# ON THE CHOICE OF COORDINATE FUNCTIONS

*R.S. Anderssen*

## 1. INTRODUCTION

Numerically, there are two independent aspects to the problem of solving (partial) differential and integral equations computationally.  On the one hand, it is necessary to have results concerning the convergence, stability and accuracy of various classes of methods such as finite difference methods for initial value problems, finite element methods for elliptic partial differential equations, shooting methods for two-point boundary value problems, etc.  The general philosophy and expertise of numerical problem solving is based on such information.  On the other hand, for a specific equation which arises in an application, it is necessary to distinguish between the various algorithms which can be constructed.  For the particular equation under examination, the aim is not simply to apply any appropriate algorithm but to use the algorithm best suited to the task in hand.  Thus, the requirements of the latter differ considerably from that of the former.

In fact, the success of any algorithm constructed for a specific problem will depend heavily on the extent to which its design exploits the mathematical characteristics of the problem under examination.  Some specific examples are:  the use of the boundary integral method to solve potential problems defined on irregularly shaped regions;  the use of the inversion formulas to solve Abel integral equations;  the numerical stability of modified Gram-Schmidt;  Fourier methods on a regular grid;  sparse matrix computations;  parallelism in algorithm construction.

In situations where the starting point for the construction of a

numerical method is an approximation of the form

$$u_n(x) = \sum_{j=1}^{n} a_j^{(n)} \phi_j^{(n)}(x) \; ,$$

this exploitation of the structure of the specific equation under examination can be coupled to the choice of the basis functions $\phi_j^{(n)}(x)$ , $j = 1,2,\ldots,n$ . It is this aspect of computational problem solving which is examined in this paper.

Any examination of the choice of the basis functions divides naturally into the following two cases:

(i) the $\phi_j^{(n)}$ , $j = 1,2,\ldots,n$ , are independent of $n$ , which typifies the situation in the application of the classical variational methodology using globally defined functions, and of spectral methods;

(ii) the $\phi_j^{(n)}$ , $j = 1,2,\ldots,n$ , depend on $n$ , which typifies the situation in the application of finite element methods.

Though there is an obvious overlap between these two situations, we limit attention to the former. A detailed discussion of the latter can be found in Arnold *et al*. [4]. In order to distinguish between these two situations, we shall refer to globally defined basis functions as *coordinate functions*.

In fact, we limit attention to the following three aspects:

1. The practical appeal of the spectral method, where the coordinate functions $\phi_j$ , $j = 1,2,\ldots,n$ , correspond to the first $n$ components of an orthonormal system.

2. The choice of the coordinate functions as the eigenfunctions of a related but simpler operator than that defining the equation to be solved.

3. The flexibility of the Petrov-Galerkin strategy.

Much of the discussion will be concerned with densely defined linear operator equations

$$(1.1) \qquad \underset{\sim}{L}u = f , \quad u = u(x) , \quad x \in \Omega , \quad \underset{\sim}{L} : \underset{=}{D}(\underset{\sim}{L}) \to \underset{=}{R}(\underset{\sim}{L}) ,$$

where $\Omega$ is a bounded region in $\mathbb{R}^q$, and the domain $\underset{=}{D}(\underset{\sim}{L})$ and the range $\underset{=}{R}(\underset{\sim}{L})$ are both contained in the same Hilbert space $\underset{=}{H}$ with inner product $(u,v)$ and norm $\|u\| = (u,u)^{\frac{1}{2}}$.

## 2. THE PRACTICAL APPEAL OF THE SPECTRAL METHOD

There is no uniform use in the terminology "spectral method". It is sometimes used to describe theoretical studies of the spectral properties of operators. It is also used loosely to describe any numerical method which exploits, in one way or another, the properties of known orthonormal systems of functions (cf. Peyret and Taylor [23], Chapter 3).

The starting point for many spectral methods is the adoption of approximations of the form

$$(2.1) \qquad u_n(x) = \sum_{j=1}^{n} a_j^{(n)} \phi_j(x) ,$$

where the coordinate (basis, trial, shape) functions $\phi_j(x)$, $j = 1,2,\ldots,n$, are chosen to be the first $n$ elements of an orthonormal system $\{\phi_j\}_1^\infty = \{\phi_j ; j = 1,2,\ldots\}$. Clearly, the qualifier "spectral" identifies this particular choice for the coordinate functions. Such methods are subclassified in terms of the procedure used to determine the unknowns $a_j^{(n)}$, $j = 1,2,\ldots,n$; i.e. for the linear operator equation (1.1), the procedure is the $n$ conditions which, in conjunction with (2.1), yield the non-singular matrix equation of order $n$

$$L_n \underset{\sim}{a}^{(n)} = \underset{\sim}{f}^{(n)} , \quad \underset{\sim}{a}^{(n)} = [a_1^{(n)}, a_2^{(n)}, \ldots, a_n^{(n)}]^T$$

for determining the $a_j^{(n)}$, $j = 1,2,\ldots,n$. A discussion of various types of spectral methods can be found in Gottlieb and Orszag (1977) and Fletcher (1984). The problems analysed there are time dependent and therefore the

approximations (2.1) now take the form

(2.3)
$$u_n(x,t) = \sum_{j=1}^{n} a_j^{(n)}(t) \, \phi_j(x) \, ,$$

where the unknown constants $a_j^{(n)}$ , $j = 1,2,\ldots,n$ , of (2.1) have now

become unknown functions of the time $t$ .

Two examples of popular forms of spectral methods are given by:

EXAMPLE 2.1. *Spectral Collocation.* When the collocation method is used to

construct approximations of the form (2.1) for the solution of a linear

operator equation (1.1), the $L_n$ and $\underset{\sim}{f}^{(n)}$ of (2.2) take the form

(2.4)
$$L_n = \begin{pmatrix} \underset{\approx}{L}\phi_1(x_1) & \cdots & \underset{\approx}{L}\phi_n(x_1) \\ \cdots\cdots\cdots\cdots\cdots\cdots\cdots \\ \underset{\approx}{L}\phi_1(x_n) & \cdots & \underset{\approx}{L}\phi_n(x_n) \end{pmatrix} \, , \quad \underset{\sim}{f}^{(n)} = (f(x_1),f(x_2),\ldots,f(x_n))^T \, ,$$

where the $x_k \in \Omega$ , $k = 1,2,\ldots,n$ , define $n$ distinct collocation points.

