

SLAC-PUB-774  
CS 141  
October 1969  
(MISC)

BOUNDS FOR THE ERROR OF LINEAR SYSTEMS  
OF EQUATIONS USING THE THEORY OF MOMENTS

Germund Dahlquist<sup>+</sup>

Stanley C. Eisenstat

Gene H. Golub<sup>\*</sup>

Reproduction in whole or in part is permitted  
for any purpose of the United States Government.

+ Royal Institute of Technology, 10044 Stockholm 70, Sweden.

\* The work of this author was in part supported by the  
Atomic Energy Commission and by the National Science  
Foundation.

(To be submitted for publication)

## Abstract

Consider the system of linear equations  $A\tilde{x} = \tilde{b}$  where  $A$  is an  $n \times n$  real symmetric, positive definite matrix and  $\tilde{b}$  is a known vector. Suppose we are given an approximation to  $\tilde{x}$ ,  $\tilde{\xi}$ , and we wish to determine upper and lower bounds for  $\|\tilde{x} - \tilde{\xi}\|$  where  $\|\dots\|$  indicates the euclidean norm. Given the sequence of vectors  $\{\tilde{r}_i\}_{i=0}^k$  where  $\tilde{r}_i = A\tilde{r}_{i-1}$  and  $\tilde{r}_0 = \tilde{b} - A\tilde{\xi}$ , it is shown how to construct a sequence of upper and lower bounds for  $\|\tilde{x} - \tilde{\xi}\|$  using the theory of moments.

In addition, consider the Jacobi algorithm for solving the system  $\tilde{x} = M\tilde{x} + \tilde{b}$  viz.  $\tilde{x}_{i+1} = M\tilde{x}_i + \tilde{b}$ . It is shown that by examining  $\delta_i = \tilde{x}_{i+1} - \tilde{x}_i$ , it is possible to construct upper and lower bounds for  $\|\tilde{x}_i - \tilde{x}\|$ .

## 1. Introduction

Consider the system of linear algebraic equations

$$\underset{\sim}{A} \underset{\sim}{x} = \underset{\sim}{b} \quad (1.1)$$

where  $A$  is an  $n \times n$  real symmetric, positive definite matrix and  $\underset{\sim}{b}$  is a given vector. Assume we have an approximation to  $\underset{\sim}{x}$  so that

$$\underset{\sim}{x} = \underset{\sim}{\xi} + \underset{\sim}{e} \quad (1.2)$$

where  $\underset{\sim}{\xi}$  is the approximation vector and  $\underset{\sim}{e}$  is the error vector. We are concerned with determining upper and lower bounds for  $\|\underset{\sim}{e}\|$  where  $\|\dots\|$  indicates the euclidean norm of the vector.

In order to compute bounds for the norm of the error vector, it is natural to compute the residual vector,

$$\underset{\sim}{r}_0 = \underset{\sim}{b} - \underset{\sim}{A} \underset{\sim}{\xi} \quad (1.3)$$

Thus since  $\underset{\sim}{r}_0 = \underset{\sim}{A} \underset{\sim}{e}$ ,

$$\frac{\|\underset{\sim}{r}_0\|}{\|\underset{\sim}{A}\|} \leq \|\underset{\sim}{e}\| \leq \|\underset{\sim}{A}^{-1}\| \|\underset{\sim}{r}_0\| \quad (1.4)$$

Here  $\|\underset{\sim}{A}\|$  indicates the spectral norm of the matrix  $A$ . Assuming that  $\|\underset{\sim}{A}\| = 1$  (this can be accomplished via a simple scaling of (1.1)), we see that even though  $\|\underset{\sim}{r}_0\|$  is "small", the bound for  $\|\underset{\sim}{e}\|$  can be quite large when  $\|\underset{\sim}{A}^{-1}\|$  is very large.

By computing additional information, it is possible to obtain more precise upper and lower bounds on the euclidean length of the error vector. In Section 2, we give an algorithm which depends upon computing an auxiliary sequence of vectors and an explicit knowledge of all the eigenvalues of the matrix  $A$ . The bounds are actually obtained as a solution to a linear programming problem.

In Section 3, we use the same sequence of vectors as described in Section 2 but we assume that the only information that the investigator has is an upper bound on the largest eigenvalue of  $A$  and a non-trivial lower bound for the smallest eigenvalue. Using the theory of moments, an algorithm is given for determining upper and lower bounds. Then in Section 4, we consider the Jacobi iterative method for solving the system of equations (1.1), and we show it is possible to establish bounds for the error by examining the difference of successive iterates. Finally, a numerical example is given in Section 5. In a future report, we shall give methods for improving the approximate solution using the techniques described in this paper.

