

Improved Condition Number for Spectral Methods

By Wilhelm Heinrichs

Abstract. For the known spectral methods (Galerkin, Tau, Collocation) the condition number behaves like $O(N^4)$ (N : maximal degree of polynomials). We introduce a spectral method with an $O(N^2)$ condition number. The advantages with respect to propagation of rounding errors and preconditioning are demonstrated. A direct solver for constant coefficient problems is given. Extensions to variable coefficient problems and first-order problems are discussed. Numerical results are presented, showing the effectiveness of our methods.

1. Introduction. Spectral methods involve representing the solution to a problem in terms of a truncated series of smooth global functions. For Dirichlet problems the Chebyshev polynomials are the trial functions. It turns out that for the standard spectral methods (Galerkin, tau, collocation) the condition number is very large and grows as $O(N^4)$ (N : maximal degree of polynomials) (see Orszag [16]).

As a consequence we observe for direct solvers a strong propagation of rounding errors. For iterative methods an effective preconditioning is necessary (see Phillips et al. [17] or Heinrichs [8], [9]). A direct solver based on an ADI-factorization technique is proposed by Haidvogel and Zang [7], where the tau method has been used for discretization.

We derive an improved spectral method with an $O(N^2)$ condition number which is also known from finite difference and finite element methods. The main idea is that we employ polynomials fulfilling the homogeneous boundary conditions. The Laplace operator (or any other elliptic operator) applied to a truncated series of these trial functions is then developed in a series of Chebyshev polynomials. The coefficients of this series are taken to be equal to the coefficients of the right-hand side expansion. The resulting spectral system has the improved condition number. Obviously, the derived approximation is identical to the approximation of Lanczos' tau method. A comparison with the approximation of [7] is made in Section 3.

The described treatment can also be applied to the Bubnov-Galerkin method where polynomials with homogeneous boundary conditions are trial functions and the usual polynomials are test functions. A disadvantage of this approach is the fact that the Galerkin systems now become nonsymmetric. Since, in addition, the use of Fast Fourier Transforms (FFT's) is no longer possible, we prefer the treatment given here.

In Sections 2 and 3 we describe the method for one- and two-dimensional constant coefficient problems. In Section 4 we propose an efficient elimination process

Received December 29, 1987; revised April 7, 1988.

1980 *Mathematics Subject Classification* (1985 Revision). Primary 65N30, 65N35.

Key words and phrases. Spectral methods, condition number, direct solver, iterative methods, elliptic problems, first-order problems.

for the corresponding spectral systems. Subsequently, the system is solved by a block-Gauss elimination procedure. In Section 5 we investigate the case of inhomogeneous boundary conditions and in Section 6 we give a short convergence analysis. Numerical results for some test examples are presented in Section 7. In Sections 8 and 9 we show how this treatment can be extended to variable coefficient and first-order problems. For variable coefficient problems we employ a collocation method for discretization and propose a suitable iterative solver for these systems. Furthermore it is indicated that a straightforward extension of these ideas to time-dependent problems is not possible. We expect that the proposed treatment will find further applications in fluid dynamics and yield similar improvements.

2. One-Dimensional Case. We consider the boundary value problem

$$(2.1) \quad -u'' = f \quad \text{on } \Omega = (-1, 1)$$

with homogeneous boundary conditions $u(-1) = u(1) = 0$. For the spectral approximation of (2.1) we choose a basis, given by $\{(1-x^2)T_k(x) : k = 0, \dots, N\}$, which fulfills the homogeneous boundary conditions. Here, T_k denotes the k th Chebyshev polynomial, i.e., $T_k(x) = \cos(k \arccos(x))$. By means of a simple calculation we get the following formula.

LEMMA 2.1. *We have, for arbitrary constants a_p ,*

$$-\left(\sum_{p=0}^N a_p(1-x^2)T_p(x)\right)'' = \sum_{p=0}^N b_p T_p(x),$$

where

$$b_p = (p+1)(p+2)a_p + \alpha_p^{-1} \sum_{\substack{q=p+2 \\ q+p \text{ even}}}^N (6q)a_q \quad \text{and} \quad \alpha_p = \begin{cases} 2 & \text{for } p = 0, \\ 1 & \text{for } p \geq 1. \end{cases}$$

