

RICE UNIVERSITY

ROUNDING ERRORS IN THE SOLUTION OF
THE ONE DIMENSIONAL HEAT EQUATION
USING A GALERKIN TECHNIQUE

by

Catherine Gardner

A THESIS SUBMITTED
IN PARTIAL FULFILLMENT OF THE
REQUIREMENTS FOR THE DEGREE OF

MASTER OF ARTS

Thesis Director's signature:

A handwritten signature in cursive script, appearing to read "H. H. Rachford". The signature is written in dark ink and is positioned above a horizontal line.

Houston, Texas

May, 1970

ABSTRACT

ROUNDING ERRORS IN THE SOLUTION OF THE ONE DIMENSIONAL HEAT EQUATION USING A GALERKIN TECHNIQUE

Catherine Gardner

In solving the one dimensional heat equation, the solution is approximated using a Galerkin technique. Then, if the approximate solution U is required to lie in a Hermite space of piecewise polynomials of degree 3 based on a rectangular grid of mesh size h and if $\Delta t = c_1 h$ or $\Delta t = c_2 h^2$ where c_1, c_2 are constants and if as $\Delta t, h \rightarrow 0$ $\frac{\nu}{\Delta t h}$ is kept constant then, in either case, it can be shown that the computed solution \hat{U} satisfies $\|\hat{U} - U\|_{L^2} \leq a$ constant. The solution is computed using floating point base N -arithmetic with a τ digit mantissa and $\nu = N^{1-\tau}$.

ACKNOWLEDGMENTS

Of the many people who deserve my gratitude, I would especially like to thank Professor Henry H. Rachford, Jr., my thesis director, whose understanding guidance encouraged me in this undertaking. I am also indebted to Professors Graeme Fairweather and Howard L. Resnikoff, who served on my thesis committee, and to Mrs. Shirley Payne, who patiently typed the very difficult manuscript. I would also like to thank the United States Air Force and Rice University whose generous support in the form of an AFOSR grant and a Rice Fellowship enabled me to continue my education on the graduate level.

ROUNDING ERRORS IN THE SOLUTION OF THE
ONE DIMENSIONAL HEAT EQUATION USING
A GALERKIN TECHNIQUE

Introduction

We shall be concerned with solutions of the following boundary value problem

$$u_{xx} = u_t \quad \text{on } (0,1) \times (0,T] \quad (1a)$$

$$u(x,0) = u_0(x) \quad (1b)$$

$$u(1,t) = u(0,t) = 0 \quad (1c)$$

and we shall approximate the solution of (1) by a function

$$U(x,t) = \sum_{\ell=1}^N \alpha_{\ell}(t) v_{\ell}(x)$$

where $\{v_{\ell}(x)\}$ is a suitable collection of functions such that $v_{\ell}(x) \in C_0^1(0,1)$ and the collection $\{v_{\ell}(x)\}$ is linearly independent. If M is the span of the $\{v_{\ell}(x)\}$ then U is defined by

$$\langle U_{xx}, v \rangle = \langle U_t, v \rangle \quad \text{for all } v \in M \quad (2a)$$

$$\langle U_0, v \rangle = \langle u_0, v \rangle \quad t=0 \quad v \in M \quad (2b)$$

where $\langle W, v \rangle = \int_0^1 Wv dx$

It can be shown [1] that under reasonable hypotheses the approximation $U(x,t)$ and its first order space derivatives tend to the solution of (1a) with the error bounded by a

multiple of $h^{2m-1} + \Delta t^2$ where the approximate solution is required to lie in a Hermite space of piecewise polynomials of degree $2m-1$ based on a rectangular grid of mesh size h .

The system is actually an initial value problem for the set of linear ordinary differential equations

$$A \alpha'(t) + B \alpha(t) = 0 \quad (3a)$$

$$A \alpha(0) = (\langle u_0 v_1 \rangle, \dots, \langle u_0 v_N \rangle)^T \quad (3b)$$

where $\alpha = (\alpha_1, \dots, \alpha_N)^T$

and A and B are the positive definite matrices given by

$$A = a_{ij} = \langle v_i, v_j \rangle$$

$$B = b_{ij} = \langle v_{ix}, v_{jx} \rangle$$

We shall analyze the following difference scheme

$$A \frac{(\alpha_{n+1} - \alpha_n)}{\Delta t} = -B \frac{(\alpha_{n+1} + \alpha_n)}{2} \quad (4a)$$

$$\alpha_{n+1} = \left(A + \frac{\Delta t}{2} B \right)^{-1} \left(A - \frac{\Delta t}{2} B \right) \alpha_n \quad (4b)$$

The direct solution of the linear algebraic system (4) can be carried out by factoring $\left(A + \frac{\Delta t}{2} B \right) = LU$ where L is unit lower triangular and U is upper triangular then solving

$$L\delta = \left(A - \frac{\Delta t}{2} B \right) \alpha_n$$

$$U\alpha_{n+1} = \delta.$$

We shall be interested in the bound on the norm of the rounding error when the above scheme is carried out in floating point computation. Throughout we shall use a symbol with a caret, e.g. \hat{U}_n to denote the computed value of the exact element U_n . Thus, we seek a bound on $||\hat{U}_n - U_n||_{L^2}$

$$\text{where } U_n = \sum_{\ell=1}^N \alpha_{\ell}(t_n) v_{\ell}(x)$$

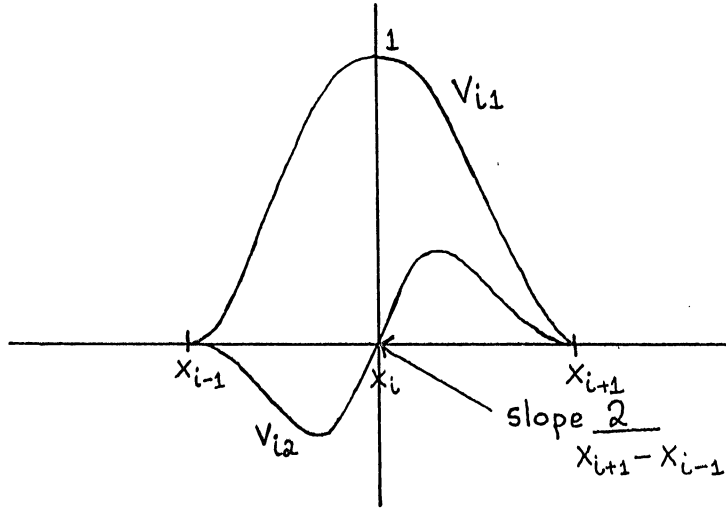
$$\text{and } \hat{U}_n = \sum_{\ell=1}^N \hat{\alpha}_{\ell}(t_n) v_{\ell}(x)$$

$$\begin{aligned} \text{Then } ||\hat{U}_n - U_n||_{L^2}^2 &= \int_0^1 \left(\sum_{\ell} \hat{\alpha}_{\ell}(t_n) v_{\ell}(x) - \sum_{\ell} \alpha_{\ell}(t_n) v_{\ell}(x) \right)^2 dx \\ &= \int_0^1 \left(\sum_{\ell} (\hat{\alpha}_{\ell}(t_n) - \alpha_{\ell}(t_n)) v_{\ell}(x) \right)^2 dx \\ &= (\hat{\alpha}_n - \alpha_n)^T A (\hat{\alpha}_n - \alpha_n) \\ &= ||\hat{\alpha}_n - \alpha_n||_A^2. \end{aligned}$$

And thus the bound we seek is on the A norm of $\hat{\alpha}_n - \alpha_n$. We shall be interested in two cases, i.e. that where we assume $\Delta t = ch$ and that where $\Delta t = ch^2$ where c is a constant, and will compare the results so obtained with those obtained in [2].