2. Bounds using linear programming

Consider the Krylov sequence [6],

$$\tilde{r}_{i+1} = A\tilde{r}_i \quad (i = 0, 1, \dots, k-1 \leq n)$$

where  $\tilde{r}_0$  is defined by (1.3). Thus

$$\tilde{r}_i = A^i \tilde{r}_0 \quad (i = 0, 1, \dots, k) .$$

We define

$$(x, y) = \sum_{i=1}^n x_i y_i$$

so that

$$\begin{aligned} (r_p, r_q) &= (A^p r_0, A^q r_0) \\ &= (A^{p+q} r_0, r_0) \\ &\equiv \mu_{p+q} \quad (p, q = 0, 1, \dots, k) . \end{aligned}$$

Since  $A$  is symmetric and positive definite, we have

$$A u_i = \lambda_i u_i \quad (i = 1, 2, \dots, n)$$

with

$$(u_i, u_j) = \begin{cases} 0 & \text{for } i \neq j \\ 1 & \text{for } i = j \end{cases}$$

and

$$0 < a \leq \lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n \leq b .$$

Now writing

$$\tilde{r}_0 = \sum_{i=1}^n \alpha_i u_i ,$$

we have

$$\mu_m = (A^m r_o, r_o) = \sum_{i=1}^n \alpha_i^2 \lambda_i^m \quad (m = 0, 1, \dots, 2k) \quad . \quad (2.1)$$

Since  $\underline{e} = A^{-1} \underline{r}_o$  ,

$$\|\underline{e}\|^2 = (A^{-2} \underline{r}_o, \underline{r}_o) = \sum_{i=1}^n \alpha_i^2 \lambda_i^{-2} = \mu_{-2} \quad . \quad (2.2)$$

Equations (2.1) and (2.2) are equivalent to

$$\mu_m = \int_a^b \lambda^m d\alpha(\lambda) \quad (m = -2, 0, 1, \dots, 2k) \quad (2.3)$$

where the weight function of the Stieltjes integral is determined as follows:

$$\begin{aligned} \alpha(\lambda) &= 0 && \text{for } a \leq \lambda \leq \lambda_1 \quad , \\ \alpha(\lambda) &= \alpha_1^2 + \alpha_2^2 + \dots + \alpha_t^2 && \lambda_t < \lambda \leq \lambda_{t+1} \quad (t = 1, 2, \dots, n-1) \quad , \\ \alpha(\lambda) &= \alpha_1^2 + \alpha_2^2 + \dots + \alpha_n^2 && \lambda_n < \lambda \leq b \quad . \end{aligned}$$

The problem of determining an upper and lower bound for  $\|\underline{e}\|$  is equivalent to the following:

Given the  $(2k+1)$  moments  $\mu_i$  , determine upper and lower bounds for  $\mu_{-2}$  .

The solution to this classical problem (cf. [7]) is dependent upon the amount of information available.

Suppose we know the eigenvalues of the matrix  $A$  . An example of this is the usual five point approximation to Poisson's equation with Dirichlet boundary conditions in a rectangular region. Thus to determine an upper bound for  $\|\underline{e}\|^2$  , we wish to maximize

$$\sum_{i=1}^n \gamma_i \lambda_i^{-2}$$

subject to the constraints

$$\left. \begin{aligned} \sum_{i=1}^n \gamma_i \lambda_i^m &= \mu_m & (m = 0, 1, \dots, 2k) \\ \gamma_i &\geq 0 & (i = 1, \dots, n) \end{aligned} \right\} \quad (2.4)$$

The numerical solution to this problem can be obtained by the simplex algorithm of G. Dantzig [3]. Special techniques may be used to take advantage of the fact that a Vandermonde system is solved at each iteration of the simplex algorithm. A lower bound for  $\|\tilde{e}\|^2$  may be obtained by determining the minimum of  $\sum_{i=1}^n \gamma_i \lambda_i^{-2}$  subject to the constraints (2.4) by the simplex algorithm.

### 3. Error bounds using the theory of moments

In the more usual situation, one has the information that

$$0 < a \leq \lambda_i \leq b \quad \text{for } i = 1, \dots, n .$$

This is a problem in the classical theory of moments which has been solved by A. A. Markov. In order to give a numerical algorithm for determining bounds for  $\|e\|$ , we review some facts from the theory of Gaussian quadrature.

