = \left[ \begin{array}{cccc|cc} 0 & & & & & 0 \\ -a_{n-2} & 1 & & & & \\ -a_{n-3} & a_{n-1} & 1 & & & \\ \vdots & & a_{n-1} & \ddots & & \\ \vdots & & & \ddots & 1 & \\ -a_1 & & & & a_{n-1} & -a_0 \\ 1 & & & & 0 & a_{n-1} \end{array} \right], \quad \text{and } p_0(F) = I.$$

The number of columns in the first block-column of  $p_{n-k-1}(F)$  is  $n - k - 1$ , and the number of columns in the second block column is  $k + 1$ . The reader may check that the  $(i, j)$  entry of  $p_{n-k-1}(F)^T$  is the coefficient of  $E_{ij}$  in (3.15).

Our last example is a pentadiagonal Fiedler matrix. For  $n \geq 3$ , there are four pentadiagonal matrices corresponding to bijections whose PCIS are  $(1, 0, 1, 0, \dots)$ ,  $(0, 1, 0, 1, \dots)$ ,  $(1, 1, 0, 1, 0, \dots)$ , and  $(0, 0, 1, 0, 1, 0, \dots)$  (see De Terán *et al.* (2010)). Formulas here become much more involved (see (De Terán *et al.*, 2014b, Th. 3.17)), and the corresponding matrices  $p_{n-k-1}(M_\sigma)$ , for  $k = 0, 1, \dots, n-1$ , do not have a simple structure. For illustrative purposes, we include here

a  $6 \times 6$  example. Let  $M_\sigma$  be the Fiedler companion matrix of the polynomial  $p(z) = z^6 + \sum_{k=0}^5 a_k z^k$  associated with a bijection  $\sigma$  such that  $\text{PCIS}(\sigma) = (1, 0, 1, 0, 1)$ . This matrix is

$$M_\sigma = \begin{bmatrix} -a_5 & 1 & 0 & 0 & 0 & 0 \\ -a_4 & 0 & -a_3 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & -a_2 & 0 & -a_1 & 1 \\ 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & -a_0 & 0 \end{bmatrix}.$$

Then, it can be seen that

$$\begin{aligned} p_0(M_\sigma) &= \begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \end{bmatrix}, & p_1(M_\sigma) &= \begin{bmatrix} 0 & 1 & 0 & 0 & 0 & 0 \\ -a_4 & a_5 & -a_3 & 1 & 0 & 0 \\ 1 & 0 & a_5 & 0 & 0 & 0 \\ 0 & 0 & -a_2 & a_5 & -a_1 & 1 \\ 0 & 0 & 1 & 0 & a_5 & 0 \\ 0 & 0 & 0 & 0 & -a_0 & a_5 \end{bmatrix}, \\ p_2(M_\sigma) &= \begin{bmatrix} 0 & 0 & -a_3 & 1 & 0 & 0 \\ -a_3 & 0 & -a_2 - a_3 a_5 & a_5 & -a_1 & 1 \\ 0 & 1 & a_4 & 0 & 0 & 0 \\ -a_2 & 0 & -a_1 - a_2 a_5 & a_4 & -a_0 - a_1 a_5 & a_5 \\ 1 & 0 & a_5 & 0 & a_4 & 0 \\ 0 & 0 & -a_0 & 0 & -a_0 a_5 & a_4 \end{bmatrix}, \\ p_3(M_\sigma) &= \begin{bmatrix} 0 & 0 & -a_2 & 0 & -a_1 & 1 \\ -a_2 & 0 & -a_1 - a_2 a_5 & 0 & -a_0 - a_1 a_5 & a_5 \\ 0 & 0 & 0 & 1 & 0 & 0 \\ -a_1 & -a_2 & -a_0 - a_1 a_5 - a_2 a_4 & a_3 & -a_0 a_5 - a_1 a_4 & a_4 \\ 0 & 1 & a_4 & 0 & a_3 & 0 \\ -a_0 & 0 & -a_0 a_5 & 0 & -a_0 a_4 & a_3 \end{bmatrix}, \\ p_4(M_\sigma) &= \begin{bmatrix} 0 & 0 & -a_1 & 0 & -a_0 & 0 \\ -a_1 & 0 & -a_0 - a_1 a_5 & 0 & -a_0 a_5 & 0 \\ 0 & 0 & 0 & 0 & -a_1 & 1 \\ -a_0 & -a_1 & -a_0 a_5 - a_1 a_4 & 0 & -a_0 a_4 - a_1 a_3 & a_3 \\ 0 & 0 & 0 & 1 & a_2 & 0 \\ 0 & -a_0 & -a_0 a_4 & 0 & -a_0 a_3 & a_2 \end{bmatrix}, \\ p_5(M_\sigma) &= \begin{bmatrix} 0 & 0 & -a_0 & 0 & 0 & 0 \\ -a_0 & 0 & -a_0 a_5 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & -a_0 & 0 \\ 0 & -a_0 & -a_0 a_4 & 0 & -a_0 a_3 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & -a_0 & -a_0 a_2 & a_1 \end{bmatrix}. \end{aligned}$$

Unlike the previous cases  $C_2$  and  $F$ , there does not seem to be a simple pattern for  $p_{n-k-1}(M_\sigma)$  for arbitrary  $n$ , with  $\sigma : \{0, 1, \dots, n-1\} \rightarrow \{1, \dots, n\}$  being the bijection such that  $\text{PCIS}(\sigma) = (1, 0, 1, 0, \dots)$ .

### 3.3 Balancing and backward errors

Balancing is a standard preprocessing technique for computing the eigenvalues of a given matrix  $A$ , which leads, very often, to more accurate results, especially when the entries of  $A$  have very different magnitudes (see Parlett & Reinsch (1969)). Actually, balancing is implemented by default as an initial step in the command `eig` for computing eigenvalues in MATLAB. Balancing consists of performing diagonal similarities  $DAD^{-1}$  (i. e., with  $D$  diagonal) to  $A$ , in order to reduce the norm of  $A$  by equilibrating as much as possible the  $\infty$ -norm of all rows and columns. In

addition, very frequently balancing reduces the eigenvalue condition numbers (see (Golub & Van Loan, 1996, §7.2.2)). We recall that we are not interested in the eigenvalue condition number, but in the condition number of the coefficients of the characteristic polynomial or, equivalently, in the backward error of the polynomial root-finding problem solved as an eigenvalue problem. However, our numerical experiments show that balancing is also, in general, a good strategy to reduce this backward error.

Balancing first computes in exact arithmetic a matrix  $DM_\sigma D^{-1}$ , which has the same characteristic polynomial as  $M_\sigma$ , namely  $p(z)$ . Then a backward stable algorithm is applied to compute the eigenvalues of  $DM_\sigma D^{-1}$ , so that we get the exact eigenvalues of  $DM_\sigma D^{-1} + \tilde{E}$ , with

$$\|\tilde{E}\| = O(u)\|DM_\sigma D^{-1}\|, \quad (3.16)$$

for some matrix norm  $\|\cdot\|$ . Now, we can get a crude formula like (3.6) for the change of the coefficients of the characteristic polynomial of  $DM_\sigma D^{-1}$  using the identity:

$$\det(zI - DM_\sigma D^{-1} - \tilde{E}) = \det(zI - M_\sigma - D^{-1}\tilde{E}D),$$

and applying Theorem 3.3 with the perturbation  $D^{-1}\tilde{E}D$  instead of  $E$ . In particular, following the arguments in the proof of Corollary 3.2, we get

$$|\tilde{a}_k - a_k| \leq n^2 \max_{1 \leq i, j \leq n} \left( \left| p_{ij}^{(\sigma, k)}(a_0, a_1, \dots, a_{n-1}) \frac{d_j}{d_i} \right| \right) \cdot \max_{1 \leq i, j \leq n} |\tilde{E}_{ij}|,$$

with  $\tilde{E}$  as in (3.16). In this way, we get a formula which provides an ‘‘a posteriori’’ (that is, once the diagonal parameters  $d_i$  are known) measure for the backward error of the polynomial root-finding problem using balanced Fiedler matrices.

Though the numerical experiments carried out in Section 4 indicate that balancing usually produces smaller backward errors, we see in Proposition 3.10 that, for any degree, there are infinitely many polynomials for which the condition numbers of all coefficients of the characteristic polynomial of any matrix  $DM_\sigma D^{-1}$  are large. This shows that, though in practice balancing Fiedler matrices may be a good strategy for the root-finding problem, there are polynomials, with any degree, for which the strategy does not lead to small backward errors.

### 3.4 Conditioning of the characteristic polynomial of a matrix $A$

The developments carried out at the beginning of this section are closely related to the conditioning of the characteristic polynomial of the matrix  $A$ . The condition number of the characteristic polynomial provides a measure of its sensitivity to perturbations of the matrix. As we have seen at the beginning of this section, this is in turn related with the gradient of the coefficients of the characteristic polynomial. In this subsection, we introduce the condition number (absolute and relative) for the coefficients of the characteristic polynomial, and we relate it with (the norm of) its gradient. In this way, we will see that the backward stability of the polynomial root-finding problem via eigenvalue methods is determined by the conditioning of the characteristic polynomial. In other words, the conditioning of the map from the matrix  $A$  to the coefficients of the (monomial basis) characteristic polynomial, that is, the absolute condition number of the vector functions  $a_k(A)$ .

Let us assume that the entries of the matrix  $E$  in (3.1) satisfy  $|E_{ij}| \leq \varepsilon \|\text{vec}(A)\|_\infty$ . Using Holder’s inequality  $|u^T v| \leq \|u^T\|_\infty \|v\|_\infty$  (with  $\| [u_1 \dots u_n] \|_\infty = |u_1| + \dots + |u_n|$ )<sup>1</sup>, from (3.1) we get, up to first order, the inequalities:

$$|a_k(A + E) - a_k(A)| \leq \|\nabla a_k(A)\|_\infty \cdot \|\text{vec}(E)\|_\infty \leq \varepsilon \|\nabla a_k(A)\|_\infty \cdot \|\text{vec}(A)\|_\infty. \quad (3.17)$$

It is straightforward to show that there exists a particular matrix  $E$  with  $\|\text{vec}(E)\|_\infty = \varepsilon \|\text{vec}(A)\|_\infty$  such that  $|\nabla a_k(A) \cdot \text{vec}(E)| = \|\nabla a_k(A)\|_\infty \|\text{vec}(E)\|_\infty$ . For this matrix the bound in (3.17) is attained to first order in  $\varepsilon$ . With this in mind, Proposition 3.7 immediately follows.

