

Scaling, sensitivity and stability in the numerical solution of quadratic eigenvalue problems

Nicholas J. Higham¹, D. Steven Mackey², Françoise Tisseur^{1,*,†}
and Seamus D. Garvey³

¹*School of Mathematics, The University of Manchester, Manchester M13 9PL, U.K.*

²*Department of Mathematics, Western Michigan University, Kalamazoo, MI 49008, U.S.A.*

³*School of Mechanical, Materials, Manufacturing Engineering and Management, University of Nottingham, University Park, Nottingham NG7 2RD, U.K.*

SUMMARY

The most common way of solving the quadratic eigenvalue problem (QEP) $(\lambda^2 M + \lambda D + K)x = 0$ is to convert it into a linear problem $(\lambda X + Y)z = 0$ of twice the dimension and solve the linear problem by the QZ algorithm or a Krylov method. In doing so, it is important to understand the influence of the linearization process on the accuracy and stability of the computed solution. We discuss these issues for three particular linearizations: the standard companion linearization and two linearizations that preserve symmetry in the problem. For illustration we employ a model QEP describing the motion of a beam simply supported at both ends and damped at the midpoint. We show that the above linearizations lead to poor numerical results for the beam problem, but that a two-parameter scaling proposed by Fan, Lin and Van Dooren cures the instabilities. We also show that half of the eigenvalues of the beam QEP are pure imaginary and are eigenvalues of the undamped problem. Our analysis makes use of recently developed theory explaining the sensitivity and stability of linearizations, the main conclusions of which are summarized. As well as arguing that scaling should routinely be used, we give guidance on how to choose a linearization and illustrate the practical value of condition numbers and backward errors. Copyright © 2007 John Wiley & Sons, Ltd.

Received 27 November 2006; Revised 13 March 2007; Accepted 21 March 2007

KEY WORDS: quadratic eigenvalue problem; sensitivity; condition number; backward error; stability; scaling; linearization; companion form; damped beam

*Correspondence to: Françoise Tisseur, School of Mathematics, The University of Manchester, Manchester M13 9PL, U.K.

†E-mail: ftisseur@ma.man.ac.uk

Contract/grant sponsor: Engineering and Physical Sciences Research Council; contract/grant numbers: GR/S31693, GR/S31679, EP/D079403

Contract/grant sponsor: Royal Society-Wolfson Research Merit Award (to the first author)

1. INTRODUCTION

The purpose of this paper is to emphasize the importance of *scaling* the coefficient matrices of second-order systems before numerically computing the eigenvalues *via* linearization. Our discussion is illustrated with a simple but nontrivial example consisting of a slender beam simply supported at both ends and damped at the midpoint, as shown in Figure 1. The equation of motion governing the transverse displacement $u(x, t)$ of the beam has the form

$$\rho A \frac{\partial^2 u}{\partial t^2} + c(x) \frac{\partial u}{\partial t} + EI \frac{\partial^4 u}{\partial x^4} = 0 \tag{1}$$

where ρA is the mass per unit length, $c(x) \geq 0$ represents the external damping and EI is the bending stiffness. The boundary conditions are $u(0, t) = u''(0, t) = 0$ and $u(L, t) = u''(L, t) = 0$, where L is the length of the beam. Note that this beam model does not include the effect of shear gradients or rotary inertia. Making the separation hypothesis $u(x, t) = e^{\lambda t} v(x)$ yields the boundary-value problem for the free vibrations

$$\begin{aligned} \lambda^2 \rho A v(x) + \lambda c(x) v(x) + EI \frac{d^4}{dx^4} v(x) &= 0 \\ v(0) = v''(0) = v(L) = v''(L) &= 0 \end{aligned} \tag{2}$$

Since the beam system is perfectly symmetric, its vibration behaviour divides into two distinct sets of modes: the symmetric and the anti-symmetric (see Appendix A.1). The symmetric modes lie in the eigenspace associated with the eigenvalue $+1$ of the reflection-across-the-midpoint symmetry of the system. At the midpoint they have zero slope and (to avoid triviality) nonzero displacement. Anti-symmetric modes are the ones lying in the eigenspace corresponding to the eigenvalue -1 of the reflection-across-the-midpoint symmetry of the system, and hence must have zero displacement and nonzero slope at the midpoint. The damper is irrelevant to anti-symmetric modes since they all have zero displacement at the midpoint. Thus these modes feel zero damping, which implies that the corresponding eigenvalues λ are pure imaginary.

We discretize the boundary-value problem (2) by finite elements using cubic Hermite polynomials as interpolation shape functions. This gives the finite-dimensional quadratic eigenvalue problem (QEP)

$$Q(\lambda)x = (\lambda^2 M + \lambda D + K)x = 0 \tag{3}$$

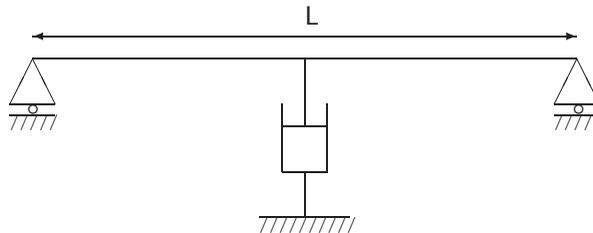


Figure 1. Beam simply supported at both ends and damped at the midpoint.

The resulting mass matrix M and stiffness matrix K are symmetric positive definite ($M > 0$, $K > 0$) by construction and D is symmetric positive semidefinite ($D \geq 0$). As a consequence, the roots of the quadratic equation $x^*Q(\lambda)x = 0$ have nonpositive real parts for all vectors x . This implies that all the eigenvalues of (3) lie in the closed left half plane and the beam problem is (weakly) stable [1]. The finite element discretization that we use preserves the property of the continuous problem of having all the modes either symmetric or anti-symmetric (see the proof of Theorem A1 in Appendix A.1).