Suppose we are given  $\{\mu_i\}_{i=0}^{2k}$ , and a function  $\varphi(\lambda)$  ( $a \leq \lambda \leq b$ ), and we wish to determine  $(L, U)$  so that

$$L \leq \int_a^b \varphi(\lambda) d\alpha(\lambda) \leq U .$$

We can determine a quadrature rule such that

$$\mu_r = \int_a^b \lambda^r d\alpha(\lambda) = \sum_{i=1}^k A_i t_i^r + \sum_{j=1}^m B_j z_j^r \quad \text{for } r = 0, 1, \dots, 2k+m-1$$

where  $\{A_i, t_i\}_{i=1}^k$  and  $\{B_j\}_{j=1}^m$  are unknown and  $\{z_j\}_{j=1}^m$  is specified.

Then

$$\int_a^b \varphi(\lambda) d\alpha(\lambda) = \sum_{i=1}^k A_i \varphi(t_i) + \sum_{j=1}^m B_j \varphi(z_j) + R[\varphi]$$

where

$$R[\varphi] = \frac{\varphi^{(2k+m)}(\eta)}{(2k+m)!} \int_a^b \prod_{j=1}^m (\lambda - z_j) \left[ \prod_{i=1}^k (\lambda - t_i) \right]^2 d\alpha(\lambda) \tag{3.1}$$

$$a < \eta < b .$$

Thus if  $\varphi(\lambda) = \lambda^{-2}$  and  $m = 1$ ,

$$R[\lambda^{-2}] = -2(k+1)\eta^{-(2k+3)} \int_a^b (\lambda - z_1) \left[ \prod_{i=1}^k (\lambda - t_i) \right]^2 d\alpha(\lambda) .$$



Hence for  $z_1 = a$ , the Gauss-Radau type quadrature rule yields an upper bound for  $\int_a^b \lambda^{-2} d\alpha(\lambda)$  and if  $z_1 = b$ , we have a lower bound. It can be shown (cf. [7, pg. 80]) that these bounds are attainable.

It is not necessary to solve the equations for the quadrature rule.

Let us note that

$$\mu_r = \sum_{i=1}^k A_i t_i^r + B_1 z_1^r \quad (r = 0, 1, \dots, 2k)$$

where  $z_1$  may be  $a$  or  $b$ . Let us write

$$\bar{\mu}_r = \sum_{i=1}^k \bar{A}_i \bar{t}_i^r + \bar{B}_1 a^r \quad \text{for all } r \text{ and for } z_1 = a \quad (3.2)$$

so that

$$\bar{\mu}_{-2} \geq \mu_{-2} \quad .$$

From (3.2), we see that  $\bar{\mu}_r$  satisfies a  $(k+1)$ <sup>th</sup> order difference equation

$$\bar{g}_0 \bar{\mu}_r + \bar{g}_1 \bar{\mu}_{r-1} + \dots + \bar{g}_k \bar{\mu}_{r-k} - \bar{\mu}_{r-(k+1)} = 0 \quad (3.3)$$

and  $\bar{t}_1, \bar{t}_2, \dots, \bar{t}_k$ , and  $a$  are the roots of the characteristic polynomial

$$p(\xi) = \bar{g}_0 \xi^{k+1} + \bar{g}_1 \xi^k + \dots + \bar{g}_k \xi - 1 \quad .$$

Since  $\bar{p}(a) = 0$ , we must have

$$\bar{g}_0 a^{k+1} + \bar{g}_1 a^k + \dots + \bar{g}_k a - 1 = 0 \quad . \quad (3.4)$$

Thus using (3.3) and (3.4) and the fact that  $\mu_r = \bar{\mu}_r$  for  $r = 0, 1, \dots, 2k$ , we have enough equations to determine  $\{\bar{g}_i\}_{i=0}^k$ . Having determined  $\{\bar{g}_i\}_{i=0}^k$ , one can solve for  $\bar{\mu}_{-2}$  by recurring twice backwards with equation (3.3).

To determine a lower bound for the error viz.  $\mu_{-2}$ , it is necessary to solve for  $\{\underline{g}_i\}_{i=0}^k$  from equations similar to (3.3) and the additional equation

$$g_0 b^{k+1} + g_1 b^k + \dots - g_k b - 1 = 0 .$$

Note to solve for  $\{g_i\}_{i=0}^k$ , it is necessary to change only one row in the matrix and one can use the devices given in [2] for solving such a modified system efficiently.

For large  $k$ , the system of linear equations which one solves for the coefficients of the difference equation may be quite ill-conditioned. For that reason it is sometimes preferable to solve explicitly for the quadrature rule. As is well known, the nodes of the quadrature rule are the roots of orthogonal polynomials. Now the orthogonal polynomials satisfy a three term recurrence relationship viz.

$$p_{j+1}(\lambda) = (\xi_{j+1} - \lambda)p_j(\lambda) - \eta_j^2 p_{j-1}(\lambda)$$

$$p_{-1}(\lambda) = 0 \quad , \quad p_0(\lambda) = 1 \quad .$$

The coefficients  $\{\xi_j\}_{j=1}^k$ ,  $\{\eta_j\}_{j=1}^{k-1}$  can be computed directly using the Lanczos algorithm [8].

