

Structured condition numbers for linear systems with parameterized quasiseparable coefficient matrices*

Froilán M. Dopico[†] Kenet Pomés[†]

March 31, 2016

Abstract

Low-rank structured matrices have attracted much attention in the last decades, since they arise in many applications and all share the fundamental property that can be represented by $\mathcal{O}(n)$ parameters, where $n \times n$ is the size of the matrix. This property has allowed the development of fast algorithms for solving numerically many problems involving low-rank structured matrices by performing operations on the parameters describing the matrices, instead of directly on the matrix entries. Among these problems the solution of linear systems of equations is probably the most basic and relevant one. Therefore, it is important to measure, via structured computable condition numbers, the relative sensitivity of the solutions of linear systems with low-rank structured coefficient matrices with respect to relative perturbations of the parameters representing such matrices, since this sensitivity determines the maximum accuracy attainable by fast algorithms and allows us to decide which set of parameters is the most convenient from the point of view of accuracy. To develop and analyze such condition numbers is the main goal of this paper. To this purpose, a general expression is obtained for the condition number of the solution of a linear system of equations whose coefficient matrix is any differentiable function of a vector of parameters with respect to perturbations of such parameters. Since there are many different classes of low-rank structured matrices and many different types of parameters describing them, it is not possible to cover all of them in a single work. Therefore, the general expression of the condition number is particularized to the important case of $\{1, 1\}$ -quasiseparable matrices and to the quasiseparable and the Givens-vector representations, in order to obtain explicit expressions of the corresponding two condition numbers that can be estimated in $\mathcal{O}(n)$ operations. In addition, detailed theoretical and numerical comparisons of these two condition numbers between themselves, and with respect to unstructured condition numbers are provided, which show that there are situations in which the unstructured condition number is much larger than the structured ones, but that the opposite never happens. The approach presented in this manuscript can be generalized to other classes of low-rank structured matrices and parameterizations.

Key words. condition numbers, linear systems, low-rank structured matrices, quasiseparable matrices, quasiseparable representation, Givens-vector representation

AMS subject classification. 65F05, 65F35, 15A06, 15A12

1 Introduction

A low-rank structured matrix is, in plain words, a matrix such that large submatrices of it have ranks much smaller than the size of the matrix. Banded matrices with small band-width are classical examples of low-rank structured matrices, but many other examples corresponding to dense matrices and appearing in many applications exist [9, 10, 16, 17]. Most of the classes of $n \times n$ low-rank structured matrices share the key property that they can be described by different sets of $\mathcal{O}(n)$ parameters, called *representations* [16, Ch. 2], which may be used in the development of *fast algorithms*, i.e., algorithms with a smaller exponent in the dependence on n of the computational cost than classical matrix algorithms. Many fast algorithms for low-rank structured matrices have been developed in the last years and they are often very sophisticated. However, all of them are based on the same simple and fundamental idea: they operate on the parameters describing the matrices instead of directly on the entries of the matrix. We refer the reader to the recent monographs [9, 10, 16, 17] and the survey paper [3], as well as to the huge amount

*Partially supported by Ministerio de Economía y Competitividad of Spain through grant MTM2012-32542.

[†]Departamento de Matemáticas, Universidad Carlos III de Madrid, Avda. Universidad 30, 28911 Leganés, Spain (dopico@math.uc3m.es, kpomes@math.uc3m.es).

of references therein, for a detail account of algorithms, properties, and applications of many different classes of low-rank structured matrices and their different parameterizations.

Despite the large number of references available on fast algorithms for computations with low-rank structured matrices, there exist very few references on the corresponding *a priori* rounding error analyses [1, 2, 4, 14, 18, 19]. Some reasons of this might be that such fast algorithms are frequently complicated and that some of them are potentially unstable, although work well in practice most of the times. This situation makes necessary the development of structured condition numbers with respect to the parameters on which fast algorithms operate and of reliable methods to estimate *a posteriori* the backward errors on such parameters from the residuals of the computed outputs of the algorithms, because in this way the forward errors committed by fast algorithms for low-rank structured matrices may be reliably and optimally estimated *a posteriori* as the product of the structured condition numbers times the backward errors on the parameters. This is a challenging research plan that has been initiated recently in [5], where for the first time in the literature some structured condition numbers (for eigenvalues in that case) with respect to some parameterizations of a certain class of low-rank structured matrices were introduced and analyzed. Among many other results, it was proved in [5] that simple eigenvalues of $\{1, 1\}$ -quasiseparable matrices [6, 16] may be much less sensitive to perturbations of the parameters describing the matrices than to perturbations of the entries of the matrix. In this paper, we extend the development of structured condition numbers with respect to perturbations of the parameters to the fundamental case of the solutions of linear systems of equations with low-rank structured coefficient matrices, and we will prove that, also in this case, the solutions may be much better conditioned with respect to perturbations of the parameters than with respect to perturbations of the entries of the matrix, but that the opposite can not happen. This work is influenced by the recent references [11, 5], which deal with the sensitivity of eigenvalues of some low-rank structured matrices, but also by the classical reference [13], in which the use of differential calculus for developing condition numbers was introduced.

Since there are many classes of low-rank structured matrices and many possible parameterizations, or representations, describing them, we focus in this paper on the particular, but very relevant, subclass of low-rank structured matrices known as $\{1, 1\}$ -quasiseparable matrices (see Definition 3.1 in Section 3) and on two of their most important representations, the general *quasiseparable representation*, which is non unique, and the essentially unique *Givens-vector representation* [6, 8, 15, 16]. One of the goals of considering two different representations is to illustrate another application of condition numbers with respect to different parameterizations for the solution of linear systems. Such application is to determine which representation is better to use, from the point of view of accuracy, for developing a fast algorithm for solving a linear system, since the most sensible choice is the one with smallest condition number. In this work, we will prove that, as it happens in the case of the eigenvalues [5], the condition number with respect to any quasiseparable representation is the same, and can not be too much larger than the condition number with respect to the Givens-vector representation, which is always the smallest. Moreover, we will show how these two condition numbers can be reliably estimated in $O(n)$ flops.

We emphasize that, although the results in this paper are particularized for $\{1, 1\}$ -quasiseparable matrices, the general framework established in Section 2 can be used to study structured condition numbers of solutions of linear systems with respect to different representations for many other classes of low-rank structured matrices.

The rest of the paper is organized as follows. Section 2 presents general results on condition numbers for the solutions of linear systems whose coefficient matrices are differentiable functions of some parameters with respect to perturbations of such parameters. Section 3 refreshes very briefly the definition of quasiseparable matrices and some of their properties. Sections 4, 5, and 6 include the most important results in this paper on the condition numbers for the solutions of linear systems of equations whose coefficient matrices are $\{1, 1\}$ -quasiseparable with respect to the quasiseparable and the Givens-vector representations and on the comparison between them. Section 7 discusses how these condition numbers can be estimated in $O(n)$ flops. Numerical experiments are presented in Section 8 and conclusions and lines of future research are established in Section 9.

Notation. Following a common notation in Numerical Linear Algebra, we will use capital Roman letters A, B, \dots , for matrices and lower case bold Roman letters $\mathbf{x}, \mathbf{y}, \dots$ for column vectors. We will only consider real vectors and matrices. Given a real column vector \mathbf{y} of size n , its transpose is denoted by \mathbf{y}^T . The vector ∞ -norm is used very often, so we recall its definition: $\|\mathbf{y}\|_\infty := \max_{1 \leq i \leq n} |y_i|$, where y_i denotes the i -th coordinate of the vector \mathbf{y} . The reader is referred to [12] for additional information on vector and matrix norms. Standard MATLAB notation for submatrices is used, i.e., given a matrix $A \in \mathbb{R}^{m \times n}$, the expression $A(i : j, k : l)$, where $1 \leq i \leq j \leq m$ and $1 \leq k \leq l \leq n$, denotes the submatrix of A consisting of rows i up to and including j of A and of columns k up to and including l of A .

2 Basics on condition numbers for linear systems

We start this section by presenting some well-known results about condition numbers for the solution of a linear system of equations. Note first that any perturbation of a matrix $A \in \mathbb{R}^{n \times n}$ can be expressed as a sum $A + \delta A$, where $\delta A \in \mathbb{R}^{n \times n}$ is called the *perturbation matrix*.

Associated with a *normwise backward error* we have the condition number in Definition 2.1 [12, Sec. 7.1], valid for any vector norm and the corresponding subordinate matrix norm.

Definition 2.1. *Let $A\mathbf{x} = \mathbf{b}$, where $A \in \mathbb{R}^{n \times n}$ is nonsingular, and $0 \neq \mathbf{x} \in \mathbb{R}^n$. Then, for $E \in \mathbb{R}^{n \times n}$ and $\mathbf{f} \in \mathbb{R}^n$, we define*

$$\kappa_{E,\mathbf{f}}(A, \mathbf{x}) := \limsup_{\eta \rightarrow 0} \left\{ \frac{\|\delta \mathbf{x}\|}{\eta \|\mathbf{x}\|} : (A + \delta A)(\mathbf{x} + \delta \mathbf{x}) = \mathbf{b} + \delta \mathbf{b}, \|\delta A\| \leq \eta \|E\|, \|\delta \mathbf{b}\| \leq \eta \|\mathbf{f}\| \right\}.$$

Observe that $\kappa_{E,\mathbf{f}}(A, \mathbf{x})$ is a normwise relative condition number, i.e., it measures the relative sensitivity of the solution \mathbf{x} of the linear system $A\mathbf{x} = \mathbf{b}$ with respect to relative normwise perturbations of the matrix and the right-hand side (note that, in this case, the perturbations are measured against the tolerances E and \mathbf{f}). This condition number has the expression presented in the following theorem proved in [12, Sec. 7.1].

Theorem 2.2. *Under the same hypotheses of Definition 2.1,*

$$\kappa_{E,\mathbf{f}}(A, \mathbf{x}) = \frac{\|A^{-1}\| \|\mathbf{f}\|}{\|\mathbf{x}\|} + \|A^{-1}\| \|E\|.$$

Recall that the usual matrix condition number is given by $\kappa(A) := \|A\| \|A^{-1}\|$ and note that if we take $E = A$ and $\mathbf{f} = \mathbf{b}$, then we have $\kappa(A) \leq \kappa_{E,\mathbf{f}}(A, \mathbf{x}) \leq 2\kappa(A)$, and therefore they are numerically equivalent. On the other hand, it is well-known that considering normwise perturbations of the matrix A and the vector \mathbf{b} may lead to pessimistic bounds on the forward errors, since there are matrices and vectors for which we may have a small relative normwise perturbation that produces some large relative perturbations over their small entries, and that may affect the zero pattern of the matrix or the vector (see, for instance, the numerical example in [12, pp. 121-124]). Therefore, it makes sense to consider componentwise perturbations and the corresponding componentwise condition number. We denote by $|A|$ the matrix whose entries are the absolute values of the entries of A (i.e., $|A|_{ij} := |A_{ij}|$) and we adopt a similar notation for vectors. In addition, inequalities $|A| \leq |B|$ mean $|A_{ij}| \leq |B_{ij}|$ for all i, j . Definition 2.3 and Theorem 2.4 can both be found in [12, Sec. 7.2], together with a brief discussion on how to choose the tolerances E and \mathbf{f} .

Definition 2.3. *Let $A\mathbf{x} = \mathbf{b}$, where $A \in \mathbb{R}^{n \times n}$ is nonsingular, and $0 \neq \mathbf{x} \in \mathbb{R}^n$. Then, for $0 \leq E \in \mathbb{R}^{n \times n}$ and $0 \leq \mathbf{f} \in \mathbb{R}^n$, we define the relative componentwise condition number as*

$$\text{cond}_{E,\mathbf{f}}(A, \mathbf{x}) := \limsup_{\eta \rightarrow 0} \left\{ \frac{\|\delta \mathbf{x}\|_\infty}{\eta \|\mathbf{x}\|_\infty} : (A + \delta A)(\mathbf{x} + \delta \mathbf{x}) = \mathbf{b} + \delta \mathbf{b}, |\delta A| \leq \eta E, |\delta \mathbf{b}| \leq \eta \mathbf{f} \right\}.$$

Theorem 2.4. *Under the same hypotheses of Definition 2.3,*

$$\text{cond}_{E,\mathbf{f}}(A, \mathbf{x}) = \frac{\| |A^{-1}| E |\mathbf{x}| + |A^{-1}| \mathbf{f} \|_\infty}{\|\mathbf{x}\|_\infty}.$$

A proof of this theorem is provided in [12, Sec. 7.2], but it can be seen as a consequence of the more general Theorem 2.9 that we introduce below, and, so, we will present a proof of Theorem 2.4 at the end of this section.