Proof. Since

$$-((1-x^2)T_p(x))'' = -(1-x^2)T_p''(x) + 4xT_p'(x) + 2T_p(x),$$

we have by standard formulas (see [6]):

$$(1-x^2) \sum_{p=0}^N a_p T_p = \sum_{p=0}^{N+2} e_p T_p, \quad \text{where } e_p = -\frac{1}{4} \alpha_{p-2} a_{p-2} \\ + \left(1 - \frac{1}{4}(\alpha_p + \alpha_{p-1})\right) a_p - \frac{1}{4} a_{p+2},$$

$$x \sum_{p=0}^N a_p T_p = \sum_{p=0}^{N+1} f_p T_p, \quad \text{where } f_p = \frac{1}{2}(\alpha_{p-1} a_{p-1} + a_{p+1})$$

and

$$\sum_{p=0}^N a_p T_p'' = \sum_{p=0}^{N-2} g_p T_p, \quad \text{where } g_p = \alpha_p^{-1} \sum_{\substack{k=p+2 \\ k+p \text{ even}}}^N k(k^2 - p^2) a_k,$$

$$\sum_{p=0}^N a_p T_p' = \sum_{p=0}^{N-1} h_p T_p, \quad \text{where } h_p = \alpha_p^{-1} \sum_{\substack{k=p+1 \\ k+p \text{ odd}}}^N (2k) a_k.$$

The operator $-\Delta$ applied to a basis function gives

$$-\Delta(1-x^2)(1-y^2)T_k(x)T_l(y) = -(((1-x^2)T_k(x))''(1-y^2)T_l(y) + ((1-y^2)T_l(y))''(1-x^2)T_k(x)).$$

The Chebyshev expansion of $-\Delta u_N$ (u_N is the spectral approximation) is known from the formulas in the proof of Lemma 2.1. It becomes obvious that now the Chebyshev polynomials T_{N+1} and T_{N+2} occur in the representation. In order to get an equation system for Chebyshev polynomials of degree $\leq N$, we neglect these terms and approximate

$$(1-x^2) \sum_{p=0}^N a_p T_p$$

by

$$\sum_{p=0}^N e_p T_p, \quad \text{where } e_p = -\frac{1}{4}\alpha_{p-2}a_{p-2} + \left(1 - \frac{1}{4}(\alpha_p + \alpha_{p-1})\right)a_p - \frac{1}{4}a_{p+2},$$

α_p as in Lemma 2.1. Let $E \in \mathbf{R}^{N+1, N+1}$ denote the matrix which represents the connection between a and e , i.e., $e = Ea$. Obviously, E is a tridiagonal matrix with positive diagonal entries and nonpositive off-diagonal elements. E is irreducible and diagonally dominant, hence an M -matrix, i.e., $E^{-1} \geq 0$ (all elements of E^{-1} are greater or equal zero) (see also Meis and Marcowitz [14, Theorem 13.16]).

Now the spectral matrix A in the coefficient space can be written by means of the tensor product \otimes as follows:

$$A = B \otimes E + E \otimes B, \quad B \otimes E = (Be_{i,j})_{i,j=0,\dots,N},$$

where B is defined in (2.2). The Fourier coefficients $\hat{f}_{p,q}$ of f can be calculated by FFT's and the corresponding linear system is

$$Aa = \hat{f}.$$

The spectral approximation, given by

$$u_N = (1-x^2)(1-y^2) \sum_{p,q=0}^N a_{p,q} T_p(x) T_q(y)$$

can now be easily evaluated by FFT's and scaling. The derived approximation u_N is identical to the approximation of Lanczos' tau method. It fulfills the homogeneous boundary conditions and is determined by matching the Chebyshev coefficients. In comparison with [7], we observe better approximation results (see Table 4) since we employ an approximation with polynomials of degree $N + 2$ (in each direction). In contrast to [7], matching is done using all Chebyshev coefficients $\hat{f}_{p,q}$ ($p, q = 0, \dots, N$).

By using Gerschgorin's estimate it becomes obvious that the largest eigenvalues grow as $O(N^2)$. The smallest eigenvalues are calculated numerically by means of the QR-factorization (see Table 1). We note that they are strictly positive and have a fixed positive lower bound. Hence the condition number behaves like $O(N^2)$.

TABLE 1

Smallest and largest eigenvalues λ_{\min} and λ_{\max} of A

N	λ_{\min}	λ_{\max}	condition number	condition number/ N^2
4	3.12188	$2.99319 \cdot 10^1$	$9.39146 \cdot 10^0$	$5.86967 \cdot 10^{-1}$
8	3.13694	$9.45206 \cdot 10^1$	$3.01332 \cdot 10^1$	$4.70831 \cdot 10^{-1}$
16	3.13767	$3.74279 \cdot 10^2$	$1.19286 \cdot 10^2$	$4.65961 \cdot 10^{-1}$
32	3.13770	$1.61179 \cdot 10^3$	$5.13685 \cdot 10^2$	$5.01646 \cdot 10^{-1}$

Further insight into what happens may be gained by a local consideration of the analytical behavior of $-\Delta(T_k T_l)$, resp. $-\Delta((1-x^2)(1-y^2)T_k T_l)$. We get

$$(3.2) \quad -\Delta(T_k T_l) = \left(\frac{k^2}{1-x^2} + \frac{l^2}{1-y^2} \right) T_k T_l + \text{lower-order terms in } k, l$$

and

$$(3.3) \quad -\Delta((1-x^2)(1-y^2)T_k T_l) = (k^2(1-y^2) + l^2(1-x^2))T_k T_l + \text{lower-order terms in } k, l.$$