The standard approach to the numerical solution of the QEP is to convert the quadratic $Q(\lambda) = \lambda^2 M + \lambda D + K$ into a linear polynomial

$$L(\lambda) = \lambda X + Y$$

of twice the dimension of Q but with the same spectrum. The resulting generalized eigenproblem $L(\lambda)z = 0$ is usually solved by the QZ algorithm for small- to medium-size problems or by a Krylov method for large sparse problems [2, 3]. A common choice of L in practice is the first companion form, given by

$$C_1(\lambda) = \lambda \begin{bmatrix} M & 0 \\ 0 & I \end{bmatrix} + \begin{bmatrix} D & K \\ -I & 0 \end{bmatrix} \quad (4)$$

It can be shown that $C_1(\lambda)$ is always a linearization in the sense that it satisfies

$$E(\lambda)C_1(\lambda)F(\lambda) = \begin{bmatrix} Q(\lambda) & 0 \\ 0 & I \end{bmatrix}$$

for some $E(\lambda)$ and $F(\lambda)$ with constant, nonzero determinants [4, Section 7.2]. This implies that $\alpha \det(C_1(\lambda)) = \det(Q(\lambda))$ for some nonzero constant α , so that C_1 and Q have the same spectrum. When K and M , respectively, are nonsingular the two pencils

$$L_1(\lambda) = \lambda \begin{bmatrix} M & 0 \\ 0 & -K \end{bmatrix} + \begin{bmatrix} D & K \\ K & 0 \end{bmatrix} \quad (5)$$

and

$$L_2(\lambda) = \lambda \begin{bmatrix} 0 & M \\ M & D \end{bmatrix} + \begin{bmatrix} -M & 0 \\ 0 & K \end{bmatrix} \quad (6)$$

are other possible linearizations [1, 5, 6]. Note that when M , D and K are symmetric, the pencils L_1 and L_2 are symmetric.

We computed the eigenvalues of the discretized beam problem in (3) for the following geometric and material properties:

$$E = 7 \times 10^{10} \text{ N m}^{-2}, \quad I = \frac{0.05 \times 0.005^3}{12} \text{ m}^4 \quad (7a)$$

$$L = 1 \text{ m}, \quad \rho AL = 0.674 \text{ kg}, \quad c = 5 \text{ kg s}^{-1} \quad (7b)$$

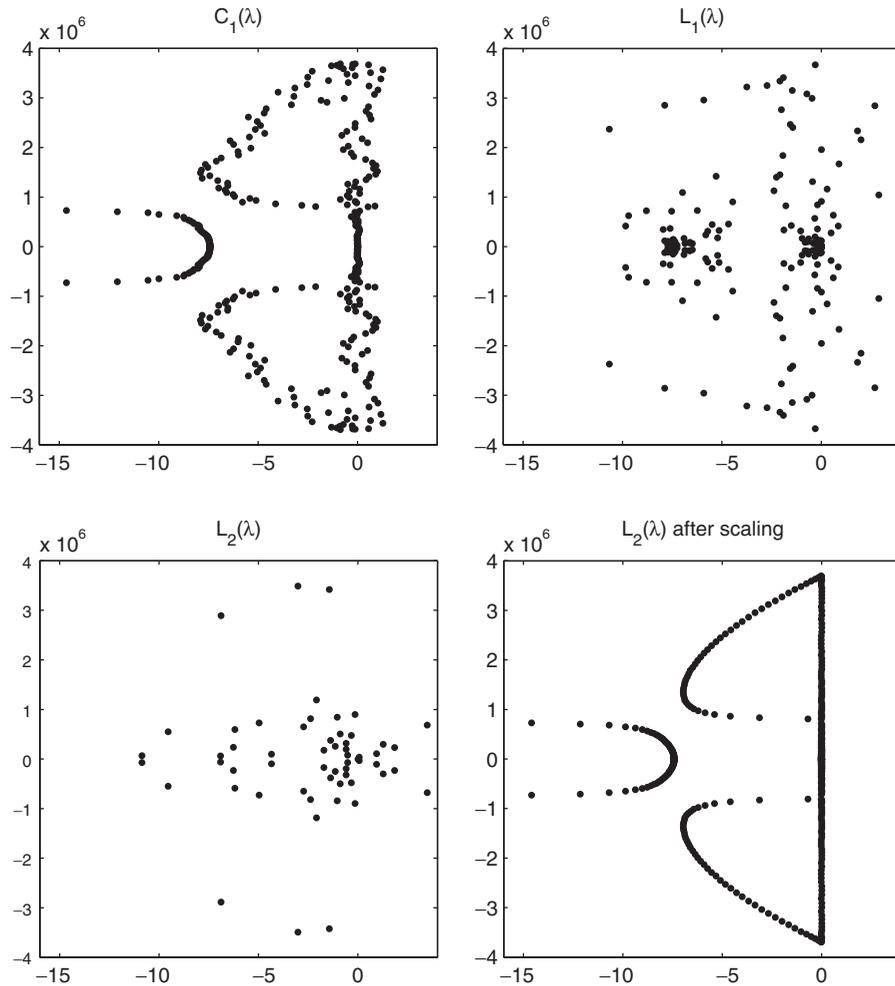


Figure 2. Beam problem discretized with 100 finite elements. Computed eigenvalues λ with $\text{Re}(\lambda) \in [-16, 4]$ of the linearizations $C_1(\lambda)$, $L_1(\lambda)$ and $L_2(\lambda)$ defined in (4)–(6).

We used 100 beam elements, which result in matrices M , D and K of dimension $n = 200$. The eigenvalues were computed by calling MATLAB's function `eig`, which implements the QZ algorithm, on each of the three linearizations (4)–(6). The first three plots in Figure 2 display those computed eigenvalues having real parts in the interval $[-16, 4]$; it follows from the analysis below that the real parts of the exact eigenvalues all lie in this range. Since (4)–(6) are all linearizations of Q , the first three plots should be identical, but in fact they are very different, and none of them correctly display the real part of the spectrum of Q to visual accuracy. These three plots have one feature in common: some eigenvalues lie in the right half plane; therefore, implying that the discretized beam problem is unstable and seems to contradict the theory.

