

## Shooting methods for a PT-symmetric periodic eigenvalue problem.

Lidia Aceto · Cecilia Magherini · Marco Marletta

Received: date / Accepted: date

**Abstract** We present a rigorous analysis of the performance of some one-step discretization schemes for a class of PT-symmetric singular boundary eigenvalue problem which encompasses a number of different problems whose investigation has been inspired by the 2003 article of Benilov, O'Brien and Sazonov [3]. These discretization schemes are analyzed as initial value problems rather than as discrete boundary problems, since this is the setting which ties in most naturally with the formulation of the problem which one is forced to adopt due to the presence of an interior singularity. We also devise and analyze a variable step scheme for dealing with the singular points. Numerical results show better agreement between our results and those obtained from small- $\epsilon$  asymptotics than has been shown in results presented hitherto.

**Keywords** Shooting methods for eigenvalues, one-step schemes, periodic eigenvalue problems, PT-symmetric, interior singularity

**Mathematics Subject Classification (2000)** 65L15 · 65L10 · 34L15 · 34L16

### 1 Introduction

Recently, Benilov, O'Brien and Sazonov [3] considered the formal ODE eigenvalue problem

$$i\epsilon \frac{d}{dx} \left( \sin(x) \frac{dy}{dx} \right) + i \frac{dy}{dx} = \lambda y, \quad x \in (-\pi, \pi), \quad (1)$$

with periodic boundary conditions. They speculated that the eigenvalues of this problem are purely real, in spite of the fact that any reasonable definition of the underlying operator is highly non-selfadjoint. This problem has now been extensively studied, both theoretically and numerically (see [5], [6], [8]). The fact that its eigenvalues are purely real for  $0 < \epsilon < 2$  has been proved rigorously by

---

L. Aceto  
Dipartimento di Matematica Applicata "U.Dini",  
Università di Pisa, Italy  
Tel.: +39-050 2213836  
Fax: +39-050 2213802  
E-mail: l.aceto@dma.unipi.it

C. Magherini  
Dipartimento di Matematica Applicata "U.Dini",  
Università di Pisa, Italy

M. Marletta  
School of Mathematics, Cardiff University,  
Cardiff CF24 4AG, United Kingdom

Weir [11] who used an approach inspired by the complex scaling method for resonances [2] to show that the eigenvalues coincide with those of an explicitly given selfadjoint Sturm-Liouville operator. An important aspect of the problem is to understand what precisely is meant by periodic boundary conditions. It turns out that one means only that  $y(-\pi) = y(\pi)$ , with no restrictions on  $y'(-\pi)$  and  $y'(\pi)$ . This second order problem therefore appears to have just one boundary condition. In fact the role of second boundary condition is played by the condition that solutions must lie in  $L^2(-\pi, \pi)$ , which restricts their behaviour at  $x = 0$ .

A more general version is the following:

$$i\epsilon \frac{d}{dx} \left( f(x) \frac{dy}{dx} \right) + i \frac{dy}{dx} = \lambda y, \quad x \in (-\pi, \pi), \quad (2)$$

in which  $f$  is a  $2\pi$ -periodic odd Lipschitz function, positive on  $(0, \pi)$  with  $f(0) = 0 = f(\pi)$ , twice differentiable on neighbourhoods of integer multiples of  $\pi$ , with non-zero first derivatives at those points. This version was considered by entirely different methods in [4], where it is also proved that the eigenvalues are real for all sufficiently small  $\epsilon > 0$ .

Because of the way that the domain for the operator underlying the equation (2) is defined, and in particular the interior singularity at 0, it is difficult to devise discretization methods in the form of matrix pencils. Instead, one is drawn to use a characterization of the eigenvalues as the zeros of the imaginary part of an analytic function, given in [4] and in Chugunova and Volkmer [6], see Theorem 1. The calculation of the values of this function requires the solution of initial value problems. Unfortunately the initial value problems are singular, which is problematic when it comes to devising a rigorous error analysis. It would be much better if we could devise a discretization amenable to the approach in [1], where we studied the discretization of singular problems by methods which yield algebraic eigenvalue problems  $A_N \mathbf{Y} = \lambda \mathbf{Y}$ , in which  $A_N$  is an  $N \times N$  matrix. These often have the nice property that the error in an isolated eigenvalue  $\lambda$  depends only on the quantity  $\tau = \|A_N \mathbf{y} - \lambda \mathbf{y}\|$ , in which the vector  $\mathbf{y}$  is the discretization-scheme-representation of the exact eigenfunction  $y$ . Since eigenfunctions are usually well behaved near singular points this gives an easy route to understanding why such schemes perform a lot better than one would naively expect from an error analysis of the underlying initial value problems. In fact the error bounds which usually come from the analysis of an initial value problem solver tend to be rather poor, because the Lipschitz constants for the discretized problems blow up as the stepsize tends to zero.

In this paper, however, a study of the initial value problems is inevitable. We avoid as much of the nastiness as possible by using logarithmic semi-norms (see Proposition 4) and by two-sided shooting (Theorem 2). We study some one-step schemes, including some Runge-Kutta collocation methods, and we show that whether or not a scheme exhibits an effective loss of order depends on several criteria, including

- whether or not a fixed stepsize is used;
- the size of  $\epsilon$ , with the error being of higher order for smaller  $\epsilon$  when fixed stepsize is used.

The number of derivatives which eigenfunctions possess at  $x = \pm\pi$  turns out to be  $\lfloor (\epsilon |f'(\pi)|)^{-1} \rfloor$ , so eigenfunctions may not be particularly regular at  $x = \pm\pi$  if  $\epsilon$  is not small. Nevertheless we devise a variable step grid which avoids loss of order when  $\epsilon$  is not small.

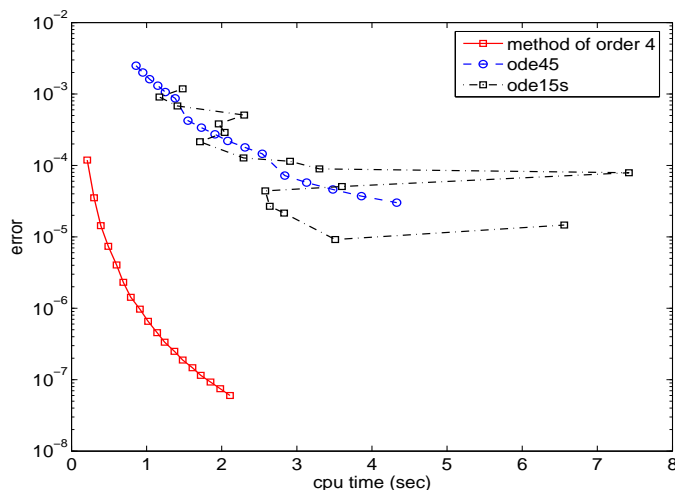
There is also an unexpected benefit of the shooting approach compared to the ‘large matrix  $A_N$ ’ approach. The original problem in [3] is actually a time-dependent PDE

$$\epsilon \frac{\partial}{\partial x} \left( f(x) \frac{\partial y}{\partial x} \right) + \frac{\partial y}{\partial x} = \frac{\partial y}{\partial t}, \quad x \in (-\pi, \pi), \quad t > 0, \quad (3)$$

for the special case  $f(x) = \sin(x)$ . A microlocal analysis shows that for  $x < 0$  it behaves like a backward heat equation and all solutions blow up instantly. Coupled with the fact that the eigenvalues are all real, this immediately establishes that the eigen- and associated functions cannot form a Riesz basis. This argument appears in Chugunova, Karabash and Pyatkov [7] for  $f(x) = \sin(x)$  and

can clearly be generalized<sup>1</sup> to a much wider class of coefficients. The failure of Riesz basisness of the eigen- and associated functions is generally associated with ill-conditioning of the spectrum and wild pseudospectral behaviour, due to the angles between eigenfunctions not being bounded away from zero. This would make any argument on eigenvalue approximation based on the measurement of local truncation errors  $A_N \mathbf{y} - \lambda \mathbf{y}$  invalid. However the shooting approach, which reduces the problem to one half of the interval, namely  $(0, \pi)$ , seems to capture enough of the structure of the problem to avoid these issues; our error analysis gives an indirect proof of this. Of course the fact that the spectrum of (1) reveals nothing of the (in)stability of the PDE (3) does mean that the accurate calculation of eigenvalues and eigenfunctions is an academic exercise. Nevertheless our demonstration that it is possible gives a rare example of a problem where one can realize a characterization of the eigenvalues which allows their accurate computation, despite the presence of unruly pseudospectra.

In the next section we prove some results which give ‘miss distances’, i.e. smooth functions of the spectral parameter whose zeros are the eigenvalues. We also give various reformulations of the problem and describe the asymptotic behaviour of solutions near singular points. Section 3 gives a rigorous error analysis for certain one-step formulae, in particular for those having ‘nice’ properties. Section 4 presents numerical results; however, since the reader may want to be convinced that existing library codes with step-size control do find these problems tough, we take the unusual step of presenting one set of results immediately. Figure 1 compares our approach based on a method of order 4 with two very respected MATLAB routines, showing the time taken to achieve a given accuracy for the eigenvalue closest to 4.



**Fig. 1** Errors in the estimates of one eigenvalue of (1) with  $\epsilon = 1$  versus the cpu time for their computation.