The difficulty with the collocation method, which limits its applicability

and is a trade-off against its simplicity, is guaranteeing the non-

singularity of $L_n$ . In fact, we know (cf. Davis [7]) that the linear

independence of the $\phi_j(x)$ , $j = 1,2,\ldots,n$ , does not guarantee for

arbitrary $x_k$ , $k = 1,2,\ldots,n$ , the non-singularity of $L_n$ when either the

dimension $q$ of $\mathbb{R}^q$ is greater than $1$ , or $\Omega$ has a branch if it is an

$R^1$-curve in $\mathbb{R}^q$ with $q \geq 2$ .

Though there are certain advantages associated with choosing the $\phi_j$ ,

$j = 1,2,\ldots,n$ , to be the first $n$ components of an orthonormal system

$\{\phi_j\}_1^{\infty}$ (e.g. the discrete Fourier transform), such a choice does not remove

the above mentioned difficulty.                                             #

EXAMPLE 2.2. *The Pseudospectral Method.* This is the name given to the

spectral collocation method when it is applied to time dependent problems

in conjunction with the approximation (2.3). In such situations, the

choice of collocation points must be such as to yield a simple structure

for the system of ordinary differential equations which must be solved for

the $a_j(t)$ , $j = 1,2,\ldots,n$ . For example, if the Chebyshev polynomials

$T_j(x)$ are used, then the collocation points are chosen to be $x_k = \cos(\tau_k)$

for appropriately chosen $\tau_k$ so that the fact that $T_n(\cos\theta) = \cos n\theta$ can

be exploited.                                                                                    #

The motivation for the use of spectral methods is two-fold:

(a) The existence of extensive mathematical properties for particular

orthonormal systems, such as the Legendre and Chebyshev polynomials, which

can be exploited in various ways to manipulate the structure of numerical

methods based on the use of orthonormal functions.  There is an extensive

literature on this aspect.  It ranges from general studies of the utility

of specific orthonormal systems such as Legendre, Chebyshev and Jacobi

polynimials, in the numerical solution of ordinary and partial differential

equations as well as integral equations (cf. Delves and Freeman [8]);  to

specific studies of how one special class of orthonormal functions such as

the Chebyshev polynomials can be used to study a variety of problems

numerically by specifically exploiting the essential properties of the

orthonormal functions chosen (cf. Elliott [10] and Horner [15]).

(b) The knowledge that, in the numerical performance of variational

methods, the choice of the coordinate functions $\phi_j(x)$ , $j = 1,2,\ldots,n$ ,

appears to play a more crucial role than the  n  conditions chosen to

define (2.2);  and thereby, the heuristic conclusion that in some sense an

orthonormal system must be better than a non-orthonormal.

Though the success of spectral methods for the approximate solution of

a wide class of practical problems (cf. Gottlieb and Orszag [14], Peyret

and Taylor [23], Fletcher [12]) yields verification for this conclusion, it

is well known (cf. Gottlieb and Orszag [14] and Anderssen and Omodei [3])

that the choice of an orthonormal system does not guarantee unconditionally that a spectral method will perform well computationally. Limitations on the utility of taking arbitrary orthonormal systems to construct approximations of the form (2.1) have been examined in Anderssen [1].

We conclude this section with a more detailed discussion of the points in (a) above. In particular, our aim is to illustrate why, in some areas of computational mathematics and physics, spectral methods are viewed with considerable respect and even reverence. However, the aim is more than simply showing how spectral methods have been applied. The idea is to identify the mathematical reasons why the use of orthogonality allows something special or advantageous to be achieved numerically or pragmatically for the solution of some practical problem.

In order to emphasize the mathematical aspects, the discussion is organized to highlight such reasons.

## 2.1 Diagonalization of Matrices Which Must be Inverted

In most numerical procedures based on approximations of the form (2.1) or (2.3), it is necessary to invert at least one matrix. If the $\phi_j(x)$ , $j = 1,2,\ldots,n$ , are chosen so that one of the matrices involved is diagonal, then the computational process is greatly simplified.

### 2.1.1 The Rayleigh-Ritz-Galerkin Method for Eigenvalue Problems

Consider the general eigenvalue problem

$$(2.5) \qquad\qquad \underset{\sim}{A} u = \lambda \underset{\sim}{B} u ,$$

where $\underset{\sim}{A}$ and $\underset{\sim}{B}$ define time independent linear operators which map a known Hilbert space $\underset{=}{H}$ , with inner product $(\cdot,\cdot)$ and norm $\|\cdot\|$ , into itself. Using the approximations (2.1), the Rayleigh-Ritz-Galerkin method replaces (2.5) by the algebraic eigenvalue problem

(2.6) $\quad A_n \underset{\sim}{a}^{(n)} = \lambda^{(n)} B_n \underset{\sim}{a}^{(n)}$ , $\quad \underset{\sim}{a}^{(n)} = [a_1^{(n)}, a_2^{(n)}, \ldots, a_n^{(n)}]^T$ ,

where the matrices $A_n$ and $B_n$ take the form

$$A_n = \begin{pmatrix} (A\phi_1,\phi_1) & \cdots & (A\phi_n,\phi_1) \\ \hline (A\phi_1,\phi_n) & \cdots & (A\phi_n,\phi_n) \end{pmatrix} , \quad B_n = \begin{pmatrix} (B\phi_1,\phi_1) & \cdots & (B\phi_n,\phi_1) \\ \hline (B\phi_1,\phi_n) & \cdots & (B\phi_n,\phi_n) \end{pmatrix} .$$

Computationally, the solution of this eigenvalue problem is greatly
simplified if the coordinate functions correspond to the orthonormal
eigenfunctions of $\underset{\approx}{B}$ . In the more common situation where $\underset{\approx}{B} = \underset{\sim}{I}$ , the
identity operator, any orthogonal system in $\underset{\equiv}{H}$ will diagonalize $B_n$ ,
though it is more appropriate computationally to work with its orthonormal
counterpart.