Now let us convert $Q(\lambda) = \lambda^2 M + \lambda D + K$ to $\tilde{Q}(\mu) = \delta Q(\gamma\mu) = \mu^2 \tilde{M} + \mu \tilde{D} + \tilde{K}$, where

$$\lambda = \gamma\mu, \quad \tilde{M} = \gamma^2 \delta M, \quad \tilde{D} = \gamma \delta D, \quad \tilde{K} = \delta K \quad (8a)$$

$$\gamma = \sqrt{k/m}, \quad \delta = 2/(k + d\gamma) \quad (8b)$$

and

$$m = \|M\|_2, \quad d = \|D\|_2, \quad k = \|K\|_2 \quad (9)$$

(the 2-norm is defined in Section 2). This scaling, with its two parameters γ and δ , was proposed by Fan *et al.* [7]. Note that the scaling does not affect any sparsity of M , D and K . The eigenvalues μ of $\tilde{Q}(\mu)$ are then computed by calling `eig` on each of the three linearizations (4)–(6) with the scaled matrices \tilde{M} , \tilde{D} and \tilde{K} in place of M , D and K . The eigenvalues of $Q(\lambda)$ are recovered from those of $\tilde{Q}(\mu)$ via $\lambda = \gamma\mu$. The last plot in Figure 2 shows the spectrum of $Q(\lambda)$ computed using the linearization L_2 after scaling; all three linearizations (4)–(6) yield similar plots after scaling. Note that all the eigenvalues are now in the left half plane and that many of the eigenvalues appear to be pure imaginary, as expected by the theory. Indeed we prove in Appendix A.1 that for the discretized beam problem half of the eigenvalues of Q are pure imaginary and that they coincide with half of the eigenvalues of the undamped quadratic $\lambda^2 M + K$. It is therefore reasonable to believe—and we will be able to conclude from our analysis—that the fourth plot in Figure 2 is a good approximation of the spectrum of Q , unlike the first three plots.

In the rest of this paper we give a theoretical explanation of these somewhat surprising numerical results by making use of some recently developed theory concerning the sensitivity and stability of linearizations. We hope to convince the engineering community of

- the importance of scaling QEPs before computing the eigenvalues *via* linearization;
- the practical value of condition numbers and backward errors for understanding the quality of the computed results.

The results herein are not confined to the beam problem but apply to any QEP.

2. SENSITIVITY AND STABILITY OF LINEARIZATIONS

Backward errors and condition numbers play an important role in modern numerical linear algebra. A condition number measures the sensitivity of the solution of a problem to perturbations in the data, whereas a backward error measures how far a problem has to be perturbed for an approximate solution to be an exact solution of the perturbed problem. Thus conditioning is a property of the problem, while backward error characterizes the stability of a method for solving the problem. Backward error and conditioning are complementary concepts: when combined with a backward error estimate, a condition number provides an approximate upper bound on the error in a computed solution. Indeed with consistent definitions we have the rule of thumb that

$$\text{error in solution} \lesssim \text{condition number} \times \text{backward error} \quad (10)$$

For a given quadratic Q , infinitely many linearizations exist (of which (4)–(6) are just three) [6]. They can have widely varying eigenvalue condition numbers [8], and approximate eigenpairs of

$Q(\lambda)$ computed *via* linearization can have widely varying backward errors [9]. Ideally, we would like the linearization L that we use to be as well conditioned as the original quadratic Q and for it to lead, after recovering an approximate eigenpair of Q from one of L , to a backward error of the same order of magnitude as that for L . In the following subsections we will define the terms condition number and backward error more precisely and investigate their size for different linearizations of a quadratic. Then we show how the scaling (8) can improve both conditioning and backward error.

We will use the 2-norm, defined for a vector x by $\|x\|_2 = (x^*x)^{1/2}$ and for a matrix A by $\|A\|_2 = \max_{x \neq 0} \|Ax\|_2 / \|x\|_2$.

2.1. Eigenvalue condition number

A normwise relative condition number of a simple, finite, nonzero eigenvalue λ of Q with corresponding right eigenvector x and left eigenvector y can be defined by

$$\kappa_Q(\lambda) = \limsup_{\varepsilon \rightarrow 0} \left\{ \frac{|\Delta\lambda|}{\varepsilon|\lambda|} : (Q(\lambda + \Delta\lambda) + \Delta Q(\lambda + \Delta\lambda))(x + \Delta x) = 0, \right. \\ \left. \|\Delta M\|_2 \leq \varepsilon m, \|\Delta D\|_2 \leq \varepsilon d, \|\Delta K\|_2 \leq \varepsilon k \right\}$$

where $\Delta Q(\lambda) = \lambda^2 \Delta M + \lambda \Delta D + \Delta K$ and, as in (9), $m = \|M\|_2, d = \|D\|_2, k = \|K\|_2$. The eigenvalue condition number $\kappa_L(\lambda)$ for the pencil $L(\lambda) = \lambda X + Y$ is defined in a similar way. Explicit formulae for these condition numbers are given by [10, Theorem 5]

$$\kappa_Q(\lambda) = \frac{(|\lambda|^2 m + |\lambda|d + k)\|y\|_2\|x\|_2}{|\lambda|\|y^*(2\lambda M + D)x|}, \quad \kappa_L(\lambda) = \frac{(|\lambda|\|X\|_2 + \|Y\|_2)\|w\|_2\|z\|_2}{|\lambda|\|w^*Xz|} \tag{11}$$

Here, w and z denote left and right eigenvectors of L corresponding to λ . (For a very readable introduction to eigenvalue condition numbers for the standard eigenvalue problem $Ax = \lambda x$ see Davis and Moler [11].)

Higham *et al.* [8] have recently investigated the conditioning of a large class of linearizations for matrix polynomials of arbitrary degree. Define $\phi_L(\lambda)$ by

$$\kappa_L(\lambda) = \phi_L(\lambda)\kappa_Q(\lambda)$$

ϕ_L can be regarded as a growth factor in the translation from conditioning for Q to conditioning for L . Ideally, we would like $\phi_L(\lambda) \approx 1$. From the analysis in [8] we obtain for the three linearizations (4)–(6) explicit approximations for $\phi_L(\lambda)$ (see Appendix A.2.1). From these, we give in Table I conditions that guarantee $\kappa_L(\lambda) \approx \kappa_Q(\lambda)$. Here, the quantity

$$\rho := \frac{\max(m, d, k)}{\min(m, k)} \tag{12}$$

measures the scaling of the problem. Table I shows that for well-scaled problems (i.e. $\rho \approx 1$) the two symmetric linearizations L_1 and L_2 are optimally conditioned for large $|\lambda|$ and small $|\lambda|$, respectively. Note that the sufficient condition $m \approx d \approx k \approx 1$ for C_1 to be optimally conditioned is more stringent than the requirement $\rho \approx 1$.

Table I. Sufficient conditions for $\kappa_Q \approx \kappa_L$; m, d, k are as in (9), ρ is defined in (12).

Linearization	Eigenvalue	Condition
C_1 in (4)	No restriction	$m \approx d \approx k \approx 1$
L_1 in (5)	$ \lambda \gtrsim 1$ $ \lambda \ll 1$	$\rho \approx 1$ 'Not available'
L_2 in (6)	$ \lambda \gtrsim 1$ $ \lambda \ll 1$	'Not available' $\rho \approx 1$

Table II. Approximations (A4)–(A6) to growth factors $\phi_L(\lambda)$ for the beam problem.