## 2 Eigenvalues and initial value problems

**Proposition 1** *Suppose that  $f$  is twice differentiable in a neighbourhood of  $x = 0$ . Then there exist solutions  $\psi_1(x, \lambda)$  and  $\psi_2(x, \lambda)$  having the asymptotic behaviours*

$$\psi_1(x, \lambda) \sim 1, \quad \psi_2(x, \lambda) \sim x^{-\nu}, \quad \nu = (\epsilon f'(0))^{-1},$$

<sup>1</sup> One needs to know that the eigen- and associated functions are complete before even considering basisness! Completeness is established in [7] for  $f(x) = \sin(x)$  but in general is still an open problem.

as  $x \rightarrow 0$ . When  $f$  is analytic in a neighbourhood of  $x = 0$  then  $\psi_1$  is analytic in a neighbourhood of  $x = 0$  and the asymptotic expression represents the first term of its Taylor expansion, which is therefore term-by-term differentiable. When  $\nu$  is non-integer,  $\psi_2$  is of the form  $x^{-\nu}g(x, \lambda)$  where  $g$  also has a Taylor expansion, and the resulting asymptotic expression is term-by-term differentiable.  $\square$

**Proposition 2** Suppose that  $f$  is twice differentiable in a neighbourhood of  $x = \pi$ . Then there exist solutions  $\zeta_1(x, \lambda)$  and  $\zeta_2(x, \lambda)$  having the asymptotic behaviours

$$\begin{aligned}\zeta_1(x, \lambda) &\sim (\pi - x)^\eta, & \eta &= -(\epsilon f'(\pi))^{-1}, \\ \zeta_2(x, \lambda) &\sim 1,\end{aligned}\tag{4}$$

as  $x \nearrow \pi$ . When  $f$  is analytic in a neighbourhood of  $x = \pi$ ,  $\zeta_1$  has the form  $(\pi - x)^\eta g(x, \lambda)$  where  $g$  has a Taylor expansion about  $x = \pi$  and the asymptotic formula is term-by-term differentiable. If, in addition,  $\eta$  is non-integer, then  $\zeta_2$  is analytic in a neighbourhood of  $x = \pi$  and the asymptotic formula represents the first term of its Taylor series, and is therefore term-by-term differentiable.  $\square$

These results follow by a standard Frobenius-type analysis; alternatively they may be proved with rigorous error bounds using a change of variables and an application of Levinson's Theorem, as in the appendix to [4].

An important observation here is that since  $f'(0) > 0$  and  $f'(\pi) < 0$ , there is only one solution bounded near  $x = 0$ , whereas all solutions are bounded near  $x = \pi$ . If  $0 < \epsilon f'(0) < 2$  then there is only one solution, namely  $\psi_1(x, \lambda)$ , which lies in  $L^2(-\pi, \pi)$ . Any eigenfunction must be a multiple of this solution, for the corresponding value of  $\lambda$ .

Proposition 2 is not quite sufficient for local truncation error analysis in the case where  $\eta = -(\epsilon f'(\pi))^{-1}$  is integer. The standard Frobenius analysis then yields the following.

**Proposition 3** If  $f$  is analytic in a neighbourhood of  $\pi$  and  $\eta = -(\epsilon f'(\pi))^{-1} \in \mathbb{N}$ , then the solution  $\zeta_2(x, \lambda)$  has the form

$$\zeta_2(x, \lambda) = \hat{\zeta}_2(x, \lambda) + C' \zeta_1(x, \lambda) \log(\pi - x),$$

in which  $C' \neq 0$  is constant and  $\hat{\zeta}_2$  is analytic in a neighbourhood of  $x = \pi$ , with  $\hat{\zeta}_2(\pi, \lambda) = 1$ .  $\square$

**Theorem 1** The differential equation (2) has a non-trivial solution  $y \in L^2(-\pi, \pi)$  satisfying  $y(-\pi) = y(\pi)$  if and only if  $\psi_1(-\pi, \lambda) = \psi_1(\pi, \lambda)$ . Moreover the solution  $\psi_1$  has the symmetries

$$\psi_1(-x, \lambda) = \psi_1(x, -\lambda), \quad \psi_1(-x, \bar{\lambda}) = \overline{\psi_1(x, \lambda)}\tag{5}$$

and so  $\lambda$  is a (real) eigenvalue if and only if

$$\Im(\psi_1(\pi, \lambda)) = 0.\tag{6}$$

*Proof* Most of the results in this theorem have appeared in several recent articles including [4] and [7]. A proof is included here for the convenience of the reader, and applies equally well to the more general equation

$$i \frac{d}{dx} \left( (1 + a g(x)) y + \epsilon f(x) \frac{dy}{dx} \right) = \lambda y,\tag{7}$$

in which  $g$  is even and  $f$  is odd with  $f > 0$  on  $(0, \pi)$ ,  $f'(0) > 0$ ,  $f$  and  $g$  both  $2\pi$ -periodic.

Since, up to scalar multiples,  $\psi_1(x, \lambda)$  is the only solution of the differential equation in  $L^2(-\pi, \pi)$ , the statement that a non-zero  $L^2$ -solution  $y$  exists with  $y(-\pi) = y(\pi)$  if and only if  $\psi_1(-\pi, \lambda) = \psi_1(\pi, \lambda)$  is obvious.

The symmetries of  $\psi_1$  in eqn. (5) are proved in [4] and [6]. The proofs are elementary. For instance, taking complex conjugates in (2) one immediately proves that  $\psi_1(x, \lambda)$  satisfies the same ODE as  $\psi_1(-x, \bar{\lambda})$ ; the fact that these two functions coincide is then immediate from the fact that they both have the value 1 at  $x = 0$ , and the solution determined by this condition is unique.

Finally, the condition (6) follows by combining the eigenvalue condition  $\psi_1(-\pi, \lambda) = \psi_1(\pi, \lambda)$  with the fact that eigenvalues are real ( $\bar{\lambda} = \lambda$ ) and the symmetry  $\psi_1(-\pi, \bar{\lambda}) = \overline{\psi_1(\pi, \lambda)}$ , to obtain  $\psi_1(\pi, \lambda) = \overline{\psi_1(\pi, \lambda)}$ .  $\square$

Theorem 1 is good if one is able to solve the eigenvalue problem by forward shooting. Starting from (asymptotic) initial conditions given by Proposition 1 one integrates forward and, if the integration method is exact, one finds, by Proposition 2, a finite value  $\psi_1(\pi, \lambda)$ . The eigenvalues are then the roots of the equation (6). However if one wishes to use two-sided shooting with a matching point inside the interval of integration, a different miss distance is required.

**Theorem 2** *Let  $\zeta_1(x, \lambda)$  be the solution of (2) determined by Proposition 2. Let  $q$  be any positive function which satisfies the ODE*

$$\frac{d}{dx} \log(q/f) = 1/(\epsilon f). \quad (8)$$

*Then  $\lambda$  is an eigenvalue if and only if  $\Im(W(x, \lambda)) = 0$ , where*

$$W(x, \lambda) = \begin{vmatrix} \psi_1(x, \lambda) & \zeta_1(x, \lambda) \\ q(x)\psi_1'(x, \lambda) & q(x)\zeta_1'(x, \lambda) \end{vmatrix} = W(\lambda) \quad (9)$$

*is, in fact, independent of  $x$ .*

*Proof* A direct calculation shows that if  $q$  satisfies (8) then the differential equation (2) is equivalent to  $-(q(x)y')' = (i\lambda q/(\epsilon f))y$  and standard calculations then show that the  $q$ -weighted Wronskian  $W(x, \lambda)$  of any two solutions is, in fact, independent of  $x$ . Thus, for the particular solutions  $\psi_1$  and  $\zeta_1$ , we have  $W(x, \lambda) = \lim_{t \nearrow \pi} W(t, \lambda)$ . We know the asymptotic behaviour of  $\zeta_1$  near  $x = \pi$ , but we also need to know the behaviour of  $q$ . Since  $f(x) \sim -f'(\pi)(\pi - x)$  for  $x \nearrow \pi$ , an integration of (8) shows that

$$q(x) \sim C(\pi - x)^{1-\eta}, \quad (10)$$

for some real constant  $C$  which we may take to be 1. Now the definition of  $W$  gives immediately  $-W = q\zeta_1^2 \left(\frac{\psi_1}{\zeta_1}\right)'$ , and so using the asymptotics

$$\zeta_1(x) \sim (\pi - x)^\eta [1 + \dots], \quad q(x) \sim (\pi - x)^{1-\eta},$$

we obtain

$$\left(\frac{\psi_1}{\zeta_1}\right)' \sim -\frac{W}{(\pi - x)^{1+\eta}[1 + \dots]}.$$

This yields, upon integrating,

$$\frac{\psi_1}{\zeta_1} \sim \frac{1}{\eta} W(\pi - x)^{-\eta} [1 + \dots] + d,$$

where  $d$  is a constant of integration. Multiplying by  $\zeta_1$  and letting  $x \nearrow \pi$  gives

$$\psi_1(\pi, \lambda) = \frac{1}{\eta} W.$$

Thus  $\Im(\psi_1(\pi, \lambda)) = 0$  if and only if  $\Im W(\lambda) = 0$ , and, from the characterization of eigenvalues in Theorem 1, the result is proved.  $\square$

### 3 Error evolution for one-step discretization schemes

We observe first that the differential equation (2) can be cast in the form of a first order system

$$\mathbf{y}'(x) = J(x)\mathbf{y}(x), \quad (11)$$

where

$$\mathbf{y}(x) = \begin{pmatrix} y \\ \epsilon f y' \end{pmatrix}, \quad J(x) = \begin{pmatrix} 0 & \frac{1}{\epsilon f(x)} \\ -i\lambda & -\frac{1}{\epsilon f(x)} \end{pmatrix}. \quad (12)$$

This is the form in which we shall integrate the differential equation. It is also worth noting that Propositions 1, 2 and 3 give corresponding asymptotic information on the solutions of (11), which may also be found in [4], and which we shall use later.

### 3.1 Forward shooting with one-step methods and fixed stepsize

This section is devoted to the analysis of convergence of a one-step method used for solving (11) subject to  $\mathbf{y}(0) = (1, 0)^T$  over the uniform grid

$$x_n = nh, \quad n = 0, 1, \dots, N, \quad h = \frac{\pi}{N}.$$

This forward integration is required for each trial value of  $\lambda$  involved in the forward shooting procedure. The analysis is performed by considering the local error of the formula and its propagation across the integration steps. In particular, a careful study of both these aspects is required near the singularities at  $x = 0, \pi$ .

The proof that the error propagation does not cause an order reduction of the formula can be deduced from the following result.