**PROPOSITION 3.7** Let  $A \in \mathbb{C}^{n \times n}$  and  $a_k : \mathbb{C}^{n^2} \rightarrow \mathbb{C}$  be the  $k$ th coefficient of the characteristic polynomial of  $X \in \mathbb{C}^{n \times n}$ , considered as a function of  $X$ . We define the condition numbers  $\kappa(a_k, A)$  and  $\kappa_{\text{rel}}(a_k, A)$  as

$$\kappa(a_k, A) := \limsup_{\varepsilon \rightarrow 0} \left\{ \frac{|a_k(A + E) - a_k(A)|}{\varepsilon} : \|\text{vec}(E)\|_\infty \leq \varepsilon \|\text{vec}(A)\|_\infty \right\} \quad (3.18)$$

<sup>1</sup>Note that, according to the definition of  $\|\cdot\|_\infty$  for  $m \times n$  matrices, see (Higham, 2002, p. 108), the expressions for  $\|u\|_\infty$  and  $\|u^T\|_\infty$  are different.

and

$$\kappa_{\text{rel}}(a_k, A) := \limsup_{\varepsilon \rightarrow 0} \left\{ \frac{|a_k(A+E) - a_k(A)|}{\varepsilon |a_k(A)|} : \|\text{vec}(E)\|_\infty \leq \varepsilon \|\text{vec}(A)\|_\infty \right\}. \quad (3.19)$$

Then

$$\kappa(a_k, A) = \|\nabla a_k(A)\|_\infty \cdot \|\text{vec}(A)\|_\infty \quad \text{and} \quad \kappa_{\text{rel}}(a_k, A) = \frac{\|\nabla a_k(A)\|_\infty \cdot \|\text{vec}(A)\|_\infty}{|a_k(A)|}.$$

The definition of condition number introduced in (3.18) and (3.19) may look non-standard, because of the inclusion of vectorizations. However, the presence of  $\text{vec}(E)$  is motivated by (3.1). We have included also  $\text{vec}(A)$  in the definition to make it more natural. Moreover, due to the identity

$$\|\text{vec}(M_\sigma)\|_\infty = \|p\|_\infty, \quad (3.20)$$

valid for any Fiedler matrix  $M_\sigma$ , this choice will allow us to get a simpler formula for  $\kappa(a_k, M_\sigma)$  (see (3.22) below).

Now, Proposition 3.7, together with (3.11), give us the following formulas for  $\kappa(a_k, A)$  and  $\kappa_{\text{rel}}(a_k, A)$ .

**COROLLARY 3.3** Let  $A \in \mathbb{C}^{n \times n}$  and let  $\kappa(a_k, A)$  and  $\kappa_{\text{rel}}(a_k, A)$  be the condition numbers defined in (3.18) and (3.19), respectively. Then, for  $k = 0, 1, \dots, n-1$ ,

$$\kappa(a_k, A) = \|\text{vec}(p_{n-k-1}(A))\|_1 \cdot \|\text{vec}(A)\|_\infty \quad \text{and} \quad \kappa_{\text{rel}}(a_k, A) = \frac{\|\text{vec}(p_{n-k-1}(A))\|_1 \cdot \|\text{vec}(A)\|_\infty}{|a_k(A)|}, \quad (3.21)$$

where  $p_{n-k-1}(z)$  is the degree  $n-k-1$  Horner shift of the polynomial  $p(z) := \det(zI - A)$ .

According to (3.21), the relative and absolute condition numbers depend on the norms of  $A$  and the degree  $n-k-1$  Horner shift of the characteristic polynomial of  $A$ . This Horner shift depends in turn on the coefficients  $a_{k+1}, \dots, a_{n-1}$  of the characteristic polynomial evaluated at  $A$ , namely:  $p_{n-k-1}(A) = A^{n-k-1} + a_{n-1}(A)A^{n-k-2} + \dots + a_{k+1}(A)I$ .

In particular, when  $A = M_\sigma$ , formula (3.21) together with Theorem 3.3 and (3.20), give

$$\kappa(a_k, M_\sigma) = \|p\|_\infty \sum_{i,j=1}^n |p_{ij}^{(\sigma,k)}(a_0, a_1, \dots, a_{n-1})|, \quad (3.22)$$

where  $p_{ij}^{(\sigma,k)}(a_0, a_1, \dots, a_{n-1})$  are given in Theorem 3.3, and they are polynomials of degree at most 2 in  $a_0, \dots, a_{n-1}$ .

By considering the maximum condition numbers of all coefficients  $a_k$  we arrive to the following notion.

**DEFINITION 3.8** Let  $A \in \mathbb{C}^{n \times n}$  and set  $p(z) = \det(zI - A)$ . Let  $\kappa(a_k, A)$  and  $\kappa_{\text{rel}}(a_k, A)$  be the condition numbers defined in (3.18) and (3.19), respectively. We define the condition number and the relative condition number of the characteristic polynomial of  $A$  with respect to perturbations of  $A$  as

$$\kappa(p, A) = \max_{k=0,1,\dots,n-1} \kappa(a_k, A) \quad \text{and} \quad \kappa_{\text{rel}}(p, A) = \max_{k=0,1,\dots,n-1} \kappa_{\text{rel}}(a_k, A). \quad (3.23)$$

The following result provides bounds for the absolute and relative condition numbers of the characteristic polynomial when  $A$  is a Fiedler matrix.

**PROPOSITION 3.9** Let  $p(z) = z^n + \sum_{k=0}^{n-1} a_k z^k$  be a monic polynomial, and  $M_\sigma$  be a Fiedler companion matrix of  $p(z)$ . Let  $\kappa(p, M_\sigma)$  and  $\kappa_{\text{rel}}(p, M_\sigma)$  be as in (3.23). Then,

$$\|p\|_\infty^2 \leq \kappa(p, M_\sigma) \leq n^3 \|p\|_\infty^3 \quad \text{and} \quad \frac{\|p\|_\infty^2}{\max\{|a_0|, |a_1|, \dots, |a_{n-1}|\}} \leq \kappa_{\text{rel}}(p, M_\sigma) \leq \frac{n^3 \|p\|_\infty^3}{\min\{|a_0|, |a_1|, \dots, |a_{n-1}|\}}.$$

Moreover, if  $C = C_1, C_2$  denotes both the first and second Frobenius companion matrices, then

$$\|p\|_\infty^2 \leq \kappa(p, C) \leq n^3 \|p\|_\infty^2 \quad \text{and} \quad \frac{\|p\|_\infty^2}{\max\{|a_0|, |a_1|, \dots, |a_{n-1}|\}} \leq \kappa_{\text{rel}}(p, C) \leq \frac{n^3 \|p\|_\infty^2}{\min\{|a_0|, |a_1|, \dots, |a_{n-1}|\}}.$$

*Proof.* The bound  $\kappa(a_k, M_\sigma) \leq n^3 \|p\|_\infty^3$  follows immediately from (3.22) and the bound  $|p_{ij}^{(\sigma,k)}(a_0, a_1, \dots, a_{n-1})| \leq n \|p\|_\infty^2$  (see Corollary 3.1), valid for all  $i, j = 1, \dots, n$ .

From (3.22) and Lemma 3.4, it follows that  $\kappa(a_k, M_\sigma) \geq (k+1)|a_{k+1}| \cdot \|p\|_\infty$ , for  $k = 0, 1, \dots, n-1$ , and  $\kappa(a_0, M_\sigma) \geq |a_0| \cdot \|p\|_\infty$ . Therefore

$$\kappa(p, M_\sigma) = \max_{k=0,1,\dots,n-1} \kappa(a_k, M_\sigma) \geq \|p\|_\infty^2.$$

Finally, from

$$\frac{\kappa(a_k, M_\sigma)}{\max\{|a_0|, |a_1|, \dots, |a_{n-1}|\}} \leq \frac{\kappa(a_k, M_\sigma)}{|a_k|} \leq \frac{\kappa(a_k, M_\sigma)}{\min\{|a_0|, |a_1|, \dots, |a_{n-1}|\}}$$

we get the bounds for  $\kappa_{\text{rel}}(p, M_\sigma)$  in the statement.

For the Frobenius companion matrices, Corollary 3.1 gives  $|p_{ij}^{(\sigma,k)}(a_0, a_1, \dots, a_{n-1})| \leq n\|p\|_\infty$ , where  $\sigma$  is the permutation corresponding to either the first or the second Frobenius companion matrix.  $\square$

**REMARK 3.4** The factor  $n^3$  appearing in all upper bounds in Proposition 3.9 usually overestimates the condition numbers. It is due to an  $n^2$  factor coming from the maximum possible number of nonzero polynomials  $p_{ij}^{(\sigma,k)}(a_0, a_1, \dots, a_{n-1})$  in the sum of the right-hand side in (3.22). This number is usually much less than  $n^2$ . For instance, it is equal to  $(k+1)(2n-2k-1)$  for the first and second Frobenius companion matrices, as can be seen from (3.14). It is also  $(k+1)(2n-2k-1)$  for the coefficients  $a_k$  with  $k = 2, \dots, n-1$ , equal to  $3n-4$  for  $a_1$  and equal to  $n$  for  $a_0$ , for the Fiedler matrix  $F$  in Theorem 3.6, as can be seen by looking at the matrices  $p_{n-k-1}(F)$  in Section 3.2.2.