From the expression in Theorem 2.4, if we consider $E = |A|$ and $\mathbf{f} = |\mathbf{b}|$, then it is straightforward to prove that the condition number $\text{cond}_{|A|,|\mathbf{b}|}(A, \mathbf{x})$ is invariant under row scaling. This useful property is stated in the following proposition.

Proposition 2.5. *Let $A\mathbf{x} = \mathbf{b}$, where $A \in \mathbb{R}^{n \times n}$ is nonsingular and $0 \neq \mathbf{x} \in \mathbb{R}^n$, and let $K \in \mathbb{R}^{n \times n}$ be an invertible diagonal matrix. Then, for $K A \mathbf{x} = K \mathbf{b}$, we have*

$$\text{cond}_{|A|,|\mathbf{b}|}(A, \mathbf{x}) = \text{cond}_{|K A|,|K \mathbf{b}|}(K A, \mathbf{x}).$$

Since many interesting classes of matrices can be represented by sets of parameters different from their entries (see Theorem 4.1, for example), we generalize the definitions above to these representations and, following the ideas in [5], we will focus on componentwise relative condition numbers for representations.

Definition 2.6. Let $A\mathbf{x} = \mathbf{b}$, where $A \in \mathbb{R}^{n \times n}$ is a nonsingular matrix whose entries are differentiable functions of a vector of parameters $\Omega = (\omega_1, \omega_2, \dots, \omega_m)^T \in \mathbb{R}^m$, this is denoted by $A(\Omega)$, and $0 \neq \mathbf{x} \in \mathbb{R}^n$. Let $0 \leq \mathbf{f} \in \mathbb{R}^n$ and $E = (e_1, e_2, \dots, e_m)^T \in \mathbb{R}^m$ with nonnegative entries. Then, we define

$$\text{cond}_{E, \mathbf{f}}(A(\Omega), \mathbf{x}) := \limsup_{\eta \rightarrow 0} \left\{ \frac{\|\delta \mathbf{x}\|_\infty}{\eta \|\mathbf{x}\|_\infty} : (A(\Omega + \delta \Omega))(\mathbf{x} + \delta \mathbf{x}) = \mathbf{b} + \delta \mathbf{b}, |\delta \Omega| \leq \eta E, |\delta \mathbf{b}| \leq \eta \mathbf{f} \right\}.$$

The main goal of this section is to find an explicit expression for the componentwise relative condition number with respect to a general representation introduced in Definition 2.6. For such a purpose we will use differential calculus and we will need Lemma 2.7. In Lemma 2.7, \mathbf{e}_i denotes the i th vector of the canonical basis of \mathbb{R}^n .

Lemma 2.7. Let $A\mathbf{x} = \mathbf{b}$, where $A \in \mathbb{R}^{n \times n}$ is an invertible matrix whose entries are differentiable functions of a vector of real parameters $\Omega = (\omega_1, \omega_2, \dots, \omega_m)^T$, and $0 \neq \mathbf{x} \in \mathbb{R}^n$. Then, the following equalities hold:

- (a) $\frac{\partial A^{-1}}{\partial \omega_k} = -A^{-1} \frac{\partial A}{\partial \omega_k} A^{-1}$, for $k \in \{1, 2, \dots, m\}$,
- (b) $\frac{\partial \mathbf{x}}{\partial b_i} = A^{-1} \mathbf{e}_i$, for $i \in \{1, 2, \dots, n\}$,
- (c) $\frac{\partial \mathbf{x}}{\partial \omega_k} = -A^{-1} \frac{\partial A}{\partial \omega_k} \mathbf{x}$, for $k \in \{1, 2, \dots, m\}$.

Proof. (a) Derivating in both sides of $AA^{-1} = I_n$, we get

$$\frac{\partial A}{\partial \omega_k} A^{-1} + A \frac{\partial A^{-1}}{\partial \omega_k} = 0.$$

(b) It follows trivially from derivating $\mathbf{x} = A^{-1}\mathbf{b}$.

(c) From derivating $\mathbf{x} = A^{-1}\mathbf{b}$ and using (a), we obtain

$$\frac{\partial \mathbf{x}}{\partial \omega_k} = \frac{\partial A^{-1}}{\partial \omega_k} \mathbf{b} = -A^{-1} \frac{\partial A}{\partial \omega_k} A^{-1} \mathbf{b} = -A^{-1} \frac{\partial A}{\partial \omega_k} \mathbf{x}.$$

□

Remark 2.8. In Lemma 2.7, we have used that the entries of A^{-1} are also differentiable functions of $(\omega_1, \dots, \omega_m)$. This follows from the facts that (1) each entry of A^{-1} is a quotient of a cofactor of A divided by $\det(A)$ and that (2) products, sums, and quotients of differentiable functions are differentiable whenever the denominators are not zero.

In Theorem 2.9, we provide the desired explicit expression of the componentwise relative condition number introduced in Definition 2.6.

Theorem 2.9. Let $A\mathbf{x} = \mathbf{b}$, where $A \in \mathbb{R}^{n \times n}$ is an invertible matrix whose entries are differentiable functions of a vector of real parameters $\Omega = (\omega_1, \omega_2, \dots, \omega_m)^T$ and $0 \neq \mathbf{x} \in \mathbb{R}^n$. Let $0 \leq \mathbf{f} \in \mathbb{R}^n$ and $E = (e_1, e_2, \dots, e_m)^T \in \mathbb{R}^m$ with nonnegative entries. Then,

$$\text{cond}_{E, \mathbf{f}}(A(\Omega), \mathbf{x}) = \frac{\left\| |A^{-1}| \mathbf{f} + \sum_{k=1}^m \left| A^{-1} \frac{\partial A}{\partial \omega_k} \mathbf{x} \right| e_k \right\|_\infty}{\|\mathbf{x}\|_\infty}.$$

Proof. Since the entries of the matrix A are differentiable functions of the parameters in Ω and, from $\mathbf{x} = A^{-1}\mathbf{b}$, it is clear that \mathbf{x} is a function of Ω and \mathbf{b} , we can use differential calculus to obtain the following result:

$$\delta \mathbf{x} = \sum_{i=1}^n \frac{\partial \mathbf{x}}{\partial b_i} \delta b_i + \sum_{k=1}^m \frac{\partial \mathbf{x}}{\partial \omega_k} \delta \omega_k + \mathcal{O}(\|(\delta \Omega, \delta \mathbf{b})\|^2),$$

where $\|(\delta\Omega, \delta\mathbf{b})\| := \max\{\|\delta\Omega\|_\infty, \|\delta\mathbf{b}\|_\infty\}$. Using (b) and (c) from Lemma 2.7 in the previous equation, we obtain:

$$\delta\mathbf{x} = \sum_{i=1}^n (A^{-1}\mathbf{e}_i) \delta b_i + \sum_{k=1}^m \left(-A^{-1} \frac{\partial A}{\partial \omega_k} \mathbf{x} \right) \delta \omega_k + \mathcal{O}(\|(\delta\Omega, \delta\mathbf{b})\|^2). \quad (2.1)$$

From (2.1), using standard properties of the ∞ -norm and the inequalities $|\delta\mathbf{b}| \leq \eta \mathbf{f}$ and $|\delta\Omega| \leq \eta E$, we get

$$\begin{aligned} \|\delta\mathbf{x}\|_\infty &\leq \eta \left\| \sum_{i=1}^n |A^{-1}\mathbf{e}_i| f_i + \sum_{k=1}^m \left| A^{-1} \frac{\partial A}{\partial \omega_k} \mathbf{x} \right| e_k \right\|_\infty + \mathcal{O}(\|(\delta\Omega, \delta\mathbf{b})\|^2) \\ &= \eta \left\| |A^{-1}| \mathbf{f} + \sum_{k=1}^m \left| A^{-1} \frac{\partial A}{\partial \omega_k} \mathbf{x} \right| e_k \right\|_\infty + \mathcal{O}(\|(\delta\Omega, \delta\mathbf{b})\|^2). \end{aligned} \quad (2.2)$$

Then, if η tends to zero, from (2.2) and from Definition 2.6, it is straightforward to get

$$\text{cond}_{E,\mathbf{f}}(A(\Omega), \mathbf{x}) \leq \frac{\left\| |A^{-1}| \mathbf{f} + \sum_{k=1}^m \left| A^{-1} \frac{\partial A}{\partial \omega_k} \mathbf{x} \right| e_k \right\|_\infty}{\|\mathbf{x}\|_\infty}. \quad (2.3)$$

On the other hand, if we consider the perturbations:

$$\delta\mathbf{b} = \eta D \mathbf{f},$$

where D is a diagonal matrix such that $D(j, j) = \text{sign}(A^{-1}(l, j))$, for $j = 1, 2, \dots, n$, and

$$\delta \omega_k = -\eta \left[\text{sign} \left(A^{-1} \frac{\partial A}{\partial \omega_k} \mathbf{x} \right)_l \right] e_k, \quad \text{for } k = 1, \dots, m,$$

where l is such that

$$\left\| |A^{-1}| \mathbf{f} + \sum_{k=1}^m \left| A^{-1} \frac{\partial A}{\partial \omega_k} \mathbf{x} \right| e_k \right\|_\infty = \left(|A^{-1}| \mathbf{f} + \sum_{k=1}^m \left| A^{-1} \frac{\partial A}{\partial \omega_k} \mathbf{x} \right| e_k \right)_l,$$

we can obtain, from (2.1) and from Definition 2.6, the desired equality in (2.3). \square

As a consequence of Theorem 2.9 we can deduce the very well-known expression in Theorem 2.4 for the condition number $\text{cond}_{E,\mathbf{f}}(A, \mathbf{x})$ by considering Ω as the entries of A . Therefore, we conclude this section by providing its proof.

Proof. (of Theorem 2.4) Note first that, in this case, we can rewrite the expression in Theorem 2.9 as:

$$\text{cond}_{E,\mathbf{f}}(A, \mathbf{x}) = \frac{\left\| |A^{-1}| \mathbf{f} + \sum_{j,k=1}^n \left| A^{-1} \frac{\partial A}{\partial a_{jk}} \mathbf{x} \right| e_{jk} \right\|_\infty}{\|\mathbf{x}\|_\infty}.$$

Then, since it is obvious that $\partial A / \partial a_{jk} = \mathbf{e}_j \mathbf{e}_k^T$, we have:

$$\begin{aligned} \sum_{j,k=1}^n \left| A^{-1} \frac{\partial A}{\partial a_{jk}} \mathbf{x} \right| e_{jk} &= \sum_{j,k=1}^n |A^{-1}(:, j)| |x_k| e_{jk} = \sum_{k=1}^n \left(\sum_{j=1}^n |A^{-1}(:, j)| e_{jk} \right) |x_k| \\ &= \sum_{k=1}^n |A^{-1}| E(:, k) |x_k| = |A^{-1}| \sum_{k=1}^n E(:, k) |x_k| = |A^{-1}| E |\mathbf{x}|, \end{aligned}$$

and the proof follows trivially. \square

or, in a more compact notation,

$$A = \begin{bmatrix} d_1 & & & & \\ & d_2 & & & \\ & & \ddots & & \\ p_i a_{ij}^\times q_j & & & & \\ & & & & d_n \end{bmatrix},$$

where $a_{ij}^\times = a_{i-1}a_{i-2}\cdots a_{j+1}$, for $i-1 \geq j+1$, $b_{ij}^\times = b_{i+1}b_{i+2}\cdots b_{j-1}$, for $i+1 \leq j-1$, $a_{j+1,j}^\times = 1$, and $b_{j,j+1}^\times = 1$ for $j = 1, \dots, n-1$.

We call the vector of parameters Ω_{QS} in Theorem 4.1 a *quasiseparable representation* of the quasiseparable matrix A . Let us see the following 5×5 example for a more clear view of this representation.