Choice of Basis Functions and Calculation
of Entries for the Matrices A & B

For our basis functions $\{v_\ell(x)\}$ we shall choose Hermite cubics of the following form: —



And if we order the basis functions as follows

$$v_{0,2} \quad v_{1,1} \quad v_{1,2} \quad v_{2,1} \quad v_{2,2} \quad \dots \quad v_{\frac{N}{2}-1,1} \quad v_{\frac{N}{2}-1,2} \quad v_{\frac{N}{2},2}$$

then the boundary conditions (1c) are satisfied and the matrix A is an $N \times N$ matrix with the following entries.

$$A = a_{jk} = \begin{cases} \left\langle v_{\frac{j-1}{2},2}, v_{\frac{k-1}{2},2} \right\rangle & j, k \text{ odd} \\ \left\langle v_{\frac{j}{2},1}, v_{\frac{k}{2},1} \right\rangle & j, k \text{ even} \\ \left\langle v_{\frac{j-1}{2},2}, v_{\frac{k}{2},1} \right\rangle & j \text{ odd, } k \text{ even} \\ \left\langle v_{\frac{j}{2},1}, v_{\frac{k-1}{2},2} \right\rangle & j \text{ even, } k \text{ odd} \end{cases}$$

v_{i1} and v_{i2} are given by the equations

$$v_{i1} = -2 \frac{(x - x_{i-1})^3}{(x_i - x_{i-1})^3} + 3 \frac{(x - x_{i-1})^2}{(x_i - x_{i-1})^2} \quad x \in [x_{i-1}, x_i]$$

$$= -2 \frac{(x_{i+1} - x)^3}{(x_{i+1} - x_i)^3} + 3 \frac{(x_{i+1} - x)^2}{(x_{i+1} - x_i)^2} \quad x \in [x_i, x_{i+1}]$$

$$v_{i2} = \frac{2}{x_{i+1} - x_{i-1}} \left[\frac{(x - x_{i-1})^3}{(x_i - x_{i-1})^2} - \frac{(x - x_{i-1})^2}{(x_i - x_{i-1})} \right] \quad x \in [x_{i-1}, x_i]$$

$$= \frac{-2}{x_{i+1} - x_{i-1}} \left[\frac{(x_{i+1} - x)^3}{(x_{i+1} - x_i)^2} - \frac{(x_{i+1} - x)^2}{(x_{i+1} - x_i)} \right] \quad x \in [x_i, x_{i+1}]$$

Since the support of v_{i1} , v_{i2} is contained in $[x_{i-1}, x_{i+1}]$ we

note immediately that $\langle v_{ik}, v_{jk} \rangle = 0 \quad |i-j| > 1 \quad k = 1, 2.$

And that $\langle v_{ik}, v_{jk} \rangle = \langle v_{jk}, v_{ik} \rangle \quad k = 1, 2.$

If we further assume that $x_{i+1} - x_i = h$ $\forall i$ then since v_{i1} is

even about x_i and v_{i2} is odd about x_i then $\langle v_{i1}, v_{i2} \rangle = 0$

and $\langle v_{i1}, v_{j2} \rangle = - \langle v_{i2}, v_{j1} \rangle$.

We also need to calculate $\| |A - \frac{\Delta t}{2} B| \|_2 \leq |a - \frac{\Delta t}{2} a'|$

$$+ 2|f - \frac{\Delta t}{2} f'| + 2|e - \frac{\Delta t}{2} e'|$$

$$\leq \frac{5\Delta t}{2h} + \frac{223}{210} h.$$

Calculation of a Bound for $\|\hat{\alpha}_n - \alpha_n\|_A$.

Recall that we are trying to solve equation (4b)

$$\left(A + \frac{\Delta t}{2} B\right) \alpha_{n+1} = \left(A - \frac{\Delta t}{2} B\right) \alpha_n$$

if $\Delta t = k$ and $\left(A + \frac{k}{2} B\right)^{-1} = Q$, $\left(A - \frac{k}{2} B\right) = M$

then (4b) reduces to $\alpha_{n+1} = QM \alpha_n$.

We seek a uniform bound for $\|\hat{\alpha}_n - \alpha_n\|_A$ where $\{\hat{\alpha}_n\}$ is the computed sequence.

The recursion for $\{\hat{\alpha}_n\}$ is given by

$$\hat{\alpha}_{n+1} = \hat{Q} \hat{d}_n \text{ where } \hat{d}_n = fl(M\hat{\alpha}_n) \text{ and } \hat{Q} = (1+R) Q$$

for some R which depends on the method of solving $Q^{-1}\alpha = \hat{d}$.

Now $\hat{d}_n = fl(M\hat{\alpha}_n) = M\hat{\alpha}_n + e_n$

$$\hat{\alpha}_{n+1} = (1+R) Q (M\hat{\alpha}_n + e_n)$$

$$\alpha_{n+1} = QM\alpha_n$$

Let $w_n = \hat{\alpha}_n - \alpha_n$ and by subtracting the recursions for $\hat{\alpha}_n$ and α_n we obtain

$$\hat{\alpha}_{n+1} - \alpha_{n+1} = (1+R) Q (M\hat{\alpha}_n + e_n) - QM\alpha_n$$

and simplifying we obtain

$$w_{n+1} = R\alpha_{n+1} + (1+R) Q (Mw_n + e_n) \tag{5}$$

In using $fl(\cdot)$ we have extended Wilkinson's notation [4] to include matrix and vector operations. Briefly, if a and b are scalars $fl(a \square b)$ means the floating point realization of the operation \square between them, where \square is $+$, $-$, \times , or \div . The relevant facts are $fl(a \square b) = (a \square b) (1 + \mu)$ for some μ satisfying $|\mu| \leq v = N^{1-\tau}$, where τ is the number of floating base N -digits in the mantissa, and

$$fl \left(\sum_{i=1}^p a_i b_i \right) = \sum_{i=1}^p a_i b_i (1 + \mu_i)$$

where $|\mu_i| \leq 1.06pv$.

We also assume that we can ignore any terms of order v^2 .

$$\begin{aligned} \text{Now } |e_n| &= 1.06 \times 6v \times |M| |\hat{\alpha}_n| \\ &= 6.36 v |M| |w_n + \alpha_n| \end{aligned}$$

since at most six products suffice for each element of $M\hat{\alpha}_n$.

$$\therefore e_n = 6.36 v \check{M} (w_n + \alpha_n) \text{ where } \check{m}_{ij} \leq \theta_{ij} m_{ij}.$$

$$|\theta_{ij}| \leq 1$$

$$\text{and } \|\check{M}\|_2 \leq \|M\|_2$$

Therefore substituting into (5) we have

$$\begin{aligned} w_{n+1} &= R\alpha_{n+1} + (1+R) Q (Mw_n + 6.36v \check{M} (w_n + \alpha_n)) \\ &= (1+R) (QM + 6.36v \check{QM})w_n + R\alpha_{n+1} \\ &\quad + (1+R) Q 6.36v \check{M}\alpha_n \end{aligned}$$

Thus

$$\begin{aligned} \|w_{n+1}\|_A \leq (1 + \|R\|_A) (\|QM\|_A + 6 \cdot 36\nu \|Q\check{M}\|_A) \|w_n\|_A \\ + \|R\|_A \|\alpha_{n+1}\|_A + (1 + \|R\|_A) 6 \cdot 36\nu \|Q\check{M}\|_A \|\alpha_n\|_A \end{aligned} \quad (6)$$

We shall now show that $\|QM\|_A \leq 1$ and that $\|\alpha_n\|_A \leq \|\alpha_0\|_A \leq$ a constant. Now $\|QM\|_A = \|A^{\frac{1}{2}}QMA^{-\frac{1}{2}}\|_2$ and if we let $\frac{k}{2}B = C$ then

$$\begin{aligned} A^{\frac{1}{2}}QMA^{-\frac{1}{2}} &= A^{\frac{1}{2}} (1+A^{-1}C)^{-1} (1-A^{-1}C) A^{-\frac{1}{2}} \\ &= A^{\frac{1}{2}} (1+A^{-1}C)^{-1} A^{-\frac{1}{2}} A^{\frac{1}{2}} (1-A^{-1}C) A^{-\frac{1}{2}} \\ &= (A^{\frac{1}{2}} (1+A^{-1}C) A^{-\frac{1}{2}})^{-1} (A^{\frac{1}{2}} (1-A^{-1}C) A^{-\frac{1}{2}}) \\ &= (1+A^{-\frac{1}{2}}CA^{-\frac{1}{2}})^{-1} (1-A^{-\frac{1}{2}}CA^{-\frac{1}{2}}) \end{aligned}$$