$ \lambda $	$\phi_{C_1}(\lambda)$	$\phi_{L_1}(\lambda)$	$\phi_{L_2}(\lambda)$
10^2	1×10^2	1×10^4	1×10^4
10^4	1×10^4	1×10^8	1×10^8
10^6	2×10^5	2×10^{11}	2×10^{11}

For the particular instance of the beam problem described in Section 1 the matrices vary widely in norm: the norms of the mass, damping and stiffness matrices are

$$m = 6.7 \times 10^{-3}, \quad d = 5, \quad k = 1.7 \times 10^9$$

Thus $\rho = 2.6 \times 10^{11}$, and so the beam problem is badly scaled. The moduli of the eigenvalues lie in the interval $(7 \times 10^1, 4 \times 10^6)$. Approximate values of $\phi_L(\lambda)$ are given in Table II for $|\lambda|$ ranging from 10^2 to 10^6 ; they indicate that the eigenvalues of all three linearizations (4)–(6) are much more sensitive to perturbations than those of the original quadratic $Q(\lambda)$.

To understand how the theory relates to the results in Figure 2, note that from the definition of $\kappa_L(\lambda)$ a relative perturbation of order ε in $L(\lambda)$ can perturb λ by

$$|\Delta\lambda| \lesssim \varepsilon |\lambda| \kappa_L(\lambda) = \varepsilon |\lambda| \phi_L(\lambda) \kappa_Q(\lambda)$$

From Table II, for $L = L_1$ and $|\lambda| = 10^6$, and with $\varepsilon \approx 10^{-16}$, which represents the perturbation introduced by the QZ algorithm in MATLAB's floating point arithmetic, we have

$$|\Delta\lambda| \lesssim 10^{-16} \times 10^6 \times 10^{11} \times \kappa_Q(\lambda) = 10 \kappa_Q(\lambda)$$

Since $\kappa_Q(\lambda) \approx 10^6$ for the largest eigenvalues of Q , we see that the eigenvalues on the imaginary axis can be perturbed by a large distance into the right half plane (and the perturbations in Figure 2 are far from the worst case).

The values in Table II show that the first companion linearization C_1 is better conditioned than the two symmetric linearizations L_1 and L_2 . This explains the slightly better looking plot for C_1 in Figure 2.

In particular applications some eigenvalues may be more important than others. For the beam problem the small eigenvalues are, like the large ones, quite ill conditioned, with $\kappa_Q(\lambda) \approx 10^8$ for $|\lambda| \approx 10^2$. Table II shows that the corresponding eigenvalue condition numbers for the linearized problem are magnified by a factor 10^2 for $L = C_1$ and 10^4 for $L = L_1$ or L_2 , which is a nonnegligible increase.

2.2. Backward error

Consider an approximate eigenpair $(\hat{x}, \hat{\lambda})$ of $Q(\lambda)$ with $\hat{\lambda}$ finite. We can interpret $(\hat{x}, \hat{\lambda})$ as the exact eigenpair of a perturbed quadratic $(Q + \Delta Q)(\lambda) = \lambda^2(M + \Delta M) + \lambda(D + \Delta D) + K + \Delta K$, where there are many possible choices of ΔM , ΔD and ΔK . We define the backward error of $(\hat{x}, \hat{\lambda})$ to be the size of the smallest of all such perturbations

$$\eta_Q(\hat{x}, \hat{\lambda}) = \min\{\varepsilon : (Q + \Delta Q)(\hat{\lambda})\hat{x} = 0, \|\Delta M\|_2 \leq \varepsilon m, \|\Delta D\|_2 \leq \varepsilon d, \|\Delta K\|_2 \leq \varepsilon k\}$$

An analogous definition holds for the backward error $\eta_L(\hat{z}, \hat{\lambda})$ of an approximate eigenpair $(\hat{z}, \hat{\lambda})$ of the pencil $L(\lambda)$. The explicit formulae [10, Theorem 1]

$$\eta_Q(\hat{x}, \hat{\lambda}) = \frac{\|Q(\hat{\lambda})\hat{x}\|_2}{(|\hat{\lambda}|^2 m + |\hat{\lambda}|d + k)\|\hat{x}\|_2}, \quad \eta_L(\hat{z}, \hat{\lambda}) = \frac{\|L(\hat{\lambda})\hat{z}\|_2}{(|\hat{\lambda}|\|X\|_2 + \|Y\|_2)\|\hat{z}\|_2} \tag{13}$$

show that the backward errors of $(\hat{x}, \hat{\lambda})$ and $(\hat{z}, \hat{\lambda})$ are scaled residuals.

Two main factors affect the backward error. First, because L is usually highly structured (see, for instance, (4)–(6)), perturbations to L cannot directly be interpreted as equivalent perturbations to Q . Second, the ‘short’ eigenvectors of Q can be recovered from the ‘long’ eigenvectors of L in many ways, with differing implications on the backward error for Q . For all the linearizations in (4)–(6) the right eigenvectors z of the linearization have the form

$$z = \begin{bmatrix} \lambda x \\ x \end{bmatrix} \tag{14}$$

where x is a right eigenvector of Q , hence we can recover eigenvectors of Q from either of the components

$$z_1 = z(1:n) \text{ (if } \lambda \neq 0), \quad z_2 = z(n+1:2n)$$

In the analysis below we will make the reasonable assumption that (14) remains approximately true for the approximate eigenpair, to the extent that the ratio $\|\hat{z}\|_2/\|\hat{z}_i\|_2$ is of order 1 for $i = 1$ when $|\hat{\lambda}| \geq 1$ and for $i = 2$ when $|\hat{\lambda}| \leq 1$.

The backward error properties of a large class of linearizations have recently been investigated by Higham *et al.* [9]. For the linearizations (4)–(6) they obtain bounds of the form

$$\frac{1}{2} \leq \frac{\eta_Q(\hat{z}_i, \hat{\lambda})}{\eta_L(\hat{z}, \hat{\lambda})} \leq c \psi_L^{(i)} \frac{|\lambda|^2 + 1}{|\lambda|^2 m + |\lambda|d + k} \frac{\|\hat{z}\|_2}{\|\hat{z}_i\|_2}, \quad i = 1, 2 \tag{15}$$

(see Appendix A.3 for explicit expressions for ψ_L), where c is a constant of order 1 that depends on the linearization, and they derive the sufficient conditions for $\eta_Q \approx \eta_L$ that are summarized

Table III. Sufficient conditions for $\eta_Q \approx \eta_L$; ρ is defined in (12) and m, d, k in (9).