## 2.1.2 The Numerical Solution of Parabolic and Hyperbolic Partial Differential Equations

Consider the time dependent partial differential equations which take
the form

$$\underset{\sim}{L}_t u = \underset{\sim}{L} u + f$$

where $\underset{\sim}{L}$ is a linear time-independent operator and $\underset{\sim}{L}_t$ denotes a partial
differential operator only involving derivatives with respect to $t$ . If
it is solved using the Ritz-Galerkin method in conjunction with the
approximation (2.1), then the system of ordinary differential equations
which determine the $a_j(t)$ are given by

(2.7) $\quad A_n \underset{\sim}{L}_t \underset{\sim}{a}^{(n)}(t) = L_n \underset{\sim}{a}^{(n)} + \underset{\sim}{f}$ , $\quad \underset{\sim}{a}^{(n)}(t) = [a_1^{(n)}(t), a_2^{(n)}(t), \ldots, a_n^{(n)}(t)]$ ,

where the matrices $A_n$ and $L_n$ take the form

$$A_n = \begin{pmatrix} (\phi_1,\phi_1) & \cdots & (\phi_n,\phi_1) \\ \hline (\phi_1,\phi_n) & \cdots & (\phi_n,\phi_n) \end{pmatrix} , \quad L_n = \begin{pmatrix} (L\phi_1,\phi_1) & \cdots & (L\phi_n,\phi_1) \\ \hline (L\phi_1,\phi_n) & \cdots & (L\phi_n,\phi_n) \end{pmatrix} .$$

Computationally, the problem of integrating (2.7) is greatly simplified if

$A_n$ is diagonal, which occurs if the $\{\phi_j\}_1^\infty$ are orthogonal in $\underline{\underline{H}}$ , though it is more appropriate to work with its orthonormal counterpart.

## 2.2 Decoupling

The complexity of many numerical procedures is a direct consequence of the cross-coupling which the n conditions, which define the numerical process to be solved for the $a_j^{(n)}$ , $j = 1,2,\ldots,n$ , forces between the different terms which define the approximation (2.1). This is implicit in the above discussion about diagonalization. Different numerical techniques have been proposed which explicitly exploit the decoupling inherent in an orthogonal (orthonormal) system.

### 2.2.1 The Tau Method

Consider the linear operator equation (1.1); i.e.

$$\underline{L}u = f .$$

The tau method of Lanczos [17] is essentially an analytic application of the backwards error analysis argument. The approximation $u_n$ is interpreted as the exact solution of

(2.8) $$\underline{L}u_n = f_n .$$

If $\underline{L}$ has a bounded inverse, then the difference $f - f_n$ yields a characterization of the quality of the approximation $u_n$ . Lanczos [17] proposed that this difference be modelled as

(2.9) $$f - f_n = \sum_{k=1}^{\infty} \tau_k \phi_{n+k}(x) ,$$

and thereby reduced the problem of estimating $f - f_n$ to the problem of determining the $\tau$ parameters in (2.9). For general systems $\{\phi_j\}_1^\infty$ , this reduces to a complex problem computationally. If however, the system $\{\phi_j\}_1^\infty$

is orthonormal in $\underline{\underline{H}}$ and its elements satisfy the boundary conditions

associated with $\underset{\sim}{L}$ , then it is a simple matter to show that

(2.10)  $\qquad\qquad \tau_k = (\underset{\sim}{L}u_n - f, \phi_{n+k})$ , $k = 1,2,\ldots$ .

If the elements of the system $\{\phi_j\}_1^\infty$ do not satisfy the boundary

conditions, it is a simple matter to modify the above argument (cf.

Gottlieb and Orszag [14], Section 2).

## 2.2.2 Manipulating Non-Linearities

In many applications, the non-linearity which makes the relevant

equations non-linear are only quadratic. For example, the Navier-Stokes

equations and Burger's equation in which the non-linearity takes the form

(2.11)  $\qquad\qquad u\, u_x$ , $\quad u = u(x,t)$ , $\quad u_x = \partial u/\partial x$ .

Various procedures are employed which exploit orthogonality in order to

decouple the cross-coupling inherent in the non-linearity. To a certain

extent, the methodology used is problem dependent and complex, because the

problems themselves are complex and the non-linear terms cannot be treated

in isolation from the other terms in the equations being solved. Never-

theless, the essence of the mathematics behind what is being done can be

described. We give two illustrations:

EXAMPLE 2.3. If the Ritz-Galerkin method is applied to the non-linearity

(2.11) using the approximations (2.1), it is necessary to construct a

matrix $B_n$ with elements $b_{jk}^{(n)}$ defined by

(2.12)  $\qquad\qquad b_{jk}^{(n)} = \sum_{i=1}^{n} a_i^{(n)} \left( \phi_j \frac{d\phi_i}{dx}, \phi_k \right)$ ,

where we have assumed that the functions $u = u(x)$ are contained in $\underline{\underline{H}}$ .

The non-linearity manifests itself through the dependence of the

coefficients $b_{jk}^{(n)}$ on the $a_i^{(n)}$ , $i = 1,2,\ldots,n$ . The manipulation of
this non-linearity is reduced to finding an orthonormal system for the
$\{\phi_j\}_1^\infty$ (e.g. the Legendre polynomials) which allows the evaluation of the
inner products in (2.12) to be greatly simplified.                                    #

EXAMPLE 2.4.  Using either a Petrov-Galerkin or method of integral
relations framework, and after appropriate changes of variable and other
manipulations, the approximation of the non-linearity is reduced to an
examination of integrals of the form

(2.13)
$$\int w(v) \, f_k(v) \, dx \, , \quad v = v(x) \, ,$$

where the resulting coordinate functions $f_k(v)$ end up being functions of
some transformed unknown $v$ rather than $x$ . The key step is then to
reverse the roles of $v$ and $x$ with $v$ becoming the independent
variable and $x$ a function of $v$ given by

$$x(v) = \int^v \eta(\tau) \, d\tau \, .$$

This has proved to be an incredibly successful way of coping with
non-linearities in various situations, since the non-linear aspect is
transferred to a linear.  Here, (2.13) becomes

(2.14)
$$\int w(v) \, f_k(v) \, \eta(v) \, dv \, .$$

If the unknown $\eta(v)$ is now approximated by

$$\eta_n(v) = \sum_{j=1}^n b_j^{(n)} \, f_k(v) \, ,$$

the task of manipulating (2.13) is greatly simplified if the $f_k(v)$ are
constructed to be orthonormal with respect to the weight function $w(x)$ .
The implementation of such manipulations is in general quite difficult;
but has considerable computational advantages when achieved as, for example,

the discussion of Fleet and Fletcher [11] shows.                    #

Fuller details for specific examples such as the solution of Burger's equation can be found in Gottlieb and Orszag [14], Orszag [21] and Fletcher [12].