**Proposition 4** *For any  $\bar{x} \in (0, \pi)$ , let  $\mathbf{y}(x; \bar{x}, \bar{\mathbf{y}})$  and  $\mathbf{z}(x; \bar{x}, \bar{\mathbf{z}})$  be the solution of (11)-(12) with initial values  $\mathbf{y}(\bar{x}) = \bar{\mathbf{y}}$  and  $\mathbf{z}(\bar{x}) = \bar{\mathbf{z}}$ , respectively. Then, for any  $t \in [\bar{x}, \pi)$ ,*

$$\|\mathbf{y}(t; \bar{x}, \bar{\mathbf{y}}) - \mathbf{z}(t; \bar{x}, \bar{\mathbf{z}})\| \leq e^{|\lambda|\pi} \|\bar{\mathbf{y}} - \bar{\mathbf{z}}\|.$$

*Proof* Let  $I_2$  be the identity matrix of size 2. Since  $J(x)$  is continuous in  $[\bar{x}, \pi)$  and since the logarithmic semi-norm

$$\mu_1(J(x)) = \lim_{\theta \rightarrow 0^+} \frac{\|I_2 + \theta J(x)\|_1 - 1}{\theta} = |\lambda|$$

for each  $x$ , the result follows immediately from Theorem 11.1 in [10, page 64].  $\square$

The immediate consequence of this Proposition is that if we denote by  $\mathbf{e}_n$  the local error at the  $n$ -th integration step, defined by

$$\mathbf{e}_n = \mathbf{y}(x_n; x_{n-1}, \mathbf{y}_{n-1}) - \mathbf{y}_n, \quad (13)$$

then

$$\|\mathbf{y}(x_N; x_0, \mathbf{y}_0) - \mathbf{y}_N\| \leq e^{|\lambda|\pi} \sum_{n=1}^N \|\mathbf{e}_n\|, \quad (14)$$

where, we recall,  $\mathbf{y}_0 = (1, 0)^T$ . The discussion of the behaviour of the local errors is closely related to the method used. In the sequel we shall consider the case where the one-step scheme is an  $s$ -stage Runge-Kutta method which, when applied to solve (11)-(12), leads to the discrete problem

$$\mathbf{y}_{n+1} = \mathbf{y}_n + h (\mathbf{b}^T \otimes I_2) D_n \mathbf{Y}_n, \quad (15)$$

$$\mathbf{Y}_n = \mathbf{E} \otimes \mathbf{y}_n + h (A \otimes I_2) D_n \mathbf{Y}_n. \quad (16)$$

Here  $A \in \mathbb{R}^{s \times s}$  and  $\mathbf{b} \in \mathbb{R}^s$  are the coefficient matrix and the vector of weights of the method, respectively;  $\mathbf{E} = (1, 1, \dots, 1)^T \in \mathbb{R}^s$ ; and  $\mathbf{Y}_n \in \mathbb{R}^{2s}$  contains the internal stages of the method. Finally, (see (12))

$$D_n = \text{blockdiag}(J(x_n + c_1 h), \dots, J(x_n + c_s h)), \quad (17)$$

$c_1, \dots, c_s$  being the abscissae of the scheme.

We shall assume that the Runge-Kutta method is a collocation method of order  $p$  satisfying the simplifying assumptions  $C(s)$  and  $B(p)$ , with  $p \geq s$  (see, for example, [10]). In this case, if we insert  $\mathbf{y}(x_{n+1}; x_n, \mathbf{y}_n)$  and

$$\mathbf{Y}(x_n, \mathbf{y}_n) = (\mathbf{y}(x_n + c_1 h; x_n, \mathbf{y}_n)^T, \dots, \mathbf{y}(x_n + c_s h; x_n, \mathbf{y}_n)^T)^T$$

into (15)-(16), we get

$$\begin{aligned}\mathbf{y}(x_{n+1}; x_n, \mathbf{y}_n) &= \mathbf{y}_n + h (\mathbf{b}^T \otimes I_2) D_n \mathbf{Y}(x_n, \mathbf{y}_n) + \delta(h, n), \\ \mathbf{Y}(x_n, \mathbf{y}_n) &= \mathbf{E} \otimes \mathbf{y}_n + h (A \otimes I_2) D_n \mathbf{Y}(x_n, \mathbf{y}_n) + \Delta(h, n),\end{aligned}$$

where, if  $\gamma_j$  ( $j = 0, 1, \dots, s$ ) are the principal error coefficients of the formula, then

$$\delta(h, n) = \gamma_0 h^{p+1} \mathbf{y}^{(p+1)}(\vartheta_{n_0}; x_n, \mathbf{y}_n), \quad (18)$$

$$\Delta(h, n) = h^{s+1} \begin{pmatrix} \gamma_1 \mathbf{y}^{(s+1)}(\vartheta_{n_1}; x_n, \mathbf{y}_n) \\ \vdots \\ \gamma_s \mathbf{y}^{(s+1)}(\vartheta_{n_s}; x_n, \mathbf{y}_n) \end{pmatrix}, \quad (19)$$

for suitable  $\vartheta_{n_j} \in [x_n, x_{n+1}]$ ,  $j = 0, 1, \dots, s$ . By virtue of this fact, from (15)-(16) and (18)-(19), one easily obtains (see (13))

$$\mathbf{e}_{n+1} = h (\mathbf{b}^T \otimes I_2) D_n (I_{2s} - h (A \otimes I_2) D_n)^{-1} \Delta(h, n) + \delta(h, n). \quad (20)$$

The classical theory of convergence for Runge-Kutta methods applies away from the singularities. Near the singular points a careful analysis of the behaviours of  $D_n$ ,  $\Delta(h, n)$  and  $\delta(h, n)$  is needed.

By using (12) and (17) it is not difficult to verify that  $D_n = O((f'(0)(n+1)h)^{-1})$ , if  $(n+1)h$  is sufficiently small. In particular this implies that  $\|(I_{2s} - h(A \otimes I_2) D_n)^{-1}\|$  is bounded with respect to  $h$  for each  $n$ .

Recalling the solutions  $\psi_1$  and  $\psi_2$  introduced in Proposition 1, we can represent the local solution  $\mathbf{y}(x; x_n, \mathbf{y}_n)$  as a linear combination

$$\mathbf{y}(x; x_n, \mathbf{y}_n) = (1 - \xi_{1,n})\Psi_1(x) - \xi_{2,n}\Psi_2(x) \quad (21)$$

where  $\xi_{1,n}$  and  $\xi_{2,n}$  are suitable coefficients and

$$\Psi_j(x) = \begin{pmatrix} \psi_j(x, \lambda) \\ \epsilon f(x) \psi_j'(x, \lambda) \end{pmatrix}, \quad j = 1, 2.$$

When  $n = 0$ , the fact that  $\mathbf{y}_0 = (1, 0)^T$  implies that  $\xi_{1,0} = \xi_{2,0} = 0$  so that  $\mathbf{y}(x; 0, \mathbf{y}_0)$  is analytic,  $\Delta(h, 0) = O(h^{s+1})$ ,  $\delta(h, 0) = O(h^{p+1})$  and, consequently  $\mathbf{e}_1 = O(h^{s+1})$ .

**Proposition 5** *The coefficients  $\xi_{1,n}$  and  $\xi_{2,n}$  in (21) satisfy bounds*

$$\xi_{1,n} = O(h^{s+1} \log(n+1)), \quad \xi_{2,n} = O(h^{s+1} (nh)^\nu),$$

while the local error  $\mathbf{e}_{n+1}$  satisfies

$$\mathbf{e}_{n+1} = O\left(\frac{h^{s+1}}{n+1}\right), \quad (22)$$

for  $n = 0, \dots, \bar{n}$ , where  $\bar{n}h = O(1)$ .

*Proof* The proof is by induction. The initial case  $n = 0$  has already been checked, so we need only verify that the results hold for  $n+1$  if they hold for  $n$ . Recall that, by definition, see (13) and (21),

$$\begin{aligned}\mathbf{y}_n &= -\mathbf{e}_n + \mathbf{y}(x_n; x_{n-1}, \mathbf{y}_{n-1}) \\ &= O\left(\frac{h^{s+1}}{n}\right) + (1 - \xi_{1,n-1})\Psi_1(x_n) - \xi_{2,n-1}\Psi_2(x_n) \\ &= (1 - \xi_{1,n})\Psi_1(x_n) - \xi_{2,n}\Psi_2(x_n),\end{aligned}$$

where the second equality uses (22) with  $n$  replaced by  $n - 1$ , and the third equality is immediate from the definitions of the coefficients  $\xi_{1,n}$  and  $\xi_{2,n}$ . We subtract the third equation from the second to obtain an equation for the differences of the coefficients:

$$(\Psi_1(x_n); \Psi_2(x_n)) \begin{pmatrix} \xi_{1,n} - \xi_{1,n-1} \\ \xi_{2,n} - \xi_{2,n-1} \end{pmatrix} = \frac{h^{s+1}}{n} \mathbf{w}, \quad (23)$$

where  $\mathbf{w}$  is a vector whose norm is  $O(1)$  uniformly in  $h$  and  $n$ . We now invert the matrix on the left hand side of (23); to do this we require its determinant. Recall from (9) that if  $q$  is the coefficient given by (8) then the determinant  $\begin{vmatrix} \psi_1 & \psi_2 \\ q\psi_1' & q\psi_2' \end{vmatrix}$  is constant, for any solutions  $\psi_1$  and  $\psi_2$  of the differential equation; in particular, for the solutions which we use here. We may therefore assume without loss of generality that  $\det(\Psi_1(x_n); \Psi_2(x_n)) = \epsilon f(x_n)/q(x_n)$ . A simple Frobenius analysis shows that in a neighbourhood of the origin,

$$q(x) \sim x^{1+\nu}, \quad \frac{\epsilon f(x)}{q(x)} \sim \nu^{-1} x^{-\nu}. \quad (24)$$

Solving the linear system (23) we obtain

$$\begin{pmatrix} \xi_{1,n} - \xi_{1,n-1} \\ \xi_{2,n} - \xi_{2,n-1} \end{pmatrix} = \frac{h^{s+1}}{n} \frac{q(x_n)}{\epsilon f(x_n)} \begin{pmatrix} \epsilon f(x_n)\psi_2'(x_n) - \psi_2(x_n) \\ -\epsilon f(x_n)\psi_1'(x_n) - \psi_1(x_n) \end{pmatrix} \mathbf{w}. \quad (25)$$

We shall now estimate the sizes of all the terms using the information on  $\xi_{1,n-1}$  and  $\xi_{2,n-1}$  from the previous step, together with the known asymptotic behaviours of the solutions  $\psi_1$  and  $\psi_2$  and the estimate of the determinant in (24).

