# THE ROLE OF CONDITIONING IN THE
# NUMERICAL SOLUTION OF BOUNDARY VALUE PROBLEMS

*Frank de Hoog*

## ABSTRACT

The manner in which small perturbations in the data give rise to a perturbation in the solution of a two point boundary value problem will determine how effectively the solution can be approximated numerically. However this 'conditioning' also has a number of other implications that are less obvious. These include the effect on the stability of numerical schemes, the conditioning of shooting methods and the existence of a dichotomy. This paper aims to examine some of these inter-relationships.

## 1. INTRODUCTION

We shall examine the system of ordinary differential equations

$$(1.1) \qquad Ly := y' - Ay = f , \qquad 0 < t < 1$$

where $A \in [L_p(0,1)]^{n \times n}$ and $f \in [L_p(0,1)]^n$ for $1 \le p \le \infty$. In order to specify the solution uniquely, we require some restrictions on the solution and in our case we limit attention to conditions involving only values of $y$ at the ends of the interval. Specifically, we consider

$$(1.2) \qquad By = B_0 y(0) + B_1 y(1) = b$$

for $B_0, B_1 \in \mathbb{R}^{n \times n}$ satisfying

$$\max\{|B_0|, |B_1|\} = 1$$

where $|\cdot|$ denotes the Euclidean norm.

Any solution to (1.1) (see for example Keller [3]) can be written in the form

(1.3) $$y(t) = Y(t)c + \tilde{y}_p(t)$$

where $Y \in [L_p^1(0,1)]^{n \times n}$ and $y_p \in [L_p^1(0,1)]^n$ are solutions of the initial value problems

$$LY = 0, \qquad Y(0) = I,$$

$$L\tilde{y}_p = f, \qquad \tilde{y}_p(0) = 0,$$

respectively. On substituting (1.3) into the boundary conditions (1.2) we obtain the system of equations

$$\mathcal{B}Yc = (B_0 + B_1Y(1))c = b - B_1\tilde{y}_p(1)$$

and it follows that the boundary value problem (1.1)-(1.2) has a unique solution iff the matrix $B_0 + B_1Y(1)$ is nonsingular. Furthermore, on writing

$$\Phi(t) := Y(t)(B_0 + B_1Y(1))^{-1}$$

we obtain

(1.4) $$y(t) = \Phi(t)b + \int_0^1 G(t,s) \, f(s) \, ds$$

$$= \Phi(t)b + y_p(t)$$

where

$$y_p(t) = \int_0^1 G(t,s) \, f(s) \, ds$$

and $G(t,s)$ is the Green's function defined by

$$(1.5) \qquad G(t,s) = \begin{cases} \Phi(t) \ B_0 \ \Phi(0) \ \Phi^{-1}(s) & t > s \\ \\ -\Phi(t) \ B_1 \ \Phi(0) \ \Phi^{-1}(s) & t < s . \end{cases}$$

Thus, a knowledge of the fundamental solution allows us to write down (at least in principle) the solution to the boundary value problem (1.1), (1.2).

When solving boundary value problems of the form (1.1), (1.2) numerically, it is important to know how small perturbations in the right hand sides of (1.1) and (1.2) will affect the solution. In order to characterise this mathematically we need to define some norms. With $b \in \mathbb{R}^n$, $E \in \mathbb{R}^{n \times n}$ we denote the Euclidean norm by $|b|$ and the norm of $A$ by $|A| = \sup_b |Ab|/|b|$ respectively. For $f \in [L_p(0,1)]^n$ we define

$$\|f\|_p := [\int_0^1 |f(t)|^p \ dt]^{1/p} , \qquad 1 \le p < \infty$$

with its limiting value

$$\|f\|_\infty = \sup_t |f(t)|$$

and for $A \in [L_p(0,1)]^{n \times n}$ we take the induced norm

$$\|A\|_p = \sup_{f \in [L_p(0,1)]^n} \left\{ \frac{\|Af\|_p}{\|f\|_p} \right\} .$$