Example 4.2. Let A be a $\{1, 1\}$ -quasiseparable matrix of size 5×5 and consider a quasiseparable representation of A :

$$\Omega_{QS} = (\{p_i\}_{i=2}^5, \{a_i\}_{i=2}^4, \{q_i\}_{i=1}^4, \{d_i\}_{i=1}^5, \{g_i\}_{i=1}^4, \{b_i\}_{i=2}^4, \{h_i\}_{i=2}^5).$$

Then,

$$A = \begin{bmatrix} d_1 & g_1 h_2 & g_1 b_2 h_3 & g_1 b_2 b_3 h_4 & g_1 b_2 b_3 b_4 h_5 \\ p_2 q_1 & d_2 & g_2 h_3 & g_2 b_3 h_4 & g_2 b_3 b_4 h_5 \\ p_3 a_2 q_1 & p_3 q_2 & d_3 & g_3 h_4 & g_3 b_4 h_5 \\ p_4 a_3 a_2 q_1 & p_4 a_3 q_2 & p_4 q_3 & d_4 & g_4 h_5 \\ p_5 a_4 a_3 a_2 q_1 & p_5 a_4 a_3 q_2 & p_5 a_4 q_3 & p_5 q_4 & d_5 \end{bmatrix}.$$

From Theorem 4.1 we see how any square $\{1, 1\}$ -quasiseparable matrix of size $n \times n$ can be represented with $\mathcal{O}(n)$ parameters instead of its n^2 entries. This fact is crucial for developing fast algorithms for performing computations with these matrices. There exist several fast algorithms, with cost $\mathcal{O}(n)$ operations, and working with different representations, for performing computations such as matrix-vector multiplication, solution of linear systems, and matrix inversion for quasiseparable matrices [6, 7, 9, 10, 16], and even for computing structured eigenvalue condition numbers [5, Secs. 4.3, 4.4].

On the other hand note that the quasiseparable representation is not unique as we can see in the following 5×5 example.

Example 4.3. Let A be a 5×5 quasiseparable matrix with a quasiseparable representation Ω_{QS} as in Example 4.2 and consider a real value $\alpha \in \mathbb{R}/\{0, 1\}$, then $\Omega'_{QS} = (\{\alpha p_i\}_{i=2}^5, \{a_i\}_{i=2}^4, \{q_i/\alpha\}_{i=1}^4, \{d_i\}_{i=1}^5, \{g_i\}_{i=1}^4, \{b_i\}_{i=2}^4, \{h_i\}_{i=2}^5)$ is also a quasiseparable representation of A .

Taking into account that our goal is to obtain explicit expressions of structured condition numbers for the solution of a linear system involving a quasiseparable matrix in the quasiseparable representation by using differential calculus, the next lemma will become useful. The easy proof is omitted since it can be found inside the proof of [5, Theorem 4.4].

Lemma 4.4. Let $A \in \mathbb{R}^{n \times n}$ be a $\{1, 1\}$ -quasiseparable matrix and $A = A_L + A_D + A_U$, with A_L strictly lower triangular, A_D diagonal, and A_U strictly upper triangular. Let Ω_{QS} be a quasiseparable representation of A , where $\Omega_{QS} = (\{p_i\}_{i=2}^n, \{a_i\}_{i=2}^{n-1}, \{q_i\}_{i=1}^{n-1}, \{d_i\}_{i=1}^n, \{g_i\}_{i=1}^{n-1}, \{b_i\}_{i=2}^{n-1}, \{h_i\}_{i=2}^n)$. Then the entries of A are differentiable functions of the parameters in Ω_{QS} and:

- a) $\frac{\partial A}{\partial d_i} = e_i e_i^T$, for $i = 1, \dots, n$.
- b) $p_i \frac{\partial A}{\partial p_i} = e_i A_L(i, :)$, for $i = 2, \dots, n$.
- c) $a_i \frac{\partial A}{\partial a_i} = \begin{bmatrix} 0 & 0 \\ A(i+1:n, 1:i-1) & 0 \end{bmatrix}$, for $i = 2, \dots, n-1$.
- d) $q_i \frac{\partial A}{\partial q_i} = A_L(:, i) e_i^T$, for $i = 1, \dots, n-1$.
- e) $g_i \frac{\partial A}{\partial g_i} = e_i A_U(i, :)$, for $i = 1, \dots, n-1$.

$$f) b_i \frac{\partial A}{\partial b_i} = \begin{bmatrix} 0 & A(1 : i - 1, i + 1 : n) \\ 0 & 0 \end{bmatrix}, \text{ for } i = 2, \dots, n - 1.$$

$$g) h_i \frac{\partial A}{\partial h_i} = A_U(:, i) e_i^T, \text{ for } i = 2, \dots, n.$$

4.2 The condition number for $\{1, 1\}$ -quasiseparable matrices in the quasiseparable representation: expression and properties

Since from Lemma 4.4 we have that the entries of a quasiseparable matrix A are differentiable functions of the parameters in a quasiseparable representation of A , we can deduce relative-relative component-wise condition numbers of the solution of linear systems with respect to these representations by using Theorem 2.9. This leads to Theorem 4.5.

Theorem 4.5. *Let $A\mathbf{x} = \mathbf{b}$, where $0 \neq \mathbf{x} \in \mathbb{R}^n$, $A \in \mathbb{R}^{n \times n}$ is a nonsingular $\{1, 1\}$ -quasiseparable matrix with a quasiseparable representation Ω_{QS} , and such that $A = A_L + A_D + A_U$, with A_L strictly lower triangular, A_D diagonal, and A_U strictly upper triangular. Let $0 \leq \mathbf{f} \in \mathbb{R}^n$ and $0 \leq E_{QS} \in \mathbb{R}^{7n-8}$. Then*

$$\begin{aligned} \text{cond}_{E_{QS}, \mathbf{f}}(A(\Omega_{QS}), \mathbf{x}) &= \frac{1}{\|\mathbf{x}\|_\infty} \left\| |A^{-1}| \mathbf{f} + |A^{-1}| |Q_d| |\mathbf{x}| + |A^{-1}| |Q_p| |A_L \mathbf{x}| + |A^{-1}| |A_L| |Q_q| |\mathbf{x}| \right. \\ &\quad + |A^{-1}| |Q_g| |A_U \mathbf{x}| + |A^{-1}| |A_U| |Q_h| |\mathbf{x}| \\ &\quad + \sum_{i=2}^{n-1} \left| A^{-1} \begin{bmatrix} 0 & 0 \\ A(i+1 : n, 1 : i-1) & 0 \end{bmatrix} \mathbf{x} \right| \left| \frac{e_{a_i}}{a_i} \right| \\ &\quad \left. + \sum_{j=2}^{n-1} \left| A^{-1} \begin{bmatrix} 0 & A(1 : j-1, j+1 : n) \\ 0 & 0 \end{bmatrix} \mathbf{x} \right| \left| \frac{e_{b_j}}{b_j} \right| \right\|_\infty, \end{aligned}$$

where:

$$\begin{aligned} \Omega_{QS} &= (\{p_i\}_{i=2}^n, \{a_i\}_{i=2}^{n-1}, \{q_i\}_{i=1}^{n-1}, \{d_i\}_{i=1}^n, \{g_i\}_{i=1}^{n-1}, \{b_i\}_{i=2}^{n-1}, \{h_i\}_{i=2}^n), \\ E_{QS} &= (\{e_{p_i}\}_{i=2}^n, \{e_{a_i}\}_{i=2}^{n-1}, \{e_{q_i}\}_{i=1}^{n-1}, \{e_{d_i}\}_{i=1}^n, \{e_{g_i}\}_{i=1}^{n-1}, \{e_{b_i}\}_{i=2}^{n-1}, \{e_{h_i}\}_{i=2}^n), \\ Q_d &= \text{diag}(e_{d_1}, \dots, e_{d_n}), \quad Q_p = \text{diag}\left(1, \frac{e_{p_2}}{p_2}, \dots, \frac{e_{p_n}}{p_n}\right), \quad Q_q = \text{diag}\left(\frac{e_{q_1}}{q_1}, \dots, \frac{e_{q_{n-1}}}{q_{n-1}}, 1\right), \\ Q_g &= \text{diag}\left(\frac{e_{g_1}}{g_1}, \dots, \frac{e_{g_{n-1}}}{g_{n-1}}, 1\right), \quad Q_h = \text{diag}\left(1, \frac{e_{h_2}}{h_2}, \dots, \frac{e_{h_n}}{h_n}\right), \end{aligned}$$

and each quotient whose denominator is zero must be understood as zero if the numerator is also zero and, otherwise, the zero parameter in the denominator should be formally cancelled out with the same parameter in the corresponding piece of A .

Proof. We will proceed by calculating the contribution of each subset of parameters to the expression for $\text{cond}_{E_{QS}, \mathbf{f}}(A(\Omega_{QS}), \mathbf{x})$ given in Theorem 2.9 as follows.

Derivatives with respect to $\{d_i\}_{i=1}^n$. By using a) in Lemma 4.4 we get:

$$\kappa_d := \sum_{i=1}^n \left| A^{-1} \frac{\partial A}{\partial d_i} \mathbf{x} \right| e_{d_i} = \sum_{i=1}^n |A^{-1} \mathbf{e}_i \mathbf{e}_i^T \mathbf{x}| e_{d_i} = \sum_{i=1}^n |A^{-1}(:, i)| |x_i| e_{d_i} = |A^{-1}| |Q_d| |\mathbf{x}|.$$

Derivatives with respect to $\{p_i\}_{i=2}^n$. From b) in Lemma 4.4 we have:

$$\begin{aligned} \kappa_p &:= \sum_{i=2}^n \left| A^{-1} \frac{\partial A}{\partial p_i} \mathbf{x} \right| |e_{p_i}| = \sum_{i=2}^n \left| A^{-1} p_i \frac{\partial A}{\partial p_i} \mathbf{x} \right| \left| \frac{e_{p_i}}{p_i} \right| = \sum_{i=2}^n |A^{-1} \mathbf{e}_i A_L(i, :) \mathbf{x}| \left| \frac{e_{p_i}}{p_i} \right| \\ &= \sum_{i=2}^n |A^{-1}(:, i) A_L(i, :) \mathbf{x}| \left| \frac{e_{p_i}}{p_i} \right| = \sum_{i=2}^n |A^{-1}(:, i)| \left| \frac{e_{p_i}}{p_i} \right| |A_L(i, :) \mathbf{x}| = |A^{-1}| |Q_p| |A_L \mathbf{x}|. \end{aligned}$$

Derivatives with respect to $\{a_i\}_{i=2}^{n-1}$. From c) in Lemma 4.4 we have:

$$\kappa_a := \sum_{i=2}^{n-1} \left| A^{-1} a_i \frac{\partial A}{\partial a_i} \mathbf{x} \right| \left| \frac{e_{a_i}}{a_i} \right| = \sum_{i=2}^{n-1} \left| A^{-1} \begin{bmatrix} 0 & 0 \\ A(i+1 : n, 1 : i-1) & 0 \end{bmatrix} \mathbf{x} \right| \left| \frac{e_{a_i}}{a_i} \right|.$$

Derivatives with respect to $\{q_j\}_{j=1}^{n-1}$. From d) in Lemma 4.4 we have:

$$\begin{aligned}\kappa_q &:= \sum_{j=1}^{n-1} \left| A^{-1} \frac{\partial A}{\partial q_j} \mathbf{x} \right|_{|e_{q_j}|} = \sum_{j=1}^{n-1} \left| A^{-1} q_j \frac{\partial A}{\partial q_j} \mathbf{x} \right|_{\left| \frac{e_{q_j}}{q_j} \right|} = \sum_{j=1}^{n-1} \left| A^{-1} A_L(:, j) \mathbf{e}_j^T \mathbf{x} \right|_{\left| \frac{e_{q_j}}{q_j} \right|} \\ &= \sum_{j=1}^{n-1} \left| A^{-1} A_L(:, j) \right|_{|x_j|} \left| \frac{e_{q_j}}{q_j} \right| = |A^{-1} A_L| |Q_q| |\mathbf{x}|.\end{aligned}$$

Analogously, we can find the contribution to the condition number $\text{cond}_{E_{QS}, \mathbf{f}}(A(\Omega_{QS}), \mathbf{x})$ of the derivatives of A with respect to the parameters $\{g_i\}_{i=1}^{n-1}$, $\{b_i\}_{i=2}^{n-1}$, and $\{h_i\}_{i=2}^n$, which describe the strictly upper triangular part of A . The results are the following.