**3.4.1 Balancing and condition numbers.** Though similar matrices have the same characteristic polynomial, the sensitivity of its coefficients may be quite different. In other words, the condition numbers  $\kappa(a_k, A)$  and  $\kappa_{\text{rel}}(a_k, A)$  defined in (3.18) and (3.19) are not invariant under diagonal similarity. Since  $q(SAS^{-1}) = Sq(A)S^{-1}$ , for any polynomial  $q(z)$  and any invertible matrix  $S$ , formula (3.21) gives

$$\kappa(a_k, SAS^{-1}) = \|\text{vec}(Sp_{n-k-1}(A)S^{-1})\|_1 \|\text{vec}(SAS^{-1})\|_\infty \quad (3.24)$$

and

$$\kappa_{\text{rel}}(a_k, SAS^{-1}) = \frac{\|\text{vec}(Sp_{n-k-1}(A)S^{-1})\|_1 \|\text{vec}(SAS^{-1})\|_\infty}{|a_k(A)|}.$$

The norms of the vectors in the right hand side of the previous expression can be quite different for different matrices  $S$ . The optimal balancing for a given  $A$  (or, equivalently, a given polynomial  $p(z) = \det(zI - A)$ ) from the point of view of the sensitivity of the characteristic polynomial (or, equivalently, from the point of view of backward errors of the root-finding problem via eigenvalue methods) would be given by some nonsingular diagonal matrix  $D$  such that  $\kappa_{\text{rel}}(p, DAD^{-1})$  is minimal among all nonsingular diagonal matrices  $D$  (see Parlett & Reinsch (1969) for the eigenvalue problem). In the case of Fiedler matrices, the following result provides a lower bound for this minimal conditioning.

**PROPOSITION 3.10** . Let  $p(z) = z^n + \sum_{k=0}^{n-1} a_k z^k$  be a monic polynomial and set  $a_n = 1$ . Let  $\sigma : \{0, 1, \dots, n-1\} \rightarrow \{1, \dots, n\}$  be a bijection, let  $M_\sigma$  be the Fiedler companion matrix of  $p(z)$  associated with  $\sigma$  and let  $D \in \mathbb{C}^{n \times n}$  be a diagonal and nonsingular matrix. Then, for  $k = 0, 1, \dots, n-1$ ,

$$\kappa(a_k, DM_\sigma D^{-1}) \geq (k+1)|a_{n-1}| \cdot |a_{k+1}| \quad \text{and} \quad \kappa_{\text{rel}}(a_k, DM_\sigma D^{-1}) \geq \frac{(k+1)|a_{n-1}| \cdot |a_{k+1}|}{|a_k|}.$$

*Proof.* We prove the result for  $\kappa(a_k, DM_\sigma D^{-1})$ , since the bound for the relative condition number can be obtained just dividing by  $|a_k|$ . The result is a consequence of the fact that diagonal similarity does not change the diagonal entries of a matrix. From (3.24),

$$\kappa(a_k, DM_\sigma D^{-1}) \geq \|\text{diag}(Dp_{n-k-1}(M_\sigma)D^{-1})\|_1 \cdot \|\text{diag}(DM_\sigma D^{-1})\|_\infty = \|\text{diag}(p_{n-k-1}(M_\sigma))\|_1 \cdot \|\text{diag}(M_\sigma)\|_\infty.$$

Now we prove that  $\text{diag}(M_\sigma) = (-a_{n-1}, 0, \dots, 0)$  and  $\text{diag}(p_{n-k-1}(M_\sigma)) = (0, \dots, 0, a_{k+1}, \dots, a_{k+1})$ , where the coefficient  $a_{k+1}$  appears  $(k+1)$  times.

For the diagonal of  $M_\sigma$  the proof proceeds by induction in  $n$ . The case  $n = 2$  is immediate, since the only possible  $M_\sigma$  are  $\begin{bmatrix} -a_1 & -a_0 \\ 1 & 0 \end{bmatrix}$  and  $\begin{bmatrix} -a_1 & 1 \\ -a_0 & 0 \end{bmatrix}$ . We assume that the identity is true for Fiedler matrices associated with polynomials of degree  $n-1$ . For degree  $n$ , we assume that  $\sigma$  has a consecution at  $n-2$  (the case where  $\sigma$  has an inversion at  $n-2$  is similar). Then, using MATLAB notation for columns and rows,  $M_\sigma$  may be written as,

$$M_\sigma = \begin{bmatrix} -a_{n-1} & 1 & 0 \\ W(:, 1) & 0 & W(:, 2:n-1) \end{bmatrix},$$

where  $W \in \mathbb{C}^{(n-1) \times (n-1)}$  is a Fiedler companion matrix of the polynomial  $z^{n-1} + \sum_{k=0}^{n-2} a_k z^k$  (see (De Terán *et al.*, 2013, p. 949)). Therefore,  $\text{diag}(M_\sigma) = (-a_{n-1}, 0, W(2,2), W(3,3), \dots, W(n-1, n-1)) = (-a_{n-1}, 0, \dots, 0)$ , by induction.

From Lemma 3.4 and equation (3.11), the  $(i, i)$  entry of  $p_{n-k-1}(M_\sigma)$  is equal to  $p_{ii}^{(\sigma, k)}(a_0, a_1, \dots, a_{n-1}) = a_{k+1}$ , if  $n-1 \geq k \geq n-i$  (that is,  $i \geq n-k$ ), and  $p_{ii}^{(\sigma, k)}(a_0, a_1, \dots, a_{n-1}) = 0$ , otherwise. This concludes the proof.  $\square$

### 3.5 Backward stability for $\|p\|_\infty$ moderate and coefficientwise backward stability

Corollary 3.2 indicates that computing the roots of scalar polynomials as the eigenvalues of an arbitrary Fiedler matrix is not backward stable if  $\|p\|_\infty$  is large, even if we compute the eigenvalues using a backward stable algorithm. This is revealed by the presence of the factor  $\|p\|_\infty$  in (3.7) and  $\|p\|_\infty^2$  in (3.8). However, when  $\|p\|_\infty$  is moderate, (3.8) guarantees backward stability. This fact is in accordance with results in (Van Dooren & Dewilde, 1983, p. 576), where it is proven that solving matrix *Polynomial Eigenvalue Problems* by applying the QZ algorithm to the Frobenius companion matrix is backward stable, provided that the original matrix polynomial has been previously scaled so that all coefficients have norm less than or equal to 1. For scalar polynomials (not necessarily monic), this condition can be always achieved by dividing all coefficients of the original polynomial  $p(z)$  by some sufficiently large number. However, if we want to restrict ourselves to the set of monic polynomials to use the QR algorithm, this is not a valid strategy any more, since we could get a non-monic polynomial after dividing the coefficients of  $p(z)$  (monic). To keep the polynomial  $p(z)$  in (1.1) within the set of monic polynomials, we can consider another kind of scaling, like:

$$\widehat{p}(z) := \alpha^n p(z/\alpha) = z^n + \sum_{k=0}^{n-1} a_k \alpha^{n-k} z^k.$$

Now,  $\alpha$  can be chosen so that  $|a_k \alpha^{n-k}| \leq 1$ , for all  $k = 0, 1, \dots, n-1$ . The roots of  $p(z)$  can be easily recovered from those of  $\widehat{p}(z)$  just dividing by  $\alpha$ . Once all coefficients of  $\widehat{p}(z)$  have absolute value less than or equal to 1, we can apply the QR algorithm to any Fiedler companion matrix of  $\widehat{p}(z)$  to get its roots, and then recover the roots of  $p(z)$ . However, this does not guarantee that the method is backward stable. It is not difficult to find examples of quadratic polynomials  $p(z)$  such that there is a polynomial  $\widehat{q}(z)$  with  $\|\widehat{p} - \widehat{q}\| = O(u)\|\widehat{p}\|$ , but  $\|p - q\|/\|p\|$  is  $O(1)$ , with  $q(z) = (1/\alpha^2)\widehat{q}(\alpha z)$ .

We want to emphasize that we are not considering in this paper the backward errors of single roots of  $p$ , but the backward error of the set of all roots of  $p$ . Backward errors of single roots has been considered in Tisseur (2000) for the more general case of matrix Polynomial Eigenvalue Problems. In particular, the backward error of a single computed root  $\tilde{\lambda}$  considered in Tisseur (2000) is:

$$\eta(\tilde{\lambda}) = \min \left\{ \varepsilon : (p + \Delta p)(\tilde{\lambda}) = 0, \quad |\Delta a_i| \leq \varepsilon |a_i|, \quad i = 0, 1, \dots, n \right\},$$

where  $p(z) = \sum_{k=0}^n a_k z^k$ , and  $\Delta p(z) = \sum_{k=0}^n (\Delta a_k) z^k$  are not necessarily monic. It is shown in (Tisseur, 2000, Theorem 7) that, for quadratic matrix polynomials whose coefficients have 2-norm equal to 1, computing the eigenvalues of its companion pencil (defined in (Tisseur, 2000, p. 347)) with a backward stable algorithm gives a coefficientwise backward stable method for the Quadratic Eigenvalue Problem. Though we are considering different notions of backward error, this fact seems to be in accordance with Corollary 3.2 when  $\|p\|_\infty = 1$  and with the discussion right below.

We also emphasize that the backward stability of polynomial root-finding when  $\|p\|_\infty = 1$  does not guarantee small relative backward errors in each coefficient. In other words, we can not guarantee that

$$\max_{k=0,1,\dots,n-1} \frac{|\tilde{a}_k - a_k|}{|a_k|} = O(u) \tag{3.25}$$

even in the case  $\|p\|_\infty = 1$ . In Section 4 we show some numerical experiments where  $\|p\|_\infty = 1$  and (3.25) does not hold. However, when  $|a_k|$  is moderate, for all  $k = 0, 1, \dots, n-1$ , and not too close to zero (loosely speaking, of order  $\Theta(1)$ ), then (3.7)–(3.8) imply that (3.25) holds, also in accordance with Tisseur (2000).

## 4. Numerical experiments

In this section we provide numerical experiments that support our theoretical results. Our goals are: (i) to show whether or not the bounds in (3.7)–(3.8) correctly predicts the dependence on the norm of  $p(z)$  of the largest backward error that may be obtained if the roots of  $p(z)$  are computed as the eigenvalues of a Fiedler matrix with a backward stable eigenvalue algorithm; (ii) to show that if the roots of a polynomial  $p(z)$ , with moderate coefficients, are computed as

the eigenvalues of a Fiedler matrix, then this process is normwise backward stable, regardless of the Fiedler matrix that is used, which implies that, in this situation, any Fiedler matrix can be used for the root-finding problem with the same reliability as the Frobenius companion matrices; (iii) to investigate, from the point of view of backward errors, the effect of balancing Fiedler matrices; and (iv) following Edelman & Murakami (1995), to show that Theorem 3.3 may be used to predict the backward error when the roots of a monic polynomial are computed as the eigenvalues of a Fiedler matrix. Along this section we denote by  $u = 2^{-52}$  the machine epsilon in IEEE double precision arithmetic.