Consider the first row in (25). We know from Proposition 1 that  $|\epsilon f(x_n)\psi_2'(x_n)| + |\psi_2(x_n)| \leq C(nh)^{-\nu}$ ; also, (24) gives  $\frac{q(x_n)}{\epsilon f(x_n)} \sim \nu(nh)^\nu$ . Thus we obtain  $\xi_{1,n} - \xi_{1,n-1} = O(h^{s+1}/n)$ , with initial condition  $\xi_{1,0} = 0$ . This gives the estimate  $\xi_{1,n} = O(h^{s+1} \log(n+1))$  in the usual way.

Next consider the second row in (25). We know from Proposition 1 that  $|\epsilon f(x_n)\psi_1'(x_n)| + |\psi_1(x_n)| = O(1)$ , uniformly in  $n$ , and we still have  $\frac{q(x_n)}{\epsilon f(x_n)} \sim \nu(nh)^\nu$ . Combining these yields  $\xi_{2,n} - \xi_{2,n-1} = O\left((nh)^\nu \frac{h^{s+1}}{n}\right) = O(h^{s+\nu+1} n^{\nu-1})$ . Using the fact that  $\sum_{j=1}^n j^{\nu-1} = O(n^\nu)$  we obtain  $\xi_{2,n} = O(h^{s+1}(nh)^\nu)$ , as required.

Finally, we estimate the local error  $\mathbf{e}_{n+1}$ . This estimate comes from (20) using the estimates of  $\Delta(h, n)$  and  $\delta(h, n)$  in equations (18)-(19) in terms of the derivatives of the local solution  $\mathbf{y}(x; x_n, \mathbf{y}_n)$  in whose representation (21) we now have the necessary estimates for the coefficients  $\xi_{1,n}$  and  $\xi_{2,n}$ . Observe first that the term  $hD_n$  in (20) yields a factor  $O(h/((n+1)h)) = O(1/(n+1))$ . The factor  $h^{s+1}$  is clear in (19); in (18) we have a smaller factor of  $h^{p+1}$  since  $p \geq s$ . It remains only to estimate the size of the derivatives of  $\mathbf{y}(x; x_n, \mathbf{y}_n)$ . From (21) we have

$$\mathbf{y}^{(s+1)}(x; x_n, \mathbf{y}_n) = (1 - \xi_{1,n})\Psi_1^{(s+1)}(x) - \xi_{2,n}\Psi_2^{(s+1)}(x).$$

The factors  $\xi_{1,n}$  are bounded, as are the derivatives of the analytic solution  $\Psi_1$ . The singular solution  $\Psi_2$  has the property that  $\Psi_2^{(s+1)}(x_n) \sim x_n^{-\nu-s-1} = (nh)^{-\nu-s-1}$ , which cancels the factor of  $h^{s+1}(nh)^\nu$  in the coefficient  $\xi_{2,n}$  leaving just a factor of  $n^{-s-1}$ . For small  $n$  this is  $O(1)$ , while for larger  $n$  the dominant term in the estimate of  $\mathbf{y}^{(s+1)}(x; x_n, \mathbf{y}_n)$  is thus  $\Psi_1^{(s+1)}(x)$ , which is bounded.

To deal with  $\delta(h, n)$  we need to take derivatives up to order  $p+1$ . In this case the regular solution  $\Psi_1$  contributes an  $O(1)$  term to  $\mathbf{y}^{(p+1)}(x; x_n, \mathbf{y}_n)$ . The singular term  $\xi_{2,n}\Psi_2^{(p+1)}$  makes a contribution

$$O((nh)^\nu h^{s+1} (nh)^{-\nu-p-1}) = O\left(\frac{h^{s+1}}{n^{p+1} h^{p+1}}\right),$$

which could now be the dominant term for small  $n$ . However this term is multiplied by  $h^{p+1}$  (see (18)) leaving a term which is  $O\left(\frac{h^{s+1}}{n^{p+1}}\right)$ . This is dominated by the  $O(h^{s+1}/(n+1))$  term coming from  $\Delta(h, n)$  and completes the proof.  $\square$



Proposition 5 explains why forward integration does not, in general, suffer a catastrophic loss of order despite the singularity at the origin. Unfortunately pure forward integration with constant stepsize all the way to the singularity at  $x = \pi$  is not possible without a serious loss of accuracy, as we shall now argue: in fact, the order of convergence is controlled by the size of  $\eta = (\epsilon|f'(\pi)|)^{-1}$ .

Suppose that  $\eta > 0$  is non-integer. In a neighbourhood of  $x = \pi$ , from Proposition 2 and (18)-(19) one deduces that there exists a constant  $C$  such that

$$\|\Delta(h, n)\| \leq Ch^{s+1} \left( (\pi - x_{n+\frac{1}{2}})^{\eta-s-1} + 1 \right), \quad (26)$$

$$\|\delta(h, n)\| \leq Ch^{p+1} \left( (\pi - x_{n+\frac{1}{2}})^{\eta-p-1} + 1 \right). \quad (27)$$

Consequently, see (20),

$$\|\mathbf{e}_n\| \leq \bar{C} \left( h^{s+2} \left( (\pi - x_{n-\frac{1}{2}})^{\eta-s-2} + (\pi - x_{n-\frac{1}{2}})^{-1} \right) + h^{p+1} \left( (\pi - x_{n-\frac{1}{2}})^{\eta-p-1} + 1 \right) \right), \quad (28)$$

for a suitable constant  $\bar{C}$ . Therefore, if we denote by  $n_1$  the step-number such that  $x_{n_1-1} < \pi - 1 \leq x_{n_1}$ , we obtain

$$\begin{aligned} \sum_{n=n_1}^{N-1} \|\mathbf{e}_n\| &\lesssim \hat{C}h^{s+1} \int_{\pi-1}^{\pi-h} (\pi-x)^{\eta-s-2} dx + \hat{C}h^{s+1} \int_{\pi-1}^{\pi-h} (\pi-x)^{-1} dx \\ &+ \hat{C}h^p \int_{\pi-1}^{\pi-h} (\pi-x)^{\eta-p-1} dx + O(h^p) \\ &= O(h^\eta) + O(h^{s+1} \log h) + O(h^p), \end{aligned} \quad (29)$$

since  $p \geq s$ , where again  $\hat{C}$  is a suitable constant.

In the case where  $\eta$  is integer, the behaviour of solution derivatives near  $x = \pi$  is different, and is given in Proposition 3. Eqn.s (26) and (27) must be modified accordingly. Observe that

$$\frac{d^\ell}{dx^\ell} (\pi-x)^\eta \log(\pi-x) = O(C_{\eta,\ell} (\pi-x)^{\eta-\ell} \log(\pi-x) + (\pi-x)^{\eta-\ell}),$$

in which  $C_{\eta,\ell} = 0$  for  $\ell > \eta$ . The estimate (28) is now replaced by

$$\begin{aligned} \|\mathbf{e}_n\| &\leq \hat{C}h^{s+2} C_{\eta,s+1} (\pi - x_{n-\frac{1}{2}})^{\eta-s-2} |\log(\pi - x_{n-\frac{1}{2}})| \\ &+ \hat{C}h^{s+2} \left( (\pi - x_{n-\frac{1}{2}})^{\eta-s-2} + (\pi - x_{n-\frac{1}{2}})^{-1} \right) \\ &+ \hat{C}h^{p+1} C_{\eta,p+1} (\pi - x_{n-\frac{1}{2}})^{\eta-p-1} |\log(\pi - x_{n-\frac{1}{2}})| \\ &+ \hat{C}h^{p+1} \left( (\pi - x_{n-\frac{1}{2}})^{\eta-p-1} + 1 \right), \end{aligned}$$

in which additional logarithmic terms have now appeared. It follows that in (29) the right hand side now acquires an additional term, of order

$$C_{\eta,s+1} h^{s+1} \int_{\pi-1}^{\pi-h} (\pi-x)^{\eta-s-2} |\log(\pi-x)| dx + C_{\eta,p+1} h^{p+1} \int_{\pi-1}^{\pi-h} (\pi-x)^{\eta-p-1} |\log(\pi-x)| dx.$$

There are two cases to consider.

1. If  $\eta < s+1$ , then we also have  $\eta < p+1$  since  $p \geq s$ . In this case  $C_{\eta,s+1} = 0 = C_{\eta,p+1}$  and there are no new contributions to the error found in (29).

2. In general we have, for  $\ell = s + 1$  or  $\ell = p$ ,

$$\begin{aligned} h^\ell \int_{\pi-1}^{\pi-h} (\pi-x)^{\eta-\ell-1} |\log(\pi-x)| dx &= h^\eta \int_{\log(h)}^0 t e^{(\eta-\ell)t} dt \\ &= \begin{cases} O(h^\eta \log^2(h)), & \ell = \eta; \\ O(h^\eta |\log(h)|) + O(h^\eta) + O(h^\ell), & \ell \neq \eta. \end{cases} \end{aligned}$$

A careful re-examination of the different cases shows that the error is unchanged from (29) in the cases  $\eta \neq s + 1$  but that in the case  $\eta = s + 1$  then we obtain

$$\|\mathbf{e}_n\| = O(h^{s+1} \log^2(h) + h^p). \quad (30)$$

**Theorem 3** *Figures 2 and 3 summarize the orders of convergence in different parts of the interval  $[0, \pi]$  for forward integration.*

$$\begin{array}{ccc} \mathbf{h}^{s+1} \log(\mathbf{h}) & \mathbf{h}^p & \max(\mathbf{h}^\eta, \mathbf{h}^{s+1} \log(\mathbf{h}), \mathbf{h}^p) \\ \hline 0 & & \pi \end{array}$$

**Fig. 2** Orders of convergence for collocation RK formulae for  $\eta \neq s + 1$ .

$$\begin{array}{ccc} \mathbf{h}^{s+1} \log(\mathbf{h}) & \mathbf{h}^p & \max(\mathbf{h}^{s+1} \log^2(\mathbf{h}), \mathbf{h}^p) \\ \hline 0 & & \pi \end{array}$$

**Fig. 3** Orders of convergence for collocation RK formulae for  $\eta = s + 1$ .

*Proof* The results follow from the estimates (29) and (30), the observations concerning the behaviour of the local error near the origin in Proposition 5, and the error propagation estimate (14).  $\square$

Clearly a significant order reduction occurs when  $\eta$  is close to 1. For example, if  $f(x) = \sin(x)$  and  $\epsilon = 1$ , then  $\eta = 1$  and the shooting procedure with constant stepsize has order 1 independently of  $s$ . This is why we consider the alternative approach described in the next section.

