

STIFF SYSTEMS OF ORDINARY DIFFERENTIAL EQUATIONS. PART 1. COMPLETELY STIFF, HOMOGENEOUS SYSTEMS

J. J. MAHONY and J. J. SHEPHERD

(Received 8 November 1979),

(Revised 10 December 1980)

Abstract

For the completely stiff real homogeneous system

$$\epsilon \dot{\mathbf{x}} = A(t, \epsilon)\mathbf{x},$$

where ϵ is a small positive parameter, a method is given for the construction of a basis for the solution space.

If A has n linearly independent eigenvector functions, then there exists a choice of these, $\{s_i\}$, with corresponding eigenvalue functions $\{\lambda_i\}$, such that there is a local basis for solution, that takes the form

$$\left\{ [s_i + v_i] \exp \left[\epsilon^{-1} \int^t \lambda_i \right] \right\},$$

where v_i is a vector that tends to zero with ϵ . In general, a basis of this form exists only on an interval in which the distinct eigenvalues have their real parts ordered. A construction is provided for continuing any solution across the boundaries of any such interval. These results are proved for a finite or infinite interval for which there are only a finite number of points at which the ordering of the real parts of eigenvalues changes.

1. Introduction

Stiff is an adjective, frequently used by people interested in numerical solution procedures, to describe a class of differential equations of the form

$$\dot{\mathbf{z}} = \mathbf{h}(t, \mathbf{z}). \tag{1.1}$$

Here, \mathbf{z} and \mathbf{h} are real n -vector valued functions, and the dot denotes differentiation with respect to the real variable t . Standard numerical integration techniques, based on a step size h , may be applied to such a system, but are suspect

whenever the quantity $|\lambda h|$ is not small. In the above, λ is the eigenvalue of maximum modulus of the Jacobian derivative of \mathbf{h} evaluated at any point (t, \mathbf{z}) relevant to the solution. In these circumstances, the term stiff is used when practical considerations place a lower bound on h such that $|\lambda h|$ fails to be small over a range of t -values too large to be covered by the use of local analytic approximations.

For systems satisfying this criterion, useful algorithms for the location of the solutions of the relevant boundary value problems are not readily available. This is due, in part, to a lack of knowledge of the mathematical properties of these exact solutions. The standard system that has been studied to provide such knowledge is

$$\left. \begin{aligned} \epsilon \dot{\mathbf{x}} &= \mathbf{f}(t, \mathbf{x}, \mathbf{y}, \epsilon), \\ \dot{\mathbf{y}} &= \mathbf{g}(t, \mathbf{x}, \mathbf{y}, \epsilon), \end{aligned} \right\} \quad (1.2)$$

where ϵ is a small positive parameter. Although this assumed known division into stiff component \mathbf{x} and nonstiff component \mathbf{y} is not always available for the more general system (1.1), the form (1.2) is convenient for generating a system having Jacobian derivative for which some of the eigenvalues are large. There is an extensive literature for such a system (see, for instance, Vasil'eva [9]), but it does not deal with many of the cases which cause difficulties in numerical procedures as the real parts of large eigenvalues are required to be one-signed over the interval.

A heuristic argument frequently used in connection with such systems is based on the idea, or hope, that $\epsilon \dot{\mathbf{x}}$ will be small almost everywhere, so that the system may be approximated by the system

$$\left\{ \begin{aligned} \mathbf{0} &= \mathbf{f}(t, \mathbf{x}_0, \mathbf{y}_0, \epsilon), \\ \dot{\mathbf{y}}_0 &= \mathbf{g}(t, \mathbf{x}_0, \mathbf{y}_0, \epsilon), \end{aligned} \right\} \quad (1.3)$$

termed the reduced system. This assumption reduces the order of the system to be studied and hence, in general, the number of boundary conditions which may be satisfied. The ability to meet the original boundary conditions is regained by the appending of suitable localized corrections to the solution of (1.3). The literature reviewed in [9] provides a set of sufficient conditions under which this method may be used.

For natural initial value problems, where all large eigenvalues have negative real parts, there are well proven numerical procedures for the integration of stiff systems: see, for example, the review article by Shampine and Gear [6]. The criticisms one hears of these methods appear to arise from attempts to apply them to boundary value problems.

The heuristic basis of most approximation arguments for dealing with stiff systems is the idea that, if a differential equation is satisfied with a small error,

the trial solution proposed will be close to an exact solution. That this is an overly optimistic view can be shown by considering the simple example

$$\varepsilon \dot{x} = -t^2(x^2 - t^2), \quad |t| < 1. \quad (1.4)$$

The behaviour of solutions of this equation can be deduced by phase plane methods, the details of which will not be presented here. The general behaviour is exhibited by the representative curves displayed in Figure 1. Note that, for small values of ε , curves not close to the curves $x_0(t)$ satisfying $t^2(x_0^2 - t^2) = 0$ are extremely steep. Thus any solution which remains bounded lies close to $x_0^2 - t^2 = 0$ for most values of t . Further, it is for only a very restricted range of values of $x(-1, \varepsilon)$ that solutions remain close to this pair of all t in $[-1, 1]$. It may be observed that while $x_0 = -t$ approximately satisfies the differential equation everywhere, there is no exact solution close to it even on the open interval $(-1, 1)$. On the other hand, there *are* exact solutions close to the other three heuristically generated trials $x_0 = t$ and $x_0 = \pm |t|$.

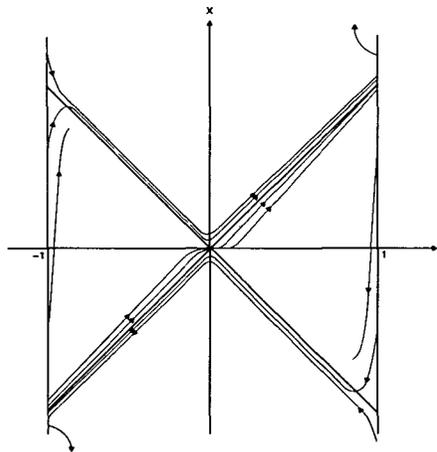


FIG. 1. Displaying solutions of $\varepsilon \dot{x} = -t^2(x^2 - t^2)$.

The basis of most of the effective existence proofs for systems like (1.2) is the establishment of a trial solution \mathbf{X}, \mathbf{Y} satisfying (1.2) in some appropriate approximate sense, a linearization of the differential equation about \mathbf{X}, \mathbf{Y} , and proof of the existence, with appropriate estimates, of suitable corrections $(\mathbf{x} - \mathbf{X}), (\mathbf{y} - \mathbf{Y})$. The methods almost invariably utilize the Contraction Mapping Theorem, and are thus constructive by nature. However, the example (1.4) provides two distinct cases where such methods are unlikely to work. There is a small set of solutions which lies close to $x = -|t|$ whose range of values of $x(-1, \varepsilon)$ is so small that it is unlikely that one could get close to them to contract onto them. However, dichotomy arguments in association with the continuous dependence of $x(t, \varepsilon)$ on $x(-1, \varepsilon)$ permit one to show that such solutions do

indeed exist. Mathematically, such proofs would employ a more general fixed point theorem with a loss of uniqueness, in general. This feature causes no problem to a shooting method approach which merely seeks to straddle the interesting range; but in many solutions the potential nonuniqueness of the fixed point can cause difficulties. In the example above, there are two families of solutions, one like $x_0 = t$ and the other like $x_0 = -|t|$, whose boundary values $x(-1, \epsilon)$ lie extremely close together. It is clear that, to obtain a practical hold on both these types of solution, one will have to formulate the problem in such a way that this over sensitive dependence on $x(-1, \epsilon)$ is suppressed. Even further, it is clear that, for systems of larger order, extremely delicate questions regarding behaviour of solutions can arise.

It is the principal objective of the present series of papers, of which this is the first, to provide a mathematical framework which offers some possibility of investigating the difficulties that may arise in a given stiff system. We hope to provide guidance regarding the type of difficulties likely to be encountered by any robust software package in applications to boundary value problems for stiff systems. To this end, and in view of the above linearization procedure for (1.2), we see an examination of the basic system

$$\left. \begin{aligned} \epsilon \dot{\mathbf{x}} &= A\mathbf{x} + B\mathbf{y} + \mathbf{r}_1, \\ \dot{\mathbf{y}} &= C\mathbf{x} + D\mathbf{y} + \mathbf{r}_2 \end{aligned} \right\} \quad (1.5)$$

to be of fundamental interest, and will seek to characterize those equations and boundary value problems for which difficulties will be encountered. For suitably strong results about (1.5), we might hope to handle nonlinear systems by applying some form of fixed-point theorem.

The papers are structured as follows. In the present paper we consider the completely stiff homogeneous system

$$\epsilon \dot{\mathbf{x}} = A\mathbf{x}, \quad (1.6)$$

and investigate the nature of the most useful basis for the solution space. In the second paper, we investigate inhomogeneous systems of the form

$$\epsilon \dot{\mathbf{x}} = A\mathbf{x} + \mathbf{r}_1, \quad (1.7)$$

and obtain conditions for the occurrence of behaviour of the kind germane to the difficulties described above for the simple example (1.4). The relationships between properties of the linear system and the failures of the nonlinear system are discussed. These two papers establish the basic methods to be used, and the third paper extends the results to the more general system (1.5).

There already exists an extensive literature concerning a suitable basis for the solution space of the homogeneous system relevant to (1.5), with (1.6) as a special case: see, for example, the review by Wasow in [10, page 3]. The greater portion of this work is based on the assumption that A is an analytic function of

t. In the context of the numerical solution of stiff systems, A will derive from $f_x(t, X, Y)$, while X will possess a finite (normally small) number of continuous derivatives, so that it is necessary to relax the analyticity assumption. In the work cited, this analyticity assumption is applied in two quite distinct ways. One use is merely to construct a formal approximation (that is, X, Y) to a solution. This step is still possible, to a limited extent, in the present context. The other use, however, involves the path independence of integrals, to obtain tight estimates on solutions. This step is not assumed to be available here, and we will show that the existing results may be substantially modified as a consequence.

There turns out to be a significant difference between the results for analytic and non-analytic systems which has serious implications for numerical integration techniques. Analytic theory is generally based on the assumption that the matrix A in (1.6) has an analytic diagonal Jordan form in a sufficiently large domain in the complex plane. Then equation (1.6) has a basis of solutions of the form $\{s_i + \epsilon v_i\} \exp\{\epsilon^{-1} \int^t \lambda_i\}$, where s_i and λ_i are independent eigenvectors and associated eigenvalues of A . If two eigenvalue functions become equal at some point in the complex plane, in general the above form of the solution basis set becomes local rather than global. From each such point there will emerge a Stokes line across which the exponentially dominant order of the corresponding solution will change. Such is the complexity of the case that no serious effort appears to have been made to study global properties of solutions on the real line when the restrictive assumption on eigenvalues in the complex plane is relaxed. But it may be observed that the global nature of the validity of the usual basis on the real line no longer necessarily holds. Testable assumptions on the real axis are not available in the analytic theory for questions of significance in numerical studies. We show here that, for non-analytic coefficients defined on the real line, there is a change of exponential dominance of some solutions wherever there is a change in the order of the real parts of the eigenvalue functions. The implications of this for numerical integration needs emphasis. Consider an analytic system (1.6) which has no points where two eigenvalue functions of the Jacobian derivative become equal on the real line but where the real parts of eigenvalues change order. Under these circumstances, a linearization about a numerical solution will lead to an approximate basis of quite different form from the exact analytic theory. However, there will certainly be cases where a Stokes line emerges from a nearby point in the complex plane and cuts the real axis. Though the non-analytic result will be in error, it could be equally wrong to assume that the correct result contains no change in exponential order of the solution basis. Numerical solution techniques for non-linear analytic systems will provide no information about solutions off the real line. Hence it is desirable to understand the implications of this phenomenon for the validity of various numerical or analytic approximations.