Again, let

$$\underline{r}_0 = \underline{b} - A\underline{x} \quad .$$

We generate a sequence of vectors  $\{\underline{z}_i\}_{i=0}^{k+1}$  such that

$$(\underline{z}_i, \underline{z}_j) = \begin{cases} 0 & \text{for } i \neq j \\ 1 & \text{for } i = j \end{cases} .$$

Let  $\underline{z}_0 = \underline{r}_0 \times (\|\underline{r}_0\|)^{-1}$  .

Then for  $j = 0, 1, \dots, k$ ,





4. Error bounds for the Jacobi method

Consider the system of equations,

$$C\tilde{y} = \tilde{f} \quad (4.1)$$

where  $C$  is a real, symmetric positive definite matrix of order  $n$ .

Let  $D = \text{diag}[(c_{11})^{-1/2}, \dots, (c_{nn})^{-1/2}]$ . We may write (4.1) in the form

$$DCDD^{-1}\tilde{y} = D\tilde{f} \quad (4.2)$$

or equivalently,

$$A\tilde{x} = \tilde{b} \quad (4.3)$$

Note the diagonal elements of the matrix  $A$  are all equal to one. Hence

$$A = I - M$$

where the diagonal elements of  $M$  are zero. We shall assume that  $M$

is convergent viz.  $\max_{1 \leq i \leq n} |\lambda_i(M)| < 1$ . The Jacobi method viz.

$$\tilde{x}_{i+1} = M\tilde{x}_i + \tilde{b} \quad (i = 0, 1, \dots)$$

is frequently used to solve (4.3).

Let

$$\tilde{e}_i = \tilde{x} - \tilde{x}_i = M^i \tilde{e}_0$$

and

$$\tilde{\delta}_i = \tilde{x}_{i+1} - \tilde{x}_i = M^i \tilde{\delta}_0 \quad :$$

The vector  $\tilde{\delta}_i$  is the difference vector. Since  $\tilde{\delta}_i = \tilde{x}_{i+1} - \tilde{x}_i = M\tilde{x}_i + \tilde{b} - \tilde{x}_i = \tilde{b} - A\tilde{x}_i$ , the difference vector is the residual vector associated with  $\tilde{x}_i$ . Note

$$e_{\sim i} = (I-M)^{-1} \delta_{\sim i} = (I-M)^{-1} M^i \delta_{\sim 0} .$$

Given  $\{\delta_{\sim i}\}_{i=0}^k$ , we compute

$$(\delta_{\sim p}, \delta_{\sim q}) = (\delta_{\sim 0}, M^{p+q} \delta_{\sim 0}) \equiv v_{p+q} , \quad (p+q = 0, \dots, 2k).$$

Thus

$$\|e_{\sim k+1}\|^2 = (e_{\sim k+1}, e_{\sim k+1}) = ((I-M)^{-1} M^{k+1} \delta_{\sim 0}, (I-M)^{-1} M^{k+1} \delta_{\sim 0}) .$$

Since  $M$  is symmetric, we have

$$M u_{\sim i} = \xi_i u_{\sim i} , \quad (i = 1, 2, \dots, n)$$

$$(u_{\sim i}, u_{\sim j}) = \begin{cases} 0 & \text{for } i \neq j \\ 1 & \text{for } i = j \end{cases}$$

and we assume

$$-1 < c \leq \xi_1 \leq \xi_2 \leq \dots \leq \xi_n \leq d < 1 .$$

Thus

$$v_m = \int_c^d \xi^m d\beta(\xi) \quad (m = 0, 1, \dots, 2k)$$

and

$$(e_{\sim k+1}, e_{\sim k+1}) = \int_c^d \frac{\xi^{2k+2}}{(1-\xi)^2} d\beta(\xi) .$$

We wish to determine upper and lower bounds for  $\|e_{\sim k+1}\|$ . This problem was first discussed by H. Weinberger [9] for  $k = 1$ .