Given a monic polynomial  $p(z)$  of degree  $n$ , we denote by  $\{\tilde{\lambda}_1, \tilde{\lambda}_2, \dots, \tilde{\lambda}_n\}$  the roots of  $p(z)$  computed as eigenvalues of a Fiedler matrix  $M_\sigma$  using a backward stable eigenvalue algorithm. If  $\tilde{p}(z)$  denotes the monic polynomial of degree  $n$  whose roots are  $\{\tilde{\lambda}_1, \tilde{\lambda}_2, \dots, \tilde{\lambda}_n\}$ , namely,  $\tilde{p}(z) = \prod_{k=0}^n (z - \tilde{\lambda}_k) = z^n + \sum_{k=0}^{n-1} \tilde{a}_k z^k$ , then we are interested in:

- the normwise backward error (NBE):  $\|\tilde{p} - p\|_\infty / \|p\|_\infty$ , and
- the coefficientwise backward error (CBE):  $\max_{k=0,1,\dots,n-1} (|\tilde{a}_k - a_k| / |a_k|)$ .

In the numerical experiments, we consider monic polynomials of degree 20 and the following Fiedler companion matrices associated with degree-20 polynomials:

- the second Frobenius companion matrix  $M_{\sigma_1} = C_2$ ,
- the Fiedler matrix  $M_{\sigma_2}$  with  $\text{PCIS}(\sigma_2) = (1, 0, 1, 0, 1, 0, 1, 0, 1, 0, 1, 0, 1, 0, 1, 0, 1, 0, 1, 0, 1)$ , which is a pentadiagonal matrix,
- the Fiedler matrix  $M_{\sigma_3}$  with  $\text{PCIS}(\sigma_3) = (0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1)$ , and
- the Fiedler matrix  $M_{\sigma_4}$  with  $\text{PCIS}(\sigma_4) = (1, 1, 1, 0, 0, 1, 0, 0, 1, 1, 0, 0, 0, 1, 1, 0, 1, 1, 1, 1)$ .

Recall that  $M_{\sigma_2}$  is the Fiedler matrix considered in the last example of Section 3.2.2, and  $M_{\sigma_3}$  is the Fiedler matrix in Theorem 3.6.

Given a monic polynomial  $p(z)$  of degree 20, to compute the polynomial  $\tilde{p}(z)$  we proceed as follows. First, we compute the eigenvalues of  $M_\sigma$  using the function `eig` in MATLAB (with and/or without balancing, see comments below); then, if  $\{\tilde{\lambda}_1, \tilde{\lambda}_2, \dots, \tilde{\lambda}_{20}\}$  denote the computed eigenvalues, we compute the polynomial  $\tilde{p}(z) = \prod_{k=1}^{20} (z - \tilde{\lambda}_k) = z^{20} + \sum_{k=0}^{19} \tilde{a}_k z^k$  using the function `vpa` (variable precision arithmetic) followed by the command `poly` on a diagonal matrix whose diagonal entries are  $\{\tilde{\lambda}_1, \tilde{\lambda}_2, \dots, \tilde{\lambda}_{20}\}$ , in MATLAB with 32 decimal digits of accuracy.

#### 4.1 Numerical experiments that show the dependence of the normwise backward error with $\|p\|_\infty$

In this subsection, we perform numerical experiments to determine whether or not the largest normwise backward errors that may be obtained if the roots of monic polynomials are computed as the eigenvalues of a Fiedler matrix  $M_\sigma$  with a backward stable eigenvalue algorithm, behave like  $\|\tilde{p} - p\|_\infty / \|p\|_\infty = O(u) \|p\|_\infty^2$ , when  $M_\sigma$  is a Fiedler matrix other than the Frobenius ones, or like  $\|\tilde{p} - p\|_\infty / \|p\|_\infty = O(u) \|p\|_\infty$ , when  $M_\sigma$  is one of the Frobenius companion matrices, as it is predicted by Corollary 3.2. We perform numerical experiments with and without balancing the Fiedler matrices. Our results show that if we do not balance the Fiedler matrices the bound in Corollary 3.2, although in a lot of cases is very pessimistic, predicts well the dependence with  $\|p\|_\infty$  of the largest backward errors. If the Fiedler matrices are balanced, our results show that there is still a dependence with  $\|p\|_\infty$  of the largest normwise backward errors, and that this dependence is similar for all Fiedler matrices. Also we show that the backward errors that are usually obtained when the Fiedler matrices are balanced are almost independent of the norm of the polynomials, and that polynomial root-finding algorithms using Fiedler matrices are usually normwise backward stable.

In order to see the dependence of the backward error with  $\|p\|_\infty$  we proceed as follows. For each  $k = 0, 1, \dots, 10$  we generate 500 random degree-20 polynomials with coefficients of the form  $a \cdot 10^c$ , where  $a$  is drawn from the uniform distribution on the interval  $[-1, 1]$  and  $c$  is drawn from the uniform distribution on  $[-k, k]$ . We set  $a_0 = 10^k$  to fix the infinity norm of the 500 random polynomials to be  $10^k$ . For each of these 11 samples of 500 random polynomials, we compute the normwise backward errors, as it is explained at the beginning of Section 4, when their roots are computed as the eigenvalues of the four Fiedler matrices  $M_{\sigma_1}, M_{\sigma_2}, M_{\sigma_3}, M_{\sigma_4}$ , with and without balancing them.

In Figures 1 (a)–(d) we plot the decimal logarithms of the maximum and the minimum normwise backward errors obtained for each of the 11 samples of 500 random polynomials against the logarithms of the norm of the polynomials, when their roots are computed as the eigenvalues of  $M_{\sigma_1}, M_{\sigma_2}, M_{\sigma_3}, M_{\sigma_4}$ , respectively, without balancing. We also plot a linear fitting for the logarithms of the maximum normwise backward errors to get the dependence with  $\|p\|_\infty$ . As

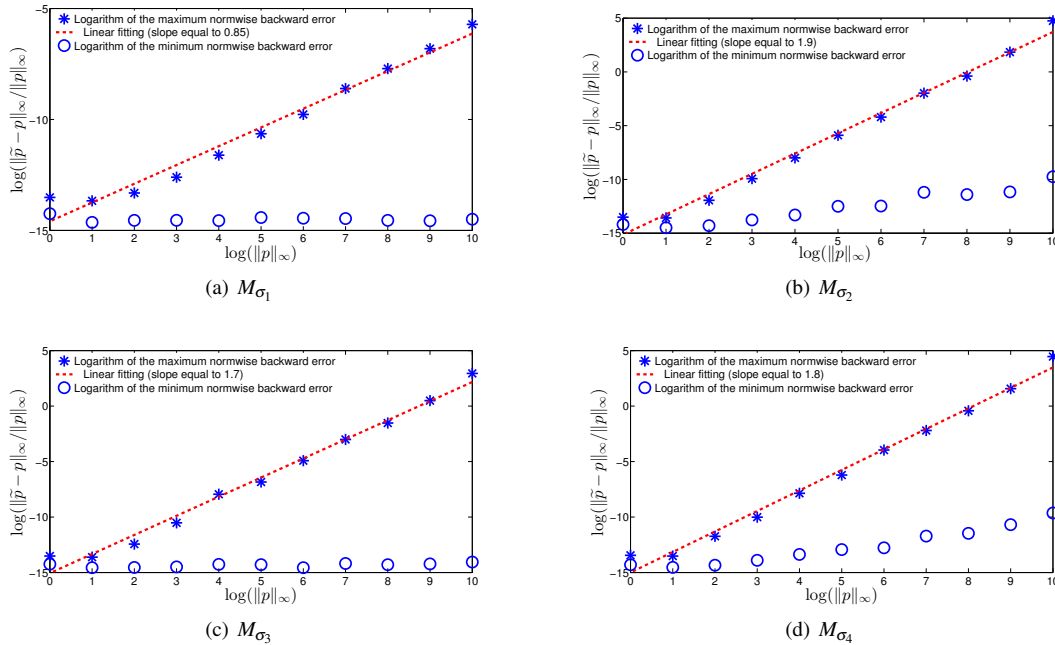


FIG. 1. Decimal logarithms of the maximum and minimum normwise backward errors obtained for each of the 11 samples of 500 random degree-20 polynomials, for  $k = 0, 1, \dots, 10$ , with a fixed infinite norm equal to  $10^k$  and with coefficients of the form  $a \cdot 10^c$ , where  $a$  is drawn from the uniform distribution on  $[-1, 1]$  and  $c$  is drawn from the uniform distribution on  $[-k, k]$ , and where we set  $a_0 = 10^k$ , when their roots are computed as the eigenvalues of the Fiedler matrices  $M_{\sigma_1}, M_{\sigma_2}, M_{\sigma_3}, M_{\sigma_4}$ , without balancing them.

may be seen in Figures 1 (a)–(d), there is a dependence with  $\|p\|_\infty$  of the largest normwise backward errors of the form  $\|p\|_\infty^\alpha$ . From the linear fittings we obtain  $\alpha = 0.85$  for  $M_{\sigma_1} = C_2$ ,  $\alpha = 1.9$  for  $M_{\sigma_2}$ ,  $\alpha = 1.7$  for  $M_{\sigma_3}$ , and  $\alpha = 1.8$  for  $M_{\sigma_4}$ . This is consistent with the bound in Corollary 3.2, which predicts  $\alpha = 1$  for the Frobenius companion matrices  $C_1$  and  $C_2$ , and  $\alpha = 2$  for Fiedler matrices other than the Frobenius ones. Also note that in Figures 1 (a)–(d) it may be seen that the bound in Corollary 3.2 is in some cases very pessimistic, since there are polynomials for which we get small normwise backward errors, regardless of their norms.