However, we will be invoking the smallness of ϵ to establish the validity of approximations via the use of the contraction mapping theorem. The process envisaged is as follows. The numerical approximation has been obtained for a given small value ϵ . The required coefficients can be estimated and tested to see whether the contraction mapping theorem could be invoked. It is possible that this method would never lead to the establishment of the existence of a solution. But from a numerical viewpoint it is equally impossible to establish uniform bounds on coefficients for sufficiently small ϵ . We take the view that, in attempting a practical judgement as to whether to accept the result of a numerical integration procedure, it is appropriate to test for the validity of the use of contraction mapping theorem for a given ϵ . We shall insist that the bounds established for the particular value ϵ_0 , if applied as uniform in ϵ , must imply that the result would be true in the limit ϵ tending to zero.

In Part 1 we shall largely present the arguments in terms of appropriate bounds where substantive results are involved but revert to the standard $O(\epsilon)$ and $O(\epsilon)$ notation where the presentation would otherwise become too complex.

There is one further aspect of the present problem that enforces a substantial departure from the standard asymptotic analysis discussion. For large order systems, and possibly even for small order systems, the trial solution for the nonlinear problem will almost certainly be derived by numerical means. This will happen for a discrete set of values of ϵ , and interest will be centred on what happens, for a given ϵ , in the limit of step size h tending to zero, rather than as ϵ tends to zero.

2. Background to the basic assumptions

Suppose that $A(t)$ is a smooth matrix function on a given interval of t -values possessing n eigenvalues λ_k , and corresponding eigenvectors s_k , that are distinct at each value of t . Then it is well known that the differential equation $\epsilon \dot{x} = A(t)x$ has n formal solutions of the form $\{s_k + o(1)\} \exp(\epsilon^{-1} \int^t \lambda_k)$, where $o(1)$ is a term that vanishes with ϵ . Nayfeh [5, page 332] demonstrates how to generate further terms in such a formal expansion. When A is analytic in a region of the complex plane, Wasow [9] demonstrates methods by which these formal solutions may be shown to be asymptotic approximations to exact solutions of the above system.

Central to these considerations is the assumption that the eigenvalues of A are distinct; for then a simple diagonal canonical form for A is available. Sibuya [7] has considered the case where this fails, and has shown that, given an $A(t)$ analytic in the above sense, a smooth basis exists that reduces A to a canonical

form on any open interval in which the multiplicity signature of the eigenvalues remains constant. However, this canonical form is not diagonal in general, and the smoothness of the basis cannot be guaranteed at the end points of such an interval.

In the following sections, we propose to consider the solutions of the above system on intervals in which changes of this multiplicity signature may occur, but on which a smooth eigenvector basis is maintained throughout. Our motivation for such a choice of a smaller class of problems lies in a desire to consider the implications of changes in the multiplicity of eigenvalues, but to avoid the involved analysis that results from a direct application of the canonical forms of A constructed by Sibuya to the system above.

With these ideas in mind, we are led to make our first basic assumption regarding the matrix A occurring in the equation (1.6).

ASSUMPTION 1. *For each t on an interval of interest I (which could be $[0, T]$ for some $T > 0$, or $[0, \infty)$), and for any given positive ϵ lying in a suitable neighbourhood of zero, $A(t, \epsilon)$ is a linear map from R^n to R^n , to which correspond n linearly independent eigenvectors $s_k(t, \epsilon)$ and eigenvalues $\lambda_k(t, \epsilon)$, with both sets of functions being continuously differentiable as functions of t on I .*

Note that we have allowed $A(t, \epsilon)$ to depend on ϵ ; this should cause no real alteration to the above arguments regarding motivation for the present approach.

Assumption 1 assures us that $A(t, \epsilon)$ may be written in the canonical form

$$A(t, \epsilon) = S(t, \epsilon)\Lambda(t, \epsilon)S^{-1}(t, \epsilon) \quad (2.1)$$

for each such t and ϵ , where the $n \times n$ matrices S and Λ are defined by

$$S = [s_1 s_2 \cdots s_n] \quad (2.2)$$

and

$$\Lambda = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_n). \quad (2.3)$$

REMARKS. 1. Complex eigenvalues and eigenvectors are admissible, but they may be arranged in conjugate pairs. The assumption of n linearly independent differentiable eigenvectors rules out the possibility of an eigenvector function's being complex on part of I and not on the rest since, at the changeover point, the continuous imaginary parts of the pair of conjugate complex eigenvectors vanish, so that the set of eigenvectors is not linearly independent at such a point.

2. Systems generated from an n th-order differential equation by the standard means can satisfy this assumption only if all eigenvalues have multiplicity one.

3. At any point t_0 where the multiplicity signature of the eigenvalues changes, Assumption 1 is restrictive on the behaviour, but only mildly so. Some indication of the nature of this restriction may be obtained from the following analysis.

Consider the case of every eigenvalue's being distinct on $t > 0$, with two real eigenvalues λ_1 and λ_2 , say, being equal at the endpoint $t = 0$. Then Sibuya's results referred to the above show that, in $t > 0$, there is a basis of eigenvector functions such that

$$A = S\Lambda S^{-1}, \tag{2.4}$$

with S and Λ as defined above.

Differentiation and a little manipulation yields, on $t > 0$,

$$\dot{A}S - S\dot{\Lambda} = \dot{S}\Lambda - A\dot{S}. \tag{2.5}$$

If, for each $t > 0$, we make the substitution

$$\dot{s}_i = \sum_j m_{ij}s_j, \tag{2.6}$$

the i th column in the matrix representation of the right-hand side of (2.5) is

$$\sum_j m_{ij}(\lambda_i - \lambda_j)s_j. \tag{2.7}$$

For given continuously differentiable A and Λ , the corresponding column of the left-hand side of (2.5) is given by

$$\sum_j D_{ij}s_j, \tag{2.8}$$

where D_{ij} are coefficients defined by \dot{A} and $\dot{\Lambda}$ from any set of eigenvectors at a given value of t . Under the assumptions on A and Λ , these must be bounded and continuous on $t \geq 0$. Moreover, since

$$m_{ij}(\lambda_i - \lambda_j) = D_{ij}, \tag{2.9}$$

it follows that the m_{ij} are bounded on the closed interval $t \geq 0$, unless $\lambda_i = \lambda_j$. Thus, under the assumptions above, only m_{12} and m_{21} (and *perhaps* m_{11} and m_{22}) may be unbounded at $t = 0$.

If we now restrict attention to the first two solutions of (2.5), and note (2.6), we may readily derive differential equations for s_1 and s_2 that take the form

$$(\lambda_1 - \lambda_2)\dot{s}_1 = m_{11}s_1 + B_1 + C_1s_2 \tag{2.10}$$

and

$$(\lambda_2 - \lambda_1)\dot{s}_2 = m_{22}s_2 + B_2 + C_2s_1, \tag{2.11}$$

where B_1 and B_2 are bounded maps whose domains lie in the space spanned by s_3, \dots, s_n , while C_1 and C_2 are bounded maps that depend on the properties of A and $\dot{\Lambda}$.

This equation shows that s_1 and s_2 satisfy a pair of differential equations that display a singular point at $t = 0$, because of the zero of $\lambda_1 - \lambda_2$ there. The nature of this singular point depends on the nature of this zero, and those of m_{11} and m_{22} . If m_{11} and m_{22} are bounded away from zero at $t = 0$, and the zero of $\lambda_1 - \lambda_2$ there is simple, this system may be expected to be dominated by a regular singular point system, and to behave like one. Thus s_1 and s_2 can behave like positive or negative powers of t near $t = 0$, and whether or not Assumption 1 is met there is highly dependent on the detailed structure of the right-hand sides of (2.10) and (2.11).

The above argument demonstrates that, at such a point, many different cases can arise, but Assumption 1, as well as others to be made subsequently, leave the number of cases covered by the material of this paper far from vacuous.

The cases where \dot{A} is not bounded, even though A is, add further to the great variety of cases. Thus the form of our assumption is quite unsuited to the general case. In Section 6 we shall look at appropriate ways of dealing with the cases not covered.

3. Existence of a basis for the solution space

We now turn to the task of constructing a useful basis for the solution space of the system (1.6) on a closed interval $J \subseteq I$, where I is the interval described in Assumption I. Our choice of J is motivated by the observation that Assumption 1 alone is not sufficient to construct a basis having the properties we desire, and so we imagine $A(t, \epsilon)$ to be further restricted on J . We further envisage that J could depend on ϵ , although it will become clear that we settle on a particular value of this parameter for which certain estimates hold.

Our first assumption concerns the behaviour of the eigenvalues and eigenvectors of A on J .

ASSUMPTION 2. For any given $\epsilon > 0$, a set of $m \leq n$ distinct differentiable eigenvalue functions $\lambda_\alpha(t, \epsilon)$ of A may be defined on J such that, for each eigenvalue $\lambda_i(t, \epsilon)$, we have $\lambda_i(t, \epsilon) \equiv \lambda_\alpha(t, \epsilon)$ for some α and all $t \in J$.

DEFINITION 1. For any given $\epsilon > 0$, and corresponding to each eigenvalue function λ_α , we define, at each $t \in J$, a vector subspace R_α , spanned by all eigenvectors $s_k(t, \epsilon)$ of A corresponding to eigenvalues $\lambda_k(t, \epsilon)$ such that $\lambda_k(t, \epsilon) = \lambda_\alpha(t, \epsilon)$ at the given value of t .

REMARK. We note that, under Assumptions 1 and 2, $\dim R_\alpha$ is constant on J , for each α , and $\cup_\alpha R_\alpha = R^n$ throughout J .

Our next assumption concerns the relative properties of the various eigenvalue functions λ_α defined on J .

ASSUMPTION 3. *For any given $\epsilon > 0$, and each pair α and β , the zeros of $\text{Re}(\lambda_\alpha - \lambda_\beta)$ and $\text{Im}(\lambda_\alpha - \lambda_\beta)$ on J are finite in number and of finite order, while either*

(a) $\text{Re}(\lambda_\alpha - \lambda_\beta)$ does not change sign on J ,

or

(b) $\text{Re}(\lambda_\alpha - \lambda_\beta) \equiv 0$ while $\text{Im}(\lambda_\alpha - \lambda_\beta) \not\equiv 0$ on J .

REMARK. When, for some $\epsilon > 0$, Assumptions 1, 2 and 3 define such an interval J , the end-points of J may be identified as the zeros of $\text{Re}(\lambda_\alpha - \lambda_\beta)$, for some α and β at which there is a sign change, or the boundaries of an interval on which $\text{Re}(\lambda_\alpha - \lambda_\beta)$ vanishes identically, or points at which the multiplicity of an eigenvalue changes, or an original end point of I . However, although we have defined J to be closed, we envisage that the above assumptions could hold on a semi-infinite interval (when I were infinite), with the relevant properties holding for the given ϵ as $t \rightarrow \infty$.

Related to the behaviour of eigenvalue functions as given above, we will find it useful to define some terms that will have considerable significance later in this section.

DEFINITION 2. *For each α, β and given $\epsilon > 0$, we define*

(i) $\lambda_\alpha < \lambda_\beta$ on J if and only if $\text{Re}(\lambda_\beta - \lambda_\alpha) > 0$ and $\text{Re}(\lambda_\beta - \lambda_\alpha) \not\equiv 0$ for all $t \in J$;

(ii) $\lambda_\alpha \simeq \lambda_\beta$ on J if and only if $\text{Re}(\lambda_\beta - \lambda_\alpha) \equiv 0$ there.

Further, we define, for the eigenvectors s_k of A ,

(i) $s_p < s_q$ if and only if $s_p \in R_\alpha, s_q \in R_\beta$ and $\lambda_\alpha < \lambda_\beta$;

(ii) $s_p \simeq s_q$ if and only if $s_p \in R_\alpha$ and $s_q \in R_\beta$ for some α and β ; for which $\lambda_\alpha \simeq \lambda_\beta$.