Linearization	Eigenvalue	Right eigenvector	Condition
C_1 in (4)	$ \hat{\lambda} \geq 1$ $ \hat{\lambda} \leq 1$	\hat{z}_1 \hat{z}_2	$d \leq m \approx k \approx 1$
L_1 in (5)	$ \hat{\lambda} \geq 1$ $ \hat{\lambda} \leq 1$	\hat{z}_1 \hat{z}_2	$\rho \approx 1$ $\rho \max(1, (m+d)\ K^{-1}\ _2) \approx 1$
L_2 in (6)	$ \hat{\lambda} \geq 1$ $ \hat{\lambda} \leq 1$	\hat{z}_1 \hat{z}_2	$\rho \max(1, (d+k)\ M^{-1}\ _2) \approx 1$ $\rho \approx 1$

in Table III. We see that for well-scaled problems (i.e. $\rho \approx 1$) the two symmetric linearizations L_1 and L_2 are optimally stable for large $|\hat{\lambda}|$ and small $|\hat{\lambda}|$, respectively. This is entirely consistent with the sufficient conditions for optimal conditioning given in Table I. Note that which portion of z must be used to recover the eigenvector of Q depends on the size of $|\hat{\lambda}|$, as specified in the third column of the table.

For our beam problem example, $\psi_L^{(i)} \gg 1$, $i = 1, 2$, for all the linearizations (4)–(6), indicating that these three linearizations are potentially unstable. We found that for all three linearizations $\eta_L(\hat{z}, \hat{\lambda}) \leq 10^{-15} \approx nu$ for all computed eigenpairs, where $u \approx 1.1 \times 10^{-16}$ is the unit roundoff. This is not surprising since the QZ algorithm for the generalized eigenvalue problem is backward stable. However, we found that the ratio $\eta_Q(\hat{z}_1, \hat{\lambda})/\eta_L(\hat{z}, \hat{\lambda})$ can be as large as 10^7 for the companion linearization and as large as 10^{12} for the other linearizations. Therefore, none of the three linearizations produces computed eigenpairs with satisfactory backward error for the beam example, which again is consistent with the poor results observed in Figure 2.

2.3. Scaled quadratics

In view of the sufficient conditions for $\kappa_L \approx \kappa_Q$ and $\eta_L \approx \eta_Q$ given in Tables I and III it is natural to scale the problem to try to bring the 2-norms of M , D , and K close to 1 in order to ameliorate any sensitivity or instability induced by the linearization process. The scaling of Fan *et al.* [7] defined in (8) has precisely this aim. Note that $\kappa_Q(\lambda) = \kappa_{\tilde{Q}}(\mu)$ and $\eta_Q(x, \lambda) = \eta_{\tilde{Q}}(x, \mu)$, where $\tilde{Q}(\mu) = \mu^2 \tilde{M} + \mu \tilde{D} + \tilde{K}$ is the scaled quadratic in (8) and $\lambda = \gamma\mu$, so this scaling has no effect on the condition number and backward error for the quadratic; its purpose is to improve the condition number of the linearization L and backward error of the eigenpairs of Q obtained from L . That it might do so is clear from the fact that scaling generally decreases ρ : it can be shown [9] that $\tilde{\rho}$ defined as in (12) for the scaled problem satisfies $\tilde{\rho} = \max(1, \tau) \leq \rho$, where

$$\tau := \frac{d}{\sqrt{mk}}$$

For the scaled quadratic $\tilde{Q}(\mu)$ the analyses of Sections 2.1 and 2.2 simplify. We define the key quantity

$$\omega(\mu) := \frac{1 + \tau}{1 + (|\mu|/(1 + |\mu|^2))\tau} \quad (16)$$

Let $(\widehat{z}, \widehat{\mu})$ be an approximate eigenpair of a linearization for \widetilde{Q} with \widehat{z} partitioned as $\widehat{z} = \begin{bmatrix} \widehat{z}_1 \\ \widehat{z}_2 \end{bmatrix}$ and $\mu, \widehat{\mu}$ are assumed to be finite. For the linearizations (4)–(6) we can show that

$$\frac{\kappa_L(\mu)}{\kappa_{\widetilde{Q}}(\mu)} \approx \widetilde{\phi}_L(\mu) = \begin{cases} \omega(\mu), & L = C_1 \\ \frac{1 + |\mu|}{|\mu|} \omega(\mu), & L = L_1 \\ (1 + |\mu|) \omega(\mu), & L = L_2 \end{cases} \tag{17}$$

and

$$\frac{\eta_{\widetilde{Q}}(\widehat{z}_k, \widehat{\mu})}{\eta_L(\widehat{z}, \widehat{\mu})} \lesssim \widetilde{\psi}_L(\widehat{z}_k) \frac{\|\widehat{z}\|_2}{\|\widehat{z}_k\|_2} \tag{18}$$

where

$$\begin{aligned} \widetilde{\psi}_{C_1}(\widehat{z}_k) &= \omega(\widehat{\mu}), \quad k = 1:2 \\ \widetilde{\psi}_{L_1}(\widehat{z}_k) &= \begin{cases} \omega(\widehat{\mu}), & k = 1 \\ \|\widetilde{K}^{-1}\|_2 \omega(\widehat{\mu}), & k = 2 \end{cases} \\ \widetilde{\psi}_{L_2}(\widehat{z}_k) &= \begin{cases} \|\widetilde{M}^{-1}\|_2 \omega(\widehat{\mu}), & k = 1 \\ \omega(\widehat{\mu}), & k = 2 \end{cases} \end{aligned} \tag{19}$$