Thus $A^{\frac{1}{2}}QMA^{-\frac{1}{2}}$ is symmetric and therefore if λ_i is an eigenvalue of $A^{-\frac{1}{2}}CA^{-\frac{1}{2}}$ then since $A^{-\frac{1}{2}}CA^{-\frac{1}{2}}$ is similar to $A^{-1}C$, λ_i is both real and positive.

$$\therefore \|A^{\frac{1}{2}}QMA^{-\frac{1}{2}}\|_2 = \frac{1 - \lambda_i}{1 + \lambda_i} \leq 1$$

$$\|QM\|_A \leq 1$$

Now we note that if in equation (3a), $V = U_{n+1} + U_n$ and the same difference scheme is used, then

$$-\frac{1}{2} \int_0^1 \frac{\partial}{\partial x} \left[\left(\frac{U_{n+1} + U_n}{2} \right)^2 \right] = \frac{1}{\Delta t} \int_0^1 (U_{n+1} - U_n) (U_{n+1} + U_n)$$

or

$$\left\| U_{n+1} \right\|_{L^2}^2 - \left\| U_n \right\|_{L^2}^2 + \frac{\Delta t}{2} \left\| U_{n+1} + U_n \right\|_{H_0^1}^2 = 0$$

thus

$$\left\| U_{n+1} \right\|_{L^2}^2 \leq \left\| U_n \right\|_{L^2}^2$$

Now

$$U_{n+1} = \sum_{\ell=1}^N \alpha_{\ell}(t_{n+1}) v_{\ell}(x)$$

$$\begin{aligned} \left\| U_{n+1} \right\|_{L^2}^2 &= \int_0^1 \left(\sum_{\ell} \alpha_{\ell}(t_{n+1}) v_{\ell}(x) \right)^2 \\ &= \alpha_{n+1}^T A \alpha_{n+1} \end{aligned}$$

and similarly

$$\left\| U_{n+1} + U_n \right\|_{H_0^1}^2 = (\alpha_{n+1} + \alpha_n)^T B (\alpha_{n+1} + \alpha_n)$$

$$\left\| \alpha_{n+1} \right\|_A^2 - \left\| \alpha_n \right\|_A^2 + \frac{\Delta t}{2} \left\| \alpha_{n+1} + \alpha_n \right\|_B^2 = 0$$

$$\left\| \alpha_{n+1} \right\|_A \leq \left\| \alpha_n \right\|_A$$

and thus $\left\| \alpha_n \right\|_A \leq \left\| \alpha_0 \right\|_A \quad \forall n$

Now utilizing equation (2b) we can write

$$\langle U_0, U_0 \rangle = \langle u_0, U_0 \rangle$$

$$\begin{aligned} ||U_0||_{L_2}^2 &= \frac{1}{0} \left(\sum_{\ell} \alpha_{\ell}(t_0) v_{\ell}(x) \right)^2 = (\alpha_0^T A \alpha_0) = ||\alpha_0||_A^2 \\ &= \langle u_0, U_0 \rangle \leq ||u_0||_{L^2} ||U_0||_{L^2} \end{aligned}$$

$$||\alpha_0||_A \leq ||u_0||_{L^2} = \text{a constant which we shall}$$

denote β .

Thus, if P is a bound for $||R||_A$ then substituting into equation (6)

$$\begin{aligned} ||w_{n+1}||_A &\leq (1+P) (1+6 \cdot 36v ||\check{Q}\check{M}||_A) ||w_{n+1}||_A \\ &\quad + P\beta + (1+P) 6 \cdot 36v ||\check{Q}\check{M}||_A \beta \\ &\leq ||w_n||_A \phi_1 + \beta \phi_2 \end{aligned}$$

where $\phi_1 = (1+P) (1+6 \cdot 36v ||\check{Q}\check{M}||_A)$

$$\phi_2 = P + (1+P) 6 \cdot 36v ||\check{Q}\check{M}||_A$$

$$\therefore ||w_{n+1}||_A \leq \beta \phi_2 \frac{(\phi_1^n - 1)}{(\phi_1 - 1)}$$

Since $||\check{Q}\check{M}||_A = ||A^{\frac{1}{2}} \check{Q} M A^{-\frac{1}{2}}||_2$

$$\leq ||A^{\frac{1}{2}}||_2 ||A^{-\frac{1}{2}}||_2 ||Q||_2 ||\check{M}||_2$$

$$\leq ||A^{\frac{1}{2}}||_2 ||A^{-\frac{1}{2}}||_2 ||Q||_2 ||M||_2$$

it merely remains to obtain bounds on $||A^{-\frac{1}{2}}||_2$, $||Q||_2$

and the bound P on $||R||_A$

Calculation of $\|Q\|_2$ and $\|A^{-1}\|_2$

Recall that $Q = \left(A + \frac{k}{2} B\right)^{-1}$
 A^{-1} and Q^{-1} are positive definite and therefore

$$\|Q\|_2 = \frac{1}{\lambda_{\min} Q^{-1}} \quad \text{and} \quad \|A^{-1}\|_2 = \frac{1}{\lambda_{\min} A}$$

Consider first the matrix A which has the form discussed earlier. We shall show that if λ_1 is the minimum eigenvalue of A then

$$\lambda_1 \geq \frac{h}{5 \times 420}.$$

We shall define a series of matrices

$$A^{(i)} = a_{jk}^{(i)} = \begin{cases} \left\langle v_{\frac{j-1}{2}, 2}, v_{\frac{k-1}{2}, 2} \right\rangle^{(i)} & j, k \text{ odd} \\ \left\langle v_{\frac{j}{2}, 1}, v_{\frac{k}{2}, 1} \right\rangle^{(i)} & j, k \text{ even} \\ \left\langle v_{\frac{j-1}{2}, 2}, v_{\frac{k}{2}, 1} \right\rangle^{(i)} & j \text{ odd, } k \text{ even} \\ \left\langle v_{\frac{j}{2}, 1}, v_{\frac{k-1}{2}, 2} \right\rangle^{(i)} & j \text{ even, } k \text{ odd} \end{cases}$$

Where $\left\langle v_{eg}, v_{mg} \right\rangle^{(i)} = \frac{ih}{(i-1)h} \int v_{eg}, v_{mg}$.

Then clearly $\sum_{i=1}^{N/2} A^{(i)} = A$.

Now for $i \neq 1, N/2$ each $A^{(i)}$ is an $N \times N$ matrix with non-zero entries only in the $(i, i+3 : i, i+3)$ principal submatrix,

and for v_i the entries in the non-zero principal submatrix are identical. Thus, we can denote by D the $(i, i+3 : i, i+3)$ principal submatrix of each $A^{(i)}$. It is clear that D is positive definite, since it is a grammian and the basis functions are linearly independent over $[(i-1)h, ih]$. For boundary conditions satisfying (1c), $A^{(1)}$ is an $N \times N$ matrix with non-zero entries only in the $(1, 3 : 1, 3)$ principal submatrix, we shall denote this principal submatrix D' , and note that it is positive definite and also a principal submatrix of D . $A^{N/2}$ is also an $N \times N$ matrix with non-zero entries only in the $\left(\frac{N}{2} - 2, \frac{N}{2} : \frac{N}{2} - 2, \frac{N}{2}\right)$ principal submatrix, we shall denote this principal submatrix D'' and note that it also is positive definite and a principal submatrix of D .