It is clear that the symbol $<$ defined above defines an ordering on the set of eigenvalue functions. Note also that the symbols $<$ and \simeq have been used to define a similar relation on the set of eigenvectors. No confusion should result from this usage.

We may now use the relations of Definition 2 to define some vector subspaces, generated by the eigenvectors, that will greatly assist our notation in later results.

DEFINITION 3. Corresponding to any eigenvector $\mathbf{s}_k(t, \epsilon)$, and a given interval J , we define, for a given $\epsilon > 0$, a partition of R^n into two subspaces:

R_{kl} , spanned by \mathbf{s}_k and all \mathbf{s}_j such that $\mathbf{s}_j \prec \mathbf{s}_k$, and

R_{kg} , its complement in R^n .

We also define the subspace $R_{k\approx}$ to be the subspace spanned by \mathbf{s}_k and all \mathbf{s}_j such that $\mathbf{s}_j \simeq \mathbf{s}_k$.

REMARK. In view of Definition 2, we see that, for any given k , the elements \mathbf{s}_j of R_{kl} and R_{kg} are in one-to-one correspondence with terms $\exp(\epsilon^{-1} f^t \lambda_j)$ of exponential order lower, greater or equal to that of $\exp(\epsilon^{-1} f^t \lambda_k)$.

We now return to the problem of constructing a basis for the system (1.6). Defining the variables

$$\zeta_i(t, \epsilon) = \epsilon^{-1} \int^t \lambda_i(s, \epsilon) ds \quad (3.1)$$

for each i and $t \in J$, where the lower limit can be taken as some point not strictly interior to J , we see that the heuristics described in [5] imply that SZ would be a suitable approximate fundamental matrix for this system, where Z is the diagonal matrix

$$Z(t, \epsilon) = \exp \operatorname{diag}[\zeta_1, \zeta_2, \dots, \zeta_n]. \quad (3.2)$$

Prompted by this observation, and the methods employed in a recent paper by Chapman and Mahony [2], we attempt to establish the existence of solutions to (1.6) that have the form

$$\mathbf{x} = SZ\mathbf{B}, \quad (3.3)$$

where $\mathbf{B}(t, \epsilon)$ is a vector-valued function. Because both S and Z are nonsingular for some $\epsilon > 0$ and all $t \in J$, there is a one-to-one correspondence between functions \mathbf{x} and \mathbf{B} so related.

It follows from the definition of Λ that

$$\epsilon \dot{Z} = \Lambda Z = Z \Lambda, \quad (3.4)$$

and it is an easy calculation to verify the existence of a one-to-one correspondence between solutions of (1.6) and those of

$$\dot{\mathbf{B}} = Z^{-1} H Z \mathbf{B}, \quad (3.5)$$

where $H(t, \epsilon)$ is the $n \times n$ matrix defined by

$$H(t, \epsilon) = [h_{ik}(t, \epsilon)] = -S(t, \epsilon)^{-1} \dot{S}(t, \epsilon), \quad (3.6)$$

for given $\epsilon > 0$ and all $t \in J$.

The equation (3.5) is readily converted to an integral equation on J , of the form

$$\mathbf{B} = \mathbf{B}_0 + \int' Z^{-1} H Z \mathbf{B} \, ds, \tag{3.7}$$

where \mathbf{B}_0 , a constant vector (that could vary with ϵ), and one terminal of integration are, as yet, unspecified.

Solving this equation would generate solutions

$$\mathbf{x} = S Z \mathbf{B} = S Z \mathbf{B}_0 + S Z \int' Z^{-1} H Z \mathbf{B} \, ds. \tag{3.8}$$

This equation makes it clear that \mathbf{B} occurs only through the combination

$$\mathbf{b} = Z \mathbf{B}, \tag{3.9}$$

while, since all terms are premultiplied by S , we may regard \mathbf{b} as the vector specifying the components of the solution \mathbf{x} relative to the basis provided by the eigenvectors of A . Thus we will find it more convenient to work with \mathbf{b} than \mathbf{B} , and to write the system (3.8) as

$$\mathbf{b} = Z \mathbf{b}_0 + Z \int' Z^{-1} H \mathbf{b} \, ds, \tag{3.10}$$

with

$$\mathbf{x} = S \mathbf{b}, \tag{3.11}$$

for an arbitrary constant vector \mathbf{b}_0 . Obviously solutions to this system provide solutions to the original system (1.6).

The vector \mathbf{b}_0 and the terminals of integration in (3.10) may be chosen in different ways, and each such suitable choice generates a solution \mathbf{x} of (1.6). When the $\text{Re } \lambda_i$ are one-signed on J for all $0 < \epsilon < \epsilon_0$, the simple choice $\mathbf{b}_0 = O(1)$ enables us to show, by means of the contraction mapping theorem, the existence of solutions $\mathbf{x}(t, \epsilon)$ of the form

$$\mathbf{x} = S Z \mathbf{b}_0 + O(\epsilon), \tag{3.12}$$

provided H is assumed $O(1)$. This result agrees with the predictions of heuristic theory as given in [5], to the extent that boundary layers may exist at one or both ends of J , and the solution is small in the interior. However, (3.12) is not sufficiently precise for many purposes; and, in fact, if $\text{Re } \lambda_i$ changes sign interior to J for some i , the estimate is even coarser. To deal with these more general cases, and to obtain a more useful estimate of the form of (3.12), we will need to examine the properties of the kernel H in (3.10) more carefully, and to make a more useful choice of the vector \mathbf{b}_0 .

Regarding this last point, if we choose

$$\mathbf{b}_0 = \mathbf{e}_k, \tag{3.13}$$

the usual basis vector in R^n , we might expect from heuristic arguments that the solution we would construct would look like $\mathbf{s}_k(t, \epsilon) \exp(\zeta_k(t, \epsilon))$. With this in

mind, we adopt the following organization of (3.10) which, with the choice (3.12), permits us to establish the existence of a solution \mathbf{b}_k for which the solution \mathbf{x}_k of (1.6), given by (3.11), has this behaviour, in the sense that $\mathbf{b}_k = Z\{\mathbf{e}_k + o(1)\}$, which is a considerable improvement on estimates of the type (3.12).

Choosing, for definiteness,

$$J = [t_0, t_1], \tag{3.14}$$

and making the transformation

$$\mathbf{b} = \boldsymbol{\beta} \exp \zeta_k \tag{3.15}$$

for a given value of k , we may write (3.10) in terms of $\boldsymbol{\beta}$, using the choice (3.13) for \mathbf{b}_0 , as

$$\beta_k = 1 + \int_{t_0}^t \sum_j h_{kj} \beta_j, \tag{3.16}$$

$$\beta_l = \exp(\zeta_l - \zeta_k) \int_{t_0}^t \exp(\zeta_k - \zeta_l) \sum_j h_{lj} \beta_j \quad \text{when } l \neq k \text{ and } s_l \in R_{kl}, \tag{3.17}$$

and

$$\beta_l = \exp(\zeta_l - \zeta_k) \int_t^{t_1} \exp(\zeta_k - \zeta_l) \sum_j h_{lj} \beta_j \quad \text{when } l \neq k \text{ and } s_l \in R_{kg}.$$

It is apparent from the above that a different set of integral equations is to be chosen for each solution of the form (3.15), that is, for each choice of \mathbf{e}_k . The choice of the terminals of integration in the above has been motivated by an observation of the asymptotic result that, for small values of ϵ , and for suitable ϵ -dependence of eigenvectors of A , the integrals in (3.17) and (3.18) are dominated by the values of the integrand at t . Further examination of these equations shows that, under the simple hypothesis that the h_{kj} are bounded, the contraction mapping theorem is not directly applicable, because the right-hand sides of (3.16) to (3.18) do not define a contraction on the space $R_{k\approx}$ for each k , in terms of the usual supremum norm for functions continuous on J . We are thus prompted to make a suitable choice of the eigenvectors \mathbf{s}_k of A that gives S , and hence H , the properties necessary for our construction to proceed. This choice is defined in the following lemma.

LEMMA 1. *Let Assumptions 1, 2 and 3 hold for a given $\epsilon > 0$ and some interval J . Then, for each $t \in J$ and each α , there exists a choice of eigenvectors of A giving the map H the property that*

$$R_\alpha \cap H(R_\alpha) = \{0\}. \tag{3.19}$$

PROOF. We begin by noting that, if s_k and s_j are eigenvectors corresponding to the same eigenvalue function, then the linear combination $c_k(t, \epsilon)s_k(t, \epsilon) + c_j(t, \epsilon)s_j(t, \epsilon)$ is also an eigenvector corresponding to this eigenvalue, for any scalar functions c_k and c_j . We exploit this property by choosing a suitably combined set of eigenvectors $s_k(t, \epsilon)$ that will give H the desired property.

Define the $n \times n$ matrix M by

$$M = \text{diag}[M_\alpha, M_\beta, \dots], \tag{3.20}$$

where M_α, M_β, \dots are square submatrices that map the subspaces R_α, R_β, \dots into themselves for given ϵ and each $t \in J$. Further, define the matrix $S_1(t, \epsilon)$ by

$$S_1(t, \epsilon) = S(t, \epsilon)M(t, \epsilon). \tag{3.21}$$

Differentiating, we obtain

$$S_1^{-1}\dot{S}_1 = M^{-1}[S^{-1}\dot{S}M + \dot{M}]. \tag{3.22}$$

If we define $(S^{-1}\dot{S})_\alpha$ to be the restriction of $S^{-1}\dot{S}$ from R_α to R_α , we will obtain the desired results for each $t \in J$ on choosing

$$\dot{M}_\alpha + (S^{-1}\dot{S})_\alpha M_\alpha = 0. \tag{3.23}$$

This is an initial-value problem for M_α (which may be a singleton), that may be solved over all of J by an appropriate choice of initial value. Obviously, a nonsingular choice ensures the nonsingularity of M , and hence the validity of the above procedure. Thus, we obtain the result as stated.

REMARKS. 1. Note that the above construction implies that we can choose the eigenvectors arbitrarily within the subspace at the initial point, but thereafter, within J , the eigenvectors are uniquely determined. In the subsequent parts of this section, we will assume that the above choice of eigenvectors has been made, and we will use the symbol $S(t, \epsilon)$ to denote the matrix given by $S_1(t, \epsilon)$ in (3.21), while $H(t, \epsilon)$ will denote the kernel defined by (3.6) in terms of this choice of S .

2. Apart from clearly delineating our choice of eigenvectors, the matrix M defined by (3.21) is significant in that, for any given k corresponding to $s_k \in R_\alpha$, the map defined by the integral operators on the right-hand sides of (3.16) to (3.18) has the property of mapping R_α off itself. As we will see, this property is sufficient, along with Assumption 4 below, for our contraction proof to proceed.

3. It may also be observed that this process of the selection of M , when viewed against the background of the heuristic theory for the system (1.6), provides the appropriate generalization of Green's interpretation of the Liouville formula for the amplitude of waves in a slowly varying medium.

Our last assumption describes the properties of the kernel H , and is motivated by the asymptotic results, as $\epsilon \rightarrow 0$, mentioned above. However, as presented here, we regard it as being valid for small, but not necessarily vanishing, values of ϵ .