Next, we investigate the effect of balancing the Fiedler matrices in the backward errors. In Figures 2 (a)–(d), we plot the decimal logarithms of the maximum and the minimum normwise backward errors obtained for each of the 11 samples of 500 random polynomials against the logarithms of the norm of the polynomials, when their roots are computed as the eigenvalues of  $M_{\sigma_1}, M_{\sigma_2}, M_{\sigma_3}, M_{\sigma_4}$ , respectively, but in this case the Fiedler matrices are balanced before we compute their eigenvalues. As in the previous experiment, we plot a linear fitting for the logarithms of the maximum normwise backward errors in order to get the dependence with  $\|p\|_\infty$ . We also plot the ninth decile of the normwise backward error for each of the 11 samples. Figures 2 (a)–(d), show that there is a dependence of the largest backward errors with the norm of the polynomials of the form  $\|p\|_\infty^\alpha$ , but this dependence is similar for all four Fiedler matrices. In particular, from the linear fittings, we get  $\alpha = 0.59$  for  $M_{\sigma_1} = C_2$ ,  $\alpha = 0.71$  for  $M_{\sigma_2}$ ,  $\alpha = 0.67$  for  $M_{\sigma_3}$ , and  $\alpha = 0.71$  for  $M_{\sigma_4}$ . Note that 90% of the backward errors obtained when the roots of the polynomials are computed as the roots of  $M_{\sigma_1}, M_{\sigma_2}, M_{\sigma_3}, M_{\sigma_4}$  are excellent, since they are more or less between  $10^{-12}$  and  $10^{-16}$ , even for polynomials with norms as large as  $10^{10}$ .

#### 4.2 Numerical experiments with polynomials of moderate coefficients

We have obtained numerical evidence that supports what we claim in Section 3.5, namely, that computing the roots of a monic polynomial  $p(z)$  as in (1.1), with  $|a_i|$  moderate, for  $i = 0, 1, \dots, n-1$ , as the eigenvalues of a Fiedler matrix using a backward stable eigenvalue algorithm is normwise backward stable, regardless of the Fiedler matrix that is used. In addition, we show that to have  $|a_i|$  moderate, for  $i = 0, 1, \dots, n-1$ , is not enough to guarantee coefficientwise backward stability. Finally, we provide numerical evidence that supports the last sentence in Section 3.5, namely, that (3.25) holds when  $|a_i| = \Theta(1)$ , for  $i = 0, 1, \dots, n-1$ , regardless of the Fiedler matrix that is used. For this, we have run two sets of numerical experiments. Each set consists of random samples of 1000 degree-20



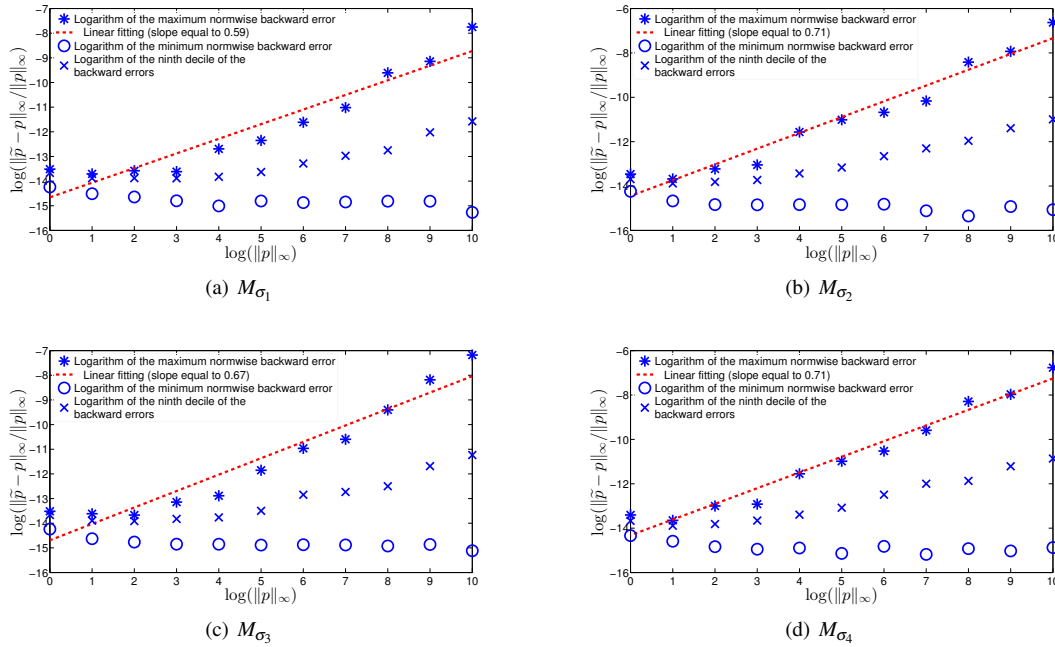


FIG. 2. Decimal logarithms of the maximum and minimum normwise backward errors obtained for the 11 samples of 500 random degree-20 polynomials with, for  $k = 0, 1, \dots, 10$ , a fixed infinite norm equal to  $10^k$  and with coefficients of the form  $a \cdot 10^c$ , where  $a$  is drawn from the uniform distribution on  $[-1, 1]$  and  $c$  is drawn from the uniform distribution on  $[-k, k]$ , and where we set  $a_0 = 10^k$ , when their roots are computed as the eigenvalues of the Fiedler matrices  $M_{\sigma_1}, M_{\sigma_2}, M_{\sigma_3}, M_{\sigma_4}$ , balancing them before computing their eigenvalues.

polynomials. In both cases, the coefficients of the polynomial have moderate norm (i. e., not too large). But in the first one, we set  $a_{19} = 10^{-10}$ , which is very close to zero. The numerical experiments in this case provide excellent backward error, but poor coefficientwise backward error. In the second set all coefficients have moderate norm not close to zero (their absolute value differ in at most four orders of magnitude). The experiments in this case give excellent normwise backward error and very good coefficientwise backward error. For more information on this, we refer the reader to De Terán *et al.* (2014b) (see, in particular, Tables 1 and 2).

### 4.3 Numerical experiments balancing Fiedler matrices

In this subsection we perform numerical experiments to study, from the point of view of backward errors, the effect of balancing Fiedler matrices. We show that, when a Fiedler matrix  $M_\sigma$  is balanced before computing its eigenvalues, the backward error obtained if we compute the roots of  $p(z)$  as the eigenvalues of  $M_\sigma$  may be much smaller than the backward error that is obtained when  $M_\sigma$  is not balanced, regardless of the Fiedler matrix that is used. We show also that balancing a Fiedler matrix is usually enough to guarantee that the process of computing the roots of a polynomial as the eigenvalues of a Fiedler matrix is normwise backward stable, even if the polynomial has large coefficients. Finally, we investigate the effect of the size of the coefficient  $a_{n-1}$ , since Proposition 3.10 suggests that it plays a key role in getting or not backward stability after balancing Fiedler matrices. To be precise, Proposition 3.10 shows that, for large values of  $|a_{n-1}|$ , the condition number of any coefficient of the characteristic polynomial of any Fiedler matrix will be large, regardless of the balancing. This leads us to expect large backward errors when  $|a_{n-1}|$  is large.

We consider a random sample of 1000 degree-20 polynomials with coefficients of the form

$$a_1 \cdot 10^{c_1} + ia_2 \cdot 10^{c_2}, \quad (4.1)$$

where  $i$  denotes the imaginary unit,  $a_1, a_2$  are drawn from the uniform distribution on the interval  $[-1, 1]$ , and  $c_1$  and  $c_2$  are drawn from the uniform distribution on  $[-10, 10]$ . These polynomials, considered in Toh & Trefethen (1994), allow us to measure the normwise backward errors with varying orders of magnitude in the coefficients of  $p(z)$ . We consider a second sample of 1000 degree-20 polynomials with coefficients of the form (4.1), but we fix  $a_{19} = 1$ .

For the first sample of random polynomials, in Tables 1-(a) and 1-(b) we give the mean, the maximum and the minimum of the decimal logarithms of the normwise backward errors (Log-Mean NBE, Log-Maximum NBE,

Log-Minimum NBE, respectively) obtained when the roots of the polynomials are computed as the eigenvalues of  $M_{\sigma_1}, M_{\sigma_2}, M_{\sigma_3}, M_{\sigma_4}$ , when these Fiedler matrices are not or are balanced, respectively.

(a) The Fiedler matrices are not balanced.

	$M_{\sigma_1}$	$M_{\sigma_2}$	$M_{\sigma_3}$	$M_{\sigma_4}$
Log-Mean NBE	-10.5	-2.4	-9.9	-3.0
Log-Maximum NBE	-5.8	3.2	0.1	3.5
Log-Minimum NBE	-14.7	-8.9	-14.7	-10.0

(b) The Fiedler matrices are balanced.

	$M_{\sigma_1}$	$M_{\sigma_2}$	$M_{\sigma_3}$	$M_{\sigma_4}$
Log-Mean NBE	-13.1	-13.1	-13.1	-12.9
Log-Maximum NBE	-8.1	-7.5	-8.0	-7.8
Log-Minimum NBE	-14.7	-14.9	-15.1	-14.8

Table 1. Mean, maximum, and minimum of the decimal logarithms of the normwise backward errors obtained for a sample of 1000 random degree-20 polynomials, with coefficients of the form (4.1), when their roots are computed as the eigenvalues of the Fiedler matrices  $M_{\sigma_1}, M_{\sigma_2}, M_{\sigma_3}, M_{\sigma_4}$ , without balancing and with balancing.

Several observations may be drawn from the data in Tables 1-(a) and 1-(b). First note, from the data in Log-Maximum NBE in Table 1-(a), that if the Fiedler matrices are not balanced, the backward errors may be very large. Note also that the largest of these backward errors is consistent with (3.7) for the Frobenius companion matrices, and with (3.8) for Fiedler matrices other than the Frobenius ones. Second, note that the process of balancing the Fiedler matrices makes that the backward errors after balancing may be much smaller than the backward errors obtained when the Fiedler matrices are not balanced (this is especially evident for  $M_{\sigma_2}$  and  $M_{\sigma_3}$ ). Finally, note, from the data in Log-Maximum NBE in Table 1-(b), that there are polynomials for which balancing the Fiedler matrices does not guarantee that the process of computing their roots as the eigenvalues of Fiedler matrices is normwise backward stable.