## 2.3 Transformation

Each approximation $u_n(x)$ has two representations: Its *physical* or *continuous* which corresponds to $u_n(x)$ itself as a function of $x$; and its *vector* or *finite dimensional* which corresponds to $\underset{\sim}{a}^{(n)} = [a_1^{(n)}, a_2^{(n)}, \ldots, a_n^{(n)}]^T$. When the system $\{\phi_j\}_1^\infty$ is orthonormal, the latter representation is often called the *spectral*. Thus, any computation involving the approximation $u_n(x)$ can be performed using either the physical or the vector (spectral) representation. The advantages of computing in one rather than the other can only be exploited if the transformations from one to the other can be evaluated economically. This is only possible if the system $\{\phi_j\}_1^\infty$ is orthonormal for then the inverse of the forward transformation

$$(2.15) \qquad u_n(x_\ell) = \sum_{j=1}^{n} a_j \, \phi_j(x_\ell) \, , \qquad \ell = 1, 2, \ldots, n \, ,$$

is given by

$$(2.16) \qquad a_k = \int u_n \, \phi_k(x) \, dx \, , \qquad k = 1, 2, \ldots, n \, .$$

Ordinarily, the evaluation of either (2.15) or (2.16) will involve $O(n^2)$ operations. Using fast methods, this can often be reduced to $O(n \log n)$ operations (cf. Orszag [21]). As explained by Orszag [21], and illustrated convincingly for multidimensional calculations by McCrory and Orszag [18] the aim is to choose the representation most appropriate for the computation required.

## 2.3.1 Evaluation of Non-Linear Terms

Consider the problem of evaluating the convective term $u\,u_x$ in Burger's equation (cf. Fletcher [12]). If at time $t_m$ it is known that

$$(2.17) \qquad u_n(x,t_m) = \sum_{j=1}^{n} a_{j,m}^{(n)} \phi_j(x)$$

then $u\,u_x$ at $x = x_1, x_2, \ldots, x_n$, can be evaluated as

$$(2.18) \quad [u(x,t_m)\,u_x(x,t_m)]_{x=x_\ell} = \left( \sum_{j=1}^{n} a_{j,m}^{(n)} \phi_j(x_\ell) \right) \left( \sum_{j=1}^{n} a_{j,m}^{(n)} \left( \frac{\partial \phi_j}{\partial x} \right)_{x=x_\ell} \right),$$

for $\ell = 1, 2, \ldots, n$. This clearly involves $O(n^2)$ operations. Alternatively, $u_x$ could be evaluated as

$$(2.19) \qquad \frac{\partial u}{\partial x} = \sum_{j=1}^{n} b_{j,m}^{(n)} \phi_j(x) \ .$$

If the $\{\phi_j\}_1^\infty$ correspond to the Fourier components, the Legendre polynomials or the Chebyshev polynomials, the recurrence relations which specify the $b_{j,m}^{(n)}$, $j = 1, 2, \ldots, n$, in terms of the $a_{j,m}^{(n)}$, $j = 1, 2, \ldots, n$, in $O(n)$ operations are known for a variety of situations which include Burger's equation. Thus, the $O(n^2)$ operations involved with evaluating $u\,u_x$ at $x = x_1, x_2, \ldots, x_n$, can be reduced to $O(n \log n)$ operations through the judicious use of (2.19) and the physical and spectral representations:

(i) transform (using $O(n \log n)$ operations) to the physical representation

$$\hat{u}_n(x_\ell, t_m) = \sum_{j=1}^{n} a_{j,m}^{(n)} \phi_j(x_\ell) \ , \qquad \ell = 1, 2, \ldots, n \ ;$$

(ii) evaluate (using $O(n)$ operations)

$$b_{j,m}^{(n)} \quad \text{from} \quad a_{j,m}^{(n)}$$

using the recurrence relations;

(iii) transform (using $O(n \log n)$ operations) to the physical
representation

$$\frac{\partial \hat{u}_n(x_\ell, t_m)}{\partial x} = \sum_{j=1}^{n} b_{j,m}^{(n)} \phi_j(x_\ell) , \quad \ell = 1, 2, \ldots, n ;$$

(iv) evaluate (using $O(n)$ operations)

$$w_n(x_\ell, t_m) = \hat{u}_n(x_\ell, t_m) \frac{\partial \hat{u}_n(t_\ell, t_m)}{\partial x} , \quad \ell = 1, 2, \ldots, n ;$$

(v) transform (using $O(n \log n)$ operations) back to the spectral
representation

$$c_{\ell,m}^{(n)} = \int w_n(x, t_m) \phi_\ell(x) \, dx .$$

## 2.4 Mimicking the Eigenfunction Solution

Again, consider the operator equation (1.1). If $\underset{\sim}{L}$ has a discrete
spectrum $\{\lambda_j\}_1^\infty$ and its (normalized) eigenfunctions $\{\psi_j\}_1^\infty$ are known
and form a basis in $\underset{=}{H}$ then the solution $u_f$ of (1.1) is automatically
given by

(2.20)
$$u_f = \sum_{j=1}^{\infty} \frac{(f, \psi_j) \psi_j}{\lambda_j} .$$

In many ways, the use of the approximations $u_n$ of (2.1) can be seen as an
attempt to mimic this eigenfunction representation of the solution $u_f$.
In addition, heuristic and intuitive considerations lead naturally to the
idea that, even if the $\psi_j(x)$ are not known exactly, every attempt should
be made to choose systems $\{\phi_j\}_1^\infty$ which approximate the $\{\psi_j\}_1^\infty$ in some
appropriate way (which will depend on the context of the problem being
solved).

This point has been discussed for spectral methods by Anderssen [1]
and is motivated from the more general point of view of the numerical
stability of variational methods in the next section of this paper. From a

practical point of view, it is instrumental in the use of Legendre

polynomials in weather prediction (cf. Fletcher [12], §5.6.1).

## 3. EIGENFUNCTIONS OF SIMILAR BUT SIMPLER OPERATORS

For a given choice of coordinate functions $\phi_j^{(n)}$ , $j = 1,2,\ldots,n$ , the

associated matrix equation (2.2) is solved computationally for the unknown

vector $\underset{\sim}{a}^{(n)}$ defining $u_n$ in (2.1). Appealing to the backward error

analysis argument, the computed solution $\underset{\sim}{b}^{(n)}$ is interpreted as the exact

solution of the perturbed system

$$(3.1) \qquad (L_n + \Delta L_n)\underset{\sim}{b}^{(n)} = \underset{\sim}{f}^{(n)} + \Delta\underset{\sim}{f}^{(n)} \ .$$

Under the restriction that

$$(3.2) \qquad r = \|L_n^{-1} \Delta L_n\| < 1 \ ,$$

a standard argument (cf. Mikhlin [19], §9) shows that

$$(3.3) \quad \|\underset{\sim}{a}^{(n)} - \underset{\sim}{b}^{(n)}\| \leq \|L_n^{-1} \Delta L_n\| \ \|\underset{\sim}{a}^{(n)}\|/(1-r) + \|L_n^{-1}\| \ \|\Delta\underset{\sim}{f}^{(n)}\|/(1-r) \ .$$

It follows automatically that a sufficient condition for the error

$\|\underset{\sim}{a}^{(n)} - \underset{\sim}{b}^{(n)}\|$ to remain bounded is that $\|L_n^{-1}\|$ remain bounded. If the

spectral norm of $L_n^{-1}$ is taken, then the boundedness of $L_n^{-1}$ can be

related to the behaviour of the smallest eigenvalue of the Gram matrix $L_n$.