Derivatives with respect to $\{g_i\}_{i=1}^{n-1}$. By using e) in Lemma 4.4 we obtain:

$$\kappa_g := \sum_{i=1}^{n-1} \left| A^{-1} \frac{\partial A}{\partial g_i} \mathbf{x} \right|_{|e_{g_i}|} = |A^{-1}| |Q_g| |A_U \mathbf{x}|.$$

Derivatives with respect to $\{b_i\}_{i=2}^{n-1}$. By using f) in Lemma 4.4 we obtain:

$$\kappa_b := \sum_{i=2}^{n-1} \left| A^{-1} \frac{\partial A}{\partial b_i} \mathbf{x} \right|_{|e_{b_i}|} = \sum_{i=2}^{n-1} \left| A^{-1} \begin{bmatrix} 0 & A(1 : i-1, i+1 : n) \\ 0 & 0 \end{bmatrix} \mathbf{x} \right|_{\left| \frac{e_{b_i}}{b_i} \right|}.$$

Derivatives with respect to $\{h_j\}_{j=2}^n$. By using g) in Lemma 4.4 we obtain:

$$\kappa_h := \sum_{j=2}^n \left| A^{-1} \frac{\partial A}{\partial h_j} \mathbf{x} \right|_{|e_{h_j}|} = |A^{-1} A_U| |Q_h| |\mathbf{x}|.$$

This proof is completed by observing that according to Theorem 2.9 we have:

$$\text{cond}_{E_{QS}, \mathbf{f}}(A(\Omega_{QS}), \mathbf{x}) = \frac{\| |A^{-1}| \mathbf{f} + \kappa_d + \kappa_p + \kappa_a + \kappa_q + \kappa_g + \kappa_b + \kappa_h \|_{\infty}}{\|\mathbf{x}\|_{\infty}}.$$

□

As it is easy to see from its explicit expression in Theorem 4.5, $\text{cond}_{E_{QS}, \mathbf{f}}(A(\Omega_{QS}), \mathbf{x})$ depends in general on the quasiseparable representation Ω_{QS} of the matrix A . More specifically, that condition number depends on the ratios between the parameters in the representation and the corresponding tolerances in E_{QS} . Next, we will restrict ourselves to the case $E_{QS} = |\Omega_{QS}|$ in Theorem 4.6, which is the most natural election for E_{QS} . In this situation we adopt for brevity in the rest of the paper, the following notation:

$$\text{cond}_{\mathbf{f}}(A(\Omega_{QS}), \mathbf{x}) \equiv \text{cond}_{|\Omega_{QS}|, \mathbf{f}}(A(\Omega_{QS}), \mathbf{x}),$$

since the parameters Ω_{QS} are already shown in $A(\Omega_{QS})$.

Theorem 4.6. *Let $A\mathbf{x} = \mathbf{b}$, where $0 \neq \mathbf{x} \in \mathbb{R}^n$ and $A \in \mathbb{R}^{n \times n}$ is a nonsingular $\{1, 1\}$ -quasiseparable matrix such that $A = A_L + A_D + A_U$, with A_L strictly lower triangular, A_D diagonal, and A_U strictly upper triangular. Let $0 \leq \mathbf{f} \in \mathbb{R}^n$. Then*

$$\begin{aligned}\text{cond}_{\mathbf{f}}(A(\Omega_{QS}), \mathbf{x}) &= \frac{1}{\|\mathbf{x}\|_{\infty}} \left\| \left| A^{-1} \right| \mathbf{f} + |A^{-1}| |A_D| |\mathbf{x}| + |A^{-1}| |A_L| |\mathbf{x}| + |A^{-1} A_L| |\mathbf{x}| + |A^{-1}| |A_U| |\mathbf{x}| \right. \\ &\quad \left. + |A^{-1} A_U| |\mathbf{x}| + \sum_{i=2}^{n-1} \left| A^{-1} \begin{bmatrix} 0 & 0 \\ A(i+1 : n, 1 : i-1) & 0 \end{bmatrix} \mathbf{x} \right| \right. \\ &\quad \left. + \sum_{j=2}^{n-1} \left| A^{-1} \begin{bmatrix} 0 & A(1 : j-1, j+1 : n) \\ 0 & 0 \end{bmatrix} \mathbf{x} \right| \right\|_{\infty}.\end{aligned}$$

Proof. It follows directly from the expression in Theorem 4.5 for $\text{cond}_{E_{QS}, \mathbf{f}}(A(\Omega_{QS}), \mathbf{x})$ by observing that, in this case, we are considering $E_{QS} = |\Omega_{QS}|$ and, therefore, using the notation in Theorem 4.5, the following equalities hold: $Q_d = |A_D|$, $|Q_p| = |Q_q| = |Q_g| = |Q_h| = I$, and $|e_{a_i}/a_i| = |e_{b_i}/b_i| = 1$. □

Proposition 4.7 proves that $\text{cond}_{\mathbf{f}}(A(\Omega_{QS}), \mathbf{x})$ depends only on A , \mathbf{x} and \mathbf{f} , but not on the particular choice of quasiseparable parameters.

Proposition 4.7. *Let $A\mathbf{x} = \mathbf{b}$, where $0 \neq \mathbf{x} \in \mathbb{R}^n$ and $A \in \mathbb{R}^{n \times n}$ is a nonsingular $\{1, 1\}$ -quasiseparable matrix. Let $0 \leq \mathbf{f} \in \mathbb{R}^n$. Then, for any two vectors Ω_{QS} and Ω'_{QS} of quasiseparable parameters of A ,*

$$\text{cond}_{\mathbf{f}}(A(\Omega_{QS}), \mathbf{x}) = \text{cond}_{\mathbf{f}}(A(\Omega'_{QS}), \mathbf{x}).$$

Proof. It is obvious from the fact that the expression in Theorem 4.6 does not depend on the parameters of the representation Ω_{QS} but on the entries of the matrix A and the entries of the vectors \mathbf{x} and \mathbf{f} . \square

Proposition 4.8 states another important property of this relative componentwise condition number that arises from the natural comparison with the unstructured relative entrywise condition number for the solution of linear systems defined in Definition 2.3 and further developed in Theorem 2.4.

Proposition 4.8. *Let $A\mathbf{x} = \mathbf{b}$, where $0 \neq \mathbf{x} \in \mathbb{R}^n$ and $A \in \mathbb{R}^{n \times n}$ is a nonsingular $\{1, 1\}$ -quasiseparable matrix, and let Ω_{QS} be a quasiseparable representation of A . Then, for $0 \leq \mathbf{f} \in \mathbb{R}^n$, the following relation holds,*

$$\text{cond}_{\mathbf{f}}(A(\Omega_{QS}), \mathbf{x}) \leq n \text{cond}_{|A|, \mathbf{f}}(A, \mathbf{x}).$$

Proof. From Theorem 4.6 and using standard properties of absolute values and norms we obtain:

$$\begin{aligned} \text{cond}_{\mathbf{f}}(A(\Omega_{QS}), \mathbf{x}) &\leq \frac{1}{\|\mathbf{x}\|_{\infty}} \left\| |A^{-1}| \mathbf{f} + |A^{-1}| |A_D| |\mathbf{x}| + |A^{-1}| |A_L| |\mathbf{x}| \right. \\ &\quad \left. + |A^{-1}| |A_L| |\mathbf{x}| + |A^{-1}| |A_U| |\mathbf{x}| + |A^{-1}| |A_U| |\mathbf{x}| \right. \\ &\quad \left. + \sum_{i=2}^{n-1} |A^{-1}| |A_L| |\mathbf{x}| + \sum_{i=2}^{n-1} |A^{-1}| |A_U| |\mathbf{x}| \right\|_{\infty} \\ &= \frac{1}{\|\mathbf{x}\|_{\infty}} \left\| |A^{-1}| \mathbf{f} + |A^{-1}| |A_D| |\mathbf{x}| \right. \\ &\quad \left. + n |A^{-1}| |A_L| |\mathbf{x}| + n |A^{-1}| |A_U| |\mathbf{x}| \right\|_{\infty} \\ &\leq \frac{n}{\|\mathbf{x}\|_{\infty}} \left\| |A^{-1}| \mathbf{f} + |A^{-1}| |A| |\mathbf{x}| \right\|_{\infty} = n \text{cond}_{|A|, \mathbf{f}}(A, \mathbf{x}). \end{aligned}$$

\square

According to this proposition, the structured condition number $\text{cond}_{\mathbf{f}}(A(\Omega_{QS}), \mathbf{x})$ is smaller than the unstructured condition number $\text{cond}_{|A|, \mathbf{f}}(A, \mathbf{x})$, except for a factor n . In addition, we will see in the numerical experiments presented in Section 8 that it can be much smaller.

From Proposition 2.5 we know that the unstructured componentwise condition number is invariant under row scaling, which is a very convenient property (see [12, Secs. 7.2 and 7.3]). Therefore, it makes sense to study the behavior of the structured condition number under row scaling for the natural choice $\mathbf{f} = |\mathbf{b}|$ as well. This is done in Proposition 4.10, for which we will need Lemma 4.9. Proposition 4.10 proves that the structured componentwise condition number is also invariant under row scaling.

Lemma 4.9. *Let $K = \text{diag}(k_1, k_2, \dots, k_n)$ be an invertible diagonal matrix and $A \in \mathbb{R}^{n \times n}$ be a $\{1, 1\}$ -quasiseparable matrix with a quasiseparable representation $\Omega_{QS} = (\{p_i\}_{i=2}^n, \{a_i\}_{i=2}^{n-1}, \{q_i\}_{i=1}^{n-1}, \{d_i\}_{i=1}^n, \{g_i\}_{i=1}^{n-1}, \{b_i\}_{i=2}^{n-1}, \{h_i\}_{i=2}^n)$, as in Theorem 4.1. Then, the matrix KA is also a $\{1, 1\}$ -quasiseparable matrix and $\Omega'_{QS} = (\{k_i p_i\}_{i=2}^n, \{a_i\}_{i=2}^{n-1}, \{q_i\}_{i=1}^{n-1}, \{k_i d_i\}_{i=1}^n, \{k_i g_i\}_{i=1}^{n-1}, \{b_i\}_{i=2}^{n-1}, \{h_i\}_{i=2}^n)$ is a quasiseparable representation of KA .*

Proof. It follows from Theorem 4.1 and from $(KA)(i, j) = k_i A(i, j)$. \square

Proposition 4.10. *Let $A\mathbf{x} = \mathbf{b}$, where $0 \neq \mathbf{x} \in \mathbb{R}^n$ and $A \in \mathbb{R}^{n \times n}$ is a nonsingular $\{1, 1\}$ -quasiseparable matrix with a quasiseparable representation Ω_{QS} , such that $A = A_L + A_D + A_U$, with A_L strictly lower triangular, A_D diagonal, and A_U strictly upper triangular. Let $K \in \mathbb{R}^{n \times n}$ be an invertible diagonal matrix. Then*

$$\text{cond}_{|K\mathbf{b}|}((KA)(\Omega'_{QS}), \mathbf{x}) = \text{cond}_{|\mathbf{b}|}(A(\Omega_{QS}), \mathbf{x}),$$

where Ω'_{QS} is any quasiseparable representation of KA .

Proof. Since K is a diagonal matrix, all the following equalities are straightforward:

- 1) $|(KA)^{-1}|K\mathbf{b}| = |A^{-1}||K^{-1}||K||\mathbf{b}| = |A^{-1}||\mathbf{b}|,$
- 2) $|(KA)^{-1}||KA_D||\mathbf{x}| = |A^{-1}||K^{-1}||K||A_D||\mathbf{x}| = |A^{-1}||A_D||\mathbf{x}|,$
- 3) $|(KA)^{-1}||KA_L\mathbf{x}| = |A^{-1}||K^{-1}||K||A_L\mathbf{x}| = |A^{-1}||A_L\mathbf{x}|,$
- 4) $|(KA)^{-1}(KA_L)||\mathbf{x}| = |A^{-1}A_L||\mathbf{x}|,$
- 5) $|(KA)^{-1}||KA_U\mathbf{x}| = |A^{-1}||K^{-1}||K||A_U\mathbf{x}| = |A^{-1}||A_U\mathbf{x}|,$
- 6) $|(KA)^{-1}(KA_U)||\mathbf{x}| = |A^{-1}A_U||\mathbf{x}|,$
- 7) $(KA)^{-1} \begin{bmatrix} 0 & 0 \\ (KA)(i+1:n, 1:i-1) & 0 \end{bmatrix} \mathbf{x} = A^{-1} \begin{bmatrix} 0 & 0 \\ A(i+1:n, 1:i-1) & 0 \end{bmatrix} \mathbf{x},$
- 8) $(KA)^{-1} \begin{bmatrix} 0 & (KA)(1:j-1, j+1:n) \\ 0 & 0 \end{bmatrix} \mathbf{x} = A^{-1} \begin{bmatrix} 0 & A(1:j-1, j+1:n) \\ 0 & 0 \end{bmatrix} \mathbf{x}.$

The result follows trivially from 1)-8), $KA = KA_L + KA_D + KA_U$, Theorem 4.6, and the fact that $\text{cond}_{\mathcal{F}}(A(\Omega_{QS}), \mathbf{x})$ does not depend on the particular quasiseparable parameterization used. \square

5 Condition number of the solution of $\{1, 1\}$ -quasiseparable linear systems in the Givens-vector representation

Another important representation for quasiseparable matrices is the Givens-vector representation, which was introduced for the first time in [15] and that will be described in Section 5.1, along with its minor variant called tangent-Givens-vector representation, which has been introduced in [5]. The original contributions of this paper to the study of structured condition numbers for the solution of linear systems with respect to this representation are presented in Sections 5.2 and 6.