The values for $\widetilde{\phi}_L(\mu)$ are obtained from (A4)–(A6) using $|\mu|^2 \widetilde{m} + |\mu| \widetilde{d} + \widetilde{k} = 2(1 + |\mu|^2) / \omega(\mu)$, where $\widetilde{m} = \|\widetilde{M}\|_2$, $\widetilde{d} = \|\widetilde{D}\|_2$ and $\widetilde{k} = \|\widetilde{K}\|_2$ can be shown to satisfy $\widetilde{m} = \widetilde{k} = 2 / (1 + \tau)$, $\widetilde{d} = 2\tau / (1 + \tau)$. The values for $\widetilde{\psi}_L(\widehat{\mu})$ come from [9, Theorem 5.1]. Hence, we conclude that for the companion linearization C_1 , and for L_1 (if $|\mu| \geq 1$) or L_2 (if $|\mu| \leq 1$), both backward error and conditioning are essentially optimal for the scaled problem if $\omega(\mu) = O(1)$. The quantity $\omega(\mu)$ satisfies the bounds

$$1 \leq \omega(\mu) \leq \min \left\{ 1 + \tau, \frac{1 + |\mu|^2}{|\mu|} \right\} \tag{20}$$

Hence, $\omega(\mu) = O(1)$ if $\tau = O(1)$ or if $|\mu| = O(1)$. Note that if

$$d \lesssim \sqrt{mk} \tag{21}$$

then $\tau \lesssim 1$ and hence we are guaranteed that $\omega = O(1)$. In the terminology of systems originating from mechanical systems with damping, condition (21) holds for systems that are not too heavily damped. A class of problems for which (21) is satisfied is the elliptic Q [12, 13]: those for which M is Hermitian positive definite, D and K are Hermitian, and $(x^* D x)^2 < 4(x^* M x)(x^* K x)$ for all nonzero $x \in \mathbb{C}^n$.

The following bounds on $|\mu|$ (expressed in terms of the unscaled matrices) can be derived from [14, Lemma 3.1]:

$$\frac{1}{2} \tau \left(-1 + \sqrt{1 + 4 / (\tau^2 \kappa_2(K))} \right) \leq |\mu| \leq \frac{1}{2} \tau \kappa_2(M) \left(1 + \sqrt{1 + 4 / (\tau^2 \kappa_2(M))} \right) \tag{22}$$

These bounds depend on the condition numbers $\kappa_2(M) = \|M\|_2 \|M^{-1}\|_2$ and $\kappa_2(K) = \|K\|_2 \|K^{-1}\|_2$. In particular, if M is well conditioned and $\tau = O(1)$, the upper bound is close to 1 and we can safely use the symmetric linearization L_2 .

Turning to the beam problem, the Fan, Lin and Van Dooren scaling (8) yields

$$\|\tilde{M}\|_2 = \|\tilde{K}\|_2 \approx 2, \quad \|\tilde{D}\|_2 \approx 3 \times 10^{-3}$$

For this problem the scaling does not quite reach its aim of bringing the norms of all three matrices to 1. Nevertheless, since the problem is not too heavily damped, (21) holds, so that $\omega(\mu) = O(1)$. Thus, the theory above guarantees optimal conditioning and optimal stability for the companion linearization. Because M is ill conditioned with $\kappa_2(M) = 2.6 \times 10^6$, the upper bound in (22) is not very informative. But if we apply (22) to the diagonally scaled quadratic $H\tilde{Q}H$ (which has the same eigenvalues as \tilde{Q}), where $H = \text{diag}(\tilde{m}_{11}^{1/2}, \dots, \tilde{m}_{22}^{1/2})$, then (22) produces the surprisingly good bounds

$$1.4 \times 10^{-5} \leq |\mu| \leq 7.25$$

the correct interval being $[1.43 \times 10^{-4}, 7.2461]$. This tells us that the symmetric linearization L_2 in (6) is optimal in terms of both conditioning and stability. Finally, we can use the rule of thumb (10) to bound the errors in our computed eigenvalues. The condition numbers $\kappa_Q(\lambda)$ in (11) are readily computed and are found to be at most 10^8 . Since $\eta_Q \approx \eta_L \approx 10^{-15}$ for C_1 and L_2 on the scaled problem, we know that the relative error is at most about $10^8 \times 10^{-15} = 10^{-7}$. Thus, without knowing the exact eigenvalues we can be sure that all the computed eigenvalues have relative error at most about 10^{-7} , confirming that the last plot in Figure 2 is correct to (much more than) visual accuracy.

3. CONCLUSIONS

Through the specific example of the beam problem we have shown that solving a QEP by linearization can yield unsatisfactory computed eigenvalues, both quantitatively and qualitatively. To improve the quality of computed eigenvalues in general we recommend an initial scaling of the quadratic using the scaling (8) of Fan, Lin and Van Dooren. Although this scaling employs only two parameters we have shown that it can produce essentially optimal backward error and conditioning properties of the linearization. Indeed this scaling usually reduces (and never increases) the scaling factor ρ in (12), which is a key quantity in measuring the sensitivity and stability of the linearization process. For the beam problem we have shown that while the computed eigenvalues of the unscaled problem have large errors in the ‘eyeball norm’, scaling achieves optimal backward error and conditioning for C_1 and L_2 and leads to computed eigenvalues with at least seven correct significant digits, the improvements being clear from Figure 2. Scaling requires approximations of the 2-norms of the mass, damping and stiffness matrices; these can be obtained by the power method or the Lanczos method [2], or simply by computing a different norm, such as the Frobenius norm $\|A\|_F = (\sum_{i,j} |a_{ij}|^2)^{1/2}$.

Turning to the choice of linearization once the QEP is scaled, the theory summarized in (17)–(19) shows that the companion form for all eigenvalues, L_1 for large eigenvalues, and L_2 for small eigenvalues all satisfy essentially the same conditioning and backward error bounds, and

that these bounds are optimal for QEPs that are not too heavily damped. The companion form has the advantage that it is always a linearization. If preserving symmetry is an issue for reasons of storage or computational cost then L_1 and L_2 are attractive, and the choice between them should be guided by the fact that they favour large and small eigenvalues, respectively. However, if the condition number of K or M is small (these condition numbers can be estimated without explicitly computing the inverses [15, Chapter 15; 16]) then L_1 or L_2 , respectively, can safely be used to stably obtain all the eigenpairs, as shown by the conditions in (19).

Finally, it is worth emphasizing that when only a subset of the spectrum is required we should try to choose a linearization for which the physically important eigenvalues are optimally conditioned. Imposing this requirement on just a subset of the spectrum provides more flexibility in the choice of linearization and the analysis summarized here provides help in making such a choice.

APPENDIX A

A.1. On the spectrum of the beam problem

The fourth plot in Figure 2 strongly suggests that many of the eigenvalues of the QEP (3) for the beam problem lie on the imaginary axis. We prove in this Appendix that this is indeed the case, thereby providing further evidence that this plot accurately portrays the true eigenvalues of (3). To do this we need a more precise description of the structure of the matrices in (3).