Now
$$x^T A x = x^T \left(\sum_{i=1}^{N/2} A^{(i)} \right) x$$

and if we choose x to be the eigenvector corresponding to the minimum eigenvalue of A and partition x so that

$$x^T = \left(x_1, x_2, \dots, x_{\frac{N}{2}-1}, x_{\frac{N}{2}} \right)$$

where $x_i, i \neq 1, \frac{N}{2}$ is the vector such that $x_i^T = \begin{pmatrix} x_i^{(1)} & x_i^{(2)} \end{pmatrix}$ and

$x_1, x_{\frac{N}{2}}$ are vectors consisting of a single element. Then provided $x^T x = 1$

$$\lambda_{\min} A = (Ax, x) = \left(D' \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}, \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \right) + \left(D \begin{bmatrix} x_2 \\ x_3 \end{bmatrix}, \begin{bmatrix} x_2 \\ x_3 \end{bmatrix} \right) + \dots$$

$$+ \left(D'' \begin{bmatrix} x_{\frac{N}{2}-1} \\ x_{\frac{N}{2}} \end{bmatrix}, \begin{bmatrix} x_{\frac{N}{2}-1} \\ x_{\frac{N}{2}} \end{bmatrix} \right)$$

$$\geq \lambda_{\min D'} \left\| \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \right\|^2 + \lambda_{\min D} \left\| \begin{bmatrix} x_2 \\ x_3 \end{bmatrix} \right\|^2 + \dots$$

$$\dots + \lambda_{\min D''} \left\| \begin{bmatrix} x_{\frac{N}{2}-1} \\ x_{\frac{N}{2}} \end{bmatrix} \right\|^2$$

but $\lambda_{\min D'} \geq \lambda_{\min D}$ and $\lambda_{\min D''} \geq \lambda_{\min D}$ since both are principal submatrices of D .

$$\lambda_{\min A} \geq \lambda_{\min D} \left(\left\| \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \right\|^2 + \left\| \begin{bmatrix} x_2 \\ x_3 \end{bmatrix} \right\|^2 + \dots + \left\| \begin{bmatrix} x_{\frac{N}{2}-1} \\ x_{\frac{N}{2}} \end{bmatrix} \right\|^2 \right)$$

$$\geq \lambda_{\min D} \|\mathbf{x}\|^2 + \lambda_{\min D} \left(x_2^2 + x_3^2 + \dots + x_{\frac{N}{2}-1}^2 \right)$$

$$\geq \lambda_{\min D}$$

Now D is the following 4x4 matrix

$$D = \begin{bmatrix} \frac{a}{2} & -d & f & +e \\ -d & \frac{b}{2} & -e & g \\ f & -e & \frac{a}{2} & +d \\ +e & g & d & \frac{b}{2} \end{bmatrix}$$

where a, f, e, b, g have the values given previously and

$$d = \langle v_{11}, v_{12} \rangle^1 = \frac{h}{0} \begin{bmatrix} v_{11} & v_{12} \end{bmatrix} = -\frac{11h}{210}$$

$$\therefore D = \frac{h}{420} \begin{bmatrix} 156 & 22 & 54 & -13 \\ 22 & 4 & 13 & -3 \\ 54 & 13 & 156 & -22 \\ -13 & -3 & -22 & 4 \end{bmatrix}$$

Calculation shows that the characteristic polynomial of 420D is

$$\lambda^4 - 320\lambda^3h + 22617\lambda^2h^2 - 29792\lambda h^3 + 5145h^4$$

which, by calculations based on the change in sign of the

characteristic polynomial gives $\lambda_{\min} D \geq \frac{h}{5 \times 420}$

and $\lambda_{\min} A \geq \frac{h}{5 \times 420}$. Therefore, $\|A^{-1}\|_2 \leq \frac{2100}{h}$

We also seek a bound for the minimum eigenvalue of $A + \frac{k}{2}B$.

We could look at

$$A + \frac{k}{2}B = \sum_{i=1}^{N/2} A^{(i)} + \frac{k}{2} \sum_{i=1}^{N/2} B^{(i)}$$

However the non-zero principal submatrix of $B^{(i)}$ $i \neq 1, N/2$, which we shall denote E is singular since the basis functions when differentiated are linearly dependent over $[(i-1)h, ih]$. Thus the bound we would obtain in this fashion is no better than that we obtained for A .

We can, however, look at

$$A + \frac{k}{2}B = \sum_{i=1}^{N/2} \left(A^{(i)} + \frac{k}{2} B^{(i)} \right)$$

where $D+E$ is the following 4×4 matrix

$$D+E = \begin{bmatrix} \frac{a}{2} + \frac{ka'}{4} & d + \frac{k}{2} d' & f + \frac{k}{2} f' & e + \frac{k}{2} e' \\ d + \frac{k}{2} d' & \frac{b}{2} + \frac{kb'}{4} & -\left(e + \frac{k}{2} e'\right) & g + \frac{k}{2} g' \\ f + \frac{k}{2} f' & -\left(e + \frac{k}{2} e'\right) & \frac{a}{2} + \frac{ka'}{4} & -\left(d + \frac{k}{2} d'\right) \\ e + \frac{k}{2} e' & g + \frac{k}{2} g' & -\left(d + \frac{k}{2} d'\right) & \frac{b}{2} + \frac{kb'}{4} \end{bmatrix}$$

where $a, a', f, f', e, e', b, b', g, g'$ and d have the same values as before and $d' = -e'$. Calculation of the characteristic polynomial of $420 (D+E)$ gives

$$\begin{aligned} & \lambda^4 - \lambda^3 \left(320h + 560 \frac{k}{h} \right) + \lambda^2 \left(22,617h^2 + 126,770k + 27,195 \frac{k^2}{h^2} \right) \\ & - \lambda \left(29,792h^2 + 1,296,540kh + 5,938,840 \frac{k^2}{h} + 308,700 \frac{k^3}{h^3} \right) \\ & + 5145h^4 + 143,724kh^2 + 16,053,200k^2 + 61,593,058 \frac{k^3}{h^2} \end{aligned}$$

Now if $k = c_1 h$ then provided $h < 1$ and $c_1 \geq 1$

then $\lambda_{\min(D+E)} \geq \frac{198}{420} h$

$$\lambda_{\min(A + \frac{k}{2} B)} \geq \frac{198}{420} h \quad \text{and} \quad \|Q\|_2 \leq \frac{420}{198h}$$

But if $k = c_2 h^2$ then, for $1 \leq c_2 < 10$, $A + \frac{k}{2} B$ is a strictly diagonally dominant symmetric matrix and a bound for the minimum eigenvalue is given by Gesgorin.

$$\text{Minimum eigenvalue} \left(A + \frac{k}{2} B \right) \geq \frac{28h}{420}$$

$$\therefore \|Q\|_2 \leq \frac{420}{28h} .$$

Calculation of P

Recall that P is a bound for $\|R\|_A$
where $\hat{Q} = (I+R)Q$ $Q = (A + \frac{k}{2} B)^{-1}$

Now suppose $Gw = d$ and invert G by Gaussian elimination. The Gaussian algorithm expresses the set of equations $Gw = d$ in the form

$$LUw = d \quad (8)$$

where L is unit lower triangular and U is upper triangular.

The computed $L(\hat{L})$ and $U(\hat{U})$ satisfy

$$\hat{L}\hat{U} = G + E$$

hence if we solve equation (8) without further rounding error we obtain the solution of the system

$$(G+E)w = d.$$

Now equation (8) can be solved in two stages

$$\hat{L}y = d$$

$$\hat{U}w = y$$

and in practice we obtain

$$(\hat{L} + d\hat{L})\hat{y} = d$$

$$(\hat{U} + d\hat{U})\hat{w} = \hat{y}$$

or $(G + E + \hat{L}\hat{d}\hat{U} + d\hat{L}\hat{U} + d\hat{L}d\hat{U})\hat{w} = d$

or $(G + K)\hat{w} = d$

where $K = E + d\hat{L}\hat{U} + \hat{L}\hat{d}\hat{U} + d\hat{L}d\hat{U}$

and $\|K\|_{\infty} \leq \|E\|_{\infty} + \|d\hat{L}\hat{U}\|_{\infty} + \|\hat{L}\hat{d}\hat{U}\|_{\infty} + \|d\hat{L}d\hat{U}\|_{\infty}$

Now
$$\begin{aligned}\hat{w} - w &= \left((G+K)^{-1} - G^{-1} \right) d \\ &= \left((G+K)^{-1}G - I \right) w\end{aligned}$$

but $G + K = G(I + G^{-1}K)$

and $(G+K)^{-1} = (I + G^{-1}K)^{-1} G^{-1}$

$$\left((G+K)^{-1}G - I \right) = \left(I + G^{-1}K \right)^{-1} - I$$

and $\left(I + G^{-1}K \right)$ is non-singular if $\|G^{-1}K\|_2 \leq 1$.