5.1 The Givens-vector representation for $\{1, 1\}$ -quasiseparable matrices

The Givens-vector representation for $\{1, 1\}$ -quasiseparable matrices was introduced in [15] in order to improve the numerical stability in numerical computations with respect to other representations, but the first rigorous contributions that show that the Givens-vector representation is indeed “more stable” than other representations, appear in [5] in the context of eigenvalue problems. In this paper, these “stability contributions” are extended to the area of linear systems. Theorem 5.1 (see [16, Sections 2.4 and 2.8]) shows how the class of $\{1, 1\}$ -quasiseparable matrices can be represented by using Givens-vector parameters.

Theorem 5.1. *A matrix $A \in \mathbb{R}^{n \times n}$ is a $\{1, 1\}$ -quasiseparable matrix if and only if it can be parameterized in terms of the following set of parameters,*

- $\{c_i, s_i\}_{i=2}^{n-1}$, where (c_i, s_i) is a pair of cosine-sine with $c_i^2 + s_i^2 = 1$ for every $i \in \{2, 3, \dots, n-1\}$,
- $\{v_i\}_{i=1}^{n-1}, \{d_i\}_{i=1}^n, \{e_i\}_{i=1}^{n-1}$ all of them independent real parameters,
- $\{r_i, t_i\}_{i=2}^{n-1}$, where (r_i, t_i) is a pair of cosine-sine with $r_i^2 + t_i^2 = 1$ for every $i \in \{2, 3, \dots, n-1\}$,

as follows:

$$A = \begin{bmatrix} & d_1 & e_1 r_2 & e_1 t_2 r_3 & \cdots & e_1 t_2 \dots t_{n-2} r_{n-1} & e_1 t_2 \dots t_{n-1} \\ & c_2 v_1 & d_2 & e_2 r_3 & \cdots & e_2 t_3 \dots t_{n-2} r_{n-1} & e_2 t_3 \dots t_{n-1} \\ & c_3 s_2 v_1 & c_3 v_2 & d_3 & \cdots & e_3 t_4 \dots t_{n-2} r_{n-1} & e_3 t_4 \dots t_{n-1} \\ & \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ c_{n-1} s_{n-2} \dots s_2 v_1 & c_{n-1} s_{n-2} \dots s_3 v_2 & c_{n-1} s_{n-2} \dots s_4 v_3 & \cdots & & d_{n-1} & e_{n-1} \\ s_{n-1} s_{n-2} \dots s_2 v_1 & s_{n-1} s_{n-2} \dots s_3 v_2 & s_{n-1} s_{n-2} \dots s_4 v_3 & \cdots & & v_{n-1} & d_n \end{bmatrix}.$$

This representation is denoted by Ω_{QS}^{GV} , i.e., $\Omega_{QS}^{GV} := (\{c_i, s_i\}_{i=2}^{n-1}, \{v_i\}_{i=1}^{n-1}, \{d_i\}_{i=1}^n, \{e_i\}_{i=1}^{n-1}, \{r_i, t_i\}_{i=2}^{n-1})$.

Example 5.2. Let $A \in \mathbb{R}^{5 \times 5}$ be a $\{1, 1\}$ -quasiseparable matrix, and let

$$\Omega_{QS}^{GV} := (\{c_i, s_i\}_{i=2}^4, \{v_i\}_{i=1}^4, \{d_i\}_{i=1}^5, \{e_i\}_{i=1}^4, \{r_i, t_i\}_{i=2}^4)$$

be a Givens-vector representation of A . Then,

$$A = \begin{bmatrix} d_1 & e_1 r_2 & e_1 t_2 r_3 & e_1 t_2 t_3 r_4 & e_1 t_2 t_3 t_4 \\ c_2 v_1 & d_2 & e_2 r_3 & e_2 t_3 r_4 & e_2 t_3 t_4 \\ c_3 s_2 v_1 & c_3 v_2 & d_3 & e_3 r_4 & e_3 t_4 \\ c_4 s_3 s_2 v_1 & c_4 s_3 v_2 & c_4 v_3 & d_4 & e_4 \\ s_4 s_3 s_2 v_1 & s_4 s_3 v_2 & s_4 v_3 & v_4 & d_5 \end{bmatrix}.$$

Note that if we consider the following relations between the parameters in Theorems 4.1 and 5.1, respectively, $\{p_i, a_i\}_{i=2}^{n-1} = \{c_i, s_i\}_{i=2}^{n-1}$, $\{q_i\}_{i=1}^{n-1} = \{v_i\}_{i=1}^{n-1}$, $\{d_i\}_{i=1}^n = \{d_i\}_{i=1}^n$, $\{g_i\}_{i=1}^{n-1} = \{e_i\}_{i=1}^{n-1}$, $\{b_i, h_i\}_{i=2}^{n-1} = \{t_i, r_i\}_{i=2}^{n-1}$, and $p_n = h_n = 1$, then it is obvious that the Givens-vector representation is a particular case of a quasiseparable representation for $\{1, 1\}$ -quasiseparable matrices. This fact can be better observed by comparing the expressions in Examples 4.2 and 5.2. Note also that the Givens-vector representation can be made unique if c_i and r_i are taken to be nonnegative numbers (if $c_i = 0$, take $s_i = 1$ and if $r_i = 0$, take $t_i = 1$) [16, p.76].

On the other hand, since the Givens-vector representation is a particular case of the quasiseparable representation, one might think that it makes no sense to study structured condition numbers for this representation, since we know from Proposition 4.7, that the condition number for a quasiseparable matrix is independent of the particular choice of the quasiseparable representation Ω_{QS} when the natural choice $E_{QS} = |\Omega_{QS}|$ is made. However, the subtle point here is that the Givens-vector representation has correlated parameters since the pairs $\{c_i, s_i\}$ are not independent parameters and the same happens for $\{r_i, t_i\}$. In fact, independent componentwise perturbations of Ω_{QS}^{GV} destroy in general the pairs cosine-sine. Therefore, if we only consider perturbations that preserve the pairs cosine-sine then a different condition number is obtained. In Definition 5.3, introduced in [5], an additional parameterization by using tangents is provided in order to make explicit the correlations between $\{c_i, s_i\}$ and $\{r_i, t_i\}$.

Definition 5.3. For any Givens-vector representation

$$\Omega_{QS}^{GV} = (\{c_i, s_i\}_{i=2}^{n-1}, \{v_i\}_{i=1}^{n-1}, \{d_i\}_{i=1}^n, \{e_i\}_{i=1}^{n-1}, \{r_i, t_i\}_{i=2}^{n-1})$$

of a $\{1, 1\}$ -quasiseparable matrix $A \in \mathbb{R}^{n \times n}$, we define the Givens-vector representation via tangents as

$$\Omega_{GV} := (\{l_i\}_{i=2}^{n-1}, \{v_i\}_{i=1}^{n-1}, \{d_i\}_{i=1}^n, \{e_i\}_{i=1}^{n-1}, \{u_i\}_{i=2}^{n-1}), \text{ where}$$

$$c_i = \frac{1}{\sqrt{1+l_i^2}}, \quad s_i = \frac{l_i}{\sqrt{1+l_i^2}}, \quad \text{and} \quad r_i = \frac{1}{\sqrt{1+u_i^2}}, \quad t_i = \frac{u_i}{\sqrt{1+u_i^2}}, \quad \text{for } i = 2, \dots, n-1.$$

Observe, from the expressions in Definition 5.3 for the cosine-sine parameters $\{c_i, s_i\}$ and $\{r_i, t_i\}$ in terms of the tangents l_i and u_i , respectively, that tiny relative perturbations of those tangents produce tiny relative perturbations of those cosine-sine parameters. This suggests the convenience of using the tangent-Givens-vector representation in practical numerical situations as was explained and motivated in [5].

In order to use differential calculus to deduce an explicit expression of the structured condition number of the solution of a linear system with respect to the tangent-Givens-vector representation, we will need Lemma 5.4 (see the simple proof inside the proof of [5, Theorem 5.4]).

Lemma 5.4. Let $A \in \mathbb{R}^{n \times n}$ be a $\{1, 1\}$ -quasiseparable matrix and let Ω_{GV} be the tangent-Givens-vector representation of A , where $\Omega_{GV} = (\{l_i\}_{i=2}^{n-1}, \{v_i\}_{i=1}^{n-1}, \{d_i\}_{i=1}^n, \{e_i\}_{i=1}^{n-1}, \{u_i\}_{i=2}^{n-1})$. Then the entries of A are differentiable functions of the parameters in Ω_{GV} and

$$\begin{aligned} \text{a) } l_i \frac{\partial A}{\partial l_i} &= \begin{bmatrix} 0 & 0 \\ -s_i^2 A(i, 1 : i-1) & 0 \\ c_i^2 A(i+1 : n, 1 : i-1) & 0 \end{bmatrix}, \text{ for } i = 2, \dots, n-1, \\ \text{b) } u_i \frac{\partial A}{\partial u_i} &= \begin{bmatrix} 0 & -t_i^2 A(1 : i-1, i) & r_i^2 A(1 : i-1, i+1 : n) \\ 0 & 0 & 0 \end{bmatrix}, \text{ for } i = 2, \dots, n-1. \end{aligned}$$

Remark 5.5. For the partial derivatives with respect to the parameters in $\{d_i\}_{i=1}^n$, $\{v_i\}_{i=1}^{n-1}$, and $\{e_i\}_{i=1}^{n-1}$ (see a), d), e) in Lemma 4.4 (recall that those parameters can also be respectively seen as the parameters $\{d_i\}_{i=1}^n$, $\{q_i\}_{i=1}^{n-1}$, and $\{g_i\}_{i=1}^{n-1}$ in a quasiseparable representation of A).

5.2 The condition number for $\{1, 1\}$ -quasiseparable matrices in the Givens-vector representation

Theorem 5.6 is the main result of Section 5 and it presents an explicit expression of the component-wise condition number for the solution of a linear system of equations with a quasiseparable matrix of coefficients with respect to the Givens-vector representation via tangents.

Theorem 5.6. *Let $A\mathbf{x} = \mathbf{b}$, where $0 \neq \mathbf{x} \in \mathbb{R}^n$ and $A \in \mathbb{R}^{n \times n}$ is a nonsingular $\{1, 1\}$ -quasiseparable matrix with a tangent-Givens-vector representation Ω_{GV} , and such that $A = A_L + A_D + A_U$, with A_L strictly lower triangular, A_D diagonal, and A_U strictly upper triangular. Let $0 \leq \mathbf{f} \in \mathbb{R}^n$ and $0 \leq E_{GV} \in \mathbb{R}^{5n-6}$. Then*

$$\begin{aligned} \text{cond}_{E_{GV}, \mathbf{f}}(A(\Omega_{GV}), \mathbf{x}) &= \frac{1}{\|\mathbf{x}\|_\infty} \left\| |A^{-1}| \mathbf{f} + |A^{-1}| |Q_d| |\mathbf{x}| + |A^{-1} A_L| |Q_v| |\mathbf{x}| \right. \\ &\quad + |A^{-1}| |Q_e| |A_U \mathbf{x}| + \sum_{i=2}^{n-1} \left| A^{-1} \begin{bmatrix} 0 & 0 \\ -s_i^2 A(i, 1 : i-1) & 0 \\ c_i^2 A(i+1 : n, 1 : i-1) & 0 \end{bmatrix} \mathbf{x} \right| \left| \frac{e_{l_i}}{l_i} \right| \\ &\quad \left. + \sum_{j=2}^{n-1} \left| A^{-1} \begin{bmatrix} 0 & -t_j^2 A(1 : j-1, j) & r_j^2 A(1 : j-1, j+1 : n) \\ 0 & 0 & 0 \end{bmatrix} \mathbf{x} \right| \left| \frac{e_{u_j}}{u_j} \right| \right\|_\infty, \end{aligned}$$

where

$$\Omega_{GV} = (\{l_i\}_{i=2}^{n-1}, \{v_i\}_{i=1}^{n-1}, \{d_i\}_{i=1}^n, \{e_i\}_{i=1}^{n-1}, \{u_i\}_{i=2}^{n-1}), \text{ as in Definition 5.3,}$$

$$E_{GV} = (\{e_{l_i}\}_{i=2}^{n-1}, \{e_{v_i}\}_{i=1}^{n-1}, \{e_{d_i}\}_{i=1}^n, \{e_{e_i}\}_{i=1}^{n-1}, \{e_{u_i}\}_{i=2}^{n-1}),$$

$$Q_d = \text{diag}(e_{d_1}, \dots, e_{d_n}), Q_v = \text{diag}\left(\frac{e_{v_1}}{v_1}, \dots, \frac{e_{v_{n-1}}}{v_{n-1}}, 1\right), Q_e = \text{diag}\left(\frac{e_{e_1}}{e_1}, \dots, \frac{e_{e_{n-1}}}{e_{n-1}}, 1\right),$$

and each quotient whose denominator is zero must be understood as zero if the numerator is also zero and, otherwise, the zero parameter in the denominator should be formally cancelled out with the same parameter in the corresponding piece of A .