It now follows from (1.4) that

$$\|y\|_\infty \le \beta|b| + \alpha\|f\|_p$$

or equivalently

$$(1.6) \qquad \|y\|_\infty \le \beta|\mathcal{B}y| + \alpha\|Ly\|_p$$

where

$$\beta = \|\Phi\|_{\infty}$$

$$\alpha = \sup_{t} [\int_0^1 |G(t,s)|^q \, ds]^{1/q} \, , \qquad \frac{1}{p} + \frac{1}{q} = 1 \, .$$

In the sequel, we shall be concerned mainly with the case $p = 1$. Then,

(1.7)
$$\alpha = \sup_{t,s} |G(t,s)| \, .$$

However, in some applications, other choices of $p$ may be more appropriate and it may even be necessary to use weighted norms (see de Hoog and Mattheij [2]).

Basically (1.6) says that the perturbations in $f$ and $b$ may be amplified by the factors $\alpha$ and $\beta$ respectively. Thus, when $\alpha$ and $\beta$ are large, the problem is poorly conditioned and we may expect to have difficulties in obtaining a numerical solution. Note however that the converse is not true and it is easy to construct problems with moderate condition numbers that are difficult to solve numerically.

This paper looks at some of the consequences of the conditioning of (1.1), (1.2). In section 2 we examine typical stability estimates for finite difference schemes for the numerical solution of (1.1), (1.2) while in section 3 we consider the effect of the stability constants on shooting methods. Finally in section 4 we describe some results of de Hoog and Mattheij [2] which relate dichotomy and stability in two point boundary value problems.

## 2. STABILITY OF FINITE DIFFERENCE SCHEMES

We shall demonstrate in this section how the conditioning of finite difference schemes is affected by the condition numbers $\alpha$ and $\beta$ defined in section 1. For simplicity, we take $p = 1$, $A \in [C^1[0,1]]^{n \times n}$

and $\|A\|_{\infty}$ , $\|A'\|_{\infty} \leq a$ .

Let $\Delta = \{0 = t_0 < t_1 < \ldots < t_N = 1\}$ be a partition of the interval $[0,1]$ and consider the centred Euler finite difference operator defined by

$$(L_{\Delta}\underset{\sim}{y})_j = \frac{y_{j+1} - y_j}{h_j} - A(t_{j+\frac{1}{2}}) \frac{y_{j+1} + y_j}{2} , \qquad j = 0, \ldots, N-1$$

and boundary conditions

$$\mathcal{B}_{\Delta}\underset{\sim}{y} = B_0 y_0 + B_1 y_N$$

where $\underset{\sim}{y} \in \mathbb{R}^{(N+1)n}$ , $h_j = t_{j+1} - t_j$ and $t_{j+\frac{1}{2}} = (t_{j+1} + t_j)/2$ . Let $\hat{y}(t)$ be the broken line interpolant to $\underset{\sim}{y}$ at the points $t_j$ , $j = 0, \ldots, N$ . That is

$$\hat{y}(t) := (y_{j+1} + y_j)/2 + (t - t_{j+\frac{1}{2}})(y_{j+1} - y_j)/h_j , \qquad t_j \leq t \leq t_{j+1} .$$

Then, using (1.6) we obtain

$$\|\hat{y}\|_{\infty} = \max_j |y_j| \leq \beta |\mathcal{B}\hat{y}| + \alpha \| L\hat{y} \|_1 .$$

Furthermore

$$|\mathcal{B}\hat{y}| = |\mathcal{B}_{\Delta}\underset{\sim}{y}|$$

and it is easy to verify that

$$\| L\hat{y} \|_1 \leq \frac{(a + a^2)}{2} h \|\hat{y}\|_{\infty} + (1 + \frac{ah}{2}) \| L_{\Delta}\underset{\sim}{y} \|_1$$

where

$$\| L_{\Delta}\underset{\sim}{y} \|_1 := \sum_{j=1}^{N-1} h_j |(L_{\Delta}\underset{\sim}{y})_j|$$