This fact was used by Mikhlin [19] as the basis of his definition of

stability.

Because it is developed in terms of constructive concepts which can be

tested, Mikhlin's stability theory for variational methods has the

following appealing structure.

Stability Definition. The numerical process defined by

$$(3.4) \qquad L_n \underset{\sim}{a}^{(n)} = \underset{\sim}{f}^{(n)} \ , \qquad n = 1,2,\ldots,$$

is *M-stable*, if there exist constants  p, q  and  s  independent of  n

such that for  $\|\Delta L_n\| \leq s$  and arbitrary  $\Delta f^{(n)}$  the perturbed system (3.1)

is solvable and the following error estimate

(3.5)
$$\| \underset{\sim}{a}^{(n)} - \underset{\sim}{b}^{(n)} \| \leq p \| \Delta L_n \| + q \| \Delta \underset{\sim}{f}^{(n)} \|$$

holds.

The connection between (3.4) and (3.5) is obvious.

*The Strong Minimality of the*  $\{\phi_j\}_1^\infty$  *in*  $\underset{=}{H}$ .  The coordinate functions

$\{\phi_j\}_1^\infty$  are said to be *strongly minimal* in  $\underset{=}{H}$ , if the smallest eigenvalues

$\lambda_1^{(n)}$  of the positive definite Gram matrices

(3.6)
$$G_n = \begin{pmatrix} (\phi_1,\phi_1) & \cdots & (\phi_n,\phi_1) \\ \cdots\cdots\cdots\cdots\cdots \\ (\phi_n,\phi_1) & \cdots & (\phi_n,\phi_n) \end{pmatrix} , \quad n = 1,2,\ldots,$$

are bounded away from zero;  i.e.

(3.7)
$$\inf_n \lambda_1^{(n)} \geq \lambda_0 > 0 .$$

*A necessary and Sufficient Condition for M-Stability.*  If  $\underset{\sim}{L} = \underset{\sim}{A}$ ,

where  $\underset{\sim}{A}$  is a positive definite operator, a necessary and sufficient

condition for the M-stability of the numerical process (3.4) is that the

coordinate system  $\{\phi_j\}_1^\infty$  used to generate (3.4) is strongly minimal in the

energy space  $\underset{=A}{H}$  with inner product  $[u,v] = (\underset{\sim}{A}u,v)$  and norm

$\| u \| = [u,u]^{\frac{1}{2}}$ .

Thus, guaranteeing the M-stability of the numerical process (3.4)

reduces to ensuring that the coordinate system  $\{\phi_j\}_1^\infty$  is strongly minimal

in  $\underset{=A}{H}$ .  There are a number of ways in which this can be done:

(a)  use the orthonormal eigenfunctions of an operator  $\underset{\sim}{B}$  which is

simpler but similar to  $\underset{\sim}{A}$ ;  (b)  scale  minimal  systems

in  $\underset{=A}{H}$ , which are not strongly minimal, to be strongly minimal.  A system

$\{\phi_j\}_1^\infty$  is said to be *minimal*, if the span of each (and every one) of its

subsets is a proper subspace of the span of $\{\phi_j\}_1^\infty$ .

In (a) (cf. Mikhlin [19], §3), two self-adjoint and positive definite operators $\underset{\sim}{A}$ and $\underset{\sim}{B}$ are said to be *similar* if $\underline{D}(\underset{\sim}{A}) = \underline{D}(\underset{\sim}{B})$ ; i.e. the domains of $\underset{\sim}{A}$ and $\underset{\sim}{B}$ are identical. Using this definition, Mikhlin [19] has identified a variety of circumstances which guarantee strong minimality in $\underset{=A}{H}$ , and used these results to propose a rational basis for the choice of the $\{\phi_j\}_1^\infty$ . The orthonormal eigenfunctions of the more important simpler operators which arise in practical situations are listed.

Anderssen [1] has discussed how these results can be used to examine the numerical performance of spectral methods and thereby clarify to what extent the choice of an arbitrary orthonormal system as the coordinate system $\{\phi_j\}_1^\infty$ can be justified numerically. The fact that "any ortho-normal system which lies in $\underline{\underline{H}}$ is strongly minimal in $\underset{=A}{H}$ if it is also contained in $\underset{=A}{H}$ and spans $\underset{=A}{H}$ " shows that convergent and stable approximations of the form (2.1) can be constructed using arbitrary ortho-normal systems in $\underline{\underline{H}}$ , when the procedure used to construct the numerical processes (3.4) corresponds to one of the standard methodologies such as Ritz-Galerkin, Bubnov-Galerkin or least squares.

However, the need for having an appropriate mathematical framework in which to formulate such results is more crucial than it might at first sight appear. It is not simply a matter of choosing an arbitrary ortho-normal system in $\underline{\underline{H}}$ , which a loose interpretation of the above comments might imply. As the following discussion illustrates, the interrelation-ships between the different spaces involved impose their own restrictions on how the orthonormal system in $\underline{\underline{H}}$ must be chosen.

If the orthonormal system $\{\phi_j\}_1^\infty$ in $\underline{\underline{H}}$ is also a spanning set in $\underset{=A}{H}$ then it is also a spanning set in $\underline{\underline{H}}$ . It is an automatic consequence of

the continuous imbedding of $H_{=A}$ in $H_{=}$ . However, if the orthonormal system $\{\phi_j\}_1^\infty$ is contained in $H_{=A}$ but is only chosen to be a spanning set in $H_{=}$ , this does not guarantee that $\{\phi_j\}_1^\infty$ is a spanning set in $H_{=A}$ . It is an immediate consequence of the fact that, for sequences in $H_{=A}$ , this convergence in $H_{=}$ does not imply their convergence in $H_{=A}$ . The inequality defining the continuous imbedding of $H_{=A}$ in $H_{=}$ goes the wrong way.