In Tables 2-(a) and 2-(b) we display the mean, the maximum and the minimum of the decimal logarithms of the normwise backward errors (Log-Mean NBE, Log-Maximum NBE, Log-Minimum NBE, respectively) that are obtained when the roots of the polynomials of the second sample are computed as the eigenvalues of the four Fiedler matrices  $M_{\sigma_1}, M_{\sigma_2}, M_{\sigma_3}, M_{\sigma_4}$ , when the Fiedler matrices are not or are balanced, respectively. Recall that for this sample of degree-20 random polynomials we set  $a_{19} = 1$ .

(a) The Fiedler matrices are not balanced.

	$M_{\sigma_1}$	$M_{\sigma_2}$	$M_{\sigma_3}$	$M_{\sigma_4}$
Log-Mean NBE	-6.9	-3.2	-6.9	-3.4
Log-Maximum NBE	-5.6	3.0	-3.4	3.0
Log-Minimum NBE	-9.8	-10.6	-9.9	-11.1

(b) The Fiedler matrices are balanced.

	$M_{\sigma_1}$	$M_{\sigma_2}$	$M_{\sigma_3}$	$M_{\sigma_4}$
Log-Mean NBE	-13.9	-13.9	-13.9	-13.7
Log-Maximum NBE	-11.6	-11.1	-11.6	-10.4
Log-Minimum NBE	-15.1	-14.8	-15.0	-15.0

Table 2. Mean, maximum, and minimum of the decimal logarithms of the normwise backward errors obtained when the roots of the polynomials of the second sample of random polynomials (i. e., coefficients from (4.1) and  $a_{19} = 1$ ) are computed as the eigenvalues of the four Fiedler matrices  $M_{\sigma_1}, M_{\sigma_2}, M_{\sigma_3}, M_{\sigma_4}$ , without balancing and with balancing.

As in the first sample of random polynomials, we may see in Tables 2-(a) and 2-(b) that the backward errors obtained when the Fiedler matrices are not balanced may be very large. Also, we may see that the backward errors may be much smaller after balancing the Fiedler matrices. Finally note that for this second sample the largest backward errors obtained when the Fiedler matrices are balanced are smaller than the largest ones obtained for the first sample.

#### 4.4 Using Theorem 3.3 to predict the coefficientwise backward error

Theorem 3.3 can be used to predict the coefficientwise backward error, without computing explicitly the polynomial  $\tilde{p}(z)$  (something that may not be possible for high degree polynomials, since using `vpa` makes this process very slow), and to show that this backward error is usually small for all Fiedler matrices if balancing is used. Of course, the normwise backward error can be also predicted from Theorem 3.3, but we omit it for brevity. We have checked this by exploring the same eight degree-20 polynomials as in Edelman & Murakami (1995) and Toh & Trefethen (1994).

As in Edelman & Murakami (1995), we first compute the coefficients exactly or with high precision using Mathematica. We then read these coefficients into MATLAB and take the rounded coefficients stored in MATLAB as our official test cases. Also, we consider again the four Fiedler companion matrices associated with degree-20 polynomials introduced at the beginning of Section 4, namely  $M_{\sigma_1}, M_{\sigma_2}, M_{\sigma_3}, M_{\sigma_4}$ .

We repeat the numerical experiments in Edelman & Murakami (1995). Our results show that Theorem 3.3 always predicts a small coefficientwise backward error, regardless of the Fiedler matrix that is used, and that this predicted backward error is usually pessimistic by at most one, two or three orders of magnitude, except for the monic polynomial with zeros  $2^{-10}, 2^{-9}, \dots, 2^8, 2^9$ , where the predicted backward error is pessimistic by 6 orders of magnitude. Note that in this case the ratio  $(|a_{19}| \cdot |a_1|)/|a_0|$  is of order  $2^{19}$ , so Proposition 3.10 ensures that the condition number for the coefficient  $a_0$  is large. However, the perturbations in the numerical experiments does not seem to affect this coefficient in such a severe way.

For more details on this example, we refer the reader to De Terán *et al.* (2014b) (see, in particular, Table 5).

### 5. The Sylvester space of Fiedler matrices

The study of the geometry of matrix spaces sheds light on the explanation of numerical processes involving matrices or matrix pencils. In particular, the theory of orbits has been used in the analysis of errors of the algorithms for computing eigenvalues and canonical forms (see Arnold (1971), Edelman *et al.* (1997, 1999) and Edelman & Murakami (1995)). In this section, and inspired by the motivating paper by Edelman & Murakami (1995), we analyze from a geometrical point of view the polynomial root-finding problem solved as an eigenvalue problem with Fiedler companion matrices. Our main result is Theorem 5.3, where we prove that the space of Sylvester matrices associated with a given Fiedler matrix  $M_\sigma$  is transversal to the similarity orbit of  $M_\sigma$ . This result extends the corresponding one for Frobenius companion matrices (Edelman & Murakami, 1995, Prop. 2.1).

Let  $p(z)$  be a monic polynomial as in (1.1) and let  $M_\sigma$  be a Fiedler matrix of  $p(z)$ . Let us consider the Euclidean matrix space  $\mathbb{C}^{n \times n}$  with the usual Frobenius inner product  $(A, B) = \text{tr}(AB^*)$ , where  $M^*$  denotes the conjugate transpose of  $M \in \mathbb{C}^{n \times n}$ . In this space, the set of matrices similar to a given matrix  $A \in \mathbb{C}^{n \times n}$  is a differentiable manifold in  $\mathbb{C}^{n \times n}$ . This manifold is the orbit of  $A$  under the action of similarity:  $\mathcal{O}(A) := \{SAS^{-1} : \det(S) \neq 0\}$ .

We will refer to the elements of a manifold as *points*, even though all manifolds considered in this paper are manifolds whose points are matrices.

It is known that the tangent space of  $\mathcal{O}(A)$  at  $A$  is the set

$$T_A \mathcal{O}(A) := \{AX - XA \text{ for some } X \in \mathbb{C}^{n \times n}\}.$$

The *normal space* of  $\mathcal{O}(A)$  at  $A$ , denoted by  $N_A \mathcal{O}(A)$ , is the set of matrices orthogonal to any matrix in  $T_A \mathcal{O}(A)$ :

$$N_A \mathcal{O}(A) := \{Y \in \mathbb{C}^{n \times n} \text{ such that } (Y, V) = 0, \text{ for all } V \in T_A \mathcal{O}(A)\},$$

and the *centralizer* of  $A$  is the set of matrices commuting with  $A$ :

$$C(A) := \{X \in \mathbb{C}^{n \times n} \text{ such that } AX - XA = 0\}$$

The following facts are already known:

- (a)  $C(A^*) = N_A \mathcal{O}(A)$  (see (Arnold, 1971, Lemma, p. 34)).
- (b) If  $A$  is a non-derogatory matrix, then:
  - (b1)  $C(A) = \{q(A) : q \text{ is a polynomial}\}$  (see (Horn & Johnson, 1985, Th. 3.2.4.2)).
  - (b2)  $\dim C(A) = n$  (see (Arnold, 1971, Corollary, p. 35)).
- (c)  $M_\sigma$  is a non-derogatory matrix, for all  $\sigma$ .

For claim (c), just recall that  $M_\sigma$  is similar to  $C_1$ , and that  $C_1$  is non-derogatory (see (Horn & Johnson, 1985, p. 147)).

As a consequence of claims (a)–(c) above, we have that  $\dim N_{M_\sigma} \mathcal{O}(M_\sigma) = n$ , for all  $\sigma$ , so there is a basis of  $N_{M_\sigma} \mathcal{O}(M_\sigma)$  consisting of  $n$  matrices which are polynomials in  $M_\sigma^*$ . This is stated in Proposition 5.1.

**PROPOSITION 5.1** Let  $p(z) = z^n + \sum_{k=0}^{n-1} a_k z^k$  be a monic polynomial,  $\sigma : \{0, 1, \dots, n-1\} \rightarrow \{1, \dots, n\}$  be a bijection,  $M_\sigma$  be the Fiedler matrix of  $p(z)$  associated with the bijection  $\sigma$ , and let  $p_d(z)$  be the  $d$ th Horner shift of  $p(z)$ , for  $d = 0, 1, \dots, n$ . Set  $p_0(M_\sigma) = I_n$  and

$$p_{n-k}(M_\sigma) = M_\sigma^{n-k} + a_{n-1} M_\sigma^{n-k-1} + \dots + a_{k+1} M_\sigma + a_k I, \quad \text{for } k = 1, \dots, n-1.$$

Then  $\{p_k(M_\sigma)^*\}_{k=0}^{n-1}$  is a basis for  $N_{M_\sigma} \mathcal{O}(M_\sigma)$ .

Note that the set  $\{p_k(M_\sigma)^*\}_{k=0}^{n-1}$  is linearly independent because, since  $M_\sigma$  is non-derogatory, its minimal polynomial coincides with its characteristic polynomial. Any  $n$  linearly independent polynomials in  $M_\sigma^*$  would serve as a basis for  $N_{M_\sigma} \mathcal{O}(M_\sigma)$ , but in Section 3.2.1 we have seen that the matrices  $p_k(M_\sigma)$  play an important role in determining how the coefficients of the characteristic polynomial of  $M_\sigma$  change when the matrix is perturbed (see (3.11)).