The mass and stiffness matrices are obtained by a finite element discretization of (2). The beam is divided into ℓ (even) elements. Each beam element has two end nodes and four degrees of freedom. These degrees of freedom are collected in the node displacement vector $[u_1 \ \theta_1 \ u_2 \ \theta_2]^T$, where u_1, u_2 are the transverse displacements and θ_1, θ_2 the slopes of the displacements at node 1 and node 2, respectively. Cubic Hermite polynomials are used as interpolation shape functions. The beam element stiffness matrix K_e and the beam element consistent mass matrix M_e are well known [17]

$$K_e = \frac{2(EI)}{L_e^3} \begin{bmatrix} 6 & 3L_e & -6 & 3L_e \\ 3L_e & 2L_e^2 & -3L_e & L_e^2 \\ -6 & -3L_e & 6 & -3L_e \\ 3L_e & L_e^2 & -3L_e & 2L_e^2 \end{bmatrix}$$

$$M_e = \frac{\rho A L_e}{420} \begin{bmatrix} 156 & 22L_e & 54 & -13L_e \\ 22L_e & 4L_e^2 & 13L_e & -3L_e^2 \\ 54 & 13L_e & 156 & -22L_e \\ -13L_e & -3L_e^2 & -22L_e & 4L_e^2 \end{bmatrix}$$

where L_e is the length of the finite element e . Now we assume that ℓ beam elements are all of equal length. With the unknowns ordered as

$$[\theta_1 \ u_2 \ \theta_2 \ \dots \ u_\ell \ \theta_\ell \ \theta_{\ell+1}]^T$$

Letting

$$\check{S} = \begin{bmatrix} & & & R \\ & & \ddots & \\ & R & & \\ R & & & \end{bmatrix}_{(n+2) \times (n+2)}$$

one can use relations (A2) to show that \check{S} commutes with \check{K} . Deleting the first row and column and the next-to-last row and column of \check{S} then produces an $n \times n$ matrix

$$S = \begin{bmatrix} & & & -1 \\ & & R & \\ & & \ddots & \\ R & & & \\ -1 & & & \end{bmatrix}_{n \times n}$$

that commutes with K and M . This S can be viewed as a left/right mirror image symmetry of the beam that interchanges node pairs that are symmetrically placed with respect to the midpoint of the beam.

The fact that S commutes with K and M can now be exploited to simultaneously block-diagonalize K and M . Observe that S is diagonalizable since $S = S^T$, and has only the eigenvalues $\lambda = \pm 1$ since $S^2 = I$. Thus \mathbb{R}^n decomposes as the orthogonal direct sum of the two eigenspaces of S , the $\lambda = 1$ eigenspace comprising the ‘symmetric modes’ of the discretized beam problem, and the $\lambda = -1$ eigenspace corresponding to the ‘anti-symmetric modes’. Now since $\text{trace}(S) = 0$ there are $\ell = n/2$ eigenvalues 1 and ℓ eigenvalues -1 . Let $\{u_1, u_2, \dots, u_\ell\}$ and $\{v_1, v_2, \dots, v_\ell\}$ be any orthonormal bases for the $\lambda = -1$ and 1 eigenspaces of S , respectively, and let W denote the orthogonal matrix

$$W = [u_1 \ u_2 \ \dots \ u_\ell \mid v_1 \ v_2 \ \dots \ v_\ell]_{n \times n} \tag{A3}$$

with columns u_i and v_j .

It is well known that when two matrices commute any eigenspace of one matrix is necessarily an invariant subspace of the other. Thus, the S -eigenspaces $\text{span}\{u_1, u_2, \dots, u_\ell\}$ and $\text{span}\{v_1, v_2, \dots, v_\ell\}$ are invariant subspaces for K and M , and consequently similarity by W will simultaneously block-diagonalize K and M as direct sums of $\ell \times \ell$ matrices. This gives us a way to decouple the undamped problem $Q_u(\lambda)$, but what happens to the damping matrix D in $Q(\lambda)$ under this similarity? The significant feature of $D = c e_\ell e_\ell^T$ for this question is that e_ℓ is itself an eigenvector of S (provided ℓ is even, which we have been assuming from the start). This is because for even ℓ (and hence an odd number of R blocks in S), there is a middle R -block in S that lines up along the diagonal at the intersection of the ℓ th and $(\ell + 1)$ st rows and columns

$$S([\ell, \ell + 1], [\ell, \ell + 1]) = R = \begin{bmatrix} 1 & 0 \\ 0 & -1 \end{bmatrix}$$

Thus, we see that e_ℓ is in the $\lambda = 1$ eigenspace $\text{span}\{v_1, v_2, \dots, v_\ell\}$ of S and hence orthogonal to each u_i , so that

$$W^T D W = W^T (c e_\ell e_\ell^T) W = c (W^T e_\ell) (e_\ell^T W) = c z z^T$$

where z^T has the form $[0 \ \dots \ 0 \ | \ * \ \dots \ *]$, the 0 block being of length ℓ . Hence,

$$W^T D W = \begin{bmatrix} 0 & 0 \\ 0 & D_2 \end{bmatrix}$$

where all the blocks are $\ell \times \ell$. So not only do we find that $W^T D W$ is block-diagonal, but the damping is isolated in just one block. In summary, similarity by W decouples both $Q_u(\lambda)$ and $Q(\lambda)$

$$W^T Q_u(\lambda) W = \left[\begin{array}{c|c} \lambda^2 M_1 + K_1 & 0 \\ \hline 0 & \lambda^2 M_2 + K_2 \end{array} \right]$$

$$W^T Q(\lambda) W = \left[\begin{array}{c|c} \lambda^2 M_1 + K_1 & 0 \\ \hline 0 & \lambda^2 M_2 + \lambda D_2 + K_2 \end{array} \right]$$

and these decoupled forms have a common leading principal $\ell \times \ell$ block $\lambda^2 M_1 + K_1$. Since $M_1 = W_1^T M W_1$ and $K_1 = W_1^T K W_1$, where W_1 comprises the first ℓ columns of W , M_1 and K_1 are both symmetric positive definite. Thus, $Q_u(\lambda)$ and $Q(\lambda)$ have the $n = 2\ell$ pure imaginary eigenvalues of $\lambda^2 M_1 + K_1$ in common (as well as the n corresponding eigenvectors in the semisimple case). \square

Numerical computations show that the eigenvalues of $\lambda^2 M_1 + K_1$ and those of $\lambda^2 M_2 + K_2$ interlace, as expected from the physics of the problem.