let $G^{-1}K = F$

then
$$\frac{\|\hat{w} - w\|_2}{\|w\|_2} = \|(I+F)^{-1} - I\|_2 \leq \frac{\|F\|_2}{1 - \|F\|_2}$$

$$\leq \frac{\|G^{-1}\|_2 \|K\|_2}{1 - \|G^{-1}\|_2 \|K\|_2}$$

but $R = \hat{Q}Q^{-1} - I$

substituting $Q^{-1} = G$ and $\hat{Q} = (G+K)^{-1}$

then $R = (G+K)^{-1}G - I$

$$\therefore \frac{||\hat{w}-w||_2}{||w||_2} = ||R||_2 \leq \frac{||Q||_2 ||K||_2}{1-||Q||_2 ||K||_2}$$

$$\text{now } ||K||_2 \leq \sqrt{||K||_\infty ||K||_1}$$

$$\text{and } ||R||_2 \leq \frac{||Q||_2 \sqrt{||K||_\infty ||K||_1}}{1-||Q||_2 \sqrt{||K||_\infty ||K||_1}}$$

$$\text{provided } ||Q||_2 \sqrt{||K||_\infty ||K||_1} < 1$$

$$||R||_A \leq ||A^{-\frac{1}{2}}||_2 ||R||_2 ||A^{\frac{1}{2}}||_2$$

$$\leq \frac{||A^{-\frac{1}{2}}||_2 ||Q||_2 \sqrt{||K||_\infty ||K||_1} ||A^{\frac{1}{2}}||_2}{1-||Q||_2 \sqrt{||K||_\infty ||K||_1}}$$

(9)

Error Analysis of Gaussian Elimination for a
Positive Definite Matrix $A = a_{ij}$ where

for $i \leq j$	$a_{ij} = 0$	$i, \text{ even} \quad j-i \geq 4$ $i, \text{ odd} \quad j-i \geq 3.$
$i > j$	$a_{ij} = 0$	$i, \text{ even} \quad i-j \geq 3$ $i, \text{ odd} \quad i-j \geq 4$

The solution of a set of equations and the inversion of a matrix by Gaussian elimination is based on the triangularization of the matrix. If we denote the original set of equations by $A^{(1)}x = b$

then $N-1$ equivalent sets

$$A^{(r)}x = b^{(r)} \quad (r = 2 \dots N)$$

are produced the matrix $A^{(N)}$ of the final set being of upper triangular form. In general $A^{(r)}$ is already of triangular form as regards its first r rows and it has a square matrix (which for convenience we shall denote $B^{(r)}$) of non-zero elements in the bottom right-hand corner. This square matrix is of order $N+1-r$. The matrix $A^{(r+1)}$ is derived from $A^{(r)}$ by subtracting a multiple m_{ir} of the r^{th} row from the i^{th} row for values of i from $r+1$ to N . The multipliers m_{ir} are defined by

$$m_{ir} = \frac{a_{ir}^{(r)}}{a_{rr}^{(r)}}$$

If we denote the elements of E by e_{ij} and note that E always has non-zero entries where A has non-zero entries, then the history of the (i,j) element differs according as $i \leq j$ or $i > j$.

i) $i \leq j$.

The element is modified in each transformation until $A^{(i)}$ is obtained after which it remains constant. For r even

$$\hat{a}_{r,r}^{(r)} = a_{r,r}^{(1)} - \hat{m}_{r,r-1} \hat{a}_{r-1,r}^{(r-1)} - \hat{m}_{r,r-2} \hat{a}_{r-2,r}^{(r-2)} + e_{r,r}$$

$$\hat{a}_{r,r+1}^{(r)} = a_{r,r+1}^{(1)} - \hat{m}_{r,r-1} \hat{a}_{r-1,r+1}^{(r-1)} - \hat{m}_{r,r-2} \hat{a}_{r-2,r+1}^{(r-2)} + e_{r,r+1}$$

$$\hat{a}_{r,r+2}^{(r)} = a_{r,r+2}^{(1)} + e_{r,r+2}$$

$$\hat{a}_{r,r+3}^{(r)} = a_{r,r+3}^{(1)} + e_{r,r+3}$$

For r odd

$$\hat{a}_{r,r}^{(r)} = a_{r,r}^{(1)} - \hat{m}_{r,r-3} \hat{a}_{r-3,r}^{(r-3)} - \hat{m}_{r,r-2} \hat{a}_{r-2,r}^{(r-2)} - \hat{m}_{r,r-1} \hat{a}_{r-1,r}^{(r-1)} + e_{rr}$$

$$\hat{a}_{r,r+1}^{(r)} = a_{r,r+1}^{(1)} - \hat{m}_{r,r-1} \hat{a}_{r-1,r+1}^{(r-1)} + e_{r,r+1}$$

$$\hat{a}_{r,r+2}^{(r)} = a_{r,r+2}^{(1)} - \hat{m}_{r,r-1} \hat{a}_{r-1,r+2}^{(r-1)} + e_{r,r+2}$$

The computed multiplier \hat{m}_{ij} satisfies

$$\hat{m}_{ij} = \frac{\hat{a}_{ij}(j)}{\hat{a}_{jj}} + \eta_{ij}$$

where η_{ij} is the rounding error in the division. The equations are therefore: —

For r even

$$e_{r+1,r} = a_{r+1,r}^{(1)} - \hat{m}_{r+1,r-2} \hat{a}_{r-2,r}^{(r-2)} - \hat{m}_{r+1,r-1} \hat{a}_{r-1,r}^{(r-1)} - \hat{m}_{r+1,r} \hat{a}_{r,r}^{(r)}$$

$$e_{r+2,r} = a_{r+2,r}^{(1)} - \hat{m}_{r+2,r} \hat{a}_{r,r}^{(r)}$$

$$e_{r+3,r} = a_{r+3,r}^{(1)} - \hat{m}_{r+3,r} \hat{a}_{r,r}^{(r)}$$

For r odd

$$e_{r+1,r} = a_{r+1,r}^{(1)} - \hat{m}_{r+1,r-1} \hat{a}_{r-1,r}^{(r-1)} - \hat{m}_{r+1,r} \hat{a}_{r,r}^{(r)}$$

$$e_{r+2,r} = a_{r+2,r}^{(1)} - \hat{m}_{r+2,r-1} \hat{a}_{r-1,r}^{(r-1)} - \hat{m}_{r+2,r} \hat{a}_{r,r}^{(r)}$$

Therefore we see that for r even

$$e_{r,r} = \epsilon_{r,r}^{(r)} + \epsilon_{r,r}^{(r-1)}$$

$$e_{r,r+1} = \epsilon_{r,r+1}^{(r)} + \epsilon_{r,r+1}^{(r-1)}$$

$$e_{r,r+2} = 0$$

$$e_{r,r+3} = 0$$

$$e_{r+1,r} = \epsilon_{r+1,r}^{(r-1)} + \epsilon_{r+1,r}^{(r)} + \epsilon_{r+1,r}^{(r+1)}$$

$$e_{r+2,r} = \epsilon_{r+2,r}^{(r+1)}$$

$$e_{r+3,r} = \epsilon_{r+3,r}^{(r+1)}$$

For r odd

$$e_{r,r} = \epsilon_{r,r}^{(r)} + \epsilon_{r,r}^{(r-1)} + \epsilon_{r,r}^{(r-2)}$$

$$e_{r,r+1} = \epsilon_{r,r+1}^{(r)}$$

$$e_{r,r+2} = \epsilon_{r,r+2}^{(r)}$$

$$e_{r+1,r} = \epsilon_{r+1,r}^{(r)} + \epsilon_{r+1,r}^{(r+1)}$$

$$e_{r+2,r} = \epsilon_{r+2,r}^{(r)} + \epsilon_{r+2,r}^{(r+1)}$$

Where $\epsilon_{ij}^{(k)} = \hat{a}_{ij}^{(k)} - \hat{a}_{ij}^{(k-1)} + \hat{m}_{i,k-1} \hat{a}_{k-1,j}^{(k-1)}$

$$k = 2 \dots i$$

$$\epsilon_{ij}^{j+1} = a_{jj}^{(j)} \eta_{ij}$$

Now in floating point the computed $\hat{a}_{ij}^{(k)}$ is defined by: —

$$\begin{aligned} \hat{a}_{ij}^{(k)} &= fl \left[\hat{a}_{ij}^{(k-1)} - fl \left(\hat{m}_{i,k-1} \hat{a}_{k-1,j}^{(k-1)} \right) \right] \\ &= \frac{\left[\hat{a}_{ij}^{(k-1)} - \hat{m}_{i,k-1} \hat{a}_{k-1,j}^{(k-1)} (1 + \epsilon_1) \right]}{1 + \epsilon_2} \quad |\epsilon_1|, |\epsilon_2| \leq \nu \end{aligned}$$