### 3.2 Two-sided shooting: one-step variable stepsize schemes with minimal stability hypotheses

In this section we shall be interested primarily in the analysis of backward integration from the singular point  $x = \pi$  to compute an estimate of the solution  $\zeta_1$  given in Proposition 2, as required by the two-sided shooting procedure (see Theorem 2). We consider only one-step schemes which, for backward integration, have the form

$$\mathbf{y}_n = \mathbf{y}_{n+1} - h_n \Phi_n(x_{n+1}, h_n, \mathbf{y}_{n+1}), \quad h_n = x_{n+1} - x_n. \quad (31)$$

We suppose that the grid is chosen according to the rule

$$\pi - x_n = e^{-nh}, \quad n \in \mathbb{N}, \quad (32)$$

which allows us to integrate backwards as far as  $\pi - 1$  and not all the way to the singularity at the origin. This means that the chosen matching point for the two-sided shooting procedure is  $\pi - 1$ . The parameter  $h > 0$  is fixed and determines the resolution of the discretization. Notice that with the choice (32) an infinite number of gridpoints is required to reach  $x = \pi$ . Obviously, we will integrate

for a finite number of steps, say  $N$ , meaning that we introduce a layer near  $\pi$ .

By inserting the exact solution  $\mathbf{y}(x)$  of the original problem into the discrete problem (31) we obtain

$$\mathbf{y}(x_n) = \mathbf{y}(x_{n+1}) - h_n \Phi_n(x_{n+1}, h_n, \mathbf{y}(x_{n+1})) - h_n \boldsymbol{\tau}_n, \quad (33)$$

where  $\boldsymbol{\tau}_n$  is the local truncation error of the method used. As mentioned before, backward integration always aims to estimate the solution  $\zeta_1$  given in Proposition 2. In other words, we have

$$\mathbf{y}(x) = \begin{pmatrix} \zeta_1(x) \\ \epsilon f(x) \zeta_1'(x) \end{pmatrix} \approx (\pi - x)^\eta (\mathbf{v} + O(\pi - x)), \quad x \nearrow \pi, \quad (34)$$

where  $\mathbf{v} = (1, -1)^T$ . (Estimates of this form are derived more rigorously in [4].) In the sequel we shall assume that in a neighbourhood of  $\pi$  we have

$$\|\boldsymbol{\tau}_n\|_\infty = O(h_n^p (\pi - x_{n+1})^{\eta-p-1}), \quad (35)$$

$p$  being the ‘‘classical order’’ of the method.

By using the linear differential equation satisfied by  $\mathbf{y}$ , it is possible to verify, after some computations, that Runge-Kutta collocation methods satisfy (35) provided  $h_n \|J(x_{n+1})\|$  is sufficiently small. By (32) this holds in spite of the singularity at  $x = \pi$ .

We make two further important assumptions:

- (A1)  $\Phi_n(x_{n+1}, h_n, \mathbf{y}_{n+1})$  is a linear function of  $\mathbf{y}_{n+1}$ ;
- (A2) there exists a constant  $C > 0$ , independent of  $N$  and  $h$ , such that

$$\|\Phi_n(x_{n+1}, h_n, \mathbf{y}_{n+1})\|_\infty \leq \frac{L(1 + C(\pi - x_{n+1}))}{\pi - x_{n+1}} \|\mathbf{y}_{n+1}\|_\infty \quad (36)$$

with, see (4),

$$L = L(h) \leq \eta + \alpha h, \quad \text{for some } \alpha \geq 0. \quad (37)$$

The first hypothesis is natural since our ODE (11) is linear. The second reflects the simple pole singularity in our system of ODEs at  $x = \pi$ . The restriction on the value of  $L$  turns out to be surprisingly important.

In the sequel we shall denote with  $\mathbf{e}_n = \mathbf{y}(x_n) - \mathbf{y}_n$  the global error at the  $n$ th integration step (note that  $\mathbf{e}_n$  is here renamed since in the previous section it denoted the local error).

**Theorem 4** *Let the method (31) be applied for integrating (11)-(12) away from  $\pi$  over the grid defined by (32). Assume that the method satisfies (35) and assumptions (A1)-(A2). If we set  $\mathbf{y}_N = (\pi - x_N)^\eta \mathbf{v}$  (see (34)) and  $\alpha N h^2 \leq 1$ , then there exists a constant  $\bar{C}$  such that*

$$\|\mathbf{e}_0\|_\infty \leq \bar{C} (e^{-Nh} + Nh^{p+1}). \quad (38)$$

*Proof* From (31), (33) and (A1) one immediately gets

$$\mathbf{e}_n = \mathbf{e}_{n+1} - h_n \Phi_n(x_{n+1}, h_n, \mathbf{e}_{n+1}) - h_n \boldsymbol{\tau}_n$$

and hence, on taking norms and using (A2),

$$\|\mathbf{e}_n\|_\infty \leq \left( 1 + h_n \frac{L(1 + C(\pi - x_{n+1}))}{\pi - x_{n+1}} \right) \|\mathbf{e}_{n+1}\|_\infty + h_n \|\boldsymbol{\tau}_n\|_\infty. \quad (39)$$

Now, from (32) it follows that

$$h_n = (\pi - x_{n+1})(e^h - 1), \quad (40)$$

and so

$$1 + h_n \frac{L(1 + C(\pi - x_{n+1}))}{\pi - x_{n+1}} \leq e^{L(1+C(\pi-x_{n+1}))h},$$

the last inequality being valid for  $L \geq 1$  (and hence, in our equation (37), when the parameter  $\eta$  is greater than unity). From the previous inequality and (39) we therefore obtain

$$\|\mathbf{e}_n\|_\infty \leq e^{L(1+C(\pi-x_{n+1}))h} \|\mathbf{e}_{n+1}\|_\infty + h_n \|\boldsymbol{\tau}_n\|_\infty. \quad (41)$$

Define a sequence  $(\chi_n)$  by

$$\chi_0 = 0, \quad \chi_{n+1} - \chi_n = L(1 + C(\pi - x_{n+1}))h, \quad n = 0, \dots, N-1,$$

where  $N$  is the total number of steps, i.e. the backward integration starts from  $x_N = \pi - e^{-Nh}$ . In the sequel, we shall use  $C_j$  to denote suitable positive constants independent of  $h$  and  $n$ , for each  $j$ . From (32),

$$\begin{aligned} \chi_n &= \sum_{j=0}^{n-1} (\chi_{j+1} - \chi_j) = Lnh + LCh \sum_{j=0}^{n-1} (\pi - x_{j+1}) \\ &= Lnh + LCh e^{-h} \sum_{j=0}^{n-1} e^{-jh} = Lnh + LCh \frac{e^{-h}}{1 - e^{-h}} (1 - e^{-nh}) \\ &\leq Lnh + LCh(1 - e^{-nh}) \leq Lnh + \log(C_1), \end{aligned} \quad (42)$$

with  $C_1 > 1$ . Observe that the last inequality holds if, for instance,  $C_1 \geq e^{LC\pi}$ , since the stepsize  $h$  is always less than  $\pi$ . Multiplying both sides of (41) on the left by  $e^{\chi_n}$  we derive

$$e^{\chi_n} \|\mathbf{e}_n\|_\infty \leq e^{\chi_{n+1}} \|\mathbf{e}_{n+1}\|_\infty + h_n e^{\chi_n} \|\boldsymbol{\tau}_n\|_\infty.$$

so that, bearing in mind that  $\chi_0 = 0$ , by recursion from (32), (40) and (42) we obtain

$$\begin{aligned} \|\mathbf{e}_0\|_\infty &\leq e^{\chi_N} \|\mathbf{e}_N\|_\infty + \sum_{n=0}^{N-1} h_n e^{\chi_n} \|\boldsymbol{\tau}_n\|_\infty \\ &\leq C_1 \left( e^{LNh} \|\mathbf{e}_N\|_\infty + (1 - e^{-h}) \sum_{n=0}^{N-1} e^{(L-1)nh} \|\boldsymbol{\tau}_n\|_\infty \right) \\ &\leq C_1 \left( e^{LNh} \|\mathbf{e}_N\|_\infty + h \sum_{n=0}^{N-1} e^{(L-1)nh} \|\boldsymbol{\tau}_n\|_\infty \right). \end{aligned} \quad (43)$$

Notice that the initial error  $\mathbf{e}_N$  is substantially magnified in this estimate: in fact,  $e^{LNh} = (\pi - x_N)^{-L}$ . Moreover the local truncation errors  $\boldsymbol{\tau}_n$  are also multiplied by exponentially growing factors. However, from the hypothesis (35) and by virtue of (32) and (40) it follows that

$$\|\boldsymbol{\tau}_n\|_\infty = O(h_n^p (\pi - x_{n+1})^{\eta-p-1}) = O(h^p e^{-(n+1)h(\eta-1)}) \leq C_2 h^p e^{-nh(\eta-1)};$$

substituting back into (43) and using (37) we obtain

$$\|\mathbf{e}_0\|_\infty \leq C_3 \left( e^{LNh} \|\mathbf{e}_N\|_\infty + h^{p+1} \sum_{n=0}^{N-1} e^{(L-\eta)nh} \right) \leq C_3 \left( e^{LNh} \|\mathbf{e}_N\|_\infty + h^{p+1} \sum_{n=0}^{N-1} e^{\alpha n h^2} \right).$$