First order perturbations of the coefficients of  $p(z)$ , with  $p(z) = \det(zI - C_1)$ , have been studied in Edelman & Murakami (1995). To do so, the authors decompose the perturbation matrix  $E$  as

$$E = E^{\text{tan}} + E^{\text{syl}}, \quad (5.1)$$

where  $E^{\text{tan}}$  belongs to the tangent space to  $\mathcal{O}(C_1)$  at  $C_1$  and  $E^{\text{syl}}$  is of the form

$$E^{\text{syl}} = \begin{bmatrix} E_{11} & \dots & E_{1n} \\ 0 & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & 0 \end{bmatrix}.$$

The matrix  $E^{\text{syl}}$  belongs to the tangent space (at any point) to the *Sylvester space* of  $C_1$ . We recall that the (affine) Sylvester space of  $C_1$  is the set of all matrices of the form

$$\begin{bmatrix} E_{11} & E_{12} & \dots & E_{1n} \\ 1 & 0 & \dots & 0 \\ & \ddots & \ddots & \vdots \\ & & & 1 & 0 \end{bmatrix},$$

that is, the set of “all first Frobenius companion matrices”<sup>2</sup>. It may be proved that, to first order in  $E$ , the matrix  $E^{\text{tan}}$  does not affect the coefficients of  $p(z)$ . Below, we prove an equivalent result for any Fiedler matrix  $M_\sigma$ . For this, we first define the Sylvester space of any Fiedler matrix, which is a natural generalization of the Sylvester space of  $C_1$ .

**DEFINITION 5.2** (Sylvester space of a Fiedler matrix) Let  $\sigma : \{0, 1, \dots, n-1\} \rightarrow \{1, \dots, n\}$  be a bijection. Then, the (affine) Sylvester space associated with the bijection  $\sigma$ , denoted by  $\text{Syl}(\sigma)$ , is the set of Fiedler matrices associated with  $\sigma$ , that is,

$$\text{Syl}(\sigma) := \left\{ M_\sigma(p) : p(z) = z^n + \sum_{k=0}^{n-1} c_k z^k, \quad c_k \in \mathbb{C} \right\},$$

where  $M_\sigma(p)$  is the matrix in (2.2).

For example, the Sylvester space associated with the bijection  $\sigma$ , such that  $\text{PCIS}(\sigma) = (1, 1, 1, 0, 0, 0)$ , is the set of matrices of the form

$$\begin{bmatrix} c_6 & c_5 & c_4 & c_3 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & c_2 & 0 & 1 & 0 \\ 0 & 0 & 0 & c_1 & 0 & 0 & 1 \\ 0 & 0 & 0 & c_0 & 0 & 0 & 0 \end{bmatrix},$$

<sup>2</sup>We note that the companion matrix considered in Edelman & Murakami (1995) is not exactly  $C_1$ , but the companion matrix obtained from  $C_1$  in (1.2) after performing a symmetry through the main anti-diagonal, and accordingly with the Sylvester space.

where  $c_k \in \mathbb{C}$ , for  $k = 0, 1, \dots, 6$ , may take any value. The tangent space of  $\text{Syl}(\sigma)$  at a given point, denoted by  $\text{TSyl}(\sigma)$ , is the set of matrices that we get if we remove the entries identically equal to 1 in the matrix above. In other words, the underlying vector space to the affine space. Observe that the tangent space of  $\text{Syl}(\sigma)$  in any matrix  $M \in \text{Syl}(\sigma)$  is independent of  $M$ . This is the reason why we just write  $\text{TSyl}(\sigma)$  without specifying the base point.

In order to extend the transversality identity (5.1) to the Sylvester space of any Fiedler matrix, we first need the following result, which is in turn an extension of (Edelman & Murakami, 1995, Eq. (5), p. 768).

**LEMMA 5.1** Let  $E^{\text{sy}l}$  be a matrix in  $\text{TSyl}(\sigma)$  with nonzero entries equal to  $E_0^{\text{sy}l}, E_1^{\text{sy}l}, \dots, E_{n-1}^{\text{sy}l}$ , where the entry  $E_k^{\text{sy}l}$ , for  $k = 0, 1, \dots, n-1$ , is in the same position as the coefficient  $-a_k$  in  $M_\sigma$ . Then, for  $k = 0, 1, \dots, n-1$ ,

$$\text{tr}(E^{\text{sy}l} p_{n-k-1}(M_\sigma)) = -E_k^{\text{sy}l}. \quad (5.2)$$

*Proof.* Let  $\tilde{p}(z) = z^n + \sum_{k=0}^{n-1} \tilde{a}_k z^k$  be the characteristic polynomial of  $M_\sigma + E^{\text{sy}l}$ . We know, by Propositions 3.1 and 3.4, that  $\tilde{a}_k = a_k - \text{tr}(E^{\text{sy}l} p_{n-k-1}(M_\sigma)) + O(\|E^{\text{sy}l}\|^2)$ . But  $M_\sigma + E^{\text{sy}l}$  is a Fiedler matrix of the polynomial  $z^n + \sum_{k=0}^{n-1} (a_k + E_k^{\text{sy}l}) z^k$ , therefore we have  $\tilde{a}_k = a_k + E_k^{\text{sy}l}$ . From these two formulas we get

$$\text{tr}(E^{\text{sy}l} p_{n-k-1}(M_\sigma)) + O(\|E^{\text{sy}l}\|^2) = -E_k^{\text{sy}l}.$$

Since this last equation is true regardless of the value of  $E_0^{\text{sy}l}, E_1^{\text{sy}l}, \dots, E_{n-1}^{\text{sy}l}$ , (5.2) follows.  $\square$

**THEOREM 5.3** Let  $p(z) = z^n + \sum_{k=0}^{n-1} a_k z^k$  be a monic polynomial,  $\sigma : \{0, 1, \dots, n-1\} \rightarrow \{1, \dots, n\}$  be a bijection, and let  $M_\sigma$  be the Fiedler matrix of  $p(z)$  associated to the bijection  $\sigma$ . Then  $\text{Syl}(\sigma)$  is transversal to  $\mathcal{O}(M_\sigma)$  at  $M_\sigma$ , i.e., every matrix  $E \in \mathbb{C}^{n \times n}$  can be expressed as

$$E = E^{\text{tan}} + E^{\text{sy}l}, \quad (5.3)$$

where  $E^{\text{sy}l} \in \text{TSyl}(\sigma)$  and  $E^{\text{tan}} \in T_{M_\sigma} \mathcal{O}(M_\sigma)$ .

*Proof.* Let  $E^{\text{sy}l}$  be a matrix in  $\text{TSyl}(\sigma)$  with nonzero entries  $E_k^{\text{sy}l} := -\text{tr}(E p_{n-k-1}(M_\sigma))$ , for  $k = 0, 1, \dots, n-1$ , where the entry  $E_k^{\text{sy}l}$  is in the same position as  $-a_k$  in  $M_\sigma$ . We may write the matrix  $E$  as  $E^{\text{sy}l} + E^{\text{tan}}$ , where  $E^{\text{tan}} = E - E^{\text{sy}l}$ . We have to check that  $E^{\text{tan}} \in T_{M_\sigma} \mathcal{O}(M_\sigma)$ . Indeed, using Lemma 5.1,

$$\text{tr}(E p_{n-k-1}(M_\sigma)) = \text{tr}(E^{\text{sy}l} p_{n-k-1}(M_\sigma)) + \text{tr}(E^{\text{tan}} p_{n-k-1}(M_\sigma)) = \text{tr}(E p_{n-k-1}(M_\sigma)) + \text{tr}(E^{\text{tan}} p_{n-k-1}(M_\sigma)).$$

From this, we deduce that  $\text{tr}(E^{\text{tan}} p_{n-k-1}(M_\sigma)) = 0$ , for  $k = 0, 1, 2, \dots, n-1$ . But, from Proposition 5.1, we have that  $\{p_k(M_\sigma)^*\}_{k=0}^{n-1}$  is a basis for  $N_{M_\sigma} \mathcal{O}(M_\sigma)$ , therefore  $E^{\text{tan}} \in T_{M_\sigma} \mathcal{O}(M_\sigma)$ .  $\square$

Theorem 5.3 and (5.2) show us that the component  $E^{\text{tan}}$  of the perturbation matrix  $E$  does not contribute to the first order term of  $a_k(M_\sigma + E)$ , so that only the ‘‘transversal complement’’  $E^{\text{sy}l}$  contributes to first order. In other words:

$$a_k(M_\sigma + E) = a_k - \text{tr}(p_{n-k-1}(M_\sigma)E) + O(\|E\|^2) = a_k - \text{tr}(p_{n-k-1}(M_\sigma)E^{\text{sy}l}) + O(\|E\|^2) = a_k(M_\sigma + E^{\text{sy}l}) + O(\|E\|^2).$$

Also, from the considerations above, if  $E_k^{\text{sy}l}$  denotes, as in Lemma 5.1, the entry of  $E^{\text{sy}l}$  which is located in the same position as the coefficient  $-a_k$  in  $M_\sigma$ , then we have, up to first order in  $E$ ,

$$E_k^{\text{sy}l} = a_k(M_\sigma + E) - a_k = - \sum_{i,j=1}^n p_{ij}^{(\sigma,k)}(a_0, a_1, \dots, a_{n-1}) E_{ij}, \quad (5.4)$$

as in (3.6), with  $p_{ij}^{(\sigma,k)}(a_0, a_1, \dots, a_{n-1})$  given by Theorem 3.3. The remaining entries of  $E^{\text{sy}l}$  are zero. Hence, from (5.3) and (5.4) we may get explicit expressions for the entries of  $E^{\text{tan}} = E - E^{\text{sy}l}$  in terms of the entries of  $E$  and  $a_0, a_1, \dots, a_{n-1}$ .

In the approach followed by Edelman & Murakami (1995), the fact that  $E^{\text{sy}l}$  is transversal to  $T_{C_1} \mathcal{O}(C_1)$  is key to get the first order expression for  $a_k(C_1 + E)$ . More precisely: using this transversality (namely, equation (5.3) with  $E^{\text{sy}l}$  being the Sylvester space for  $C_1$ ), together with the identity  $\text{tr}(p_{n-k}(C_1)E^{\text{tan}}) = 0$ , and the explicit expression  $-E_{k-1}^{\text{sy}l} = \text{tr}(p_{n-k}(C_1)E)$ , both valid for  $k = 1, \dots, n$ , they get an explicit expression for  $\text{tr}(p_{n-k}(C_1)E)$ , which is the first order term of  $a_{k-1}(C_1 + E)$ . This can be done because the matrices  $p_{n-k}(C_1)$ , for  $k = 1, \dots, n$ , have a simple structure that allows to compute  $\text{tr}(p_{n-k}(C_1)E^{\text{sy}l})$  easily and explicitly, for all  $k = 1, \dots, n$ . Unfortunately, for arbitrary Fiedler matrices, to get explicit expressions of  $\text{tr}(p_{n-k}(M_\sigma)E)$  by hand is quite involved. Hence, we have obtained the first-order term of  $a_k(M_\sigma + E)$  directly from  $\text{adj}(zI - M_\sigma)$ . This approach is completely independent of the transversality of  $E^{\text{sy}l}$  and the tangent space, though, as we have seen in Theorem 5.3, this fact is still true for arbitrary Fiedler matrices.