Proof. The proof is straightforward from Theorem 2.9, Lemma 5.4, Remark 5.5, and the proof of Theorem 4.5. Therefore, we omit the proof. \square

For the most natural choice $E_{GV} = |\Omega_{GV}|$, we adopt the shorter notation

$$\text{cond}_{\mathbf{f}}(A(\Omega_{GV}), \mathbf{x}) \equiv \text{cond}_{|\Omega_{GV}|, \mathbf{f}}(A(\Omega_{GV}), \mathbf{x}),$$

and get Theorem 5.7 as a corollary of Theorem 5.6.

Theorem 5.7. *Let $A\mathbf{x} = \mathbf{b}$, where $0 \neq \mathbf{x} \in \mathbb{R}^n$ and $A \in \mathbb{R}^{n \times n}$ is a nonsingular $\{1, 1\}$ -quasiseparable matrix with a tangent-Givens-vector representation Ω_{GV} , and such that $A = A_L + A_D + A_U$, with A_L strictly lower triangular, A_D diagonal, and A_U strictly upper triangular. Let $0 \leq \mathbf{f} \in \mathbb{R}^n$. Then*

$$\begin{aligned} \text{cond}_{\mathbf{f}}(A(\Omega_{GV}), \mathbf{x}) &= \frac{1}{\|\mathbf{x}\|_\infty} \left\| |A^{-1}| \mathbf{f} + |A^{-1}| |A_D| |\mathbf{x}| + |A^{-1} A_L| |\mathbf{x}| + |A^{-1}| |A_U \mathbf{x}| \right. \\ &\quad + \sum_{i=2}^{n-1} \left| A^{-1} \begin{bmatrix} 0 & 0 \\ -s_i^2 A(i, 1 : i-1) & 0 \\ c_i^2 A(i+1 : n, 1 : i-1) & 0 \end{bmatrix} \mathbf{x} \right| \\ &\quad \left. + \sum_{j=2}^{n-1} \left| A^{-1} \begin{bmatrix} 0 & -t_j^2 A(1 : j-1, j) & r_j^2 A(1 : j-1, j+1 : n) \\ 0 & 0 & 0 \end{bmatrix} \mathbf{x} \right| \right\|_\infty. \end{aligned}$$

Note, from the expressions in Theorem 5.7 and 4.6, that there exists an important difference between the structured condition number in the Givens-vector representation and the structured condition number in the quasiseparable representation for a given $\{1, 1\}$ -quasiseparable matrix A , since $\text{cond}_{\mathbf{f}}(A(\Omega_{GV}), \mathbf{x})$ depends not only on the entries of the matrix A but on the parameters $\{c_i, s_i\}$ and $\{r_i, t_i\}$ as well, while $\text{cond}_{\mathbf{f}}(A(\Omega_{QS}), \mathbf{x})$ only depends on the matrix entries. Furthermore, since the cosine-sine parameters in the Givens-vector representation do not change trivially under diagonal scalings, we have that, for the natural choice $\mathbf{f} = |\mathbf{b}|$, $\text{cond}_{\mathbf{f}}(A(\Omega_{GV}), \mathbf{x})$ is not invariant under row scaling while $\text{cond}_{\mathbf{f}}(A(\Omega_{QS}), \mathbf{x})$ is invariant under row scaling (recall Proposition 4.10).

6 Comparison of $\text{cond}_{\mathbf{f}}(A(\Omega_{GV}), \mathbf{x})$ and $\text{cond}_{\mathbf{f}}(A(\Omega_{QS}), \mathbf{x})$

In [5, Section 6] it was proved that the Givens-vector representation via tangents is a more stable representation than the general quasiseparable representation for eigenvalue computations for $\{1, 1\}$ -quasiseparable matrices, in the sense that the Givens-vector representation leads to smaller eigenvalue condition numbers. This was a natural result to expect, as it is in the case of computing the solution of a quasiseparable linear system of equations, since any Givens-vector representation is also a quasiseparable representation and the perturbations considered in the condition numbers preserve the structure of the tangent-Givens-vector parametrization. In Theorem 6.1, the corresponding result is proved for linear systems, that is, $\text{cond}_{\mathbf{f}}(A(\Omega_{GV}), \mathbf{x})$ can not be larger than $\text{cond}_{\mathbf{f}}(A(\Omega_{QS}), \mathbf{x})$.

Theorem 6.1. *Let $A\mathbf{x} = \mathbf{b}$, where $0 \neq \mathbf{x} \in \mathbb{R}^n$ and $A \in \mathbb{R}^{n \times n}$ is a nonsingular $\{1, 1\}$ -quasiseparable matrix, and let Ω_{GV} be the vector of tangent-Givens-vector parameters of A . Then, for $0 \leq \mathbf{f} \in \mathbb{R}^n$, and for any vector Ω_{QS} of quasiseparable parameters of A , the following inequality holds:*

$$\text{cond}_{\mathbf{f}}(A(\Omega_{GV}), \mathbf{x}) \leq \text{cond}_{\mathbf{f}}(A(\Omega_{QS}), \mathbf{x}).$$

Proof. Throughout the proof, we use the decomposition $A = A_L + A_D + A_U$ introduced in Theorems 4.5 and 5.6. For the sums in the last two terms of the expression for $\text{cond}_{\mathbf{f}}(A(\Omega_{GV}), \mathbf{x})$ we have,

$$\begin{aligned} S_1 &:= \sum_{i=2}^{n-1} \left| A^{-1} \begin{bmatrix} 0 & 0 \\ -s_i^2 A(i, 1 : i-1) & 0 \\ c_i^2 A(i+1 : n, 1 : i-1) & 0 \end{bmatrix} \mathbf{x} \right| \\ &\leq \sum_{i=2}^{n-1} \left| A^{-1} \begin{bmatrix} 0 & 0 \\ A(i, 1 : i-1) & 0 \\ 0 & 0 \end{bmatrix} \mathbf{x} \right| + \sum_{i=2}^{n-1} \left| A^{-1} \begin{bmatrix} 0 & 0 \\ 0 & 0 \\ A(i+1 : n, 1 : i-1) & 0 \end{bmatrix} \mathbf{x} \right| \\ &= \sum_{i=2}^{n-1} |A^{-1}(:, i)| |A_L(i, :)\mathbf{x}| + \sum_{i=2}^{n-1} \left| A^{-1} \begin{bmatrix} 0 & 0 \\ A(i+1 : n, 1 : i-1) & 0 \end{bmatrix} \mathbf{x} \right| \\ &\leq |A^{-1}| |A_L \mathbf{x}| + \sum_{i=2}^{n-1} \left| A^{-1} \begin{bmatrix} 0 & 0 \\ A(i+1 : n, 1 : i-1) & 0 \end{bmatrix} \mathbf{x} \right| \end{aligned} \quad (6.1)$$

and, in an analogous way,

$$\begin{aligned} S_2 &:= \sum_{j=2}^{n-1} \left| A^{-1} \begin{bmatrix} 0 & -t_j^2 A(1 : j-1, j) & r_j^2 A(1 : j-1, j+1 : n) \\ 0 & 0 & 0 \end{bmatrix} \mathbf{x} \right| \\ &\leq |A^{-1} A_U| |\mathbf{x}| + \sum_{j=2}^{n-1} \left| A^{-1} \begin{bmatrix} 0 & A(1 : j-1, j+1 : n) \\ 0 & 0 \end{bmatrix} \mathbf{x} \right|. \end{aligned} \quad (6.2)$$

From (6.1) and (6.2) the proof is straightforward by comparing the expressions in Theorem 4.6 and Theorem 5.7 for $\text{cond}_{\mathbf{f}}(A(\Omega_{QS}), \mathbf{x})$ and $\text{cond}_{\mathbf{f}}(A(\Omega_{GV}), \mathbf{x})$, respectively. \square

On the other hand, as we prove in Theorem 6.3 below, the Givens-vector representation via tangents can only improve the relative condition number for the solution of a $\{1, 1\}$ -quasiseparable linear system of equations up to a factor of $3n$ with respect to the quasiseparable representation. Therefore, we conclude that, when computing solutions of $\{1, 1\}$ -quasiseparable linear systems of equations, these representations can be considered numerically equivalent in terms of expected accuracy.

In order to prove Theorem 6.3, we follow the ideas in [5, Section 6] and develop a proof based on Definition 2.6. Recall that from the Givens-vector representation via tangents Ω_{GV} of A we can obtain the Givens-vector representation Ω_{QS}^{GV} of A as in Definition 5.3, and that Ω_{QS}^{GV} is also a quasiseparable representation of A . Therefore, in order to use the componentwise relative condition numbers for representations in Definition 2.6, let us consider a quasiseparable perturbation $\delta\Omega_{QS}^{GV}$ of the parameters in Ω_{QS}^{GV} such that $|\delta\Omega_{QS}^{GV}| \leq \eta |\Omega_{QS}^{GV}|$, and the resulting quasiseparable matrix $\tilde{A} = A(\Omega_{QS}^{GV} + \delta\Omega_{QS}^{GV})$ (note that the perturbations $\delta\Omega_{QS}^{GV}$ do not respect in general the pairs cosine-sine of Ω_{QS}^{GV}). We will refer to η as the *level* of the relative perturbation of the parameters in the representation Ω_{QS}^{GV} . Moreover, note that \tilde{A} can also be represented by a vector

$$\Omega'_{GV} := (\{l'_i\}_{i=2}^{n-1}, \{v'_i\}_{i=1}^{n-1}, \{d'_i\}_{i=1}^n, \{e'_i\}_{i=1}^{n-1}, \{u'_i\}_{i=2}^{n-1})$$

of tangent-Givens-vector parameters and let us consider the perturbations $\delta' \Omega_{GV} := \Omega'_{GV} - \Omega_{GV}$. Then, Lemma 6.2, proved inside the proof of Theorem 6.3 in [5], states a bound for the level η' of the respective relative perturbation over the parameters in Ω'_{GV} produced by a relative perturbation of level η over the quasiseparable parameters in Ω_{QS}^{GV} .

Lemma 6.2. *Let A be a $\{1, 1\}$ -quasiseparable matrix with Givens-vector representation via tangents Ω_{GV} . Then, using the notation in the previous paragraph, and for η sufficiently small, we have:*

$$|\delta \Omega_{QS}^{GV}| \leq \eta |\Omega_{QS}^{GV}| \implies |\delta' \Omega_{GV}| \leq (3(n-2)\eta + \mathcal{O}(\eta^2)) |\Omega_{GV}|.$$

Theorem 6.3. *Let $A\mathbf{x} = \mathbf{b}$, where $0 \neq \mathbf{x} \in \mathbb{R}^n$ and $A \in \mathbb{R}^{n \times n}$ is a nonsingular $\{1, 1\}$ -quasiseparable matrix with tangent-Givens-vector representation Ω_{GV} . Let $0 \leq \mathbf{f} \in \mathbb{R}^n$. Then, for any quasiseparable representation Ω_{QS} of A :*

$$\frac{\text{cond}_{\mathbf{f}}(A(\Omega_{QS}), \mathbf{x})}{\text{cond}_{\mathbf{f}}(A(\Omega_{GV}), \mathbf{x})} \leq 3(n-2).$$

Proof. Note that from Definition 2.6 and from Lemma 6.2 we have

$$\text{cond}_{\mathbf{f}}(A(\Omega_{QS}), \mathbf{x}) \leq \limsup_{\eta \rightarrow 0} \left\{ \frac{\|\delta \mathbf{x}\|_{\infty}}{\eta \|\mathbf{x}\|_{\infty}} : (A(\Omega_{GV} + \delta \Omega_{GV}))(\mathbf{x} + \delta \mathbf{x}) = \mathbf{b} + \delta \mathbf{b}, \right. \\ \left. |\delta \Omega_{GV}| \leq (3(n-2)\eta + \mathcal{O}(\eta^2)) |\Omega_{GV}|, |\delta \mathbf{b}| \leq (3(n-2)\eta + \mathcal{O}(\eta^2)) \mathbf{f} \right\}.$$

By considering the change of variable $\eta' = (3(n-2)\eta + \mathcal{O}(\eta^2))$, we obtain

$$\text{cond}_{\mathbf{f}}(A(\Omega_{QS}), \mathbf{x}) \leq \limsup_{\eta' \rightarrow 0} \left\{ \frac{3(n-2)\|\delta \mathbf{x}\|_{\infty}}{\eta' \|\mathbf{x}\|_{\infty}} : (A(\Omega_{GV} + \delta \Omega_{GV}))(\mathbf{x} + \delta \mathbf{x}) = \mathbf{b} + \delta \mathbf{b}, \right. \\ \left. |\delta \Omega_{GV}| \leq \eta' |\Omega_{GV}|, |\delta \mathbf{b}| \leq \eta' \mathbf{f} \right\} = 3(n-2) \text{cond}_{\mathbf{f}}(A(\Omega_{GV}), \mathbf{x}).$$

□

7 Fast estimation of condition numbers: the effective condition number

To compute $\text{cond}_{\mathbf{f}}(A(\Omega_{QS}), \mathbf{x})$ and $\text{cond}_{\mathbf{f}}(A(\Omega_{GV}), \mathbf{x})$ fast, i.e., in $\mathcal{O}(n)$ flops, is not easy because of the two sums that appear in the expressions in Theorems 4.6 and 5.7, respectively, which require to compute a sum of n vectors, which cost $\mathcal{O}(n^2)$ flops, in addition to other computations. Then, a natural question now is to determine whether or not the contributions of these sums to the condition numbers in which they arise are significant. This question is answered in Theorems 7.2 and 7.3, in which we provide upper and lower bounds for $\text{cond}_{\mathbf{f}}(A(\Omega_{QS}), \mathbf{x})$ and $\text{cond}_{\mathbf{f}}(A(\Omega_{GV}), \mathbf{x})$ respectively, in terms of the *effective condition number* in Definition 7.1. We will show in this way that such effective condition number contains the essential terms in the expressions given in Theorems 4.6 and 5.7.