Again, if the eigenvalues of  $M$  are known then one can use linear programming for determining upper and lower bounds of  $\|e_{k+1}\|^2$ . Thus to determine an upper bound for  $\|e_{k+1}\|^2$ , we wish to maximize

$$\sum_{i=1}^n \omega_i \frac{\xi_i^{2k+2}}{(1-\xi_i)^2}$$

subject to the constraints

$$\sum_{i=1}^n \omega_i \xi_i^m = v_m \quad (m = 0, 1, \dots, 2k)$$

$$\omega_i \geq 0 \quad (i = 1, 2, \dots, n) \quad .$$

If the eigenvalues are unknown, then we are unable to use the arguments associated with the residual vector since the  $(2k+1)^{st}$  derivative of  $\varphi(\xi) = \xi^{2k+2}/(1-\xi)^2$  is not of constant sign in the interval  $(c, d)$ .

Now, if we can determine a polynomial  $p_{2k}(\xi)$  such that

$$p_{2k}(\xi) = c_0 + c_1 \xi + \dots + c_{2k} \xi^{2k} \geq \frac{\xi^{2k+2}}{(1-\xi)^2}$$

$$\text{for } c \leq \xi \leq d$$

then this will determine an upper bound for  $\|e_{k+1}\|^2$  since

$$\int_c^d p_{2k}(\xi) d\beta(\xi) = c_0 v_0 + \dots + c_{2k} v_{2k} \geq \int_c^d \frac{\xi^{2k+2}}{(1-\xi)^2} d\beta(\xi) \quad .$$

The polynomial  $p_{2k}(\xi)$  is not unique and consequently we desire that polynomial for which

$$c_0 v_0 + \dots + c_{2k} v_{2k} = \min.$$

Unfortunately, there does not seem to be any numerical algorithms which will satisfactorily solve this problem in general.

Let  $\lambda = (1-\xi)$  so that

$$\begin{aligned} \int_c^d \frac{\xi^{2k+2}}{(1-\xi)^2} d\beta(\xi) &= \int_a^b \frac{(1-\lambda)^{2k+2}}{\lambda^2} d\alpha(\lambda) \\ &= \mu_{-2} - 2(k+1)\mu_{-1} + \sum_{s=0}^{2k} (-1)^s \binom{2k+2}{s+2} \mu_s \end{aligned}$$

where

$$a = 1-d, \quad b = 1-c$$

$$\mu_s = \int_a^b \lambda^s d\alpha(\lambda) \quad (s = -2, -1, \dots, 2k) \quad .$$

It is easy to verify that

$$\mu_s = (-1)^s \Delta^s v_0$$

where

$$\Delta v_0 = v_1 - v_0$$

$$\Delta^s v_0 = \Delta(\Delta^{s-1} v_0)$$

and hence  $\mu_s = (\underline{\delta}_0, (I-M)^s \underline{\delta}_0)$ . Our problem now is to determine upper and lower bounds for

$$\mu_{-2} - 2(k+1)\mu_{-1} \quad .$$

In order that there exist a distribution function  $\alpha(\lambda)$  in the interval  $(a,b)$  associated with  $\{\mu_s\}_{s=-2}^{2k}$ , it is necessary and sufficient that



$$M = \begin{bmatrix} \mu_{-2} & \mu_{-1} & \cdot & \cdot & \cdot & \mu_{k-1} \\ \mu_{-1} & \mu_0 & \cdot & \cdot & \cdot & \mu_k \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \mu_{k-1} & \mu_k & \cdot & \cdot & \cdot & \mu_{2k} \end{bmatrix}_{(k+2) \times (k+2)}$$

and

$$G = \begin{bmatrix} \gamma_{-2} & \gamma_{-2} & \cdot & \cdot & \cdot & \gamma_{k-1} \\ \gamma_{-1} & \gamma_0 & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \gamma_{k-2} & \cdot & \cdot & \cdot & \cdot & \gamma_{2k-2} \end{bmatrix}_{(k+1) \times (k+1)}$$

where

$$\gamma_j = -[ab \mu_j - (a+b)\mu_{j+1} + \mu_{j+2}] \quad (4.4)$$

be positive semi-definite (cf. [1]). It is easy to see why  $G$  must be positive semi-definite. Note from (4.4)

$$\begin{aligned} \gamma_j &= - \int_a^b (ab\lambda^j - (a+b)\lambda^{j+1} + \lambda^{j+2}) d\alpha(\lambda) \\ &= \int_a^b \lambda^j (\lambda-a)(b-\lambda) d\alpha(\lambda) \end{aligned}$$

Hence,

$$\begin{aligned} \tilde{z}^T G \tilde{z} &= \sum_{i=-1}^{k-1} \sum_{j=-1}^{k-1} \gamma_{i+j} z_i z_j \\ &= \int_a^b \left( \sum_{i=-1}^{k-1} z_i \lambda^i \right)^2 (\lambda-a)(b-a) d\alpha(\lambda) \geq 0 \end{aligned}$$

A similar argument shows that  $M$  must be positive semi-definite. Observe that there are two elements which are unknown in  $G$  and two elements which are unknown in  $M$  and they occur in either the first row or column of the matrix.