A.2. Technical bounds

A.2.1. *Growth factor $\phi_L(\lambda)$.* The growth factor $\phi_L(\lambda)$ is defined by $\kappa_L(\lambda) = \phi_L(\lambda) \kappa_Q(\lambda)$. From the analysis in [8] we can show that for the three linearizations (4)–(6)

$$\phi_{C_1}(\lambda) \lesssim (1 + |\lambda|) \frac{(\max(m, 1)|\lambda| + \max(d, k, 1))}{m|\lambda|^2 + d|\lambda| + k} \min \left(\sqrt{1 + (|\lambda|m + d)^2}, \sqrt{1 + \frac{k^2}{|\lambda|^2}} \right) \quad (\text{A4})$$

$$\phi_{L_1}(\lambda) \approx \frac{1 + |\lambda|^2}{|\lambda|} \frac{(\max(m, k)|\lambda| + \max(d, k))}{m|\lambda|^2 + d|\lambda| + k} \quad (\text{A5})$$

$$\phi_{L_2}(\lambda) \approx (1 + |\lambda|^2) \frac{(\max(m, d)|\lambda| + \max(m, k))}{m|\lambda|^2 + d|\lambda| + k} \quad (\text{A6})$$

where the bound and the equalities are correct to within a constant factor of order 1.

A.3. Expressions for ψ_L

Let $(\widehat{z}, \widehat{\lambda})$ be an approximate eigenpair of a linearization $L(\lambda)$ of $Q(\lambda)$. Partition \widehat{z} as $\begin{bmatrix} \widehat{z}_1 \\ \widehat{z}_2 \end{bmatrix}$ where \widehat{z}_1 and \widehat{z}_2 have the same dimension. It is shown in [9] that for the linearizations (4)–(6)

$$\frac{\eta_Q(\widehat{z}_i, \widehat{\lambda})}{\eta_L(\widehat{z}, \widehat{\lambda})} \leq c \psi_L^{(i)} \frac{|\widehat{\lambda}|^2 + 1}{|\widehat{\lambda}|^2 m + |\widehat{\lambda}| d + k} \frac{\|\widehat{z}\|_2}{\|\widehat{z}_i\|_2}, \quad i = 1, 2$$

where

$$\psi_{C_1}^{(i)} = \max(1, m, d, k)^2, \quad i = 1, 2 \tag{A7}$$

$$\psi_{L_1}^{(1)} = \max(m, d, k), \quad \psi_{L_1}^{(2)} = \max(m, d, k) \max(1, (m + d)\|K^{-1}\|_2) \tag{A8}$$

$$\psi_{L_2}^{(1)} = \max(m, d, k) \max(1, (d + k)\|M^{-1}\|_2), \quad \psi_{L_2}^{(2)} = \max(m, d, k) \tag{A9}$$

and c is constant of order 1 that depends on the linearization.

ACKNOWLEDGEMENTS

We thank Nils Wagner of the Institute of Applied and Experimental Mechanics at the University of Stuttgart for bringing to our attention the instabilities in solving the unscaled beam problem by linearization. We also thank Marta Betcke and Timo Betcke for their useful comments.

REFERENCES

1. Lancaster P. *Lambda-Matrices and Vibrating Systems*. Pergamon Press: Oxford, 1966. Reprinted by Dover: New York, 2002.
2. Bai Z, Demmel JW, Dongarra JJ, Ruhe A, van der Vorst HA (eds). In *Templates for the Solution of Algebraic Eigenvalue Problems: A Practical Guide*. Society for Industrial and Applied Mathematics: Philadelphia, PA, U.S.A., 2000.
3. Tisseur F, Meerbergen K. The quadratic eigenvalue problem. *SIAM Review* 2001; **43**:235–286.
4. Gohberg I, Lancaster P, Rodman L. *Matrix Polynomials*. Academic Press: New York, 1982.
5. Garvey SD, Friswell MI, Prells U. Co-ordinate transformations for second order systems. Part I: General transformations. *Journal of Sound and Vibration* 2002; **258**:885–909.
6. Mackey DS, Mackey N, Mehl C, Mehrmann V. Vector spaces of linearizations for matrix polynomials. *SIAM Journal on Matrix Analysis and Applications* 2006; **28**:971–1004.
7. Fan H-Y, Lin W-W, Van Dooren P. Normwise scaling of second order polynomial matrices. *SIAM Journal on Matrix Analysis and Application* 2004; **26**:252–256.
8. Higham NJ, Mackey DS, Tisseur F. The conditioning of linearizations of matrix polynomials. *SIAM Journal on Matrix Analysis and Applications* 2006; **28**:1005–1028.
9. Higham NJ, Li R-C, Tisseur F. Backward error of polynomial eigenproblems solved by linearization. *MIMS EPrint 2006.137*, Manchester Institute for Mathematical Sciences, The University of Manchester, U.K., June 2006, *SIAM Journal on Matrix Analysis and Application*, submitted.
10. Tisseur F. Backward error and condition of polynomial eigenvalue problems. *Linear Algebra and its Applications* 2000; **309**:339–361.
11. Davis GJ, Moler CB. Sensitivity of matrix eigenvalues. *International Journal for Numerical Methods in Engineering* 1978; **12**:1367–1373.
12. Higham NJ, Tisseur F, Van Dooren PM. Detecting a definite Hermitian pair and a hyperbolic or elliptic quadratic eigenvalue problem, and associated nearness problems. *Linear Algebra and its Applications* 2002; **351–352**:455–474.

13. Lancaster P. Quadratic eigenvalue problems. *Linear Algebra and its Applications* 1991; **150**:499–506.
14. Higham NJ, Tisseur F. Bounds for eigenvalues of matrix polynomials. *Linear Algebra and its Applications* 2003; **358**:5–22.
15. Higham NJ. *Accuracy and Stability of Numerical Algorithms* (2nd edn). Society for Industrial and Applied Mathematics: Philadelphia, PA, U.S.A., 2002.
16. Higham NJ, Tisseur F. A block algorithm for matrix 1-norm estimation, with an application to 1-norm pseudospectra. *SIAM Journal on Matrix Analysis and Applications* 2000; **21**:1185–1201.
17. Collar AR, Simpson A. *Matrices and Engineering Dynamics*. Ellis Horwood: Chichester, 1987.