Now
$$\epsilon_{ij}^{(k)} = \hat{a}_{ij}^{(k)} - \left[\hat{a}_{ij}^{(k-1)} - \hat{m}_{i,k-1} \hat{a}_{k-1,j}^{(k-1)} \right]$$

$$\hat{a}_{ij}^{(k-1)} = \hat{a}_{ij}^{(k)} (1+\epsilon_2) + \hat{m}_{i,k-1} \hat{a}_{k-1,j}^{(k-1)} (1+\epsilon_1)$$

$$\begin{aligned} \epsilon_{ij}^{(k)} &= \hat{a}_{ij}^{(k)} - \left[\hat{a}_{ij}^{(k)} (1+\epsilon_2) + \hat{m}_{i,k-1} \hat{a}_{k-1,j}^{(k-1)} (1+\epsilon_1) \right. \\ &\quad \left. - \hat{m}_{i,k-1} \hat{a}_{k-1,j}^{(k-1)} \right] \end{aligned}$$

$$= - \hat{a}_{ij}^{(k)} \epsilon_2 + \hat{m}_{i,k-1} \hat{a}_{k-1,j}^{(k-1)} \epsilon_1$$

$$|\epsilon_{ij}^{(k)}| \leq \left| \hat{a}_{ij}^{(k)} + \hat{m}_{i,k-1} \hat{a}_{k-1,j}^{(k-1)} \right| \nu$$

Now we show later that $\left| \hat{a}_{ij}^{(k)} \right| \leq 1$ for all k

and now we consider $\hat{m}_{i,k-1} \hat{a}_{k-1,j}^{(k-1)}$.

$$\hat{m}_{i,k-1} \hat{a}_{k-1,j}^{(k-1)} = \frac{\hat{a}_{i,k-1}^{(k-1)} \hat{a}_{k-1,j}^{(k-1)}}{\hat{a}_{k-1,k-1}^{(k-1)}} \quad (1+\epsilon) \quad |\epsilon| \leq \nu$$

if $i \leq j$ then k has values $k = 2 \dots i$

if $i < j$ then k has values $k = 2 \dots j$

and for either case $\hat{a}_{i,k-1}^{(k-1)}$, $\hat{a}_{k-1,j}^{(k-1)}$ and $\hat{a}_{k-1,k-1}^{(k-1)}$ are all

elements of $\hat{B}^{(k-1)}$. Since $\hat{B}^{(k-1)}$ is positive definite, then

$$\left(\hat{a}_{i,k-1}^{(k-1)} \right)^2 \leq \hat{a}_{ii}^{(k-1)} \hat{a}_{k-1,k-1}^{(k-1)}$$

and
$$\left(\hat{a}_{k-1,j}^{(k-1)} \right)^2 \leq \hat{a}_{jj}^{(k-1)} \hat{a}_{k-1,k-1}^{(k-1)}$$

$$\hat{m}_{i,k-1} \hat{a}_{k-1,j}^{(k-1)} \leq \left(\sqrt{\hat{a}_{ii}^{(k-1)} \hat{a}_{jj}^{(k-1)}} \right) (1+v) \leq 1+v$$

$$|\epsilon_{ij}^{(k)}| \leq v + v (1+v)$$

$$\leq 2v$$

This result applies to all $\epsilon_{ij}^{(k)}$ except $\epsilon_{ij}^{(j+1)}$ ($i > j$)

For $\epsilon_{ij}^{(j+1)}$ we have

$$\hat{m}_{ij} = fl \left(\frac{\hat{a}_{ij}^{(j)}}{\hat{a}_{jj}^{(j)}} \right) = \frac{\hat{a}_{ij}^{(j)}}{\hat{a}_{jj}^{(j)}} (1+\epsilon)$$

$$\eta_{ij} = \frac{\hat{a}_{ij}^{(j)}}{\hat{a}_{jj}^{(j)}} \epsilon$$

$$\left| \epsilon_{ij}^{(j+1)} \right| = \left| \frac{\hat{a}_{jj}^{(j)} \hat{a}_{ij}^{(j)}}{\hat{a}_{jj}^{(j)}} \epsilon \right| = \left| \hat{a}_{ij}^{(j)} \epsilon \right| \leq v < 2v$$

and therefore we need not give this element special treatment.

Combining all these results we have that: —

$$|E| \leq v \begin{bmatrix} 0 & 0 & 0 & & & & & & & & \\ 1 & 1 & 1 & & & & & & & & \\ 1 & 2 & 2 & 1 & 1 & & & & & & \\ & 1 & 2 & 2 & 2 & & & & & & \\ 1 & 2 & 3 & 3 & 1 & 1 & & & & & \\ & & 1 & 2 & 2 & 2 & & & & & \\ & & & 1 & 2 & 3 & 3 & 1 & 1 & & \\ & & & & & & & \text{etc.} & & & \end{bmatrix}$$

and therefore that $\|E\|_{\infty} \leq 22v$

$$\|E\|_1 \leq 20v .$$

To show that $\hat{B}^{(r)}$ is positive definite we write

$\hat{B}^{(r)} = B^{(r)} + J^{(r)}$ where $J^{(r)}$ is the matrix of errors introduced up to the r^{th} stage of the elimination, i.e., $J^{(r)}$ has as its elements $e_{ij}^{(k)}$, $k = 2 \dots r$. Thus it can be seen that $J^{(r)}$ is symmetric and minimum eigenvalue $\hat{B}^{(r)} \geq$ minimum eigenvalue $B^{(r)} +$ minimum eigenvalue $J^{(r)}$.

Now

$$|J^{(r)}| \leq v \begin{bmatrix} \text{r}^{\text{th}} \text{ column} & & & & & & & & & & \\ & \downarrow & & & & & & & & & \\ & & 0 & & & & & & & & \\ & & & 22 & & & & & & & \\ & & & & 22 & & & & & & \\ & & & & & & & & & & \\ & & & & & & & & & & \\ 0 & & & & & & & & & & \end{bmatrix} \quad \text{or } \leq v \begin{bmatrix} \text{r}^{\text{th}} \text{ column} & & & & & & & & & & \\ & \downarrow & & & & & & & & & \\ & & 0 & & & & & & & & \\ & & & 311 & & & & & & & \\ & & & & 111 & & & & & & \\ & & & & & & & & & & \\ & & & & & & & & & & \\ 0 & & & & & & & & & & \end{bmatrix}$$

depending on whether r is even or odd. Thus, for any r , the minimum eigenvalue of $J^{(r)} \geq -5v$. Therefore, provided minimum eigenvalue $B^{(r)} > 5v$ for any r , then minimum eigenvalue of $\hat{B}^{(r)} > 0$ and thus $\hat{B}^{(r)}$ is positive definite.

$$\begin{aligned} \text{Now } \hat{a}_{ii}^{(2)} &= f_l \left(a_{ii}^{(1)} - f_l \left(m_{i1} a_{1j}^{(1)} \right) \right) \\ &= \left[a_{ii}^{(1)} - m_{i1} a_{1j}^{(1)} (1+\epsilon_1) \right] (1+\epsilon_2) \\ &\leq a_{ii}^{(1)} (1+\epsilon_2) \end{aligned}$$

$$\text{and in general } \hat{a}_{ii}^{(k)} \leq \hat{a}_{ii}^{(k-1)} (1+\epsilon_2) \quad |\epsilon_2| \leq v$$

Now the largest element in a positive definite matrix lies on the diagonal and since we operate on each row at most 3 times, if we scale originally so that $a_{ii}^{(1)} \leq (1+3v)$ then we can be sure that $\hat{a}_{ij}^{(k)} \leq 1$ for any k .