## 6. Conclusions and future work

In this paper, we have analyzed some numerical features of the polynomial root-finding problem when considered as a standard eigenvalue problem by means of Fiedler companion matrices. In particular, we have described the first-order change of the characteristic polynomial of any Fiedler matrix under small perturbations of the matrix. This description has led us to conclude that polynomial root-finding algorithms based on backward stable eigenvalue algorithms using Fiedler companion matrices, are backward stable only if  $\|p\|_\infty$  is moderate. More precisely, given a monic polynomial  $p(z)$ , if  $\tilde{p}(z)$  denotes the monic polynomial whose roots are the computed eigenvalues of a Fiedler companion matrix of  $p(z)$ , obtained with a backward stable eigenvalue algorithm, then it is not possible to guarantee, in general, that

$$\frac{\|\tilde{p} - p\|_\infty}{\|p\|_\infty} = O(u),$$

where  $u$  is the machine epsilon of the computer. Namely, the computed roots of  $p(z)$  are not necessarily the roots of a nearby polynomial. We have seen, however, that

$$\frac{\|\tilde{p} - p\|_\infty}{\|p\|_\infty} = O(u)\|p\|_\infty^2,$$

for any Fiedler companion matrix other than the first and second Frobenius companion matrices, and that

$$\frac{\|\tilde{p} - p\|_\infty}{\|p\|_\infty} = O(u)\|p\|_\infty,$$

for the first and second Frobenius companion matrices (which are particular cases of Fiedler matrices). Extensive numerical experiments have been included to illustrate these theoretical results.

One way to circumvent the inaccuracies due to the occurrence of large polynomial coefficients is to shift from companion matrices to companion pencils where normalization can be applied (see Jónsson & Vavasis (2004)). Though exactly the same techniques used in Jónsson & Vavasis (2004) for the Frobenius companion pencils can not be directly applied to other Fiedler companion pencils, some further analysis in this direction is still to be done, and will be the subject of future work.

**Acknowledgements.** This work has been supported by the Ministerio de Economía y Competitividad of Spain through grant MTM2012-32542. The authors wish to thank Yuji Nakatsukasa and Vanni Noferini for reading a draft of this manuscript, and for their insightful comments, which helped us to improve the paper. The authors also acknowledge the useful comments and suggestions of two anonymous referees. Finally, the authors wish to thank the Associate Editor handling this manuscript, Prof. Françoise Tisseur, for her advice and help when preparing a revised version.

## REFERENCES

- ARNOLD, V. I. (1971) Matrices depending on parameters. *Uspehi Mat. Nauk*, **26**, 101–114. English translation: Russian Math. Surveys 26 (1971), no. 2, 29–43.
- AURENTZ, J. L., VANDEBRIL, R. & WATKINS, D. S. (2013) Fast computation of the zeros of a polynomial via factorization of the companion matrix. *SIAM J. Sci. Comput.*, **35**, A255–A269.
- BERNSTEIN, D. S. (2009) *Matrix mathematics. Theory, facts, and formulas*, 2nd edn. Princeton University Press, Princeton, NJ.
- BHATIA, R. & JAIN, T. (2009) Higher order derivatives and perturbation bounds for determinants. *Linear Algebra Appl.*, **431**, 2102–2108.
- BINI, D. A., GEMIGNANI, L. & PAN, V. Y. (2004) Improved initialization of the accelerated and robust QR-like polynomial root-finding. *Electron. Trans. Numer. Anal.*, **17**, 195–205.
- BINI, D. A., GEMIGNANI, L. & PAN, V. Y. (2005) Fast and stable QR eigenvalue algorithms for generalized companion matrices and secular equations. *Numer. Math.*, **100**, 373–408.
- BINI, D. A., BOITO, P., EIDELMAN, Y., GEMIGNANI, L. & GOHBERG, I. (2010) A fast implicit QR eigenvalue algorithm for companion matrices. *Linear Algebra Appl.*, **432**, 2006–2031.
- BRUGNANO, L. & TRIGIANTE, D. (1995) Polynomial roots: the ultimate answer? *Linear Algebra Appl.*, **225**, 207–219.
- CALVETTI, D., KIM, S.-M. & REICHEL, L. (2002) The restarted QR-algorithm for eigenvalue computation of structured matrices. *J. Comput. Appl. Math.*, **149**, 415–422.
- CHANDRASEKARAN, S., GU, M., XIA, J. & ZHU, J. (2008) A fast QR algorithm for companion matrices. *Recent advances in matrix and operator theory*. Oper. Theory Adv. Appl., vol. 179. Birkhäuser, Basel, pp. 111–143.

- DE TERÁN, F., DOPICO, F. M. & MACKEY, D. S. (2010) Fiedler companion linearizations and the recovery of minimal indices. *SIAM J. Matrix Anal. Appl.*, **31**, 2181–2204.
- DE TERÁN, F., DOPICO, F. M. & PÉREZ, J. (2013) Condition numbers for inversion of Fiedler companion matrices. *Linear Algebra Appl.*, **439**, 944–981.
- DE TERÁN, F., DOPICO, F. M. & PÉREZ, J. (2014a) New bounds for roots of polynomials based on Fiedler companion matrices. *Linear Algebra Appl.*, **451**, 197–320.
- DE TERÁN, F., DOPICO, F. M. & PÉREZ, J. (2014b) Technical report on backward stability of polynomial root-finding using Fiedler companion matrices. *MIMS EPrint* 2014.38. UK: Manchester Institute for Mathematical Sciences, The University of Manchester.
- EDELMAN, A., ELMROTH, E. & KÅGSTRÖM, B. (1997) A geometric approach to perturbation theory of matrices and matrix pencils. I. Versal deformations. *SIAM J. Matrix Anal. Appl.*, **18**, 653–692.
- EDELMAN, A., ELMROTH, E. & KÅGSTRÖM, B. (1999) A geometric approach to perturbation theory of matrices and matrix pencils. II. A stratification-enhanced staircase algorithm. *SIAM J. Matrix Anal. Appl.*, **20**, 667–699 (electronic).
- EDELMAN, A. & MURAKAMI, H. (1995) Polynomial roots from companion matrix eigenvalues. *Math. Comp.*, **64**, 763–776.
- FIEDLER, M. (2003) A note on companion matrices. *Linear Algebra Appl.*, **372**, 325–331.
- GANTMACHER, F. R. (1959) *The theory of matrices*. Vols. 1, 2. Chelsea Publishing Co., New York.
- GEMIGNANI, L. (2007) Structured matrix methods for polynomial root-finding. *ISSAC 2007*. ACM, New York, pp. 175–180.
- GOLUB, G. H. & VAN LOAN, C. F. (1996) *Matrix computations*. Johns Hopkins Studies in the Mathematical Sciences, 3rd edn. Johns Hopkins University Press, Baltimore, MD.
- GRAUERT, H. & FRITZSCHE, K. (1976) *Several complex variables*. Springer-Verlag, New York-Heidelberg. Translated from the German, Graduate Texts in Mathematics, Vol. 38.
- HIGHAM, N. J. (2002) *Accuracy and stability of numerical algorithms*, 2nd edn. Society for Industrial and Applied Mathematics (SIAM), Philadelphia, PA.
- HORN, R. A. & JOHNSON, C. R. (1985) *Matrix analysis*. Cambridge University Press, Cambridge.
- IPSEN, I. C. F. & REHMAN, R. (2008) Perturbation bounds for determinants and characteristic polynomials. *SIAM J. Matrix Anal. Appl.*, **30**, 762–776.
- JÓNSSON, G. F. & VAVASIS, S. (2004) Solving polynomials with small leading coefficients. *SIAM J. Matrix Anal. Appl.*, **26**, 400–414.
- LAWRENCE, P. W. & CORLESS, R. M. (2014) Stability of rootfinding for barycentric Lagrange interpolants. *Numer. Algorithms*, **65**, 447–464.
- LEMMONIER, D. & VAN DOOREN, P. (2003) Optimal scaling of companion pencils for the  $qz$ -algorithm. *Proceedings SIAM Appl. Lin. Alg. Conference*, **Paper CP7-4**.
- MOLER, C. (1991) Cleve's corner: Roots of polynomials, that is. *The Mathworks Newsletter*, **5**, 8–9.
- NIU, X.-M. & SAKURAI, T. (2003) A method for finding the zeros of polynomials using a companion matrix. *Japan J. Indust. Appl. Math.*, **20**, 239–256.
- PARLETT, B. N. & REINSCH, C. (1969) Balancing a matrix for calculation of eigenvalues and eigenvectors. *Numer. Math.*, **13**, 293–304.
- TISSEUR, F. (2000) Backward error and condition of polynomial eigenvalue problems. *Linear Algebra Appl.*, **309**, 339–361.
- TOH, K.-C. & TREFETHEN, L. N. (1994) Pseudozeros of polynomials and pseudospectra of companion matrices. *Numer. Math.*, **68**, 403–425.
- VAN BAREL, M., VANDEBRIL, R., VAN DOOREN, P. & FREDERIX, K. (2010) Implicit double shift  $QR$ -algorithm for companion matrices. *Numer. Math.*, **116**, 177–212.
- VAN DOOREN, P. & DEWILDE, P. (1983) The eigenstructure of an arbitrary polynomial matrix: computational aspects. *Linear Algebra Appl.*, **50**, 545–579.
- ZHLOBICH, P. (2012) Differential  $qd$  algorithm with shifts for rank-structured matrices. *SIAM J. Matrix Anal. Appl.*, **33**, 1153–1171.