ASSUMPTION 4. Let $H(t, \epsilon)$ be defined by (3.6), with S generated by the particular choice of eigenvectors given in Lemma 1. We then suppose that, for given $\epsilon > 0$ and interval J , there exist positive constants K_{rs} , C_{rsj} , μ_{rs} and G_{rs} and constants $\alpha_{rsj} \geq 0$, for which the following hold:

(i) for all r and s such that $\lambda_r \simeq \lambda_s$,

$$\max_{u, v \in J} \left| \int_u^v h_{rs}(w, \epsilon) \exp(\pm i\eta_{rs}(w, \epsilon)) dw \right| < K_{rs} \epsilon^{\mu_{rs}}, \tag{3.24}$$

where

$$\eta_{rs} = \text{Im}(\zeta_r - \zeta_s); \tag{3.25}$$

(ii) for all r and s for which $\text{Re}(\lambda_r - \lambda_s) \neq 0$ on J ,

$$|\text{Re}(\lambda_r - \lambda_s)| \sim C_{rsj} |t - t_{rsj}|^{\alpha_{rsj}}, \tag{3.26}$$

for $|t - t_{rsj}|$ suitably small, where t_{rsj} are the zeros of $\text{Re}(\lambda_r - \lambda_s)$ on J ;

(iii) for all r and s ,

$$\int_J |h_{rs}(w, \epsilon) dw| < G_{rs}. \tag{3.27}$$

REMARKS. The last of the above might appear redundant in the light of the continuity assumptions regarding H . However, it is a necessary assumption for the situation in which J is an infinite or semi-infinite interval for the given value of ϵ .

In the proof and discussion for the existence theorem to follow below, it will become apparent that Assumption 4 provides a set of conditions that are sufficient for certain pessimistic estimates to be made, and which allow a construction based on the contraction mapping theorem. Such estimates may be obtained by applying the techniques to be found in Erdélyi [3], for example.

We are now in a position to state our basic existence result for this section.

THEOREM 1. Let $\epsilon > 0$ and $J = [t_0, t_1]$ be such that Assumptions 1 to 4 hold. Then there exist positive constants $C(\epsilon)$, $D_k(\epsilon)$ and $\mu(\epsilon)$ such that, provided ϵ satisfies the inequality

$$C(\epsilon) \epsilon^\mu < 1, \tag{3.28}$$

there exist n linearly independent solutions $x_k(t, \epsilon)$, $k = 1, \dots, n$, of the completely stiff systems (1.6) that have the particular form

$$x_k(t, \epsilon) = X_k(t, \epsilon) \exp(\zeta_k(t, \epsilon)), \tag{3.29}$$

where

$$|\mathbf{X}_k(t, \epsilon) - \mathbf{s}_k(t, \epsilon)| < D_k \epsilon^\mu \tag{3.30}$$

for all $t \in J$, with \mathbf{s}_k the specially chosen eigenvectors of Lemma 1. Further, we have

$$\mathbf{X}_k(t_0, \epsilon) - \mathbf{s}_k(t_0, \epsilon) \in R_{kg}, \tag{3.31}$$

and

$$\mathbf{X}_k(t_1, \epsilon) - \mathbf{s}_k(t_1, \epsilon) \in R_{kt}. \tag{3.32}$$

PROOF. We begin by proving that if there are n solutions of the system (1.6) generated by (3.16) to (3.18), then these are, indeed, linearly independent. That such solutions \mathbf{x}_k generate a basis of solutions to the system (1.6) is then immediately obvious. Thus suppose that the solutions so obtained were not linearly independent. Then there exists a set of constants c_i , not all zero, such that

$$\sum_{i=1}^n c_i \mathbf{x}_i(t, \epsilon) \equiv \mathbf{0} \tag{3.33}$$

for all $t \in J$ and given ϵ . Now let λ_α be the largest eigenvalue in the ordering for which the corresponding \mathbf{x}_i in (3.32) has a non-zero multiplier. Then the restriction of (3.33) to R_α at $t = t_1$ yields

$$\sum_{\mathbf{s}_i \in R_\alpha} c_i \mathbf{s}_i(t_1, \epsilon) \exp(\zeta_i(t_1, \epsilon)) \equiv \mathbf{0}, \tag{3.34}$$

so that $c_i = 0$ follows from the linear independence of the \mathbf{s}_i in the subspace R_α . This contradicts the assumed property of λ_α and establishes the desired result.

We now write the equations (3.16) to (3.18) in the form

$$\boldsymbol{\beta}_k = \mathbf{e}_k + \mathcal{H}_k \boldsymbol{\beta}_k, \tag{3.35}$$

where \mathcal{H}_k is the linear operator defined by the integrals on the right-hand sides of these. This may be arranged in the form

$$\mathcal{H}_k \equiv \mathcal{H}^0 + \mathcal{H}^1 + \mathcal{H}^2, \tag{3.36}$$

where, for convenience, we have suppressed an implied subscript k on the right-hand side, and where the operators \mathcal{H}^0 , \mathcal{H}^1 and \mathcal{H}^2 are defined by their ranges as follows. Let \mathbf{s}_k be an element of R_α for some α . Then \mathcal{H}^0 is that part of the matrix representation of \mathcal{H}_k which has range R_α . Similarly, \mathcal{H}^1 has range $R_{k \simeq} - R_\alpha$, while \mathcal{H}^2 is the complementary subspace of these in R^n . Note that the upper index 2 is to be interpreted as a power only when applied to the operator \mathcal{H}_k .

By applying Assumption 4, and the methods of [3, page 26], we may readily show that

$$\|\mathfrak{K}^0\| \leq C_0(\epsilon), \quad \|\mathfrak{K}^1\| \leq C_1(\epsilon), \quad \text{and} \quad \|\mathfrak{K}^2\| \leq C_2(\epsilon)\epsilon^\tau \tag{3.37}$$

and

$$\|\mathfrak{K}^0\mathfrak{K}^1\| \leq C_3(\epsilon)\epsilon^\nu, \quad \|\mathfrak{K}^1\mathfrak{K}^0\| \leq C_3(\epsilon)\epsilon^\nu \quad \text{and} \quad \|\mathfrak{K}^1\mathfrak{K}^1\| \leq C_3(\epsilon)\epsilon^\nu, \tag{3.38}$$

where ν, τ and the C 's are all positive constants, while $\|\cdot\|$ is the usual operator norm defined with respect to the vector norm

$$\|\mathbf{v}\| = \max_i \left\{ \sup_{t \in J} |v_i(t, \epsilon)| \right\}. \tag{3.39}$$

Of the inequalities (3.37) and (3.38), only the latter are not immediately obvious from (3.16) to (3.18). This can be proved in the following way. All entries in the matrix representation of \mathfrak{K}^1 are of the form

$$\exp(i\eta_{rk}) \int_{t_0}^t h_{rs} \exp(-i\eta_{rk}) b_s, \tag{3.40}$$

so that all entries in $\mathfrak{K}^1\mathfrak{K}^1$, for instance, take the form

$$\exp(i\eta_{uk}) \int_{t_0}^t \exp(-i\eta_{uk}) h_{ur} \exp(i\eta_{rk}) \int_{t_0}^t h_{rs} \exp(i\eta_{rk}) b_s, \tag{3.41}$$

and h_{ur} vanishes for those u and r for which $\eta_{uk} \equiv \eta_{rk}$. This is a consequence of our choice of M in Lemma 1, to make H map R_α off itself.

If the order of integration in (3.41) is changed, then the new inner integral is immediately estimable by applying Assumption 4, and then the stated result follows readily. The remaining estimates follow in a like manner.

We now note that $\mathfrak{K}^0\mathfrak{K}^0 \equiv 0$, again from the selection of M , so that we can write (3.35) as

$$\beta_k = \mathbf{e}_k + \mathfrak{K}_k \mathbf{e}_k + \mathfrak{K}_k^2 \beta_k, \tag{3.42}$$

where, by (3.38), and for ϵ satisfying the condition (3.28) with $C(\epsilon)$ and μ appropriately chosen, \mathfrak{K}_k^2 is a contractive linear operator on the space of vectors $\beta_k(t, \epsilon)$ normed by (3.39). The existence of a unique solution $\beta_k(t, \epsilon)$ of (3.42) may now be established by appealing to the contraction mapping theorem [4, page 27], and a standard result [8, page 38] ensures that this solution coincides with that of (3.35). The estimate (3.30) may be obtained by noting that

$$\|\mathfrak{K}_k \mathbf{e}_k\| \leq A(\epsilon)\epsilon^\gamma \tag{3.43}$$

for suitable positive constants $A(\epsilon)$ and γ .

We may now deduce the asymptotic result for $\epsilon \rightarrow 0$ as a corollary to the above theorem.

COROLLARY 1. *Let Assumptions 1 to 4 hold uniformly with respect to ϵ in a neighbourhood of zero. Then the constants $C(\epsilon)$ and $D_k(\epsilon)$ of Theorem 1 may be chosen independently of ϵ for all ϵ in a suitable neighbourhood of zero.*

PROOF. The hypothesis of the corollary ensures that the ordering and various vector spaces defined in this section remain unchanged in the limit as $\epsilon \rightarrow 0$, while, for suitable small ϵ , the constants G_{rs} , K_{rs} , C_{rsj} , μ_{rs} and α_{rsj} of Assumption 4 may be chosen independently of ϵ . By repeating the construction of Theorem 1, we arrive at a condition like (3.28), where $C(\epsilon)$ is bounded above independently of ϵ for all ϵ in a suitable neighbourhood of zero. Thus this inequality may be satisfied for all appropriately restricted values of ϵ , and the constant D_k is bounded above independently of ϵ .

REMARKS. 1. For either case considered above, Assumption 4 would appear to rule out cases in which A has rapidly varying terms, because of the presence of the term \dot{S} in (3.6). It is possible that the method of proof used above may be extensible to some cases where \dot{S} is large, either locally or globally on J , at the expense of more careful estimates than those applied here.

2. When we consider limiting behaviour as $\epsilon \rightarrow 0$, the inequality (3.24) is guaranteed when h_{rs} has the decomposition

$$h_{rs}(t, \epsilon) = h_{rs}^{(0)}(t) + \epsilon^\nu h_{rs}^{(1)}(t, \epsilon),$$

with $h_{rs}^{(0)}$ Lebesgue integrable on J , while $\nu > 0$ and $h_{rs}^{(1)}$ is uniformly bounded there. For then, with the earlier assumptions holding as $\epsilon \rightarrow 0$, (3.24) is guaranteed by the Riemann-Lebesgue lemma, and standard asymptotic techniques (see, for example, Erdélyi [3]) may be applied to evaluate the constants K_{rs} and μ_{rs} .

3. It is apparent that (3.28) imposes a limitation on the values of ϵ for which the theory works, except under the hypotheses of Corollary 1. Because we are concerned with small but nonvanishing values of ϵ , relations of the type of (3.28) will be significant, and will be required in the subsequent sections. To avoid burdening the text with explicit expressions for the constants $C(\epsilon)$ and μ , we will in future state that ϵ satisfies a relation of the type (3.28), and interpret this to mean that appropriate constants $C(\epsilon)$ and μ may be found for which the relationship holds, while the relevant theory is valid.

4. Effects of change of order

It is of interest to examine the suitable approximations to the solutions constructed in the previous section in a simple case, to gain a more detailed understanding of the results established. It suffices to consider a simple second

order system to illustrate these. We note that it is possible to choose the value of S at a given point in an arbitrary fashion, and to define a variety of local continuations for each independently chosen value of S . Further, matrices A may be generated by associating an independent choice of Λ and applying (2.1). Thus, we may consider H and Λ to be chosen with much freedom, at least in regard to whether they satisfy a local condition. The corresponding assumption for analytic matrices is unclear in its limitations.

Consider the simple 2×2 system which has been reduced, by means of the transformations of Section 3, to

$$\begin{bmatrix} b_1 \\ b_2 \end{bmatrix} = Z \begin{bmatrix} b_{10} \\ b_{20} \end{bmatrix} + Z \int^t Z^{-1} \begin{bmatrix} h_{12}b_2 \\ h_{21}b_1 \end{bmatrix}, \tag{4.1}$$

and let us assume that eigenvalues λ_1 and λ_2 are real and

$$\left. \begin{aligned} \lambda_1 &> \lambda_2 && \text{on } (0, 1), \\ \lambda_1 &< \lambda_2 && \text{on } (1, 2), \end{aligned} \right\} \tag{4.2}$$

and

$$\lambda'_1(1, \epsilon) - \lambda'_2(1, \epsilon) = -\mu < 0, \tag{4.3}$$

for a suitable value of ϵ , with μ independent of ϵ .