We shall now turn our attention to the back solution and consider the solution of a lower triangular set where the matrix of coefficients does not necessarily have unit elements on the diagonal, but where if we are solving $Ly = d$

$$L = \hat{m}_{ij} \text{ then } \hat{m}_{ij} = 0 \quad \begin{array}{ll} i \text{ even} & i-j \geq 3 \\ i \text{ odd} & i-j \geq 4 \end{array}$$

for r even

$$\hat{Y}_r = fl \left(\frac{d_r - \hat{m}_{r,r-2} \hat{Y}_{r-2} - \hat{m}_{r,r-1} \hat{Y}_{r-1}}{\hat{m}_{rr}} \right)$$

$$\hat{Y}_r = \frac{d_r - \hat{m}_{r,r-2} \hat{Y}_{r-2} (1+E_{r,r-2}) - \hat{m}_{r,r-1} \hat{Y}_{r-1} (1+E_{r,r-1})}{\hat{m}_{rr} (1+\eta_r) (1+\epsilon_r)}$$

for r odd

$$\hat{Y}_r = \frac{d_r - \hat{m}_{r,r-3} \hat{Y}_{r-3} (1+E_{r,r-3}) - \hat{m}_{r,r-2} \hat{Y}_{r-2} (1+E_{r,r-2}) - \hat{m}_{r,r-1} \hat{Y}_{r-1} (1+E_{r,r-1})}{\hat{m}_{rr} (1+\eta_r) (1+\epsilon_r)}$$

where $|\eta_r| \leq \nu$ $|\epsilon_r| \leq \nu$

$$E_{ri} \leq (r+1-i) 1.06\nu$$

It is more convenient to express the above in the following form: —

for r even

$$d_r = \hat{Y}_r \hat{m}_{rr} (1+E_{rr}) - \hat{m}_{r,r-2} \hat{Y}_{r-2} (1+E_{r,r-2}) - \hat{m}_{r,r-1} \hat{Y}_{r-1} (1+E_{r,r-1})$$

and for r odd

$$d_r = \hat{Y}_r \hat{m}_{rr} (1+E_{rr}) + \hat{m}_{r,r-3} \hat{Y}_{r-3} (1+E_{r,r-3}) + \hat{m}_{r,r-2} \hat{Y}_{r-2} (1+E_{r,r-2}) + \hat{m}_{r,r-1} \hat{Y}_{r-1} (1+E_{r,r-1})$$

where $1 + E_{rr} = \left(|+\eta_r| \right) \left(1 + \epsilon_r \right)$ $|E_{rr}| \leq 2v$

and in particular

$$|E_{r,r-2}| \leq \left(r+1-(r-2) \right) 1.06v = 3 \times 1.06v$$

$$|E_{r,r-1}| \leq \left(r+1-(r-1) \right) 1.06v = 2 \times 1.06v$$

$$|E_{r,r-3}| \leq \left(r+1-(r-3) \right) 1.06v = 4 \times 1.06v$$

Thus the computed vector is the exact solution of

$$(\hat{L} + d\hat{L})\hat{y} = d$$

where $d\hat{L}$ is bounded as shown below

$$|d\hat{L}| \leq 1.06v \left[\begin{array}{ccccccc} |\hat{m}_{11}| & & & & & & \\ 2|\hat{m}_{21}| & 2|\hat{m}_{22}| & & & & & \\ 3|\hat{m}_{31}| & 2|\hat{m}_{32}| & 2|\hat{m}_{33}| & & & & \\ & 3|\hat{m}_{42}| & 2|\hat{m}_{43}| & 2|\hat{m}_{44}| & & & \\ & 4|\hat{m}_{52}| & 3|\hat{m}_{52}| & 2|\hat{m}_{54}| & 2|\hat{m}_{55}| & & \\ & & & & & & \text{etc.} \end{array} \right]$$

If the diagonal elements of L are unity then the diagonal elements of $|d\hat{L}|$ are unity.

An analysis of $d\hat{U}$ for $(\hat{U}+d\hat{U})\hat{w} = y$ can be computed in an exactly analogous fashion and

$$|\hat{d}\hat{U}|_{\leq 1.06\nu} \left[\begin{array}{cccc} 2 \left| a_{11}^{(1)} \right| & 2 \left| a_{12}^{(1)} \right| & 3 \left| a_{13}^{(1)} \right| & \\ & 2 \left| a_{22}^{(2)} \right| & 2 \left| a_{23}^{(2)} \right| & 3 \left| a_{24}^{(2)} \right| & 4 \left| a_{25}^{(2)} \right| & \dots \\ & & & & & \text{etc.} \end{array} \right]$$

We now need a bound on $\|\hat{L}\hat{d}\hat{U}\|_{\infty}$

$$\hat{L}\hat{d}\hat{U} = 1.06H\nu$$

where for i even $h_{ij} = 0$ $j < i$ $i - j \geq 3$
 $i \leq j$ $j - i \geq 4$

$$\left| h_{i,i-2} \right| \leq 2 \left| \hat{m}_{i,i-2} \hat{a}_{i-2,i-2}^{(i-2)} \right|$$

$$\left| h_{i,i-1} \right| \leq 2 \left| \hat{m}_{i,i-2} \hat{a}_{i,i-1} \right| + 2 \left| \hat{m}_{i,i-1} \hat{a}_{i-1,i-1}^{(i-1)} \right|$$

$$\left| h_{ii} \right| \leq 3 \left| \hat{m}_{i,i-2} \hat{a}_{i-2,i}^{(i-2)} \right| + 2 \left| \hat{m}_{i,i-1} \hat{a}_{i-1,i}^{(i-1)} \right| + 2 \left| \hat{m}_{ii} \hat{a}_{ii}^{(i)} \right|$$

$$\left| h_{i,i+1} \right| \leq 4 \left| \hat{m}_{i,i-2} \hat{a}_{i-2,i+1}^{(i-2)} \right| + 3 \left| \hat{m}_{i,i-1} \hat{a}_{i-1,i+1}^{(i-1)} \right| + 2 \left| \hat{m}_{ii} \hat{a}_{i,i+1}^{(i)} \right|$$

$$\left| h_{i,i+2} \right| \leq 3 \left| \hat{m}_{ii} \hat{a}_{i,i+2}^{(i)} \right|$$

$$\left| h_{i,i+3} \right| \leq 4 \left| \hat{m}_{ii} \hat{a}_{i,i+3}^{(i)} \right|$$

for i odd $h_{ij} = 0$ $j < i$ $i - j \geq 4$
 $i \leq j$ $j - i \geq 3$

$$|h_{i,i-3}| \leq 2 \left| \hat{m}_{i,i-3} \hat{a}_{i-3,i-3}^{(i-3)} \right|$$

$$|h_{i,i-2}| \leq 2 \left| \hat{m}_{i,i-3} \hat{a}_{i-3,i-2}^{(i-3)} \right| + 2 \left| \hat{m}_{i,i-2} \hat{a}_{i-2,i-2}^{(i-1)} \right|$$

$$|h_{i,i-1}| \leq 3 \left| \hat{m}_{i,i-3} \hat{a}_{i-3,i-1}^{(i-3)} \right| + 2 \left| \hat{m}_{i,i-2} \hat{a}_{i-2,i-1}^{(i-2)} \right| \\ + 2 \left| \hat{m}_{i,i-1} \hat{a}_{i-1,i-1}^{(i-1)} \right|$$

$$|h_{ii}| \leq 4 \left| \hat{m}_{i,i-3} \hat{a}_{i-3,i}^{(i-3)} \right| + 3 \left| \hat{m}_{i,i-2} \hat{a}_{i-2,i}^{(i-2)} \right| \\ + 2 \left| \hat{m}_{i,i-1} \hat{a}_{i-1,i}^{(i-1)} \right| + 2 \left| \hat{m}_{ii} \hat{a}_{ii}^{(i)} \right|$$

$$|h_{i,i+1}| \leq 3 \left| \hat{m}_{i,i-1} \hat{a}_{i-1,i+1}^{(i-1)} \right| + 2 \left| \hat{m}_{ii} \hat{a}_{i,i+1}^{(i)} \right|$$

$$|h_{i,k+2}| \leq 4 \left| \hat{m}_{i,i-1} \hat{a}_{i-1,i+2}^{(i-1)} \right| + 3 \left| \hat{m}_{ii} \hat{a}_{i,i+2}^{(i)} \right|$$