Recalling the assumption  $\alpha N h^2 \leq 1$  we deduce that

$$\|\mathbf{e}_0\|_\infty \leq C_3 e^{LNh} \|\mathbf{e}_N\|_\infty + C_4 h^{p+1} N. \quad (44)$$

In order to complete the proof we need to discuss the behaviour of  $\|\mathbf{e}_N\|_\infty$ . From (34), the chosen initial value at  $x_N$ , and the variable stepsize used, one gets

$$\|\mathbf{e}_N\|_\infty \leq C_5 (\pi - x_N)^{\eta+1} = C_5 e^{-Nh(\eta+1)}$$

and hence, see (37),

$$e^{LNh} \|\mathbf{e}_N\|_\infty \leq C_5 e^{(L-\eta)Nh} e^{-Nh} \leq C_5 e^{\alpha N h^2} e^{-Nh} \leq (eC_5) e^{-Nh}.$$

Combining this with (44) and setting  $\bar{C} = \max(C_4, eC_3C_5)$  completes the proof.  $\square$

It is important to observe that in this setting the parameters  $h$  and  $N$  are independent. In order to obtain a solution with a prescribed accuracy, however, the estimate (38) can be conveniently used to determine a relationship between them.

Let us now discuss the hypothesis (36)-(37). Is this reasonable for our system? In a neighbourhood of  $x = \pi$  the system (11) can be written as

$$\mathbf{y}'(x) = \frac{\eta}{\pi - x} [J_{-1} + (\pi - x)J_0(x)] \mathbf{y}(x) \quad (45)$$

where  $\eta$  is defined in (4),

$$J_{-1} = \begin{pmatrix} 0 & 1 \\ 0 & -1 \end{pmatrix} \quad (46)$$

and  $J_0(x) = O(1)$ . The matrix  $J_{-1}$  is diagonalizable:

$$J_{-1} = V D V^{-1}, \quad V = \begin{pmatrix} 1 & 1 \\ 0 & -1 \end{pmatrix} \quad D = \begin{pmatrix} 0 & 0 \\ 0 & -1 \end{pmatrix}$$

and hence a change of variable  $\mathbf{z}(x) = V^{-1}\mathbf{y}(x)$  brings the system into the form

$$\mathbf{z}'(x) = \frac{\eta}{\pi - x} [D + (\pi - x)V^{-1}J_0(x)V] \mathbf{z}(x).$$

Observe that for this system the Lipschitz factor near  $\pi$  is asymptotically

$$\frac{\eta}{\pi - x} (1 + O(\pi - x))$$

and so a wide class of discretization schemes may be expected to have the property (36)-(37). In particular, it is worth mentioning that this is the case for Runge-Kutta methods, as one can verify with some calculations.

#### 4 Numerical results

We present some numerical results to illustrate the theorems of the previous sections. The one-step methods used for solving the initial (or final) value problems involved in the shooting are Runge-Kutta collocation methods whose internal nodes are the roots of scaled and shifted Chebyshev polynomials. In more detail, the abscissae of the  $s$ -stage Runge-Kutta method used are

$$c_i = \frac{1}{2} \left( \cos \left( \frac{(2(s-i)+1)\pi}{2s} \right) + 1 \right), \quad i = 1, \dots, s.$$

The coefficient matrix  $A$  and the vector of weights  $\mathbf{b}$  of the scheme are then determined by imposing the simplifying conditions  $C(s)$  and  $B(s)$ , respectively (see [10]). These methods are of order  $p = s$  if  $s$  is even or  $p = s + 1$  if  $s$  is odd. By virtue of this fact, the schemes we have actually used are those with an odd number of stages. As an example, when applied to forward integration of (11) over a uniform grid with stepsize  $h = \pi/N$  the method corresponding to  $s = 1$  reads

$$\mathbf{y}_{n+1} = \mathbf{y}_n + \frac{h}{2} J(x_{n+\frac{1}{2}}) (\mathbf{y}_n + \mathbf{y}_{n+1}), \quad n = 0, 1, \dots, N-1,$$

which coincides with the second order Gauss Runge-Kutta method. The forward shooting uses the secant method to find the zeroes of the miss-distance (6) starting from an input guess  $\lambda^{(0)}$  and using  $\lambda^{(1)} = \lambda^{(0)} + 0.1$  as a second approximation of the eigenvalue. Table 1 lists the experimentally

**Table 1** Order of convergence in the eigenvalue estimates obtained with forward shooting.

method of order $p = 2$ ( $s = 1$ )									
	$\epsilon = 1$		$\epsilon = 2/3$		$\epsilon = 1/4$		$\epsilon = 1/10$		
$N$	$\delta\lambda(N)$	order	$\delta\lambda(N)$	order	$\delta\lambda(N)$	order	$\delta\lambda(N)$	order	
1600	$2.9820 \cdot 10^{-3}$	—	$2.4793 \cdot 10^{-4}$	—	$9.7357 \cdot 10^{-6}$	—	$1.8400 \cdot 10^{-5}$	—	
3200	$1.4977 \cdot 10^{-3}$	0.99	$8.8393 \cdot 10^{-5}$	1.49	$2.4326 \cdot 10^{-6}$	2.00	$4.5999 \cdot 10^{-6}$	2.00	
6400	$7.5059 \cdot 10^{-4}$	1.00	$3.1411 \cdot 10^{-5}$	1.49	$6.0804 \cdot 10^{-7}$	2.00	$1.1500 \cdot 10^{-6}$	2.00	
method of order $p = 4$ ( $s = 3$ )									
	$\epsilon = 1$		$\epsilon = 2/3$		$\epsilon = 1/4$		$\epsilon = 1/10$		
$N$	$\delta\lambda(N)$	order	$\delta\lambda(N)$	order	$\delta\lambda(N)$	order	$\delta\lambda(N)$	order	
400	$1.3385 \cdot 10^{-3}$	—	$6.6978 \cdot 10^{-5}$	—	$3.9559 \cdot 10^{-8}$	—	$9.9663 \cdot 10^{-10}$	—	
800	$6.6911 \cdot 10^{-4}$	1.00	$2.3602 \cdot 10^{-5}$	1.50	$3.0671 \cdot 10^{-9}$	3.69	$6.3151 \cdot 10^{-11}$	3.98	
1600	$3.3452 \cdot 10^{-4}$	1.00	$8.3298 \cdot 10^{-6}$	1.50	$2.2833 \cdot 10^{-10}$	3.75	$3.9755 \cdot 10^{-12}$	3.99	
method of order $p = 6$ ( $s = 5$ )									
	$\epsilon = 1$		$\epsilon = 2/3$		$\epsilon = 1/4$		$\epsilon = 1/10$		
$N$	$\delta\lambda(N)$	order	$\delta\lambda(N)$	order	$\delta\lambda(N)$	order	$\delta\lambda(N)$	order	
100	$1.9278 \cdot 10^{-3}$	—	$1.1376 \cdot 10^{-4}$	—	$3.9346 \cdot 10^{-8}$	—	$5.2220 \cdot 10^{-11}$	—	
200	$9.6366 \cdot 10^{-4}$	1.00	$4.0075 \cdot 10^{-5}$	1.51	$2.3736 \cdot 10^{-9}$	4.05	$8.4288 \cdot 10^{-13}$	5.95	
400	$4.8174 \cdot 10^{-4}$	1.00	$1.4141 \cdot 10^{-5}$	1.50	$1.4375 \cdot 10^{-10}$	4.05	$1.3323 \cdot 10^{-14}$	5.98	

observed orders of convergence for the eigenvalue estimates obtained for various values of the parameter  $\epsilon$  and various orders of the method. We used the  $2\pi$ -periodic function  $f(x)$  whose restriction to  $[-\pi, \pi]$  is

$$f(x) = \frac{x(\pi^2 - x^2)(\cos x + 2)}{2\pi^2}. \quad (47)$$

In every case we chose  $\lambda^{(0)} = 5$  as initial guess. The order of convergence was computed using the classical formula: if we denote the eigenvalue approximation obtained with  $N$  gridpoints by  $\lambda(N)$  and let  $\delta\lambda(N) = |\lambda(N) - \lambda(2N)|$ , then the order  $r$  of the method is estimated by

$$r = \log_2(\delta\lambda(N)/\delta\lambda(2N)). \quad (48)$$

Table 1 shows that for  $\epsilon = 1$  and  $\epsilon = 2/3$ , the order of convergence is always 1 or  $3/2$  respectively, regardless of the classical order  $p$  of the Runge-Kutta method used. This agrees with Theorem 3 since for two such values of  $\epsilon$  using the function (47) we have  $\eta = -1/(\epsilon f'(\pi)) = 1$  and  $\eta = 3/2$ , respectively.

When  $\epsilon = 1/4$ , on the other hand, the observed order  $r$  depends on the method. In fact:

- with the 1-stage method it coincides with its classical order  $p = 2$  (the formula (48) is too unsophisticated to perceive the factor  $\log(h)$ , see Theorem 3);
- with the 3-stage method of order  $p = 4$ , the values of  $r$  obtained are less stable even though they approach  $p$  as  $N$  increases (we think that this is due to the presence of the factor  $\log^2(h)$ );
- with the method of order  $p = 6$ ,  $r \approx \eta = 4$  as expected.

Finally, for  $\epsilon = 1/10$ , the numerics confirm that  $r \approx p < \eta$  for each method.