Definition 7.1. *Let $A\mathbf{x} = \mathbf{b}$, where $0 \neq \mathbf{x} \in \mathbb{R}^n$ and $A \in \mathbb{R}^{n \times n}$ is a nonsingular $\{1, 1\}$ -quasiseparable matrix such that $A = A_L + A_D + A_U$, with A_L strictly lower triangular, A_D diagonal, and A_U strictly upper triangular. Let $0 \leq \mathbf{f} \in \mathbb{R}^n$. Then, for any quasiseparable representation Ω_{QS} of A , we define the effective relative condition number $\text{condeff}_{\mathbf{f}}(A(\Omega_{QS}), \mathbf{x})$ for the solution of $A\mathbf{x} = \mathbf{b}$ as*

$$\text{condeff}_{\mathbf{f}}(A(\Omega_{QS}), \mathbf{x}) := \frac{1}{\|\mathbf{x}\|_{\infty}} \left\| \left| A^{-1} \right| \mathbf{f} + |A^{-1}| |A_D| |\mathbf{x}| + |A^{-1}| |A_L| \mathbf{x} \right. \\ \left. + |A^{-1} A_L| |\mathbf{x}| + |A^{-1}| |A_U| \mathbf{x} + |A^{-1} A_U| |\mathbf{x}| \right\|_{\infty}.$$

Recall from Proposition 4.7 that the condition number $\text{cond}_{\mathbf{f}}(A(\Omega_{QS}), \mathbf{x})$ does not depend on the choice of a quasiseparable representation Ω_{QS} , and note from Definition 7.1 that the same holds for $\text{condeff}_{\mathbf{f}}(A(\Omega_{QS}), \mathbf{x})$. Therefore, this effective condition number is always the same for any vector of quasiseparable parameters representing the matrix.

Theorem 7.2. Let $A\mathbf{x} = \mathbf{b}$, where $0 \neq \mathbf{x} \in \mathbb{R}^n$ and $A \in \mathbb{R}^{n \times n}$ is a nonsingular $\{1, 1\}$ -quasiseparable matrix. Let $0 \leq \mathbf{f} \in \mathbb{R}^n$. Then, for any quasiseparable representation Ω_{QS} of A , the following relations hold:

$$\text{condeff}_{\mathbf{f}}(A(\Omega_{QS}), \mathbf{x}) \leq \text{cond}_{\mathbf{f}}(A(\Omega_{QS}), \mathbf{x}) \leq (n-1)\text{condeff}_{\mathbf{f}}(A(\Omega_{QS}), \mathbf{x}).$$

Proof. Throughout the proof, we use the decomposition $A = A_L + A_D + A_U$ introduced in Theorems 4.5 and 5.6. Note first that the left hand side of the inequality is trivial from the respective expressions of $\text{condeff}_{\mathbf{f}}(A(\Omega_{QS}), \mathbf{x})$ and $\text{cond}_{\mathbf{f}}(A(\Omega_{QS}), \mathbf{x})$. On the other hand, note that

$$\begin{aligned} \left| A^{-1} \begin{bmatrix} 0 & 0 \\ A(i+1:n, 1:i-1) & 0 \end{bmatrix} \mathbf{x} \right| &= \left| A^{-1} \begin{bmatrix} 0 & 0 \\ A_L(i+1:n, 1:i-1) & 0 \end{bmatrix} \mathbf{x} \right| \\ &= \left| A^{-1} \begin{bmatrix} 0 & 0 \\ A_L(i+1:n, 1:i-1) & A_L(i+1:n, i:n) \end{bmatrix} \mathbf{x} \right| \\ &\quad + \left| A^{-1} \begin{bmatrix} 0 & 0 \\ 0 & -A_L(i+1:n, i:n) \end{bmatrix} \mathbf{x} \right| \\ &\leq |A^{-1}| \left| \begin{bmatrix} 0 \\ A_L(i+1:n, :) \end{bmatrix} \mathbf{x} \right| \\ &\quad + \left| A^{-1} \begin{bmatrix} 0 & 0 \\ 0 & A_L(i+1:n, i:n) \end{bmatrix} \right| |\mathbf{x}| \\ &\leq |A^{-1}| |A_L \mathbf{x}| + |A^{-1} A_L| |\mathbf{x}|, \end{aligned}$$

from where we obtain:

$$\sum_{i=2}^{n-1} \left| A^{-1} \begin{bmatrix} 0 & 0 \\ A(i+1:n, 1:i-1) & 0 \end{bmatrix} \mathbf{x} \right| \leq (n-2) |A^{-1}| |A_L \mathbf{x}| + (n-2) |A^{-1} A_L| |\mathbf{x}|. \quad (7.1)$$

In an analogous way, it can be proved that

$$\sum_{j=2}^{n-1} \left| A^{-1} \begin{bmatrix} 0 & A(1:j-1, j+1:n) \\ 0 & 0 \end{bmatrix} \mathbf{x} \right| \leq (n-2) |A^{-1}| |A_U \mathbf{x}| + (n-2) |A^{-1} A_U| |\mathbf{x}|. \quad (7.2)$$

Finally, note that from the inequalities in (7.1) and (7.2), it is straightforward that

$$\begin{aligned} \text{cond}_{\mathbf{f}}(A(\Omega_{QS}), \mathbf{x}) &\leq (n-1) \left\| |A^{-1}| \mathbf{f} + |A^{-1}| |A_D| |\mathbf{x}| + |A^{-1}| |A_L \mathbf{x}| \right. \\ &\quad \left. + |A^{-1} A_L| |\mathbf{x}| + |A^{-1}| |A_U \mathbf{x}| + |A^{-1} A_U| |\mathbf{x}| \right\|_{\infty} / \|\mathbf{x}\|_{\infty} \\ &= (n-1) \text{condeff}_{\mathbf{f}}(A(\Omega_{QS}), \mathbf{x}). \end{aligned}$$

□

Theorem 7.3. Let $A\mathbf{x} = \mathbf{b}$, where $0 \neq \mathbf{x} \in \mathbb{R}^n$ and $A \in \mathbb{R}^{n \times n}$ is a nonsingular $\{1, 1\}$ -quasiseparable matrix with tangent-Givens-vector representation Ω_{GV} . Let $0 \leq \mathbf{f} \in \mathbb{R}^n$. Then, for any quasiseparable representation Ω_{QS} of A , the following relations hold:

$$\frac{\text{condeff}_{\mathbf{f}}(A(\Omega_{QS}), \mathbf{x})}{3(n-2)} \leq \text{cond}_{\mathbf{f}}(A(\Omega_{GV}), \mathbf{x}) \leq (n-1)\text{condeff}_{\mathbf{f}}(A(\Omega_{QS}), \mathbf{x}).$$

Proof. It follows trivially from Theorems 6.1, 6.3 and 7.2. □

Note that from Theorems 7.2 and 7.3 we can conclude that the structured condition numbers $\text{cond}_{\mathbf{f}}(A(\Omega_{QS}), \mathbf{x})$ and $\text{cond}_{\mathbf{f}}(A(\Omega_{GV}), \mathbf{x})$ can be both estimated, “up to a factor of order n ”, by computing the easier expression in Definition 7.1 for the effective condition number. Next, we prove in Proposition 7.4 that this effective condition number can be computed fast by using, for instance, some previous results from [6] and [7].

Proposition 7.4. *Let $A\mathbf{x} = \mathbf{b}$, where $0 \neq \mathbf{x} \in \mathbb{R}^n$ and $A \in \mathbb{R}^{n \times n}$ is a nonsingular $\{1, 1\}$ -quasiseparable matrix with a quasiseparable representation Ω_{QS} which is assumed to be known. Let $0 \leq \mathbf{f} \in \mathbb{R}^n$. Then, the effective condition number $\text{condeff}_{\mathbf{f}}(A(\Omega_{QS}), \mathbf{x})$ can be computed in $\mathcal{O}(n)$ operations, i.e., with linear complexity.*

Proof. This assertion is a consequence of the results in [6] and [7]. Using [7, Algorithm 5.1] we can obtain in $\mathcal{O}(n)$ flops the quasiseparable representation for the inverse of the $\{1, 1\}$ -quasiseparable matrix A which is also $\{1, 1\}$ -quasiseparable ([6, Theorem 5.2]). In addition, in [6, Algorithm 4.4] it is shown how to compute the matrix-vector product $\mathbf{y} = R\mathbf{z}$ for the general case when R is an $\{n_L, n_U\}$ -quasiseparable matrix with a given quasiseparable representation, with linear complexity in the size n of the vector. Therefore, we can use this algorithm (twice when necessary) for computing each of the products $|A^{-1}|\mathbf{f}$, $|A^{-1}|(|A_D||\mathbf{x}|)$, $|A^{-1}|(|A_L\mathbf{x}|)$, and $|A^{-1}|(|A_U\mathbf{x}|)$ (in practice, the contribution of these terms to $\text{condeff}_{\mathbf{f}}(A(\Omega_{QS}), \mathbf{x})$ is computed as $|A^{-1}|(\mathbf{f} + |A_D||\mathbf{x}| + |A_L\mathbf{x}| + |A_U\mathbf{x}|)$). On the other hand in [6, Theorem 4.1] it is proved that the product R_1R_2 of an $\{n_1, m_1\}$ -quasiseparable matrix R_1 times an $\{n_2, m_2\}$ -quasiseparable matrix R_2 is, in general, an $\{n_1 + n_2, m_1 + m_2\}$ -quasiseparable matrix, and in [6, Algorithm 4.3] it is shown how to compute with linear complexity a quasiseparable representation for this product given the representations of the factors. Therefore, since our matrix A is a $\{1, 1\}$ -quasiseparable matrix, we have that the products $A^{-1}A_L$ and $A^{-1}A_U$ are both quasiseparable matrices of orders $\{2, 1\}$ and $\{1, 2\}$, respectively, at most, and we can compute via [6, Algorithm 4.3] their quasiseparable representations. Then, once we have obtained such representations, we can use again [6, Algorithm 4.4] for calculating the products $|A^{-1}A_L||\mathbf{x}|$ and $|A^{-1}A_U||\mathbf{x}|$ also with a linear cost. Then, from Definition 7.1, it is straightforward that $\text{condeff}_{\mathbf{f}}(A(\Omega_{QS}), \mathbf{x})$ can be computed with linear complexity. \square

Finally, note that from Definition 7.1 and from the proof of Proposition 4.10, it is straightforward to prove that the effective condition number is invariant under row scaling for $\mathbf{f} = |\mathbf{b}|$, as $\text{cond}_{|\mathbf{b}|}(A(\Omega_{QS}), \mathbf{x})$. This is stated without proof in Proposition 7.5.