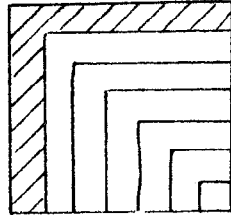


Figure I

Since  $M$  and  $G$  are positive semi-definite, it is necessary and sufficient that  $\det(M) \geq 0$  and  $\det(G) \geq 0$  in order for the values  $\mu_{-2}, \mu_{-1}$  be consistent with some distribution  $\alpha(\lambda)$  with moments  $\mu_0, \mu_1, \dots, \mu_{2k}$ . The positive semi-definite property of  $M$  and  $G$  is equivalent to the non-negativity of the sub-determinants indicated in Figure I.

We partition the Hankel matrix  $M$  as follows:

$$M = \begin{bmatrix} \mu_{-2} & \mu_{-1} & \mu_0, \dots, \mu_{k-1} \\ \mu_{-1} & \mu_0 & \mu_1, \dots, \mu_k \\ \mu_0 & \mu_1 & \mu_2, \dots, \mu_{k+1} \\ \vdots & \vdots & \vdots \\ \mu_{k-1} & \mu_k & \mu_{k+1}, \dots, \mu_{2k} \end{bmatrix} = \begin{bmatrix} A & B \\ B^T & C \end{bmatrix}$$

Since  $M$  is positive semi-definite

$$\det(M) = \det(C) \det(A - BC^{-1} B^T) \geq 0$$

so that

$$\det \begin{bmatrix} \mu_{-2} + r_1 & \mu_{-1} + r_2 \\ \mu_{-1} + r_2 & \mu_0 + r_3 \end{bmatrix} \geq 0 \quad (4.5)$$

where

$$\begin{bmatrix} r_1 & r_2 \\ r_2 & r_3 \end{bmatrix} = -BC^{-1} B^T .$$

The matrix  $-BC^{-1} B^T$  can easily be computed by applying the Choleski algorithm to the matrix

$$\left[ \begin{array}{c|c} 0 & B \\ \hline B^T & C \end{array} \right]$$

One must begin the pivoting operation, however, from the bottom diagonal element and after  $k$  eliminations, the upper  $2 \times 2$  matrix will contain  $-BC^{-1} B^T$ . In a similar fashion,

$$\det \begin{bmatrix} -ab\mu_{-2} + (a+b)\mu_{-1} - \mu_0 + s_1 & , & -ab\mu_{-1} + (a+b)\mu_0 - \mu_1 + s_2 \\ -ab\mu_{-1} + (a+b)\mu_0 - \mu_1 + s_2 & , & -ab\mu_0 + (a+b)\mu_1 - \mu_2 + s_3 \end{bmatrix} \geq 0 . \quad (4.6)$$

From equation (4.5), we see that

$$(\mu_{-2} + r_1)(\mu_0 + r_3) - (\mu_{-1} + r_2)^2 \geq 0$$

and hence since  $\mu_0 + r_3 \geq 0$  by the positive semi-definiteness of  $M$

$$\mu_{-2} \geq \frac{(\mu_{-1} + r_2)^2}{\mu_0 + r_3} - r_1 \equiv \frac{(\mu_{-1} + R_2)^2}{\mu_0 + R_3} - R_1 .$$

From (4.6), we have

$$\det \begin{bmatrix} \mu_{-2} + s_1 & \mu_{-1} + s_2 \\ \mu_{-1} + s_2 & s_3 \end{bmatrix} \geq 0$$

where

$$s_3 = \frac{ab\mu_0 - (a+b)\mu_1 + \mu_2 - s_3}{ab} < 0$$

since  $ab > 0$  ,

and hence

$$\mu_{-2} \leq \frac{(\mu_{-1} + s_2)^2}{s_3} - s_1 .$$

Therefore

$$\frac{(\mu_{-1} + R_2)^2}{R_3} - R_1 \leq \mu_{-2} \leq \frac{(\mu_{-1} + s_2)^2}{s_3} - s_1 . \quad (4.7)$$