The arguments given above indicate that we may totally disregard the matrix S , and consider only the description of the solution basis to be specified in terms of the b 's.

The solution we have constructed on $(0, 1)$ by the methods of Section 3 to have $\exp(\zeta_1)$ behaviour may be denoted by

$$\begin{bmatrix} c_1 \\ c_2 \end{bmatrix} \exp(\zeta_1),$$

where

$$c_1 = 1 + \int_0^t h_{12}c_2 \tag{4.4}$$

and

$$c_2 = e^{-\zeta} \int_0^t e^{\zeta} h_{21}c_1, \tag{4.5}$$

with $\zeta = \zeta_1 - \zeta_2$.

These equations may, by the theory of the previous section, be solved by direct iteration, for initial iterate $c_1^{(0)} = 1$ and $c_2^{(0)} = 0$. The first iterate then becomes

$$c_1^{(1)} = 1 \quad \text{and} \quad c_2^{(1)} = e^{-\zeta} \int_0^t e^{\zeta} h_{21}, \tag{4.6}$$

and simple application of the asymptotic estimates to be found, for example, in Erdélyi [3, page 37], assuming regularity properties of the h 's gives

$$c_2^{(1)} = V_1(t, \epsilon) + O(\epsilon^{3/2}), \tag{4.7}$$

uniformly on $[0, 1]$, where $V_1(t, \epsilon)$ is defined by

$$V_1(t, \epsilon) = \left\{ \begin{array}{l} \epsilon h_{21}(t, \epsilon) / (\lambda_1 - \lambda_2) \quad \text{on } [0, 1 - \epsilon^{1/2}R] \\ \epsilon^{1/2} h_{21}(1, \epsilon) \int_{\tau}^R \exp\left(-\frac{1}{2} \mu(\sigma^2 - \tau^2)\right) / d\sigma \quad \text{on } [1 - \epsilon^{1/2}R, 1] \end{array} \right\} \tag{4.8}$$

and τ is a local variable, defined by

$$t = 1 - \epsilon^{1/2}\tau, \tag{4.9}$$

while R is a given positive constant, so that $0 < \tau < R$.

For further iterates, we obtain, from (4.4),

$$c_1^{(2)} = 1 + O(\epsilon \log \epsilon) \tag{4.10}$$

uniformly on $[0, 1]$, while

$$c_2^{(2)} = c_2^{(1)} \quad \text{and} \quad c_2^{(3)} = c_2^{(2)} + O(\epsilon^{3/2} \log \epsilon), \tag{4.11}$$

again, uniformly on $[0, 1]$.

On noting error estimates for this iterative process as given in [4, page 27], for example, we may write explicitly, for the solution c_1 and c_2 ,

$$c_1 = 1 + O(\epsilon \log \epsilon) \tag{4.12}$$

and

$$c_2 = V_1(t, \epsilon) + O(\epsilon^{3/2} \log \epsilon), \tag{4.13}$$

where now order symbols are uniform on $[0, 1]$.

We may now note, for future reference, the particular result that

$$\begin{bmatrix} c_1(1, \epsilon) \\ c_2(1, \epsilon) \end{bmatrix} = \begin{bmatrix} 1 + O(\epsilon \log \epsilon) \\ V_1(1, \epsilon) + O(\epsilon^{3/2} \log \epsilon) \end{bmatrix}. \tag{4.14}$$

Similarly, one can obtain information about the second solution that takes the form

$$\begin{bmatrix} d_1 \\ d_2 \end{bmatrix} \exp(\zeta_2), \tag{4.15}$$

and it can be shown that

$$\begin{bmatrix} d_1 \\ d_2 \end{bmatrix} = \begin{bmatrix} V_2(t, \epsilon) + O(\epsilon^{3/2} \log \epsilon) \\ 1 + O(\epsilon \log \epsilon) \end{bmatrix}, \tag{4.16}$$

where $V_2(t, \epsilon)$ is defined by

$$V_2(t, \epsilon) = \begin{cases} -\epsilon h_{12}(t, \epsilon) / (\lambda_1 - \lambda_2) & \text{on } [0, 1 - \epsilon^{1/2}R] \\ -\epsilon^{1/2}h_{12}(1, \epsilon) \int_0^\tau \exp\left(\frac{1}{2}\mu(\sigma^2 - \tau^2)\right) d\sigma & \text{on } [1 - \epsilon^{1/2}R, 1]. \end{cases} \quad (4.17)$$

Similar analysis may be applied to the second interval $[1, 2]$ to obtain solutions behaving like $\bar{c} \exp(\zeta_1)$ and $\bar{d} \exp(\zeta_2)$ and which have the end-values

$$\bar{c}(1, \epsilon) = \left[\begin{array}{c} 1 \\ \bar{V}_1(1, \epsilon) + O(\epsilon^{3/2} \log \epsilon) \end{array} \right] \quad (4.18)$$

and

$$\bar{d}(1, \epsilon) = \left[\begin{array}{c} 0 \\ 1 \end{array} \right], \quad (4.19)$$

respectively, where \bar{V}_1 is defined in an analogous fashion to V_1 and V_2 and is such that $\bar{V}_1(1, \epsilon) \sim (\text{constant})\epsilon^{1/2}$ as $\epsilon \rightarrow 0$.

The end-values acquired by the first solution at $t = 1$ thus initiate a solution in the second interval that is of the form

$$\bar{c}(t, \epsilon)\exp(\zeta_1)[1 + o(1)] + [V_1(1, \epsilon) + O(\epsilon^{3/2} \log \epsilon)]\bar{d}(t, \epsilon)\exp(\zeta_2)\exp(\zeta_1(1, \epsilon) - \zeta_2(1, \epsilon)), \quad (4.20)$$

and thus it is dominated on the second interval by the solution $\bar{d}(t, \epsilon)\exp(\zeta_2)$, even when $h_{21}(1, \epsilon)$ vanishes.

On the first interval, $[0, 1]$, any solution dominated by an $\exp(\zeta_1)$ type behaviour must be of the form

$$c(t, \epsilon)\exp(\zeta_1) + A d(t, \epsilon)\exp(\zeta_2), \quad (4.21)$$

where A is a suitable constant value. If we attempt to adjust A to eliminate the unwanted term on the second interval, we require that $A \exp(\zeta_2(1, \epsilon))$ shall be able to cancel a non-zero $\epsilon^{1/2} \exp(\zeta_1(1, \epsilon))$ term, which implies that

$$A \sim \epsilon^{1/2} \exp(\zeta_1(1, \epsilon) - \zeta_2(1, \epsilon)) \quad \text{as } \epsilon \rightarrow 0, \quad (4.22)$$

and the second term of (4.21) dominates on almost all of the first interval. Thus we must conclude that there is no solution on the first interval which can match with a solution having the required properties on the second interval.

Such a result is in no way contradictory to known results for analytic systems. However, it is already indicated by the above example that similar difficulties may arise whenever a change of order occurs, even if eigenvalues do not become equal, but differ only in their imaginary parts.

An analysis of such a case involves, in the above solution, and to the accuracy of the above calculations, the estimation of integrals

$$c_2 = e^{-\zeta} \int_0^t e^{\zeta} h_{21} \quad (4.23)$$

and

$$c_{24} = -e^{-\xi} \int_{\xi}^2 e^{\xi} h_{21}. \tag{4.24}$$

For analytic functions h_{21} , these differences are exponentially small, and the problem cannot be tackled directly by considerations of this kind. Thus there is no demonstrated conflict with previous analytic results.

On the other hand, suppose we had been dealing with a function defined by a numerical calculation, which was based on an integration formula using piecewise analytic functions, but which involved finite step discontinuities in the values of the derivatives of S . This corresponds to h having step discontinuities. Then it is apparent that c_2 and \bar{c}_2 will differ sufficiently to lead to the difficulty discussed above. Note that having better regularity properties for S to some finite order of derivative serves to reduce the algebraic order of the discrepancy, but this becomes useless in comparison with exponentially larger terms.

We are thus led to expect that significantly different results pertain for non-analytic systems in contrast to the results previously obtained for analytic systems. We shall discuss the significance of this difference for the numerical solution to stiff systems later in this sequence of papers. For the moment we turn our attention to elucidating the nature of optimal bases for linear systems when changes in the ordering of the eigenvectors occur.

5. Continuation of solutions

The results of Section 4 serve to motivate the procedure we adopt for attempting to construct a single useful basis for the solution space of (1.6) on two abutting intervals J and \bar{J} , for each of which the assumptions and results of Section 3 apply.

For definiteness, we will represent J and \bar{J} by $[t_0, t_1]$ and $[t_1, t_2]$, respectively. To cope with cases where there is a multiple change of ordering of eigenvalues at t_1 , we find it convenient to define some terms as follows.

DEFINITION 4. *Let $J = [t_0, t_1]$ and $\bar{J} = [t_1, t_2]$ be two abutting subintervals of I , for which we may, for each k and a given $\epsilon > 0$, define the subspaces $R_{k\approx}, R_{kl}, R_{kg}$ and $\bar{R}_{k\approx}, \bar{R}_{kl}, \bar{R}_{kg}$ on J and \bar{J} respectively, where these quantities are as defined in Definition 3 for J and \bar{J} , respectively.*

Then, for each k , we define the crossing set at t_1 , \mathcal{C}_k , to be the set of integers k such that

$$s_k \in R_{kl} - (\bar{R}_{kl} \cup \bar{R}_{k\approx}). \tag{5.1}$$

Then elements of \mathcal{C}_k are in one-to-one correspondence with those eigenvectors which, as t passes through t_1 , go from below s_k in the ordering to above s_k .

We may now state our first main result for this section.

THEOREM 2. *Let $\epsilon > 0$ be given, and let $J = [t_0, t_1]$ and $\bar{J} = [t_1, t_2]$ be intervals for which $J \cup \bar{J} \subseteq I$ for the given ϵ , and for each of which Assumptions 1 to 4 and the hypotheses of Theorem 1 hold. Then, for such ϵ , there exists a basis $\{x_k(t, \epsilon)\}$ for the solution space of (1.6) over $J \cup \bar{J} = [t_0, t_2]$ that has the form*

$$x_k(t, \epsilon) = \begin{cases} X_k^*(t, \epsilon)\exp(\zeta_k(t, \epsilon)) & \text{for } t \in J, \\ \bar{X}_k^*(t, \epsilon)\exp(\zeta_k(t, \epsilon)) \\ + \sum_{r \in \mathcal{C}_k} d_{kr} \bar{X}_r^*(t, \epsilon)\exp\{\zeta_r(t, \epsilon) + \zeta_k(t_1, \epsilon) - \zeta_r(t_1, \epsilon)\} & \text{for } t \in \bar{J}. \end{cases} \tag{5.2}$$

Here $\{X_r^*(t, \epsilon)\exp(\zeta_r)\}$ and $\{\bar{X}_r^*(t, \epsilon)\exp(\zeta_r)\}$ are sets of n linearly independent solutions of (1.6) defined on J and \bar{J} , respectively, and satisfying the conditions

$$|X_r^*(t, \epsilon) - s_r(t, \epsilon)| < M_r(\epsilon)\epsilon^{\mu_r} \quad \text{on } J, \tag{5.3}$$

and

$$|\bar{X}_r^*(t, \epsilon) - s_r(t, \epsilon)| < \bar{M}_r(\epsilon)\epsilon^{\bar{\mu}_r} \quad \text{on } \bar{J}, \tag{5.4}$$

for positive functions $M_r(\epsilon)$ and $\bar{M}_r(\epsilon)$ and constants μ_r and $\bar{\mu}_r$. Further, we have

$$X_r^*(t_0, \epsilon) - s_r(t_0, \epsilon) \in R_{r\sigma} \tag{5.5}$$

and

$$\bar{X}_r^*(t_2, \epsilon) - s_r(t_2, \epsilon) \in \bar{R}_{r\bar{\sigma}} \cup \bar{R}_{r\approx}. \tag{5.6}$$

The $d_{kr}(\epsilon)$ are constants satisfying a condition

$$|d_{kr}(\epsilon)| < D_{kr}(\epsilon)\epsilon^{\nu_{kr}}, \tag{5.7}$$

for positive constants $D_{kr}(\epsilon)$ and ν_{kr} .