and

$$h = \max_j h_j .$$

Combining these estimates when $h < 2/(a+a^2)$ gives

(2.1) $\qquad \|\underset{\sim}{y}\|_\infty := \max_j |y_j| \leq \beta(h) \ |\mathcal{B}_\Delta y| + \alpha(h) \ \|L_\Delta y\|_1$

where

$$\beta(h) = 2\beta/(2 - (a+a^2)h) \ ,$$

$$\alpha(h) = \alpha(2+ah)/(2-(a+a^2)h) \ .$$

Equation (2.1) is the analogue of (1.6) and clearly the condition numbers $\alpha(h)$ and $\beta(h)$ of the finite difference scheme are closely related to the condition numbers $\alpha$ and $\beta$ of the continuous problem. Indeed in the present example we have $\alpha(h) \to \alpha$ and $\beta(h) \to \beta$ as $h \to 0$ .

Although the above analysis has been derived for a particularly simple finite difference scheme, the basic idea used can be generalised so that a wide variety of finite difference schemes can be analysed. Details can be found in de Boor, de Hoog and de Keller [1]. However the role played by the condition numbers $\alpha$ and $\beta$ in the above example is typical of the more general situation.


## 3. THE SHOOTING APPROACH TO BOUNDARY VALUE PROBLEMS

Suppose that $\tilde{\Phi} \in [L_p^1(0,1)]^{n \times n}$ is a fundamental solution of the differential operator $L$ defined by (1.1) and satisfies the boundary conditions

$$\tilde{\mathcal{B}}\tilde{\Phi} = \tilde{B}_0 \tilde{\Phi}(0) + \tilde{B}_1 \tilde{\Phi}(1) = I \ ,$$

$$\max\{|\tilde{B}_0|, \ |\tilde{B}_1|\} = 1 \ .$$

As in section 1 we can associate with these boundary conditions the Green's function

$$\tilde{G}(t,s) = \begin{cases} \tilde{\Phi}(t) \ \tilde{B}_0 \ \tilde{\Phi}(0) \ \tilde{\Phi}^{-1}(s) & t > s \\[2ex] -\tilde{\Phi}(t) \ \tilde{B}_1 \ \tilde{\Phi}(1) \ \tilde{\Phi}^{-1}(s) & t < s \end{cases}$$

and hence the stability constants

$$\tilde{\beta} = \|\tilde{\Phi}\|_\infty \ , \quad \tilde{\alpha} = \sup_t \left\{ \int_0^1 |\tilde{G}(t,s)|^q \ ds \right\}^{1/q} \ , \quad \frac{1}{p} + \frac{1}{q} = 1$$

so that

$$\|y\|_\infty \leq \tilde{\beta} |\tilde{B}y| + \tilde{\alpha} \|Ly\|_p \ .$$

Now let

$$\tilde{y}_p(t) = \int_0^1 \tilde{G}(t,s) \ f(s) \ ds \ .$$

Then the solution of (1.1), (1.2) can be written as

(3.1) $$y(t) = \tilde{\Phi}(t)c + \tilde{y}_p(t)$$

where  c  is determined from the linear equations

(3.2) $$(\mathcal{B}\tilde{\Phi})c = b - \mathcal{B}\tilde{y}_p \ .$$

The relevance of (3.1) and (3.2) is that if we can obtain numerical

approximations to  $\tilde{\Phi}$  and  $\tilde{y}_p$ , then the above procedure can be used to

obtain a numerical approximation to the solution of (1.1), (1.2) .  Of

course, this makes sense only if the calculation of  $\tilde{\Phi}$  and  $\tilde{y}_p$  is

substantially simpler than the calculation of  $\Phi$  and  $y_p$.  Thus the

choice almost invariably used is  $\tilde{B}_0 = I$ ,  $\tilde{B}_1 = 0$  which corresponds to

solving initial value problems for the solution of  $\tilde{\Phi}$  and  $\tilde{y}_p$ .