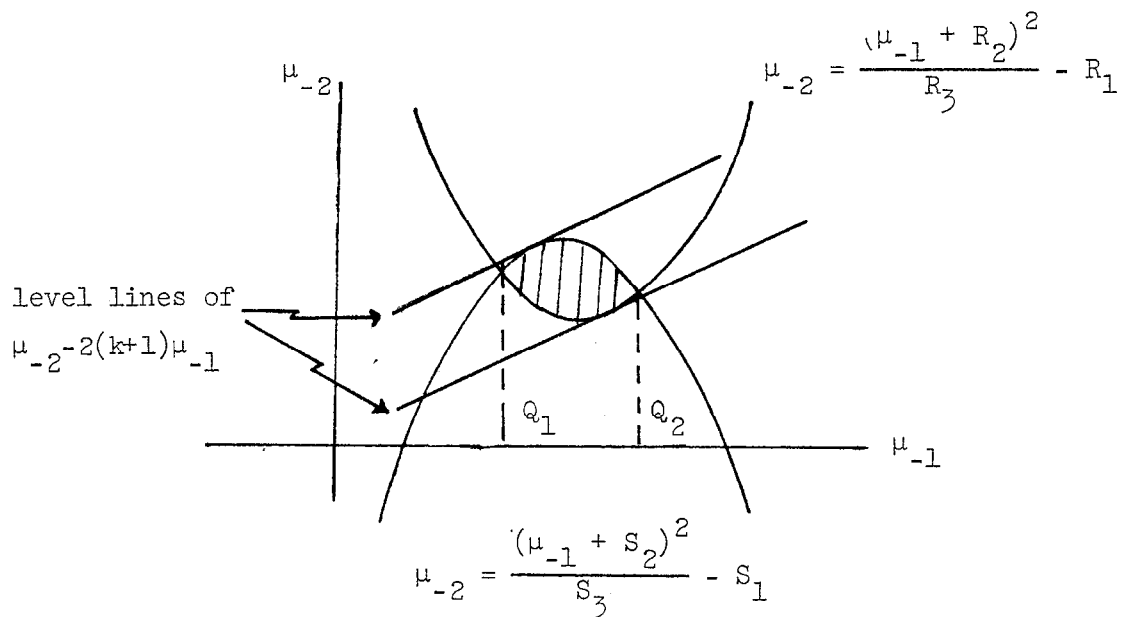


Figure II

Thus to determine the maximum of

$$\mu_{-2} - 2(k+1)\mu_{-1} \quad ,$$

it is simply necessary to examine the boundary of the shaded region in Figure II. A short calculation yields  $(Q_1, Q_2)$  for which

$$Q_1 \leq \mu_{-1} \leq Q_2 \quad .$$

Then  $\mu_{-2} - 2(k+1)\mu_{-1} = \text{maximum}$  subject to (4.7) if

$$\frac{d}{d\mu_{-1}} \frac{(\mu_{-1} + S_2)^2}{S_3} - S_1 = \frac{d}{d\mu_{-1}} 2(k+1)\mu_{-1}$$

$$\text{with } \mu_{-1}^U = -S_2 + (k+1)S_3$$

and

$$Q_1 \leq \mu_{-1}^U \leq Q_2 \quad ;$$

otherwise the maximum occurs at  $\mu_{-1}^U = Q_1$  or  $\mu_{-1}^U = Q_2$  according to

$$\mu_{-1}^U = \max\{Q_1, \min\{-S_2 + (k+1)S_3, Q_2\}\} \quad .$$

Similarly, the minimum occurs at

$$\mu_{-1}^L = \max\{Q_1, \min\{-R_2 + (k+1)R_3, Q_2\}\} \quad .$$

Thus, it is possible to determine upper and lower bounds for  $\|e_{\sim k+1}\|$ , and these bounds are attainable.

5. A numerical example

Consider the system of equations

$$\underset{\sim}{A} \underset{\sim}{x} = \underset{\sim}{b}$$

where  $A$  is a tri-diagonal matrix with elements  $(-1, 2, -1)$  and  $\underset{\sim}{b} = \underset{\sim}{0}$ , the null vector. It is well known that

$$\lambda_j(A) = 2 + 2 \cos \frac{j\pi}{n+1}, \quad j = 1, 2, \dots, n.$$

The Jacobi matrix  $M$  is also tri-diagonal and has elements  $(1/2, 0, 1/2)$ .

Here

$$\lambda_j(M) = \cos \frac{j\pi}{n+1}, \quad j = 1, 2, \dots, n.$$

The Jacobi method was used for solving the system for  $n = 20$  and  $\underset{\sim}{x}_0^T = (1, 1, \dots, 1)$ . In Tables I, II, and III, we give the error bounds associated with the error vector of  $\underset{\sim}{x}_{10}$ . To use the methods of

Section 3, we must compute in addition  $\{\underset{\sim}{r}_p\}$  for  $p = 0, 1, \dots, k$ .