By the arguments used in calculating $\|E\|_{\infty}$ we see that

$$\left| \hat{m}_{k\ell} a_{\ell m}^{(\ell)} \right| \leq 1$$

and therefore $\|\hat{L}\hat{d}\hat{U}\|_{\infty} \leq \|\hat{L}\hat{d}\hat{U}\|_{\infty} \leq 1.06 \times 35v$

$$\|\hat{L}\hat{d}\hat{U}\|_1 \leq 1.06 \times 39v$$

In exactly the same way we can calculate bounds on $||\hat{dL}\hat{U}||_\infty$ and $||\hat{dLd}\hat{U}||_\infty$ and we obtain

$$||\hat{dL}\hat{U}||_\infty \leq 1.06 \times 36v \quad ||\hat{dL}\hat{U}||_1 \leq 1.06 \times 32v$$

$$||\hat{dLd}\hat{U}||_\infty \leq (1.06v)^2 101 \quad ||\hat{dLd}\hat{U}||_1 \leq (1.06v)^2 90$$

Therefore $||K||_\infty \leq 2 \times 11v + 1.06 \times 35v + 1.06 \times 36v + (1.06v)^2 101$

$$||K||_1 \leq 2 \times 10v + 1.06 \times 39v + 1.06 \times 32v + (1.06v)^2 90$$

ignoring terms of order v^2

$$||K||_2 \leq \sqrt{||K||_\infty ||K||_1} \leq 97v$$

We can now obtain a bound P on $||R||_A$. If we write $||A^{-\frac{1}{2}}||_2 \quad ||A^{\frac{1}{2}}||_2 \leq c$ and $||Q||_2 \leq \frac{c'}{h}$ where c and c' are both constant. Then using $||K||_2 \leq 97v$ and substituting into equation (9)

$$||R||_A \leq \frac{cc'97v}{h} \left(1 - \frac{c'97v}{h} \right)^{-1}$$

$$\sim \frac{cc'97v}{h} \left(1 + \frac{c'97v}{h} \dots \dots \right)$$

$$\sim \frac{cc'97v}{h} + O\left(\frac{v}{h}\right)^2$$

$$\therefore P = \frac{cc'97v}{h}$$

We are now in a position to substitute into equation (7). We shall obtain results showing the asymptotic behaviour of the error as $k, h \rightarrow 0$ provided $\frac{v}{kh}$ is kept constant.

Consider first the case when $k = c_1 h$

then $\| |M| \|_A \leq pc_1 + p'h$ p, p' constant.

$$\phi_1 \leq \left(\frac{1 + cc'97v}{h} \right) \left(\frac{1 + 6 \cdot 36vcc'}{h} (pc_1 + p'h) \right)$$

$$\leq \frac{1 + cc'97v}{h} + \frac{6 \cdot 36cc'pc_1v}{h} + O\left(\frac{v}{h}\right)^2$$

$$\sim 1 + \frac{cc'v}{h} (97 + 6 \cdot 36pc_1)$$

$$\phi_2 \leq \frac{cc'97v}{h} + \left(\frac{1 + cc'97v}{h} \right) \frac{6 \cdot 36cc'v}{h} (pc_1 + p'h)$$

$$\leq \frac{cc'97v}{h} + \frac{6 \cdot 36cc'pc_1v}{h} + O\left(\frac{v}{h}\right)^2$$

$$\sim \frac{cc'v}{h} (97 + 6 \cdot 36pc_1)$$

Therefore

$$\phi_1 \sim 1 + \frac{k'v}{h}$$

$$\phi_2 \sim \frac{k'v}{h} \quad \text{where } k' = cc' (97 + 6 \cdot 36pc_1)$$

$$\|w_{n+1}\|_A \leq \frac{\beta k' \nu}{h} \frac{\left(\left(1 + \frac{k' \nu}{h}\right)^n - 1 \right)}{k' \frac{\nu}{h}}$$

$$\leq \beta \left(\left(1 + \frac{k' \nu}{h}\right)^n - 1 \right)$$

$$\leq \frac{\beta n k' \nu}{h} e^{\frac{n k' \nu}{h}}$$

But $n = \frac{T}{k}$ and if $\frac{\nu}{kh} = K$ a constant

$$\text{then } \|w_{n+1}\|_A = \|\hat{U}_{n+1} - U_{n+1}\|_{L^2} \leq \beta T k' K e^{TK k'} \leq K''$$

where K'' is a constant

Now consider the case where $k = c_2 h^2$, then

$$\|Q\|_2 \leq \frac{c''}{h} \quad \|M\|_2 \leq p_2 h \quad c'', p_2 \text{ constant}$$

$$\phi_1 = \left(1 + \frac{c c'' 97 \nu}{h}\right) \left(1 + \frac{6 \cdot 36 c c'' p_2 h \nu}{h}\right)$$

$$= 1 + \frac{c c'' 97 \nu}{h} + O(\nu)$$

$$\phi_2 = \frac{c c'' 97 \nu}{h} + \left(1 + \frac{c c'' 97 \nu}{h}\right) \frac{6 \cdot 36 c c'' p_2 h \nu}{h}$$

$$= \frac{c c'' 97 \nu}{h} + O(\nu)$$

$$\text{Therefore } \phi_1 = 1 + \frac{k' \nu}{h}$$

$$\phi_2 = \frac{k' \nu}{h}$$

where $k' = c c'' 97$

and
$$\|w_{n+1}\|_A \leq \beta \left(\left(1 + \frac{k'v}{h} \right)^n - 1 \right)$$

$$\leq \frac{\beta nk'v}{h} e^{\frac{nk'v}{h}}$$

But $n = \frac{T}{k}$ and if we keep $\frac{v}{kh}$ constant

then $\|w_{n+1}\|_A \leq K''$ where K'' is a constant.

We have just proved the following two theorems.

Theorem 1:

Let $v = N^{1-\tau}$ where τ is the number of floating base $-N$ digits in the mantissa. For the boundary value problem satisfying (1) whose solution is approximated by the function $U(x,t) = \sum_{\ell} \alpha_{\ell}(t) v_{\ell}(x)$ defined by (2) and using the difference scheme set out in (4) then, if the approximate solution is required to lie in a Hermite space of piecewise polynomials of degree 3, based on a rectangular grid of mesh size h , and if $\Delta t = c_1 h$, where c_1 is a constant, and if as $\Delta t, h \rightarrow 0$, $\frac{v}{\Delta t h}$ is kept constant, then the computed solution \hat{U} satisfies

$$\|\hat{U} - U\|_{L^2} \leq K' \quad \text{where } K' \text{ is a constant.}$$

Theorem 2:

For the boundary value problem described in Theorem 1, but where $\Delta t = c_2 h^2$, where c_2 is a constant and if as $\Delta t, h \rightarrow 0$, $\frac{\nu}{\Delta t h}$ is kept constant, then the computed solution \hat{U} satisfies

$$\left\| \hat{U} - U \right\|_{L^2} \leq K'' \quad \text{where } K'' \text{ is a constant.}$$

If we compare these results with those obtained in [2] for the finite difference method we see that exactly the same asymptotic behaviour is exhibited in both cases.

References

1. Douglas, J., Jr., and Dupont, T., Galerkin Methods for Parabolic Equations. To appear.
2. Rachford, H. H., Jr., Rounding Errors in Parabolic Problems I, The one space variable case. SIAM J. of Num. Analysis, 5, no.1, 156-171(1968).
3. Wilkinson, J. H., Rounding Errors in Algebraic Processes, Prentice Hall (1963).
4. Wilkinson, J. H., Error Analysis of Direct Methods of Matrix Inversion. J. Assoc. Comp. Mach. 8, 281-330(1961).