However, separable boundary conditions which correspond to the case when

rank$(\tilde{B}_0) = k$ ,  rank$(\tilde{B}_1) = n-k$  could also be used and there are a number

of good algorithms to solve such problems (see for example [4]).  The

difficulty here is the choice of appropriate separable boundary

conditions.

Before analysing the shooting method further, we require the following result.

LEMMA 3.1 *Let* $B$, $\tilde{B}$, $\Phi$ *and* $\tilde{\Phi}$ *be defined as above. Then*

$$(\tilde{B\Phi})^{-1} = \tilde{B}\Phi \ .$$

Proof  Since

$$\Phi(t) = \tilde{\Phi}(t)(\tilde{B\Phi})^{-1} = \tilde{\Phi}(t) \ \tilde{\Phi}^{-1}(0) \ \Phi(0) = \tilde{\Phi}(t) \ \tilde{\Phi}^{-1}(1) \ \Phi(1) \ ,$$

it follows that

$$(\tilde{B\Phi})^{-1} = \tilde{\Phi}^{-1}(0) \ \Phi(0) = \tilde{\Phi}^{-1}(1) \ \Phi(1) \ .$$

Thus

$$(\tilde{B\Phi})^{-1} = (\tilde{B}_0 \ \tilde{\Phi}(0) + \tilde{B}_1 \ \tilde{\Phi}(1)) \ \tilde{\Phi}^{-1}(0) \ \Phi(0)$$

$$= \tilde{B}_0 \ \Phi(0) + \tilde{B}_1 \ \tilde{\Phi}(1) \ \tilde{\Phi}^{-1}(0) \ \Phi(0)$$

$$= \tilde{B}_0 \ \Phi(0) + \tilde{B}_1 \ \tilde{\Phi}(1) \ \tilde{\Phi}^{-1}(1) \ \Phi(1)$$

$$= \tilde{B}\Phi \ . \hspace{3cm} \#$$

The above lemma enables us to obtain a simple bound on the condition number of the matrix $\tilde{B\Phi}$ . Clearly,

$$K(\tilde{B\Phi}) := |(\tilde{B\Phi})^{-1}| \ |\tilde{B\Phi}|$$

$$= |\tilde{B}\Phi| \ |\tilde{B\Phi}|$$

$$\leq (|\tilde{\Phi}(0)| + |\tilde{\Phi}(1)|) \ (|\Phi(0)| + |\Phi(1)|)$$

$$\leq 4\beta\tilde{\beta} \ .$$

Thus if $\beta\tilde{\beta}$ is large, the solution of (3.2) may be poorly conditioned. In fact, even when the condition number $K(\tilde{B\Phi})$ is not large, the calculation may be unsatisfactory because if $\|\tilde{y}_p\|_\infty$ is substantially

larger than $\|y\|_\infty$ , the addition in (3.1) must involve the subtraction

of two nearly equal members and thus loss of significant digits.  In

order to ensure that this does not occur, we require $\tilde{\alpha}$ to be of modest

size.  This requirement means that the choice $\tilde{B}_0 = I$ , $\tilde{B}_1 = 0$ is not

satisfactory as it is a simple matter to construct examples for which $\alpha$

and $\beta$ are moderate and $\tilde{\alpha}$ , $\tilde{\beta}$ are very large.

The shooting method can also be used to analyse the rate of

convergence of a numerical scheme and also its stability.  To illustrate

this, suppose that we have approximations $\tilde{\Psi}$ and $\tilde{z}_p$ to $\tilde{\Phi}$ and $\tilde{y}_p$

respectively.  Following (3.1) and (3.2) we let

$$(3.3) \qquad \hat{y}(t) := \tilde{\Psi}(t) \, (B\tilde{\Psi})^{-1} \, (b - B\tilde{z}) + \tilde{z}(t)$$

be the approximation to $y$ (we are assuming here that $B\tilde{\Psi}$ is non-

singular).  Then we find that

LEMMA 3.2  *Let* $\tilde{\Psi}$ , $\tilde{\Phi}$ , $\tilde{y}_p$ , $\tilde{z}$ , $y$ , $\hat{y}$ *be defined as above.  Further-*
*more, let*

$$\tilde{E} := \tilde{\Psi} - \tilde{\Phi} , \qquad \tilde{e} = \tilde{z} - \tilde{y}_p$$