In Tables II and III, we give bounds for the error using the difference vectors  $\{\underset{\sim}{\delta}_{g-p}\}$  for  $p = 0, 1, \dots, k$ . Note that the bound using the residual vector is slightly better than those computed using the difference vectors but it requires additional work to compute  $\{\underset{\sim}{r}_p\}_{p=0}^k$  whereas the difference vectors are computed in the natural sequence of events.

In addition, note that the lower bounds are less influenced by the interval of the eigenvalues than are the upper bounds. Furthermore, we see that in this case that a knowledge of all the eigenvalues does not provide much smaller intervals for the error.

The authors are very pleased to acknowledge the stimulating comments of Professor H. Weinberger of the University of Minnesota and Mr. David Galant of the Ames Research Center.



Error bounds after 10 iterations

$$\| \tilde{e} \|^2 = 2.700138$$

Table I

Error bounds computed from residual vector using Gauss-Radau quadrature rule

k	<u>Lower bounds</u>		<u>Upper bounds</u>	
	a=1.116917 <sub>10</sub> <sup>-2</sup> b=1.988831	a=10 <sup>-2</sup> b=1.99	a=1.116917 <sub>10</sub> <sup>-2</sup> b=1.988831	a=10 <sup>-2</sup> b=1.99
1	1.35	1.35	3.40 <sub>10</sub> <sup>1</sup>	4.21 <sub>10</sub> <sup>1</sup>
2	1.55	1.54	2.05 <sub>10</sub> <sup>1</sup>	2.52 <sub>10</sub> <sup>1</sup>
3	1.66	1.66	9.40	1.12 <sub>10</sub> <sup>1</sup>
4	1.81	1.81	6.89	8.08
5	1.89	1.89	4.84	5.50

Table II

Error bounds computed from difference vectors using determinantal inequalities

k	<u>Lower bounds</u>		<u>Upper bounds</u>	
	a = 1.116917 <sub>10</sub> <sup>-2</sup> b = 1.988831	a = 10 <sup>-2</sup> b = 1.99	a = 1.116917 <sub>10</sub> <sup>-2</sup> b = 1.988831	a = 10 <sup>-2</sup> b = 1.99
1	1.35	1.35	5.29 <sub>10</sub> <sup>1</sup>	6.59 <sub>10</sub> <sup>1</sup>
2	1.43	1.43	3.88 <sub>10</sub> <sup>1</sup>	4.82 <sub>10</sub> <sup>1</sup>
3	1.48	1.48	2.05 <sub>10</sub> <sup>1</sup>	2.51 <sub>10</sub> <sup>1</sup>
4	1.56	1.56	1.70 <sub>10</sub> <sup>1</sup>	2.08 <sub>10</sub> <sup>1</sup>
5	1.59	1.59	1.29 <sub>10</sub> <sup>1</sup>	1.57 <sub>10</sub> <sup>1</sup>

Table III

Error bounds computed from difference vectors using  
linear programming

<u>k</u>	<u>Lower bounds</u>	<u>Upper bounds</u>
1	1.35	$5.04_{10}^1$
2	1.45	$3.86_{10}^1$
3	1.55	$1.73_{10}^1$
4	1.61	$1.50_{10}^1$
5	1.62	$1.04_{10}^1$

## References

- [1] N. I. Aheizer and M. Krein, "Some Questions in the Theory of Moments," American Mathematical Society, Providence, Rhode Island, 1962.
- [2] R. Bartels and G. H. Golub, "Stable numerical methods for obtaining the Chebyshev solution to an overdetermined system of equations," Comm. A.C.M., 11 (1968), 401-406.
- [3] G. B. Dantzig, Linear Programming and Extensions, Princeton University Press, Princeton, New Jersey, 1963.
- [4] D. Galant, personal communication.
- [5] G. H. Golub and J. Welsch, "Calculation of Gauss quadrature rules," M.O.C., 23 (1969), 221-230.
- [6] A. S. Householder, The Theory of Matrices in Numerical Analysis, Blaisdell Publishing Co., New York, 1964.
- [7] S. Karlin and W. J. Studen, Tchebycheff Systems: With Application in Analysis and Statistics, Interscience Publishers, New York, 1966.
- [8] C. Lanczos, "An iteration method for the solution of the eigenvalue problem of linear differential and integral operators," J. Res. Nat. Bur. Standards, 45 (1950), 255-282.
- [9] H. Weinberger, "A posteriori error bounds in iterative matrix inversion," in Symposium on the Numerical Solution of Partial Differential Equations, James H. Bramble, ed., Academic Press, New York, 1966, 153-163.

Key words

linear algebra

theory of moments

quadrature rules