In order to obtain a strongly minimal system in $H_{=A}$ , it is not necessary to start with an orthonormal system in $H_{=}$ which lies in $H_{=A}$ . In fact, any minimal system in $H_{=}$ which lies in $H_{=A}$ can be used. Now, the alternative strategy of scaling minimal systems can be used. Dovbysh [9] has shown that any minimal system in a Hilbert space can be rescaled to yield a strongly minimal system in that space (cf. Mikhlin [19], Theorem 2.2). Either scale the minimal system in $H_{=}$ to be strongly minimal, since the resulting system will also be strongly minimal in $H_{=A}$ ; or utilize the fact that the minimal system in $H_{=}$ will also be minimal in $H_{=A}$ and therefore scale it in $H_{=A}$ to make it strongly minimal.

The advantages of this scaling from a theoretical point of view have been outlined above. However, they are based on asymptotic results which carry little information about the practical consequences of this scaling. To date such consequences have not been examined in any detail.

Since the scaling of matrix equations is known to be so mercurial, it is natural to ask what advantages result in scaling a minimal system to be strongly minimal. We know from Forsythe and Moler [13] that if two base-scaled ($\beta$-scaled) equivalent systems are used, then applying Gaussian elimination with the same pivoting sequence to both will yield the same significands in both solutions. However, if different pivoting sequences are used, then different significands will result. It is therefore

assumed that, for the matrix equations which arise from the use of minimal and strongly minimal systems, different pivoting sequences will result if Gaussian elimination with either full or partial pivoting is applied; but this assumption remains to be examined in detail.

Another important aspect of M-stability is that it is based on an examination of absolute not relative error. It is for this reason that the concept of strong minimality plays such a major role as a necessary and sufficient condition in guaranteeing M-stability. In fact, for the numerical process (3.4) and its backward error analysis counterpart (3.1), it can be shown (cf. Mikhlin [19], §9) that

$$(3.8) \qquad \| \underset{\sim}{a}^{(n)} - \underset{\sim}{b}^{(n)} \| \leq \{ \| L_n^{-1} \Delta L_n \| \, \| \underset{\sim}{a}^{(n)} \| + \| L_n^{-1} \Delta \underset{\sim}{f}^{(n)} \| \} / (1-r)$$

with $r = \| L_n^{-1} \Delta L_n \|$ . The solvability of (3.1) requires that $r < 1$ . The key role played by $L_n^{-1}$ when bounding $\| \underset{\sim}{a}^{(n)} - \underset{\sim}{b}^{(n)} \|$ is immediately apparent. Even more importantly, (3.8) shows that a sharp estimate of $\| \underset{\sim}{a}^{(n)} - \underset{\sim}{b}^{(n)} \|$ must be based on the size of $\| L_n^{-1} \Delta L_n \|$ and $\| L_n^{-1} \Delta \underset{\sim}{f}^{(n)} \|$ ; and consequently, the conservative nature of the often used alternative bound

$$(3.9) \qquad \| \underset{\sim}{a}^{(n)} - \underset{\sim}{b}^{(n)} \| \leq \| L_n^{-1} \| \, \{ \| \Delta L_n \| \, \| \underset{\sim}{a}^{(n)} \| + \| \Delta \underset{\sim}{f}^{(n)} \| \} / (1-r) \ .$$

The use of relative errors when studying the effect of rounding error has become popular for a number of reasons:

> (i) in situations where the size of $\| \underset{\sim}{a}^{(n)} \|$ is not known *a priori*, it gives a more realistic assessment of the error than the absolute.

> (ii) the forward error analysis of numerical methods often yields estimates of the form
>
> $$\| \Delta L_n \| = c(n) \, \delta \, \| L_n \| \ ,$$
>
> where $c(n)$ is a constant depending on $n$ and $\delta$ is the basic rounding error.

For such reasons, one often works with the relative error estimate

$$(3.10) \qquad \|\underset{\sim}{a}^{(n)} - \underset{\sim}{b}^{(n)}\| / \|\underset{\sim}{a}^{(n)}\| \leq \kappa(L_n) \left\{ \frac{\|\Delta L_n\|}{\|L_n\|} + \frac{\|\Delta \underset{\sim}{f}^{(n)}\|}{\|\underset{\sim}{f}^{(n)}\|} \right\} / (1-r)$$

with the condition number $\kappa(L_n) = \|L_n\| \, \|L_n^{-1}\|$ . However, because it can be derived from (3.9) (divide (3.9) by $\|\underset{\sim}{a}^{(n)}\|$ and use the inequality $\|f^{(n)}\| / \|L_n\| < \|\underset{\sim}{a}^{(n)}\|$ ), (3.10) carries no more information than (3.8). For this reason, the interpretation of (3.10) cannot be done independently of (3.8) which goes against the point made by Omodei [21] about M-stability.

For example, (3.10) yields the conclusion that the relative error is bounded if $\kappa(L_n)$ is bounded. However, this condition allows the possibility that $\lambda_n^{(n)} \uparrow \infty$ at the same rate as $\lambda_1^{(n)} \downarrow 0$ which is in conflict with the required boundedness of $\|\underset{\sim}{a}^{(n)} - \underset{\sim}{b}^{(n)}\|$ for which we need $\lambda_1^{(n)} \geq \lambda_0 > 0$ . Clearly, we require that both the absolute and relative errors are bounded. The mentioned pathology can occur when the growth of $\|\underset{\sim}{a}^{(n)} - \underset{\sim}{b}^{(n)}\|$ is dominated by the decay of $1/\|\underset{\sim}{a}^{(n)}\|$ . This however is not the full picture as the role played by the condition $\|L_n^{-1} \Delta L_n\| < 1$ has not been taken into account.

A detailed analysis of strong minimality and condition number is not appropriate here as the above argument illustrates cogently the need to first guarantee strong minimality when applying variational methods.

## 4. THE PETROV-GALERKIN FRAMEWORK

We again start with the operator equation formulation (1.1). We construct the corresponding bilinear representation

$$(4.1) \qquad (\underset{\sim}{L}u, v) = (f, v) , \quad \text{for all } v \in \underline{\underline{V}} ,$$

which is equivalent to (1.1) if $\underline{\underline{V}}$ defines a dense subset of $\underline{\underline{H}}$ such as $\underline{\underline{D}}(\underline{\underline{L}})$ .