Proposition 7.5. *Let $A\mathbf{x} = \mathbf{b}$, where $0 \neq \mathbf{x} \in \mathbb{R}^n$ and $A \in \mathbb{R}^{n \times n}$ is a $\{1, 1\}$ -quasiseparable matrix, and let $K \in \mathbb{R}^{n \times n}$ be an invertible diagonal matrix. Then*

$$\text{condeff}_{|\mathbf{b}|}(A(\Omega_{QS}), \mathbf{x}) = \text{condeff}_{|K\mathbf{b}|}((KA)(\Omega'_{QS}), \mathbf{x}),$$

where Ω_{QS} is any quasiseparable representation of A and Ω'_{QS} is any quasiseparable representation of KA .

8 Numerical experiments

This section is devoted to describe briefly some numerical experiments that have been performed in order to complete our comparison between the structured effective condition number $\text{condeff}_{|\mathbf{b}|}(A(\Omega_{QS}), \mathbf{x})$ in Definition 7.1 and the unstructured one $\text{cond}_{|A|, |\mathbf{b}|}(A, \mathbf{x})$ in Definition 2.3. We have used MATLAB for running several random numerical tests. First, the command `randn` from MATLAB has been used for generating the random parameters in a quasiseparable representation for a $\{1, 1\}$ -quasiseparable matrix of size $n \times n$, i.e., the following random vectors of parameters are generated:

$$\mathbf{p} \in \mathbb{R}^{n-1}, \mathbf{a} \in \mathbb{R}^{n-2}, \mathbf{q} \in \mathbb{R}^{n-1}, \mathbf{d} \in \mathbb{R}^n, \mathbf{g} \in \mathbb{R}^{n-1}, \mathbf{b} \in \mathbb{R}^{n-2}, \text{ and } \mathbf{h} \in \mathbb{R}^{n-1}. \quad (8.1)$$

We also generate the random right-hand side vector $\mathbf{b} \in \mathbb{R}^n$ by using the command `randn`. Then, we construct the matrix A described by the parameters in (8.1), obtain the vector of solutions \mathbf{x} using the command `A\b` from MATLAB, and compute the structured effective condition number $\text{condeff}_{|\mathbf{b}|}(A(\Omega_{QS}), \mathbf{x})$ and the unstructured condition number $\text{cond}_{|A|, |\mathbf{b}|}(A, \mathbf{x})$ by using direct matrix-vector multiplication and the `inv` command of MATLAB.

In general, when using just random parameters, we have obtained similar, often moderate, values for the effective condition number and the unstructured condition number for the solution of linear systems, i.e., $\text{condeff}_{|\mathbf{b}|}(A(\Omega_{QS}), \mathbf{x}) \approx \text{cond}_{|A|, |\mathbf{b}|}(A, \mathbf{x})$. Therefore, following the experiments in [5, Section 7], we have performed several tests using different kinds of scalings over the generated quasiseparable parameters in order to obtain unbalanced quasiseparable matrices which may be very ill conditioned.

In particular, after generating the vectors of parameters in (8.1) and the vector \mathbf{b} , we have modified \mathbf{p} and \mathbf{h} as follows

$$\mathbf{p} = k_1 * \mathbf{p} \quad \text{and} \quad \mathbf{h} = k_2^{-1} * \mathbf{h},$$

where k_1 and k_2 are fixed natural numbers not greater than 10^3 . This scaling, combined with the randomness of \mathbf{p} and \mathbf{h} and the rest of parameters, produces sometimes matrices with unbalanced lower left and upper right corners (see the matrix at the end of this section for an example), for which the unstructured condition number $\text{cond}_{|A|,|\mathbf{b}|}(A, \mathbf{x})$ can be much larger than the structured one $\text{condeff}_{|\mathbf{b}|}(A(\Omega_{QS}), \mathbf{x})$. In fact, for $n = 5$, $n = 10$, $n = 50$, and $n = 100$, we have obtained vectors of parameters generating particular matrices A and vectors \mathbf{b} such that:

n	$\frac{\text{cond}_{ A , \mathbf{b} }(A, \mathbf{x})}{\text{condeff}_{ \mathbf{b} }(A(\Omega_{QS}), \mathbf{x})}$	$\text{cond}_{ A , \mathbf{b} }(A, \mathbf{x})$	$\text{condeff}_{ \mathbf{b} }(A(\Omega_{QS}), \mathbf{x})$
5	$1.6139 \cdot 10^4$	$6.5318 \cdot 10^4$	4.0471
10	$1.9980 \cdot 10^6$	$3.1788 \cdot 10^7$	15.8768
50	$1.4107 \cdot 10^7$	$8.7762 \cdot 10^8$	62.2104
100	$1.6804 \cdot 10^9$	$6.0297 \cdot 10^{10}$	35.8823

where $\mathbf{x} = A \setminus \mathbf{b}$ in each case.

From these numerical experiments we conclude that the structured effective condition number $\text{condeff}_{|\mathbf{b}|}(A(\Omega_{QS}), \mathbf{x})$ may be indeed much smaller than the unstructured one $\text{cond}_{|A|,|\mathbf{b}|}(A, \mathbf{x})$, in other words, that there exist linear systems of equations with $\{1, 1\}$ -quasiseparable matrices of coefficients that have solutions which are very ill conditioned with respect to perturbations of the entries of the matrix, but that are very well conditioned with respect to perturbations on the quasiseparable parameters representing the matrix. The particular *structure* observed in the matrices that produced such huge differences between the structured and the unstructured condition numbers is illustrated in the following matrix and the respective vector \mathbf{b} , which are the ones that produced the results in the table above for $n = 5$:

$$A = \begin{bmatrix} -7.8876 \cdot 10^{-2} & -1.3485 \cdot 10^{-2} & -7.8066 \cdot 10^{-3} & 2.7951 \cdot 10^{-3} & 5.1089 \cdot 10^{-5} \\ 3.0423 \cdot 10^{-1} & -5.6399 \cdot 10^{-1} & 1.6206 \cdot 10^{-1} & -5.8026 \cdot 10^{-2} & -1.0606 \cdot 10^{-3} \\ 5.5451 \cdot 10^1 & -2.5873 \cdot 10^{-1} & 1.3525 \cdot 10^0 & -1.5088 \cdot 10^{-3} & -2.7578 \cdot 10^{-5} \\ -3.8947 \cdot 10^5 & 1.8172 \cdot 10^3 & -1.7047 \cdot 10^0 & 3.0944 \cdot 10^{-3} & -6.7069 \cdot 10^{-3} \\ 1.7714 \cdot 10^8 & -8.2653 \cdot 10^5 & 7.7535 \cdot 10^2 & 8.3875 \cdot 10^{-1} & -2.0998 \cdot 10^0 \end{bmatrix},$$

$$\mathbf{b} = [-8.8528 \cdot 10^{-1} \quad -1.3154 \cdot 10^{-1} \quad -1.5711 \cdot 10^0 \quad -7.8284 \cdot 10^{-1} \quad -1.0898 \cdot 10^0]^T.$$

Note that there is an obvious unbalance in the matrix entries, since the absolute values of the entries near the lower left corner are large compared with the absolute values of the entries in the opposite upper right corner of the matrix (compare the absolute values of the entries of the submatrix $A(4 : 5, 1 : 2)$ versus those from $A(1 : 2, 4 : 5)$).

9 Conclusions

A general expression for the condition number of the solution of a linear system of equations whose coefficient matrix is a differentiable function of a vector of parameters with respect to relative componentwise perturbations of such parameters has been presented. This expression involves the partial derivatives of the matrix with respect to the parameters. This result is related to results presented recently in [5] for eigenvalue condition numbers and both papers make use of differential calculus. This general expression has been used to deduce formulas for the componentwise condition numbers of the solutions of linear systems whose coefficient matrices are $\{1, 1\}$ -quasiseparable of size $n \times n$ with respect to perturbations of the parameters in any quasiseparable representation and in the tangent-Givens-vector representation of the coefficient matrices. We have compared theoretically these two structured condition numbers and we have proved that they differ at most by a factor $3n$ and, therefore, that they are numerically equivalent, though the one with respect to the tangent-Givens-vector representation is always the smallest. Moreover, it has been shown that these structured condition numbers can be estimated in $\mathcal{O}(n)$ operations via an effective condition number. We have also proved rigorously that these structured condition numbers are always smaller, up to a factor n , than the componentwise unstructured condition number. In addition, the performed numerical experiments illustrate that the structured condition numbers can be much smaller than the unstructured one in practice. This means that the structure of $\{1, 1\}$ -quasiseparable matrices may play a key role in the accuracy of computed solutions of linear systems of equations, since

these solutions can be much less sensitive to relative perturbations of the parameters representing the matrices than to relative perturbations of the matrix entries. The techniques used in this paper can be generalized to obtain structured condition numbers for the solution of linear systems involving other classes of low-rank structured matrices and they can be extended to study the structured conditioning of other problems involving low-rank structured matrices like, for instance, least squares problems.

References

- [1] AURENTZ, J.L., MACH, T., VANDEBRIL, R., WATKINS, D.S., *Fast and backward stable computation of roots of polynomials*, SIAM J. Matrix Anal. Appl., 36(3):942-973, 2015.
- [2] BELLA, T., OLSHEVSKY, V., STEWART, M., *Nested product decomposition of quasiseparable matrices*, SIAM J. Matrix Anal. Appl., 34(4):1520-1555, 2013.
- [3] BÖRM, S., GRASEDYCK, L., HACKBUSCH, W., *Introduction to hierarchical matrices with applications*, Eng. Anal. Bound. Elemen., 27:405-422, 2003.
- [4] DOPICO, F.M., OLSHEVSKY, V., ZHLOBICH, P., *Stability of QR-based fast system solvers for a subclass of quasiseparable rank one matrices*, Math. Comp., 82:2007-2034, 2013.
- [5] DOPICO, F.M., POMÉS, K., *Structured eigenvalue condition numbers for parameterized quasiseparable matrices*, published electronically in Numer. Math., DOI 10.1007/s00211-015-0779-5, 2015.
- [6] EIDELMAN, Y., GOHBERG, I., *On a new class of structured matrices*, Integral Equ. Oper. Theory, 34(3):293-324, 1999.
- [7] EIDELMAN, Y., GOHBERG, I., *Linear complexity inversion algorithms for a class of structured matrices*, Integral Equ. Oper. Theory, 35:28-52, 1999.
- [8] EIDELMAN, Y., GOHBERG, I., *On generators of quasiseparable finite block matrices*, Calcolo, 42:187-214, 2005.
- [9] EIDELMAN, Y., GOHBERG, I., HAIMOVICI, I., *Separable Type Representations of Matrices and Fast Algorithms. Volume 1. Basics. Completion Problems. Multiplication and Inversion Algorithms*, Operator Theory: Advances and Applications, 234. Birkhäuser/Springer, Basel, 2014.
- [10] EIDELMAN, Y., GOHBERG, I., HAIMOVICI, I., *Separable Type Representations of Matrices and Fast Algorithms. Volume 2. Eigenvalue Method*, Operator Theory: Advances and Applications, 235. Birkhäuser/Springer, Basel, 2014.
- [11] FERREIRA, C., PARLETT, B., DOPICO, F.M., *Sensitivity of eigenvalues of an unsymmetric tridiagonal matrix*, Numer. Math., 122(3):527-555, 2012.
- [12] HIGHAM, N.J., *Accuracy and Stability of Numerical Algorithms, 2nd ed.*, Society for Industrial and Applied Mathematics, Philadelphia, 2002.
- [13] RICE, J.R., *A theory of condition*, SIAM J. Numer. Anal., 3:287-310, 1966.
- [14] STEWART, M., *On orthogonal transformation of rank structured matrices*, submitted (available at <http://saaz.mathstat.gsu.edu/pub/pub.html>).
- [15] VANDEBRIL, R., VAN BAREL, M., MASTRONARDI, N., *A note on the representation and definition of semiseparable matrices*, Numer. Linear Algebra Appl., 12:839-858, 2005.
- [16] VANDEBRIL, R., VAN BAREL, M., MASTRONARDI, N., *Matrix Computations and Semiseparable Matrices. Volume 1. Linear Systems*, The Johns Hopkins University Press, Baltimore, 2008.
- [17] VANDEBRIL, R., VAN BAREL, M., MASTRONARDI, N., *Matrix Computations and Semiseparable Matrices. Volume 2. Eigenvalue and Singular value methods*, The Johns Hopkins University Press, Baltimore, 2008.
- [18] XI, Y., XIA, J., *On the stability of some hierarchical rank structured matrix algorithms*, submitted.
- [19] XI, Y., XIA, J., CAULEY, S., BALAKRISHNAN, V., *Superfast and stable structured solvers for Toeplitz least squares via randomized sampling*, SIAM J. Matrix Anal. Appl., 35(1):44-72, 2014.