*and*

$$|\tilde{B}\tilde{\Phi}| \, |B\tilde{E}| < 1 .$$

*Then, the approximation* (3.3) *is well defined (in the sense that* $B\tilde{\Psi}$ *is*
*nonsingular) and satisfies*

$$|\hat{y}(t) - y(t)| \le \frac{1}{1 - |B\tilde{\Phi}| \, |B\tilde{E}|} \left\{ |\Phi(t)| \, (|B\tilde{E}| \, |\tilde{B}y| + |B\tilde{e}|) \right.$$

$$\left. + |\tilde{E}(t)| \, (|\tilde{B}y| + |\tilde{B}\Phi| \, |B\tilde{e}|) \right\} + |\tilde{e}(t)| .$$

Proof  Let

$$\Phi(t) = \tilde{\Phi}(t) \, (\mathcal{B}\tilde{\Phi})^{-1} \, , \qquad \Psi(t) = \tilde{\Psi}(t) \, (\mathcal{B}\tilde{\Psi})^{-1} \, .$$

Then,

$$\left| \hat{y}(t) - y(t) \right| \leq \left| (\Phi(t) - \Psi(t)) \, (b - \mathcal{B}\tilde{y}_p) \right| + \left| \Psi(t) \right| \, \left| \mathcal{B}\tilde{e} \right| + \left| \tilde{e}(t) \right|$$

and we now proceed to estimate the terms on the right hand side.  First we obtain

$$\left| \Psi(t) \right| = \left| \tilde{\Psi}(t) \, (\mathcal{B}\tilde{\Psi})^{-1} \right|$$

$$= \left| (\tilde{\Phi}(t) + \tilde{E}(t)) \, (\mathcal{B}\tilde{\Phi})^{-1} \, (I + \mathcal{B}\tilde{E}(\mathcal{B}\tilde{\Phi})^{-1})^{-1} \right|$$

$$\leq \frac{\left| \Phi(t) \right| + \left| \tilde{E}(t) \right| \, \left| \mathcal{B}\tilde{\Phi} \right|}{1 - \left| \mathcal{B}\tilde{E} \right| \, \left| \mathcal{B}\tilde{\Phi} \right|}$$

the last being obtained using lemma 3.1.  Also

$$b - \mathcal{B}\tilde{y}_p = (\mathcal{B}\tilde{\Phi}) \, (\mathcal{B}\tilde{\Phi})^{-1} \, (b - \mathcal{B}\tilde{y}_p)$$

$$= (\mathcal{B}\tilde{\Phi}) \, (\mathcal{B}\tilde{\Phi}(b - \mathcal{B}\tilde{y}_p))$$

$$= \mathcal{B}\tilde{\Phi} \, \, \tilde{\mathcal{B}y}$$

and hence

$$\left| (\Phi(t) - \Psi(t)) \, (b - \mathcal{B}\tilde{y}_p) \right| \leq \left| \tilde{\Phi}(t) - \tilde{\Psi}(t) \, (\mathcal{B}\tilde{\Psi})^{-1} \, \mathcal{B}\tilde{\Phi} \right| \, \left| \mathcal{B}y \right|$$

$$\leq \left| \tilde{\Phi}(t) \, (1 - (\mathcal{B}\tilde{\Phi} + \mathcal{B}\tilde{E})^{-1} \, \mathcal{B}\tilde{\Phi}) \right| \, \left| \mathcal{B}\tilde{y} \right|$$

$$+ \left| \tilde{E}(t) \right| \, \left| (\mathcal{B}\tilde{\Phi} + \mathcal{B}\tilde{E})^{-1} \, \mathcal{B}\tilde{\Phi} \right| \, \left| \mathcal{B}y \right|$$