With respect to appropriately chosen $\phi_i$ , $i = 1,2,\ldots,n$ , approximations of the form (2.1) are sought which satisfy

(4.2) $\qquad\qquad (\underset{\sim}{L}u_n,v) = (f,v)$ , $\quad$ for all $v \in \underset{=}{V}_n$ ,

exactly with $\underset{=}{V}_n$ an appropriate finite dimensional subspace of $\underset{=}{V}$ . In particular, if the $\chi_j$ , $j = 1,2,\ldots,n$ , define a basis for $\underset{=}{V}_n$ , then (4.2) reduces to the following algebraic system which defines the Petrov-Galerkin method for (4.1), and hence, (1.1)

(4.3) $\qquad\qquad \sum_{i=1}^{n} (L\phi_i,\chi_j) = (f,\chi_j)$ , $\quad j = 1,2,\ldots,n$ .

Thus, (4.3) defines the Petrov-Galerkin counterpart of (2.2) and (3.4). The Galerkin method corresponds to the quite special situation where the $\chi_j = \phi_j$ , $j = 1,2,\ldots,n$ . The advantage of the Petrov-Galerkin method is the greater flexibility it gives to the construction of the algebraic system (4.3). Having made the choice of the $\phi_j$ , $j = 1,2,\ldots,n$ , for one reason, the choice of the $\chi_j$ , $j = 1,2,\ldots,n$ , can be exploited for another. For example, in the Galerkin method, the $\{\phi_j\}_1^n$ can either be chosen to optimize, with respect to $n$ , the representation used for $u_n$ or $\underset{\sim}{f}^{(n)}$ . On the other hand, the Petrov-Galerkin framework allows the $\{\phi_j\}_1^n$ to be chosen to optimize the representation used for $u_n$ leaving the $\{\chi_j\}_1^n$ to be chosen so as to optimize the representation used for $\underset{\sim}{f}^{(n)}$ .

For example, if $\{\phi_j\}_1^n$ were known to yield an efficient representation for $u_n$ , then a natural choice for the representation of $\underset{\sim}{f}^{(n)}$ would be $\{\underset{\sim}{A}\phi_j\}_1^n$ . The corresponding Petrov-Galerkin method is in fact the method of least squares which ensures that

$$\min_{u_n \in \text{span}(\{\phi_j\}_1^n)} \| \underset{\sim}{A}u_n - f \|$$

is attained. However, because we do not know  u  in advance, it is difficult to choose the  $\{\phi_j\}_1^n$  in relation to some property of  u  in advance. Nevertheless, we do know  f . It is a far simpler matter to seek the  $\{\chi_j\}_1^n$  which yield an efficient representation for  $\underset{\sim}{f}^{(n)}$ . It would then be necessary to decide what a corresponding choice of the  $\{\phi_j\}_1^n$  should be.

The obvious choice is  $\phi_j = \underset{\sim}{A}^{-1} \chi_j$  which is clearly not available except under special circumstances. When it is the corresponding Petrov-Galerkin method reduces to a pseudo-analytic method (cf. Anderssen and de Hoog [2])  for which the approximation  $u_n$  is known the moment the approximation

$$f_n = \sum_{j=1}^{n} b_j^{(n)} \chi_j$$

to  f  has been constructed;  namely,

$$u_n = \underset{\sim}{A}^{-1} f_n = \sum_{j=1}^{n} b_j^{(n)} \phi_j \ .$$

When applicable, pseudo-analytic methods are quite successful for properly posed problems in application. When the application is improperly posed, they are invariably unsuccessful as they contain no stabilization which damps out the enhancement of errors between  $f_n$  and  f  implicit in the transformation from the  $\chi_j$  to the  $\phi_j$ .

A less obvious choice is  $\phi_j = \underset{\sim}{A}^* \chi_j$ . The corresponding Petrov-Galerkin method is in fact Murray's method which has been examined in some detail by Petryshyn [22]. In fact, using the properties of  $\underset{\sim}{A}^*$ , the algebraic equations (4.3) become

$$(4.4) \qquad \sum_{i=1}^{n} (\underset{\sim}{A}^*\psi_i, \underset{\sim}{A}^*\psi_j)\, a_i^{(n)} = (f, \psi_j)\ , \qquad j = 1, 2, \ldots, n \ .$$

There are two levels at which the Petrov-Galerkin methodology is

manipulated.  On the one hand, it is exploited in a fairly explicit and practical manner to yield for specific applications numerical processes (3.4) with desired numerical properties.  This reduces to working directly with the numerical properties associated with particular choices of coordinate functions and centers on the algebraic properties of the matrices $L_n$ defining the numerical processes (3.4).  Some examples are: spline Petrov-Galerkin methods for the Korteweg-de Vries equation where the $\chi_j$ are translates of the $\phi_j$ (Schoonbie [24];  subdomain methods Fletcher [12]);  collocation methods such as spectral collocation methods (Voigt *et al.* [25]).

On the other hand, error estimates which formalize some of the observations made above about the flexibility of the Petrov-Galerkin methodology are derived and then manipulated in both a general and specific manner depending on the context in which the Petrov-Galerkin method is being examined.  For example, Canuto and Quateroni [6] derive a general error estimate which displays explicitly the dependence of the error $u - u_n$ on the form of the subspaces $\Phi_n$ and $X_n$ spanned by the $\phi_j$ and $\chi_j$ , $j = 1,2,\ldots,n$ , respectively.  They use this estimate to derive specific results about the numerical performance of the spectral method.

This error estimate is a generalization of an estimate given in Babuska and Aziz [5] which limits attention to the situation where the inner products are evaluated exactly.  Under appropriate coercivity conditions, the Babuska and Aziz estimate for $\left\vert\!\left\vert\!\left\vert u - u_n \right\vert\!\right\vert\!\right\vert$ takes the form of a best approximate estimate

$$\left\vert\!\left\vert\!\left\vert u - u_n \right\vert\!\right\vert\!\right\vert \leq K \left\{ \inf_{w \in \Phi_n} \left\vert\!\left\vert\!\left\vert u - w \right\vert\!\right\vert\!\right\vert \right\} ,$$

where $\left\vert\!\left\vert\!\left\vert \cdot \right\vert\!\right\vert\!\right\vert$ denotes an appropriate energy norm.

The role of the $\chi_j$ only enters through $K$ .  First and foremost, the

error depends on the choice of $\Phi_n$ , and then only on $X_n$ . This fits naturally the framework of the application in that the $\phi_j$ , $j = 1,2,\ldots,n$ , determine $u_n$ while the choice of the $\chi_j$ , $j = 1,2,\ldots,n$ , determines the algebraic system (3.4) to be solved. Thus, if $u$ lies in $\Phi_n$ , even if a poor choice is made for $X_n$ , the error is still zero (assuming the under-lying algebra is done exactly). Thus, the choice of $X_n$ determines the numerical properties of (3.4) (relative to the choice of $\Phi_n$ ). This confirms intuition.