In order to better contrast the clearly unsatisfactory results obtained with forward shooting for small values of  $\eta$  with the good performance of forward shooting when  $\eta$  is large, in Figure 4, we have plotted the errors in the eigenvalue estimates obtained for the function  $f(x)$  in (47) and  $\epsilon = 1, 2/3, 1/10$  versus the number of gridpoints. The errors have been evaluated with respect to reference eigenvalues computed with two-sided shooting used as described later. The results reported corresponding to  $\epsilon = 1, 2/3$  show that the rates at which the errors decrease are the same for each method so that the differences in the magnitudes of the errors can only be attributed to the principal error coefficient of the formula. On the other hand, when  $\epsilon = 1/10$  the results obtained are good, in particular for the method of order 6, especially considering that with the chosen initial guess for the root-finding procedure we obtain approximations to higher index eigenvalues than the ones computed for the other two values of  $\epsilon$ .

Motivated by these observations, we applied the two-sided shooting procedure when  $\eta$  is small. In this case, the miss-distance used was

$$\Im \left( \hat{W}(\pi - 1, \lambda) \right) = 0 \quad \text{with} \quad \hat{W}(x, \lambda) = \begin{vmatrix} \psi_1(x, \lambda) & \zeta_1(x, \lambda) \\ \epsilon f(x) \psi_1'(x, \lambda) & \epsilon f(x) \zeta_1'(x, \lambda) \end{vmatrix}. \quad (49)$$

Observe that  $\hat{W}(x, \lambda) = \frac{\epsilon f(x)}{q(x)} W(\lambda)$  with  $W(\lambda)$  defined in (9). The approximations to  $\psi_1(x, \lambda)$  are obtained by applying the relevant Runge-Kutta methods with constant stepsize for integrating (11) over  $[0, \pi - 1]$ , with initial value  $\mathbf{y}(0) = (1, 0)^T$ , while the approximation of  $\zeta_1(x, \lambda)$  is computed by backward integration with variable stepsize as described in Section 3.2. In particular, the error bound (38) has been conveniently used to choose the number of gridpoints of the discretization of the subinterval  $[\pi - 1, \pi]$  and, consequently, the position of the layer near  $\pi$ . The criterion used was the following: if  $h$  is the stepsize of the uniform partition of the left subinterval  $[0, \pi - 1]$  then the value of  $N$  in Theorem 4 has been chosen as

$$N = \lceil -(s + 1) \log(h/2)/h \rceil$$

where  $s$  is the stagelength of the method. In this way, the asymptotic behaviour of the term on the right hand-side of (38) is  $O(h^{s+1} \log(h))$  which coincides with the one of the error near the singularity at  $x = 0$  (see Figures 2-3).

In Figure 5 we report the errors in the eigenvalue estimates obtained against the total number of gridpoints  $N_{tot}$  used for discretizing the overall interval  $[0, \pi]$  for the methods of order 4 and 6,  $\epsilon = 1, 2/3$  and the function  $f(x)$  in (47). The errors were estimated by taking as reference eigenvalues the ones obtained with two-sided shooting,  $p = 6$  and  $N_{tot} = 2477$ . These are the reference eigenvalues we have used also for the results reported in Figure 4. The errors corresponding to the simple forward shooting with the method of order 6 have been also reported in order to better emphasize the great advantages arising from the use of two-sided shooting when  $\eta$  is small.

The aim of the last two examples was to compare the performance of our forward shooting procedure with that of the WKB method proposed by Benilov, O'Brien and Sazonov in [3] and the shooting procedure proposed by Chugunova and Volkmer in [6], where the results are given for the special case  $f(x) = \sin(x)$ . WKB uses an asymptotic expansion of the solutions of (1) with respect to the parameter  $\epsilon$  which, therefore, is assumed to be sufficiently small. The shooting procedure in [6], instead, solves the more general problem (7) with  $f(x) = \sin(x)$  and  $g(x) = -\cos(x)$  which clearly reduces to (1) if  $a = 0$ .

In Table 2, the eigenvalue estimates obtained for  $\epsilon = 0.13$  with our methods of orders  $p = 4, 6$  with  $N = 300$  gridpoints and forward shooting are listed. In this table  $k$  denotes the mode number while  $\lambda_{k,p}$  is the approximation of the  $k$ -th eigenvalue provided by the method of order  $p$ . The corresponding estimates  $\lambda_{k,WKB}$  and  $\mu_k$  given by the WKB method and the method in [6], respectively, are also reported. The differences reported in the last three columns of the table can be regarded as estimates of the error for the methods of order 4, WKB and the method in [6], respectively. As one can see, our method provides definitely more accurate approximations than the other two. It must

be said, however, that the accuracy of WKB is independent of the mode number of the eigenvalues while for our method it deteriorates as the index increases. This usually happens when standard discretization schemes are used for the approximation of the spectrum of differential operators; it can often be avoided by the use of specially designed Magnus methods [9] though these would have to be carefully chosen to behave well near singularities.

In view of the competitiveness of our method when  $a = 0$ , we decided to test it also for the case of a nonzero value for  $a$  even though the theory is still lacking. We reformulated (7) as a system of first order ODEs as described in (11)-(12) where the coefficient matrix is now

$$J(x) = \begin{pmatrix} 0 & (\epsilon f(x))^{-1} \\ -i\lambda - ag'(x) & -(1 + ag(x))(\epsilon f(x))^{-1} \end{pmatrix}.$$

In Table 3 we list the eigenvalues computed with  $N = 300$  by the methods of order 4 and 6, together with the estimates provided by the shooting method in [6]. The differences reported in the second last column suggest that our methods can tackle this more general case too. However the rigorous error analysis which we presented here does not immediately generalize to the case  $a \neq 0$ , with general functions  $g$  and  $f$  satisfying the properties specified in the proof of Theorem 1, because the logarithmic semi-norm  $\mu_1(J(x))$  is no longer bounded (see Proposition 4). We shall defer consideration of this case to future work.

**Table 2** Comparison of our results with those provided by the WKB method and the method in [6] for  $f(x) = \sin(x)$  and  $\epsilon = 0.13$ .

$k$	$\lambda_{k,4}$	$\lambda_{k,6}$	$\lambda_{k,WKB}$	$\mu_k$	$ \lambda_{k,4} - \lambda_{k,6} $	$ \lambda_{k,WKB} - \lambda_{k,6} $	$ \mu_k - \lambda_{k,6} $
1	1.016031	1.016031	1.016306	1.016070	$5.804468 \cdot 10^{-12}$	$2.752395 \cdot 10^{-4}$	$3.894446 \cdot 10^{-5}$
2	2.118176	2.118176	2.119515	2.118266	$1.671983 \cdot 10^{-10}$	$1.339477 \cdot 10^{-3}$	$9.039683 \cdot 10^{-5}$
3	3.359418	3.359418	3.361995	3.359580	$1.076850 \cdot 10^{-9}$	$2.576915 \cdot 10^{-3}$	$1.623123 \cdot 10^{-4}$
4	4.764312	4.764312	4.767907	4.764572	$3.648983 \cdot 10^{-9}$	$3.595254 \cdot 10^{-3}$	$2.598077 \cdot 10^{-4}$
5	6.343734	6.343734	6.348097	6.344121	$8.212884 \cdot 10^{-9}$	$4.362362 \cdot 10^{-3}$	$3.868387 \cdot 10^{-4}$
6	8.102844	8.102844	8.107778	8.103393	$1.235901 \cdot 10^{-8}$	$4.933742 \cdot 10^{-3}$	$5.487597 \cdot 10^{-4}$
7	10.044314	10.044314	10.049678	10.045061	$6.585230 \cdot 10^{-9}$	$5.364257 \cdot 10^{-3}$	$7.469191 \cdot 10^{-4}$
8	12.169637	12.169637	12.175331	12.170624	$3.378970 \cdot 10^{-8}$	$5.694361 \cdot 10^{-3}$	$9.872458 \cdot 10^{-4}$
9	14.479701	14.479701	14.485653	14.480973	$1.626462 \cdot 10^{-7}$	$5.952096 \cdot 10^{-3}$	$1.272263 \cdot 10^{-3}$
10	16.975062	16.975062	16.981219	16.976669	$4.859431 \cdot 10^{-7}$	$6.156774 \cdot 10^{-3}$	$1.606505 \cdot 10^{-3}$

We conclude by summarizing the results of the numerical experiments:

- the numerically observed orders of convergence in the eigenvalue estimates computed with forward shooting and listed in Table 1 confirm the result of Theorem 3;
- a catastrophic loss of order is observed when  $\eta = -1/(\epsilon f'(\pi))$  is close to 1 due to the singularity at  $x = \pi$ ; this severe drawback can be overcome by using the proposed two-sided shooting;
- by contrast, when  $\eta$  is small enough, forward shooting works very well, as confirmed by the last examples where our estimates have been compared with those provided by the WKB method proposed in [3] and the shooting procedure proposed in [6].

## References

1. Aceto, L.; Ghelardoni, P., Marletta, M.: Numerical solution of forward and inverse Sturm-Liouville problems with an angular momentum singularity. *Inverse Problems* **24** no. 1, 015001 (2008)
2. Aguilar, J.; Combes, J. M.: A class of analytic perturbations for one-body Schrödinger Hamiltonians. *Comm. Math. Phys.* **22**, 269–279 (1971)
3. Benilov, E. S., O'Brien, S. B. G., Sazonov, I. A.: A new type of instability: explosive disturbances in a liquid film inside a rotating horizontal cylinder. *J. Fluid Mech.* **497**, 201–224 (2003)



**Table 3** Comparison of our results with respect to those provided by the method in [6] for problem (7) for  $f(x) = \sin(x)$ ,  $g(x) = -\cos(x)$ ,  $\epsilon = 0.13$ , and  $a = 0.4$ .