PROOF. By Theorem 1, there exist linearly independent sets of solutions of (1.6), $\{X_k(t, \epsilon)\exp(\zeta_k)\}$ and $\{\bar{X}_k(t, \epsilon)\exp(\zeta_k)\}$ on J and \bar{J} , satisfying conditions of the form (3.30) as well as end conditions of the form (3.31) and (3.32) on their respective intervals.

The linear independence of these sets, together with their properties (3.3), allow us to deduce the existence of constants C_{kr} such that

$$\bar{X}_k(t_1, \epsilon) - X_k(t_1, \epsilon) = \sum_{r \in \sigma} C_{kr} X_r(t_1, \epsilon) + \sum_{r \in \bar{\sigma}} \bar{X}_r(t_1, \epsilon), \tag{5.8}$$

where

$$|C_{kr}(\epsilon)| \leq Q_{kr}(\epsilon)\epsilon^{\delta_{kr}} \tag{5.9}$$

for appropriate positive constants Q_{kr} and δ_{kr} , and where σ and $\bar{\sigma}$ is a partitioning of the integers $\{1, \dots, n\}$. We ensure that ϵ , and hence the left side of (5.8), is sufficiently small by imposing a condition like (3.28), which may be absorbed into the hypotheses of Theorem 2.

From (4.8), we obtain

$$\mathbf{X}_k(t_1, \epsilon) - \sum_{r \in \bar{\sigma}} C_{kr} \bar{\mathbf{X}}_r(t_1, \epsilon) = \mathbf{X}_k(t_1, \epsilon) + \sum_{r \in \sigma} C_{kr} \mathbf{X}_r(t_1, \epsilon), \tag{5.10}$$

which is a condition that solutions J and \bar{J} , having these end-values, should match at t_1 . The solution on the appropriate interval may be obtained by continuation of these end-values into that interval. Thus we will obtain a solution valid over $J \cup \bar{J}$ that is continuous at t_1 .

Continuing the left-hand side of (5.10) into \bar{J} , we obtain

$$\bar{\mathbf{X}}_k(t, \epsilon)\exp(\zeta_k(t, \epsilon)) - \sum_{r \in \bar{\sigma}} C_{kr} \bar{\mathbf{X}}_r(t, \epsilon)\exp(\zeta_k(t, \epsilon))\exp\left[\epsilon^{-1} \int_{t_1}^t (\lambda_r - \lambda_k)\right], \tag{5.11}$$

while the continuation of the right-hand side into J is

$$\mathbf{X}_k(t, \epsilon)\exp(\zeta_k(t, \epsilon)) + \sum_{r \in \sigma} C_{kr} \mathbf{X}_r(t, \epsilon)\exp(\zeta_k(t, \epsilon))\exp\left[\epsilon^{-1} \int_{t_1}^t (\lambda_k - \lambda_r)\right]. \tag{5.12}$$

So far, we have left the choice of σ and $\bar{\sigma}$ arbitrary. We now choose σ to be the set of integers r such that

$$\mathbf{s}_r \in R_{kg}, \tag{5.13}$$

with $\bar{\sigma}$ its complement, namely the set r such that

$$\mathbf{s}_r \in R_{kl}. \tag{5.14}$$

Two distinct cases now arise, namely, whether or not \mathcal{C}_k is empty. When, for selected k , \mathcal{C}_k is empty, we have the condition

$$R_{kl} \subseteq \bar{R}_{kl} \cup \bar{R}_{k\approx}, \tag{5.15}$$

so that the exponential factors in the expressions (5.11) and (5.12), namely

$$\exp\left(\epsilon^{-1} \int_{t_1}^t (\lambda_r - \lambda_k)\right) \quad \text{and} \quad \exp\left(\epsilon^{-1} \int_{t_1}^t (\lambda_k - \lambda_r)\right)$$

are uniformly bounded by unity on their respective intervals of application. Because we have ensured that the C_{kr} are suitably small, the additional terms may be absorbed by defining

$$\mathbf{X}_r^* = \mathbf{X}_k + \sum_{r \in \sigma} C_{kr} \mathbf{X}_r \exp\left(\epsilon^{-1} \int_{t_1}^t (\lambda_k - \lambda_r)\right) \tag{5.16}$$

and

$$\bar{\mathbf{X}}_k^* = \bar{\mathbf{X}}_k + \sum_{r \in \bar{\mathcal{C}}_k} C_{kr} \mathbf{X}_r \exp\left(\varepsilon^{-1} \int_{t_1}^t (\lambda_r - \lambda_k)\right). \tag{5.17}$$

The results (5.3) and (5.4) are immediate consequences of these definitions, since, when \mathcal{C}_k is empty, no constants d_{kr} are involved. Note that, in this case, there exists a solution of the type described in Theorem 1, that is, of the form $(\mathbf{s}_k + \mathbf{w}_k) \exp(\zeta_k)$ in both J and \bar{J} .

When \mathcal{C}_k is not empty, there are elements of R_{kl} that do not lie in $\bar{R}_{kl} \cup \bar{R}_{k\approx}$, so that there is no continuation that has the specific form of (5.12) which leads to uniformly bounded exponential functions. In this case, we rearrange (5.12) in the form

$$\bar{\mathbf{X}}_k \exp(\zeta_k) - \left\{ \sum_{\bar{\sigma} - \mathcal{C}_k} + \sum_{\mathcal{C}_k} \right\} C_{kr} \bar{\mathbf{X}}_r \exp(\zeta_k) \exp\left[\varepsilon^{-1} \int_{t_1}^t (\lambda_r - \lambda_k)\right]. \tag{5.18}$$

If we were to define

$$\bar{\mathbf{X}}_k^* = \bar{\mathbf{X}}_k - \sum_{\bar{\sigma} - \mathcal{C}_k} C_{kr} \bar{\mathbf{X}}_r \exp\left[\varepsilon^{-1} \int_{t_1}^t (\lambda_r - \lambda_k)\right],$$

and choose $d_{kr} = C_{kr}$, we would have results substantially in the form as stated in the theorem. However, for any r that was an element of \mathcal{C}_k for some k , we would have both $\bar{\mathbf{X}}_r$ and $\bar{\mathbf{X}}_r^*$ involved in describing the solution set on \bar{J} . Thus more than n solutions would be involved, and they would not be linearly independent.

We therefore proceed as follows. Arrange all those λ_k for which the \mathcal{C}_k are not empty, into sets δ_s such that any two λ_i and λ_k are in the set if and only if $\lambda_i(t_1, \varepsilon) = \lambda_k(t_1, \varepsilon)$ for the given value of ε . It follows that, for any given k , the set of eigenvalues corresponding to \mathcal{C}_k is contained in the particular set δ_s containing λ_k . Thus we can arrange those solutions with non-empty \mathcal{C}_k into distinct sets Δ_s , where elements are in one-to-one correspondence with those of the sets δ_s . Within each of these sets, there may be a number of eigenvalue functions involved and, of these, there will be one which is largest on \bar{J} . For any j such that λ_j is equal under the ordering to this eigenvalue function, the previous proof applies, and there is an $\bar{\mathbf{X}}_j^*$ defined as in equation (5.17). Consider now the next largest eigenvalue function. We now define the set $(\bar{\mathbf{X}}_k(t, \varepsilon))$ by replacing all $\bar{\mathbf{X}}_j$ by $\bar{\mathbf{X}}_j^*$, whenever λ_j is equivalent to the largest eigenvalue function in the corresponding set δ_s . All others are as before. It follows from the constructive procedure described above that we may find constants C_{kr} , for all λ_k equivalent to the second largest eigenfunction, such that

$$\bar{\mathbf{X}}_k(t_1, \varepsilon) - \sum_{\bar{\sigma} - \mathcal{C}_k} C_{kr} \bar{\mathbf{X}}_r(t_1, \varepsilon) - \sum_{\mathcal{C}_k} C_{kr} \bar{\mathbf{X}}_r^*(t_1, \varepsilon) = \mathbf{X}_k(t_1, \varepsilon) + \sum_{\sigma} C_{kr} \mathbf{X}_r(t_1, \varepsilon).$$

This equation can be continued into the respective subintervals and the results of the theorem will be established as before for all such solutions. It is obvious that any set δ_s can be exhausted by applying the above algorithm, after successively replacing \bar{X}_r by \bar{X}_r^* . The linear independence of the n solutions follow from Abel's identity, and the fact that the matrices formed with $\bar{X}_r(t_1, \epsilon)$ and $\bar{X}_r^*(t_1, \epsilon)$ are obtainable from each other by multiplication by a suitable matrix of the form $(I + \text{small terms})$.

COROLLARY 2. *Under the conditions of Theorem 2, the basis $\{x_k(t, \epsilon)\}$ may also take the form*

$$x_k(t, \epsilon) = \begin{cases} X'_k(t, \epsilon) \exp \zeta_k(t, \epsilon) \\ + \sum_{r \in \mathcal{C}_k} f_{kr} X'_r(t, \epsilon) \exp \{ \zeta_r(t, \epsilon) + \zeta_k(t_1, \epsilon) - \zeta_r(t_1, \epsilon) \} & \text{for } t \in J, \\ \bar{X}'_k(t, \epsilon) \exp \zeta_r(t, \epsilon) & \text{for } t \in \bar{J}, \end{cases} \quad (5.19)$$

where the vectors $X'_r(t, \epsilon)$ and $\bar{X}'_r(t, \epsilon)$ and the scalars $f_{kr}(\epsilon)$ satisfy conditions analogous to those satisfied by $X_r^*(t, \epsilon)$, $\bar{X}_r^*(t, \epsilon)$ and $d_{kr}(\epsilon)$, respectively.

PROOF. The proof proceeds exactly as that for Theorem 2, but this time we choose $\bar{\sigma}$ to be the set of integers r such that

$$s_r \in \bar{R}_{kl} \cup \bar{R}_{k\approx}, \quad (5.20)$$

with σ its complement.

Corresponding to Corollary 1 we have a similar result for Theorem 2, for the case where $\epsilon \rightarrow 0$.

COROLLARY 3. *Let the assumption of Theorem 2 hold uniformly with respect to ϵ in a neighbourhood of zero. Then the constants $M_r(\epsilon)$, $\bar{M}_r(\epsilon)$ and $D_{kr}(\epsilon)$ of that theorem may be chosen to be independent for ϵ for all ϵ in a suitable neighbourhood of zero.*

PROOF. The proof is an elementary application of the results of Corollary 1 to the construction procedure described above.

REMARKS. Here, we see the generalization of the phenomenon noted on Section 4, namely the occurrence of a solution that, on J , behaves like

$$\{s_k + w_k\} \exp(\zeta_k),$$

which continues into \bar{J} in such a way it is dominated, apart from a small neighbourhood of t_1 , by

$$d_{kr}\{s_r + w_r\} \exp(\zeta_k) \exp\left[\varepsilon^{-1} \int_{t_1}^t (\lambda_r - \lambda_k)\right],$$

when there is a change in ordering at t_1 .

Precisely which λ_r are involved in the exponential factors are determined by the choice of the crossing set \mathcal{C}_k corresponding to the solution considered. The choice of σ to correspond to $R_{kg} - R_{k\approx}$ leads us to define \mathcal{C}_k to correspond to $R_{kl} \cup R_{k\approx} - \bar{R}_{kl} \cup \bar{R}_{k\approx}$, which is a bigger set than that defined by (5.1). This leads us to suggest the choice (5.1) to be an improvement, though it would be premature to regard it as optimal.