Clearly, the choice of $X_n$ must not only ensure that (3.4) has appropriate numerical properties, but also that $K$ is kept small. In achieving the latter, Jenkinson [16] has shown that the size of $K$ depends crucially on the relationship between $\Phi_n$ and $X_n$ .

The Canuto and Quarteroni [6] estimate is most useful as it shows the effect of not being able to evaluate the inner products exactly. Now the form of $X_n$ used plays a more dominant role. It confirms the above point that the choice of $X_n$ is crucial in determining the numerical properties of the resulting numerical process. In fact, if the approximations used for $(\underset{\sim}{L}u,v)$ and $(f,v)$ are denoted by $B_N(u,v)$ and $f_N(v)$ , then their estimate takes the form

$$\|u - u_N\| \leq \inf_{w \in \Phi_n} \left\{ K_1(n) \|u - w\| + K_2(n) \sup_{v \in X_n} \left( \frac{|(\underset{\sim}{L}u,v) - B_N(u,v)|}{\|v\|} + \frac{|(f,v) - f_N(v)|}{\|v\|} \right) \right\}.$$

From the discussion contained in Canuto and Quarteroni [6], as well as the comments made above, it is clear that the use of such error estimates can assist greatly in tuning the choice of the $\phi_j$ and $\chi_j$ , $j = 1,2,\ldots,n$ , for specific applications.

## REFERENCES

[1]   R.S. Anderssen, 'On the numerical performance of spectral methods', *Proceedings of Workshop on Numerical Analysis and Optimization,* Proceedings of the Centre for Mathematical Analysis, Vol.6, ANU, Dec., 1983.

[2]   R.S. Anderssen and F.R. de Hoog, *Application and numerical solution of Abel-type integral equations,* Mathematics Research Report #7-1982, The Australian National University, pp.42, 1982.

[3]   R.S. Anderssen and B.J. Omodei, 'On the stability of uniformly asymptotically diagonal systems', *Math. Comp.* 28 (1974), 719-730.

[4]   D.N. Arnold, I. Babuška and J. Osborn, 'Finite element methods: principles for their selection', *Comp. Methods in Appl. Mech. and Eng.* 45 (1984), 57-96.

[5]   I. Babuška and A.K. Aziz, 'Survey lecture on the mathematical foundations of the finite element method', in *The Mathematical Foundations of the Finite Element Method with Applications to Partial Differential Equations* (editor:  A.K. Aziz), pp.3-359, Academic Press, New York, 1972.

[6]   C. Canuto and A. Quarteroni, 'Variational methods in the theoretical analysis of spectral approximations', in *Spectral Methods for Partial Differential Equations* (editors:  R.G. Voigt, D. Gottlieb and M.Y. Hussaini), SIAM, Philadelphia, 1984.

[7]   P.J. Davis, *Interpolation and Approximation*, Blaisdell, New York, 1963.

[8]   L.M. Delves and T.L. Freeman, *Analysis of Global Expansion Methods: Weakly Asymptotically Diagonal Systems*, Academic Press, London, 1981.

[9]   L.N. Dovbysh, 'A note on minimal systems', *Trudy. Matem. in-ta. im. V.A. Steklov* 96 (1968), 188-189.

[10]    D. Elliott, 'A method for the numerical integration of the one-
        dimensional heat equation using Chebyshev series', *Proc. Camb. Phil.
        Soc.* <u>57</u> (1961), 823-832.

[11]    R.W. Fleet and C.A.J. Fletcher, 'Application of the Dorodnitsyn
        boundary layer formulation to wall blowing', in *Computational
        Techniques and Applications: CTAC-83* (editors: J. Noye and
        C. Fletcher), pp.626-640, North-Holland, Amsterdam, 1984.

[12]    C.A.J. Fletcher, *Computational Galerkin Methods*, Springer-Verlag,
        Berlin, 1984.

[13]    G. Forsythe and C.B. Moler, *Computer Solution of Linear Algebraic
        Equations*, Prentice-Hall, Englewood Cliffs, N.J., 1967.

[14]    D. Gottlieb and S.A. Orszag, *Numerical Analysis of Spectral Methods:
        Theory and Applications*, SIAM, Philadelphia, Penn., 1977.

[15]    T.S. Horner, 'A double Chebyshev series method for elliptic partial
        differential equations', in *Numerical Solution of Partial
        Differential Equations* (editor: J. Noye), pp.573-590, North-
        Holland, Amsterdam, 1982.

[16]    J. Jenkinson, *The Petrov-Galerkin Method*, Ph.D Thesis - in
        preparation.

[17]    C. Lanczos, *Applied Analysis*, Prentice-Hall, Englewood Cliffs, N.J.,
        1956.

[18]    R.L. McCrory and S.A. Orszag, 'Spectral methods for multi-
        dimensional diffusion problems', *J. Comp. Phys.* <u>37</u> (1980), 93-112.

[19]    S.G. Mikhlin, *The Numerical Performance of Variational Methods*,
        Wolters-Noordhoff Publishing, Groningen, The Netherlands, 1971.

[20]    B.J. Omodei, 'On the numerical stability of the Rayleigh-Ritz
        method', *SIAM J. Numer. Anal.* <u>14</u> (1977), 1151-1171.

[21]    S.A. Orszag, 'Spectral methods for problems in complex geometries',
        *J. Comp. Phys.* <u>37</u> (1980), 70-92.

[22]    W.V. Petryshyn, 'Direct and iterative methods for the solution of
        linear operator equations in Hilbert space', *Trans. Am. Math. Soc.*
        <u>105</u> (1962), 136-175.

[23]    R. Peyret and T.D. Taylor, *Computational Methods for Fluid Flow*,
        Springer Series in Computational Physics, Springer-Verlag, New York,
        1983.

[24]    S.W. Schoombie, 'Spline Petrov-Galerkin methods for the numerical
        solution of the Korteweg-de Vries equation', *IAM J. Num. Anal.* <u>2</u>
        (1982), 95-109.

[25]    R.G. Voigt, D. Gottlieb and M.Y. Hussaini (editors), *Spectral
        Methods for Partial Differential Equations*, SIAM, Philadelphia, 1984.

[26]    J.H. Wilkinson, *Rounding Errors in Algebraic Processes*, Prentice-
        Hall, Englewood Cliffs, N.J., 1963.

Division of Mathematics and Statistics
CSIRO
CANBERRA