$$= \left| \Phi(t) \, \mathcal{B}\tilde{E}(I + (\mathcal{B}\tilde{\Phi})^{-1} \, \mathcal{B}\tilde{E})^{-1} \right| \, \left| \tilde{\mathcal{B}y} \right|$$

$$+ \left| \tilde{E}(t) \right| \, \left| (I + (\mathcal{B}\tilde{\Phi})^{-1} \, \mathcal{B}\tilde{E})^{-1} \right| \, \left| \tilde{\mathcal{B}y} \right|$$

$$\leq \frac{\left| \Phi(t) \right| \, \left| \mathcal{B}\tilde{E} \right| \, \left| \tilde{\mathcal{B}y} \right| + \left| \tilde{E}(t) \right| \, \left| \tilde{\mathcal{B}y} \right|}{1 - \left| \mathcal{B}\tilde{E} \right| \, \left| \mathcal{B}\tilde{\Phi} \right|} \, .$$

The result now follows on combining the above estimates.        #

An immediate consequence of the above lemma is

COROLLARY 3.1  *Let* $\|\tilde{E}\|_\infty < \varepsilon_1$ *and* $\|\tilde{e}\| < \varepsilon_2$ . *Then*

$$\|y - \hat{y}\|_\infty \leq \frac{2}{1 - 4\beta\varepsilon_1} \left\{ (2\beta + 1) \|y\|_\infty \varepsilon_1 + \beta(1 + 2\varepsilon_1)\varepsilon_2 \right\} + \varepsilon_2 \ . \qquad |$$

The thing to notice about the above bound is that it does not involve the condition numbers $\tilde{\alpha}$ and $\tilde{\beta}$ . They are however implicit in the quantities $\|\tilde{E}\|_\infty$ and $\|\tilde{e}\|_\infty$ because the magnitude of these when $\tilde{\psi}$ and $\tilde{z}_p$ are obtained by applying a numerical scheme to calculate $\tilde{\phi}$ and $\tilde{y}_p$ respectively will almost certainly depend on the condition numbers $\tilde{\alpha}$ and $\tilde{\beta}$ .

## 4. A RELATIONSHIP BETWEEN DICHOTOMY AND CONDITIONING

One reason why (1.1), subject to initial conditions, may be very poorly conditioned when (1.1), subject to (1.2), is well conditioned is because the fundamental solution may have both increasing and decreasing components.  In fact it has become almost traditional to assume that the solution space $S = \{Y(t)c \,|\, c \in \mathbb{R}^n\}$ can be split into $S = S_1 \oplus S_2$ such that the solutions in $S_1$ are 'decreasing' while the solutions in $S_2$ are 'increasing'.  Specifically we say that

DEFINITION  S *is dichotomic if there exist matrices* $T_1$ *and* $T_2$ *such that* rank $T_1 = k$ , rank $T_2 = n-k$ , $T_1 + T_2 = I$ *and*

$$\left| Y(t) \, T_1 \, Y^{-1}(s) \right| \leq \gamma \qquad t > s$$

(4.1)

$$\left| Y(t) \, T_2 \, Y^{-1}(s) \right| \leq \gamma \qquad t < s$$

*for some constant* $\gamma$ .

If  S  is dichotomic then on defining

$$S_1 := \{Y(t)T_1 c \mid c \in \mathbb{R}^n\} \ ,$$

$$S_2 := \{Y(t)T_2 c \mid c \in \mathbb{R}^n\}$$

and noting that $T_1$ and $T_2$ are projections (i.e. $T_1 = T_1^2$ , $T_2 = T_2^2$) we find that for $\phi \in S_1$

$$\frac{|\phi(t)|}{|\phi(s)|} \leq \max_{c \in \mathbb{R}^n} \frac{|Y(t)T_1 c|}{|Y(s)T_1 c|} \leq |Y(t)T_1 \ Y^{-1}(s)| < \gamma$$

for $t > s$ while for $\phi \in S_2$ .

$$\frac{|\phi(t)|}{|\phi(s)|} \leq \max_{c \in \mathbb{R}^n} \frac{|Y(t)T_2 c|}{|Y(s)T_2 c|} \leq |Y(t)T_2 \ Y^{-1}(s)| < \gamma$$

for $t < s$ . Thus, $S_1$ does indeed correspond to the 'decreasing' solutions (i.e. solutions that do not increase too much) while $S_2$ corresponds to the 'increasing' solutions.