Such consideration aside, we can see that, if this is a genuine effect, and not merely a consequence of the construction technique applied, it has the most important implications. On the first subinterval, a fundamental matrix for (1.6) exists, having the form $SZ(I + w)$ (where we may assert, by arguments similar to those used in Section 4, that w is only algebraically small (that is, is dominated by ε^ν for some $\nu > 0$) as $\varepsilon \rightarrow 0$), in all cases except those involving the most extreme restrictions on the h_{ij} . The continuation of this into the second interval, however, has two columns which differ, apart from a constant multiplier, by terms which are exponentially small in comparison with the dominant terms in the solution.

To investigate this further, we show below, under conditions that are quite reasonable in the light of Section 4, every fundamental matrix must exhibit such a change when there is a genuine crossing of eigenvalue functions. However, in general it is not true that the position is as bad as would be the case if the results of Theorem 2 were optimal in the sense described above. We will show that the best form of results depends on the detailed structure of the changes of order of the eigenvalues and that, therefore, any general discussion would be of excessive complexity. We are content, therefore, to show, under simple assumptions, the nature of the results to be expected. The discussion provides the basis for an algorithm to treat any specific case.

We first show that it is sometimes possible to have less change in exponential dominance than arises from the construction of Theorem 2.

THEOREM 3. *Let λ_α and λ_β be eigenvalue functions of multiplicities m_α and m_β , respectively, on $J \cup \bar{J}$ which satisfy $\lambda_\beta < \lambda_\alpha$ on J and $\lambda_\alpha < \lambda_\beta$ on \bar{J} , and let no other eigenvalue function be equivalent in the ordering on either subinterval.*

Let $D_{\alpha\beta}$ be the $m_\alpha \times m_\beta$ matrix, of rank r , whose entries are the constants $d_{ij}(\varepsilon)$ of Theorem 2, where $\lambda_i = \lambda_\alpha$ and $\lambda_j = \lambda_\beta$. Then there exist at least $(m_\alpha - r)$

solutions $\mathbf{x}_k(t, \epsilon)$ of (1.6) on $J \cup \bar{J}$ which are of order unity at t_0 , and for which

$$\mathbf{x}_k(t, \epsilon) \exp \left[-\epsilon^{-1} \int_{t_0}^t \lambda_\alpha \right]$$

is bounded in modulus uniformly on $J \cup \bar{J}$ by a term of the form $B(\epsilon)$ not exponentially large in ϵ . When the hypotheses of Corollary 3 hold, $B(\epsilon)$ may be chosen to be independent of ϵ in an appropriate neighbourhood of zero.

PROOF. Let \mathbf{g} be any m_α -vector in the (left) null space of $D_{\alpha\beta}$. Then the solution

$$\sum_{\lambda_k = \lambda_\alpha} g_k \mathbf{X}_k^* \exp(\zeta_k(t, \epsilon)) \quad \text{for } t \in J$$

has, as its continuation into \bar{J} ,

$$\sum_{\lambda_k = \lambda_\alpha} g_k \bar{\mathbf{X}}_k^* \exp(\zeta_k) + \sum_{\lambda_k = \lambda_\alpha} g_k \sum_{r \in C_k} d_{kr} \bar{\mathbf{X}}_r^* \exp\{\zeta_r(t, \epsilon) + \zeta_k(t_1, \epsilon) - \zeta_r(t_1, \epsilon)\},$$

by Theorem 2, and the assumption on \mathbf{g} implies that this is

$$\sum_{\lambda_k = \lambda_\alpha} g_k \bar{\mathbf{X}}_k^* \exp(\zeta_k(t, \epsilon)).$$

Since the null space has dimension $(m_\alpha - r)$, there are at least $(m_\alpha - r)$ linearly independent solutions of this form, obtainable in the above manner. A suitable scaling of the \mathbf{X}_k^* and $\bar{\mathbf{X}}_k^*$ gives us the result as stated (or, alternatively, we may define the ζ_k from t_0).

COROLLARY 4. *The solutions of Theorem 3 cannot behave like*

$$\{\mathbf{s}_j + \mathbf{w}_j\} \exp \left(\epsilon^{-1} \left\{ \int_{t_0}^t \lambda_\beta + \int_{t_0}^{t_1} \lambda_\alpha - \int_{t_0}^{t_1} \lambda_\beta \right\} \right)$$

on the second interval \bar{J} .

PROOF. Such terms as are given above are specifically excluded by our definition of \mathbf{g} .

REMARKS. The results of Theorem 3 demonstrate the truth of our earlier statement regarding the optimality of the results of Theorem 2 in that such results are not optimal in any case, as above, where $m_\beta < m_\alpha$.

The estimates of Section 4 indicate that the entries in $D_{\alpha\beta}$ are of the form $C(\epsilon)\epsilon^r$, and hence behave algebraically as $\epsilon \rightarrow 0$ for ϵ satisfying inequalities like (3.28). Thus it will normally be possible to find independent sets of constants \mathbf{g}_α which are of unit order or at worst algebraically small.

It may also be noted that the above selection of linear combinations of solutions with the same eigenvalue function is equivalent to a particular choice of eigenvectors $s_i(0, \epsilon)$ within that space. The discussion in the proof of Lemma 1 shows that this choice of basis was quite arbitrary, and hence is available within the present context. If, in order to remove the terms which are exponentially larger on \bar{J} , it is necessary to select a basis in R_α which has dimension less than m_α for some ϵ in a neighbourhood of zero for which all the other hypotheses hold, one is faced with a choice between two alternatives. Either one accepts on algebraic degeneracy on J or an exponential one on \bar{J} . Which difficulty is least unacceptable will be a matter for decision in a specific context.

Our last result for this section shows that this rank condition is, at least under certain assumptions, optimal, in a case for which m_α, m_β and r are all 1, and which is a direct generalization of the results of Section 4.

THEOREM 4. *Let λ_i and λ_k be simple eigenvalue functions equivalent to no other eigenvalues under the ordering on $J \cup \bar{J}$, and such that*

$$\text{and } \left. \begin{array}{l} \lambda_i < \lambda_k \text{ for } t \in J, \\ \lambda_k < \lambda_i \text{ for } t \in \bar{J}. \end{array} \right\} \tag{5.21}$$

Let $d_{ki}\epsilon^\mu$ be the non-zero constant as determined in Theorem 2, where μ is some positive constant such that $|d_{ki}|$ is bounded above and below in the sense that $d < |d_{ki}| < D$, with ϵ satisfying a relationship of the form (3.28). Then there exists no solution $x_k(t, \epsilon)$ of (1.6) on $J \cup \bar{J}$ such that

$$x_k(t, \epsilon)\exp\{-\zeta_k(t, \epsilon)\}$$

is bounded in modulus uniformly with respect to t by a term of the form $T(\epsilon)$ that is not exponentially large in ϵ .

PROOF. It is sufficient to show that there do not exist solutions $x^{(1)}$ and $x^{(2)}$, defined, with the above property, on J and \bar{J} , respectively, and such that

$$x^{(1)}(t_1, \epsilon) = x^{(2)}(t_1, \epsilon) - d_{ki}\epsilon^\mu \exp(\zeta_k(t_1, \epsilon)). \tag{5.22}$$

Any solution on J may be obtained as a linear combination of the basis solutions of Theorem 2, and we may arrange such a combination in the form

$$x^{(1)}(t, \epsilon) = \sum_{s_j \in R_{ki}} f_j X_j^* \exp(\zeta_j) + \sum_{s_j \in R_{kg}} \bar{f}_j \bar{X}_j^* \exp(\zeta_j), \tag{5.23}$$

where f_j and \bar{f}_j are suitably chosen constants. Similarly, we may introduce the representation

$$x^{(2)}(t, \epsilon) = \sum_{s_j \in \bar{R}_{ki}} g_j \bar{X}_j^* \exp(\zeta_j) + \sum_{s_j \in \bar{R}_{kg}} \bar{g}_j \bar{X}_j^* \exp(\zeta_j). \tag{5.24}$$

We note here that there is no need to include the f_k term (5.24) since it merely serves to scale $X_k^* \exp(\zeta_k)$, which is irrelevant in a linear problem. We will make this assumption implicitly in what follows and thus set $f_k = 0$.

The property of the X_j^* implied by equation (5.3) enables us to show the exponential property required by the theorem is possible only if

$$f_j = \exp(\zeta_k(t_0, \epsilon) - \zeta_j(t_0, \epsilon))f'_j \quad \text{on } R_{kl}$$

and

$$\bar{f}_j = \exp(\zeta_k(t_1, \epsilon) - \zeta_j(t_1, \epsilon))f'_j \quad \text{on } R_{kg},$$

where the constants f'_j are not exponentially large in ϵ . Similarly,

$$g_j = \exp(\zeta_k(t_1, \epsilon) - \zeta_j(t_1, \epsilon))g'_j \quad \text{on } \bar{R}_{kl}$$

and

$$\bar{g}_j = \exp(\zeta_k(t_2, \epsilon) - \zeta_j(t_2, \epsilon))g'_j \quad \text{on } \bar{R}_{kg}.$$

If these estimates are now used in equation (5.22), these results, at $t = t_1$, after cancelling a factor $\exp(\zeta_k(t_1, \epsilon))$,

$$f'_k X_k^* + \sum_{R_{kg}} f'_j X_j^* + \Phi(\mathbf{f}') = \sum_{\bar{R}_{kl}} g'_j \bar{X}_j^* + g'_k \bar{X}_k^* + \Psi(\mathbf{g}') - d_{ki} \epsilon^\mu \bar{X}_i^*, \quad (5.25)$$

where Φ and Ψ are linear operators from R^n to R^n which have the exponentially small norms

$$\|\Phi\| = \max_{R_{kl}} \left\{ \exp(-[\zeta_k(t, \epsilon) - \zeta_j(t, \epsilon)]) \right\}_{t_0}^{t_1} \quad (5.26)$$

and

$$\|\Psi\| = \max_{\bar{R}_{kg}} \left\{ \exp(-[\zeta_j(t, \epsilon) - \zeta_k(t, \epsilon)]) \right\}_{t_1}^{t_2}. \quad (5.27)$$

Thus we have n equation in terms of $2n$ unknowns, namely the constants f'_j and g'_j , so that the problem given above is not well set in the usual sense. However, if we define a vector \mathbf{v} by

$$v_j = \begin{cases} -g'_k & \text{for } j = k, \\ d_{ki} \epsilon^\mu & \text{for } j = i, \\ f'_j & \text{for } s_j \in R_{kg}, \\ g'_j & \text{for } s_j \in \bar{R}_{kl} - \{s_k\}, \end{cases} \quad (5.28)$$

the equation (6.30) may be arranged as

$$S(I + T)\mathbf{v} = \Psi(\mathbf{g}') - \Phi(\mathbf{f}'), \quad (5.29)$$

where T is a linear operator from R^n into itself, such that

$$\|T\| \leq W(\epsilon)\epsilon^\nu \quad (5.30)$$

for some positive $W(\epsilon)$ and ν ; when ϵ satisfies an inequality of the type (3.28), the above equation (5.29) may be solved iteratively, with the result that any solution \mathbf{v} must be exponentially small in ϵ . This leads to a contradiction to the above definition of v_j .