$k$	$\lambda_{k,4}$	$\lambda_{k,6}$	$\mu_k$	$ \lambda_{k,4} - \lambda_{k,6} $	$ \mu_k - \lambda_{k,6} $
1	0.911658	0.911658	0.911679	$2.073419 \cdot 10^{-11}$	$2.146221 \cdot 10^{-5}$
2	1.937155	1.937155	1.937209	$4.860525 \cdot 10^{-10}$	$5.366715 \cdot 10^{-5}$
3	3.128265	3.128265	3.128365	$2.844630 \cdot 10^{-9}$	$1.000899 \cdot 10^{-4}$
4	4.500009	4.500009	4.500174	$9.548147 \cdot 10^{-9}$	$1.647988 \cdot 10^{-4}$
5	6.056619	6.056619	6.056869	$2.371329 \cdot 10^{-8}$	$2.496469 \cdot 10^{-4}$
6	7.799412	7.799412	7.799770	$4.862277 \cdot 10^{-8}$	$3.581424 \cdot 10^{-4}$
7	9.728819	9.728819	9.729310	$8.696231 \cdot 10^{-8}$	$4.908783 \cdot 10^{-4}$
8	11.844977	11.844977	11.845628	$1.396452 \cdot 10^{-7}$	$6.511694 \cdot 10^{-4}$
9	14.147916	14.147916	14.148757	$2.039943 \cdot 10^{-7}$	$8.413897 \cdot 10^{-4}$
10	16.637629	16.637629	16.638692	$2.709244 \cdot 10^{-7}$	$1.062888 \cdot 10^{-3}$

4. Boulton, L., Levitin, M., Marletta, M., On a class of non-self-adjoint periodic eigenproblems with boundary and interior singularities. To appear in J. Differential Equations.
5. Chugunova, M., Pelinovsky, D.: Spectrum of a non-self-adjoint operator associated with the periodic heat equation. J. Math. Anal. Appl. **342** no. 2, 970–988 (2008)
6. Chugunova, M., Volkmer, H.: Spectral analysis of an operator arising in fluid dynamics. Stud. Appl. Math. **123** no. 3, 291–309 (2009)
7. Chugunova, M., Karabash, I.M., Pyatkov, S.G.: On the nature of ill-posedness of a forward-backward heat equation. Integral Equations Operator Theory **65**, 319–344 (2009)
8. Davies, E. B.: An indefinite convection-diffusion operator. LMS J. Comput. Math. **10**, 288–306 (2007)
9. *Highly oscillatory problems*. Edited by Bjorn Engquist, Athanasios Fokas, Ernst Hairer and Arieh Iserles. London Mathematical Society Lecture Note Series, 366. Cambridge University Press, Cambridge, 2009. xiv+239 pp. ISBN: 978-0-521-13443-9
10. Hairer, E., Nørsett, S.P., Wanner, G.: Solving Ordinary Differential Equations I, 2<sup>nd</sup> ed. Springer-Verlag, Berlin, 1993
11. Weir, J.: An indefinite convection-diffusion operator with real spectrum. Appl. Math. Lett. **22** no. 2, 280–283 (2009)

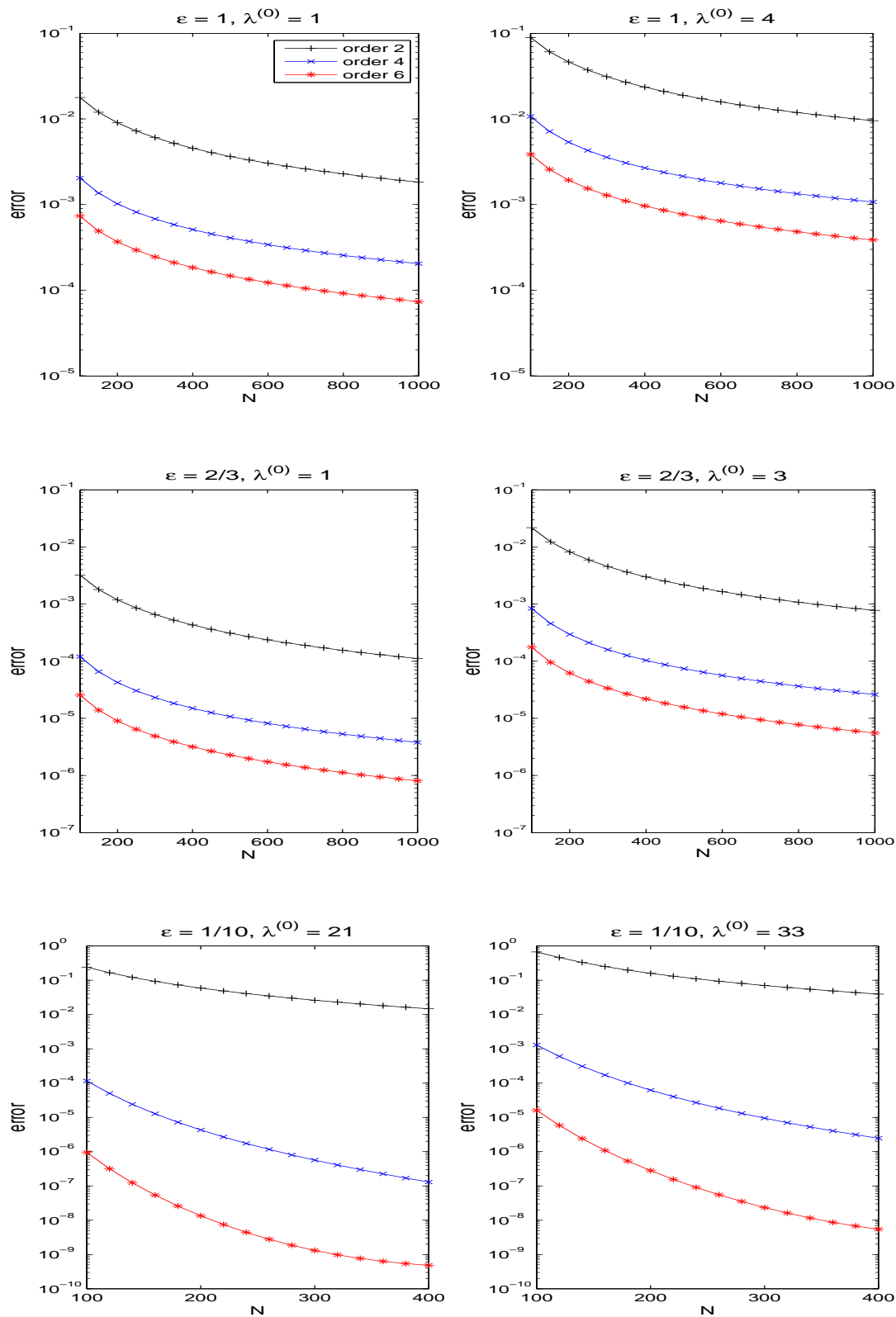


Fig. 4 Errors in the eigenvalue estimates computed with forward shooting and fixed stepsize.

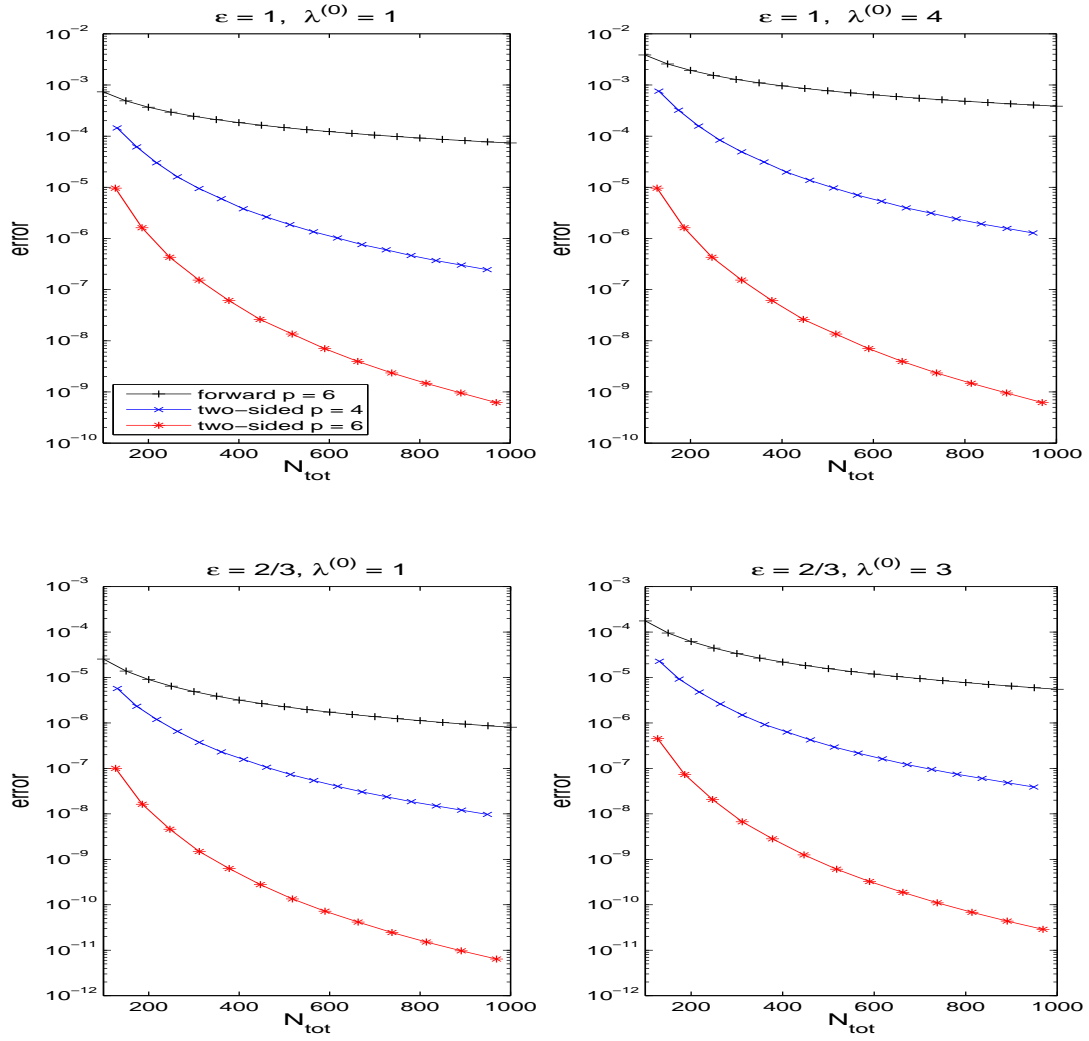


Fig. 5 Errors in the eigenvalue estimates for two-sided shooting with variable stepsize.