Of course it is always possible to find projections $T_1$ and $T_2$ such that (4.1) holds (for example $T_1 = I$ , $T_2 = 0$). However we would like to find $T_1$ and $T_2$ such that $\gamma$ is not too large when the stability constants $\alpha$ and $\beta$ are moderate. This is straightforward when the boundary conditions are separable.

LEMMA 4.1 *Let* rank $B_0 = k$ , rank $B_1 = n-k$ *and*

$$\alpha = \sup_{t,s} |G(t,s)|$$

*Then* (4.1) *holds with* $T_1 = \Phi(0) B_0$ , $T_2 = \Phi(0) B_1 \Phi(1) \Phi^{-1}(0)$ $\gamma = \alpha$ .

Proof The result follows on noting that

$$Y(t) = \Phi(t) \ \Phi^{-1}(0) \ ,$$

$$Y(t)T_1 \ Y^{-1}(s) = G(t,s) \ , \qquad t > s$$

$$Y(t)T_2 \ Y^{-1}(s) = -G(t,s) \ , \qquad t < s$$

$$\text{rank}(T_1) = \text{rank}(B_0) = k$$

$$\text{rank}(T_2) = \text{rank}(B_1) = n-k$$

and

$$T_1 + T_2 = I \ . \qquad\qquad \#$$

Thus if we can find separable boundary conditions such that the resulting problem is well conditioned we can immediately establish dichotomy. Such boundary conditions have been derived by de Hoog and Mattheij [2]. Let

$$Y(1) = \Phi(1) \ \Phi^{-1}(0) =: UDV^T$$

where U and V are orthogonal matrices and

$$D = \text{diag}(1/d_1,\ldots,1/d_k, \ d_{k+1},\ldots,d_n)$$

with $0 < d_i \le 1$ , $i = 1,\ldots,n$ . Now define

$$P_1 = \left( \begin{array}{c|c} 0 & 0 \\ \hline 0 & I_{n-k} \end{array} \right) \ , \qquad P_2 = \left( \begin{array}{c|c} I_k & 0 \\ \hline 0 & 0 \end{array} \right)$$

(4.2) $$\tilde{B}_0 = P_1 V^T \ , \qquad \tilde{B}_1 = P_2 U^T \ .$$

Then, it is shown in [2] that

LEMMA 4.2 *Let* $\tilde{B}_0$ *and* $\tilde{B}_1$ *be defined by* (4.2). *Then*

$$\|y\|_\infty \le \tilde{\beta} |\tilde{B}y| + \tilde{\alpha} \|Ly\|_1$$

*with* $\tilde{\alpha} \le \alpha + 4\alpha^2$ *and* $\tilde{\beta} < 4\alpha$ .

Thus, we find that the constant $\gamma$ in (4.1) can be bounded in terms of the stability constant for the continuous problem.

# REFERENCES

[1]    de Boor, C., de Hoog, F., and de Keller, H.B.,  "The Stability of
       one-step schemes for first-order two-point boundary value
       problems", *SIAM J. Numer. Anal.* 20 (1983), 1139-1146.

[2]    de Hoog, F., and Mattheij, R.,  "On dichotomy and well-
       conditioning in boundary value problems", submitted for
       publication.

[3]    Keller, H.B.,  "Numerical solution of two-point boundary value
       problems", *SIAM Regional Conference Series* 24, Philadelphia,
       1976.

[4]    Scott, M.R., and Watts, H.A.,  "Computational solution of linear
       two point boundary value problems via orthogonalization", *SIAM J.
       Numer. Anal.* 14 (1977), 40-70.

Department of Mathematics and Statistics,
C.S.I.R.O.,