It should be apparent that the proof extends to the case considered in Theorem 3, since we can utilize that fact that all components in R_α lead to disposable pairs $(f_j - g_j)$ to cancel the $(m_\alpha - r)$ combinations of the null space, but no other improvement is possible. Some reflection on how the above proof works provides the key to its extension. The only f'_j which arise are those for which s_j are in R_{k_g} and the only g'_j are those for which $s_j \in \bar{R}_{k_l} \cup \bar{R}_{k_{\approx}}$. Thus there will be freedom of choice within the range of those elements which occur twice. As an illustration, consider the case where, on the first subinterval, $\lambda_k < \lambda_i < \lambda_j$, whereas, on the second subinterval, $\lambda_j < \lambda_i < \lambda_k$. Then the construction forming the basis of Theorem 2 leads to a solution $\mathbf{X}_k^* \exp(\zeta_k)$ on $[t_0, t_2]$, and $\mathbf{X}_i^* \exp(\zeta_i)$ on the first interval, continuing as

$$\bar{X}_i^* \exp(\zeta_i) + d_{ik} \exp\{\zeta_k(t, \epsilon) + \zeta_i(t_1, \epsilon) - \zeta_k(t_1, \epsilon)\};$$

while $\mathbf{X}_j^* \exp(\zeta_j)$ on the first subinterval is continued as

$$\begin{aligned} \bar{X}_j^* \exp(\zeta_j) + d_{jk} \exp\{\zeta_k(t, \epsilon) + \zeta_j(t_1, \epsilon) - \zeta_k(t_1, \epsilon)\} \\ + d_{ji} \exp\{\zeta_i(t, \epsilon) + \zeta_j(t_1, \epsilon) - \zeta_i(t_1, \epsilon)\}. \end{aligned}$$

It is apparent that, if $d_{jk} \neq 0$, we can find a solution

$$\mathbf{x}, - (d_{ik}/d_{jk})\mathbf{x}_j \exp(-\{\zeta_j(t_1, \epsilon) - \zeta_i(t_1, \epsilon)\})$$

which is uniformly $\exp(\zeta_i)$ on the two subintervals. However, the order change still introduces a change of exponential order in the solution \mathbf{x}_j .

While the generalization of the above approach leads to entirely algebraic questions about systems for which answers may be calculated to sufficient accuracy to permit precise estimates, there are some points which need further consideration. In our examples, we have not considered equivalent eigenvalue functions. If one has $\zeta_1 = \epsilon^{-1} \int_0^1 (\lambda + i\mu_2)$, then it becomes a matter of definition as to whether a solution which is dominated by $s_1 \exp(\zeta_1)$ in J and $s_2 \exp(\zeta_2)$ in \bar{J} is a satisfactory continuation. If it is, then one will have more freedom to reduce the incidence of exponential order change. For the present we are content to point out that any such decision ought to be based on context in application, and not on arbitrary decision within the mathematics.

To this stage, we have considered only the continuation across one point where the ordering varies. However, we draw attention to the fact that the continuation process is essentially algebraic, and that it is possible to calculate

the continuation matrices, or at least dominant approximations to them. Further, if we start at $t = 0$, we can construct a basis $\mathbf{X}_k^* \exp(\zeta_k)$ on the first interval, a basis $\bar{\mathbf{X}}_k^* \exp(\zeta_k)$ on the second, and a connecting matrix relating them. Using $\bar{\mathbf{X}}_k^* \exp(\zeta_k)$ on the second interval and the original basis $\mathbf{X}_k \exp(\zeta_k)$ on the third, we obtain new solutions $\bar{\mathbf{x}}_k^{**} \exp(\zeta_k)$ on the second interval, obtainable from the originals by matrix multiplication. Thus the process is continuable conclusively if I is the union of a finite number of intervals J_i for which the theory of Section 3 holds. The process is routine, if tedious, the most important objection being that it may lead to excessive round-off error if computation were attempted.

It is also of interest to note that, where multiple eigenvalue functions occur, there exists the possibility of selecting initial eigenvectors to achieve further simplification. Thus, in the case of Theorem 3, one might treat cases where $m_\beta > 1$ by selecting those initial eigenvectors so that the continuation of a solution of the form $\mathbf{X}_k^* \exp(\zeta_k)$ will take the form $\bar{\mathbf{X}}_k^* \exp(\zeta_k) + d_{kr} \bar{\mathbf{X}}_r^* \exp(\zeta_r)$, where only one $\bar{\mathbf{X}}_r^*$ is involved. Alternatively, one may accept some behaviour at the first crossing to gain desirable freedom at a subsequent crossing to suppress more troublesome behaviour. In such circumstances there seems little point in attempting a general description of the way the fundamental matrix varies in behaviour.

6. Discussion

The presentation of the proof in Section 3, designed to alleviate the burden of calculation, may tend to obscure the underlying motivation for the proof itself. Further, we believe that this motivation may be of considerable use in circumstances where the assumptions of this paper do not hold. Work by Doust (private communication) on the case where the eigenvectors span a space of dimension less than n shows that a re-interpretation of this motivation is useful. It would appear that the basic idea is that on a subspace associated with an eigenvalue function λ_i , the crudest level of variation of the solution is like $\exp(\zeta_i)$. In the present paper, our method of generating the integral equations (3.16) to (3.18) led essentially to order one terms which were integrated (via the M matrix here) so that, with these terms removed, an iterative procedure could be applied to solve these equations. Where our assumptions fail, the matrix H may be unbounded, but Doust has shown that a further refinement, using dominant behaviour, may eventually lead to integral equations for which iteration is possible.

We note that if A has rapidly varying entries in the neighbourhood of an endpoint t_0 , our assumptions will fail there. Suppose the scale of these variations was ϵ^r . Then the matrix H will take the form $\epsilon^{-r} H_0$, and the integral equation for

the vector β becomes $\beta = e + \varepsilon^{-\nu} Z f' Z^{-1} H_0(\tau, t) \beta$, where $\tau = \varepsilon^{-\nu}(t - t_0)$. Now H_0 will be $O(1)$ only locally in a zone of extent ε^ν , and the form of the nonvanishing (in the limit as $\varepsilon \rightarrow 0$) contribution to β may possibly be deducible. This better estimate of the form of solution may provide the basis of generating a more useful integral equation for β . Thus there is a possibility that the method may be extended to matrices A having entries which are locally rapidly varying.

Even if this hope proves to be unbounded in such cases, we believe that the results will be useful in dealing with the problem of locating boundary layers in nonlinear systems. Thus the assumptions regarding the size of the derivatives of S should apply to zones away from the boundary layers, where the so-called outer solution is dominant. In such circumstances, the exponential variations will apply almost everywhere, and so the solutions obtained here permit an examination of the stability of an outer solution under perturbations at either endpoint, thus offering a means of separating the treatment of nonlinear boundary layers at either end of an interval.

To illustrate this, consider a system in which no change of order of eigenvalues occurs on an interval $[0, T]$, and define

$$\mu_i(t, \varepsilon) = \operatorname{Re} \int_0^t \lambda_i(s, \varepsilon) ds.$$

If μ_i achieves its maximum at $t = 0$, the corresponding solutions, as constructed here, are available to adjust to boundary conditions arising from a boundary layer at $t = 0$. On the other hand, if the solution to the system were forced by boundary conditions at $t = T$, the component corresponding to μ_i would be exponentially larger than at its 'origin', $t = T$. Analogous statements apply when μ_i achieves its maximum at $t = T$. It is of interest to note that the sign of $\operatorname{Re} \lambda_i$ may change without affecting the above statements. If μ_i has its maximum at an interior point, then the solution corresponding to μ_i will be exponentially larger in the interior than at either endpoint, and hence exponentially larger than boundary values imposed at these points. This phenomenon reminds one of the so-called resonance of Ackerberg and O'Malley [1], and is a point we shall take up in greater detail in Part II. If the maximum of μ_i is taken at both 0 and T , then data at either end may contribute to determining the solution of a boundary-value problem, and in this case, consistency may prevent the boundary layers' being treated separately.

Where a change of order takes place, accompanied by a change of sign of the real parts of the eigenvalues involved, the situation may be much more difficult. Here we content ourselves with a brief discussion of one illustrative example which indicates the nature of results to be expected. Let the (real) eigenvalue functions be λ_α of multiplicity one and λ_β of multiplicity $(n - 1)$. Let λ_β be negative on $[0, 2]$ and λ_α be less than λ_β on $[0, 1]$ and greater on $[1, 2]$. Further,

let the inequalities

$$\int_0^1 \lambda_\beta + \int_1^2 \lambda_\alpha > 0 \quad \text{and} \quad \int_0^2 \lambda_\alpha < 0$$

hold. Then it is possible to construct a basis for the solution space of vectors of the form:

$$\begin{aligned} \text{one solution:} & \quad \{s_\alpha + o(1)\} \exp\left(\varepsilon^{-1} \int_0^t \lambda_\alpha\right) \\ (n-2) \text{ solutions:} & \quad \{s_i + o(1)\} \exp\left(\varepsilon^{-1} \int_0^t \lambda_\beta\right), \end{aligned}$$

where the vectors s_i have been selected in R_β so that $d_{i\alpha} = 0$, and the remaining solution:

$$\text{and} \quad \left. \begin{aligned} & \{s_j + o(1)\} \exp\left(\varepsilon^{-1} \int_0^t \lambda_\beta\right) \quad \text{for } 0 \leq t \leq 1 \\ & d_{j\alpha} \{s_j + o(1)\} \exp\left(\varepsilon^{-1} \left\{ \int_0^1 \lambda_\beta + \int_1^t \lambda_\alpha \right\}\right) \quad \text{for } 1 < t \leq 2, \end{aligned} \right\}$$

where the vector s_j lies in R_β and is not uniquely determined. The first $(n-1)$ solutions describe boundary layer behaviour with the boundary layer located at $t = 0$, while the last solution, when rescaled, describes a boundary layer at $t = 2$. It is curious that the eigenvector function s_α is associated with two independent boundary layers, as is suggested by the crudest heuristic arguments concerning the sign λ_α , but there is an element of R_β which is not compatible with uniformly bounded solutions if invoked by the boundary condition at $t = 0$ where it appears natural. A more detailed discussion is deferred to Part II where boundary value problems for stiff systems will be discussed.

We further note that where the dimensions of the eigensubspaces change, our results indicate that there may be significant changes in the solution behaviour. However, we have already simplified matters there, by our assumptions that S is nonsingular and that \dot{S} exists. Where these fail, the behaviour is possibly much worse. The only feasible approach to the numerical treatment of general large order stiff systems would appear to be that, in the neighbourhood of such points, the integration step size should be reduced locally to a size where the system ceases to be stiff.

Acknowledgement

This research was supported in part by the Australian Research Grants Committee. The authors also wish to acknowledge Dr. P. B. Chapman's assistance in reading the manuscript, through which several errors were detected and corrected.

References

- [1] R. C. Ackerberg and R. E. O'Malley, "Boundary layer problems exhibiting resonance", *Studies in Applied Math.* 49 (1970), 277–295.
- [2] P. Chapman and J. J. Mahony, "Reflection of waves in a slowly varying medium", *S.I.A.M. J. Appl. Math.* 34 (1978), 303–319.
- [3] A. E. Erdélyi, *Asymptotic expansions* (Dover, New York 1956).
- [4] L. A. Liusternik and V. J. Sobolev, *Elements of functional analysis* (Hindustan Publishing Corp., India, 1974).
- [5] A. H. Nayfeh, *Perturbation methods* (John Wiley, New York, 1973).
- [6] L. F. Shampine and C. W. Gear, "A user's view of solving stiff ordinary differential equations", *S.I.A.M. Review* 21 (1979), 1–17.
- [7] Y. Sibuya, "Some global properties of matrices of functions of one variable", *Math. Annalen* 161 (1965), 67–77.
- [8] D. R. Smart, *Fixed point theorems* (Cambridge University, Press, London, 1974).
- [9] A. B. Vasil'eva, "The development of the theory of ordinary differential equations with a small parameter multiplying the highest derivative during the period 1966–1976", *Russian Math. Surveys* 31 (1976), 109–131.
- [10] W. Wasow, "Asymptotic expansion for ordinary differential equations: trends and problems" in Calvin H. Wilcox (Ed.), *Asymptotic solutions of differential equations and their applications*, (John Wiley, New York, 1964).

Department of Mathematics
University of Western Australia
Nedlands
Western Australia 6009