From the representation (3.2) it becomes clear that for $k \sim N$ ($l \sim N$) and $x \sim \pm 1 \mp O(N^{-2})$ ($y \sim \pm 1 \mp O(N^{-2})$) (as for the Chebyshev nodes) an N^4 term appears and leads to largest eigenvalues growing as N^4 . From (3.3) it is clear that the largest eigenvalues grow as N^2 . On the other hand, the smallest eigenvalues do not behave like N^{-2} since the corresponding influence only comes from the four edges. Hence, there is no global effect, and the smallest eigenvalues do not tend to zero.

We remark that the improved condition number can also be attained by preconditioning the usual spectral matrix by a diagonal matrix with diagonal entries equal to $(1-x_i^2)(1-y_j^2)$, $i, j = 1, \dots, N-1$. This observation becomes of great interest for iterative solvers and will be discussed in Section 8.

4. The Direct Solver. Now we describe an efficient block Gauss elimination procedure for the spectral matrix A . Obviously, A has the following block Hessenberg structure:

$$\begin{pmatrix} C & 0 & C & 0 & P & 0 & P & 0 & \dots & 0 & P \\ 0 & C & 0 & C & 0 & P & 0 & \dots & \dots & P & 0 \\ D & 0 & C & 0 & C & 0 & P & \dots & \dots & 0 & P \\ & D & 0 & C & 0 & C & 0 & \dots & \dots & P & 0 \\ & & \dots & \dots & \dots & \dots & \dots & \dots & \dots & 0 & P \\ & & & D & 0 & C & 0 & C & 0 & P & 0 \\ & & & & D & 0 & C & 0 & C & 0 & P \\ & & & & & D & 0 & C & 0 & C & 0 \\ & & & & & & D & 0 & C & 0 & C \\ & & & & & & & D & 0 & C & 0 \\ \mathbf{0} & & & & & & & & D & 0 & C & 0 \\ & & & & & & & & & D & 0 & C \end{pmatrix}.$$

Here D , P and C denote submatrices with the following structures:

- D has the same structure as B (upper triangular matrix),
- P has the same structure as E (pentadiagonal matrix),
- C has the same structure as $B + E$ (Hessenberg structure).

We now explicitly give the matrices B and E :

$$B = \begin{pmatrix} 2 & 0 & 6 & 0 & 12 & 0 & 18 & 0 & 24 & \cdots & 0 & 3N \\ & 6 & 0 & 18 & 0 & 30 & 0 & 42 & 0 & \cdots & 6(N-1) & 0 \\ & & 12 & 0 & 24 & 0 & 36 & 0 & 48 & \cdots & 0 & 6N \\ & & & 20 & 0 & 30 & 0 & 42 & 0 & \cdots & 6(N-1) & 0 \\ & & & & 30 & 0 & 36 & 0 & 48 & \cdots & 0 & 6N \\ & & & & & 42 & 0 & 42 & 0 & \cdots & 6(N-1) & 0 \\ & & & & & & 56 & 0 & 48 & \cdots & 0 & 6N \\ & & & & & & & 72 & 0 & \cdots & 6(N-1) & 0 \\ & \mathbf{0} & & & & & & & 90 & \ddots & \vdots & \vdots \\ & & & & & & & & & \ddots & 0 & 6N \\ & & & & & & & & & & N(N+1) & 0 \\ & & & & & & & & & & & (N+1)(N+2) \end{pmatrix}.$$

$$E = \begin{pmatrix} 1/2 & 0 & -1/4 & & & & & & & & & \mathbf{0} \\ 0 & 1/4 & 0 & -1/4 & & & & & & & & \\ -1/2 & 0 & 1/2 & 0 & -1/4 & & & & & & & \\ & -1/4 & 0 & 1/2 & 0 & -1/4 & & & & & & \\ & & -1/4 & 0 & 1/2 & 0 & -1/4 & & & & & \\ & & & \ddots & \ddots & \ddots & \ddots & \ddots & & & & \\ & & & & -1/4 & 0 & 1/2 & 0 & -1/4 & & & \\ \mathbf{0} & & & & & -1/4 & 0 & 1/2 & 0 & & & \\ & & & & & & -1/4 & 0 & 1/2 & & & \\ & & & & & & & -1/4 & 0 & 1/2 & & \end{pmatrix}.$$

Direct elimination for the matrix A requires storage of about $O(N^3)$ elements; for Gauss elimination about $O(N^6)$ arithmetic operations are necessary. By utilizing the block Hessenberg structure of A this can be reduced to about $O(N^5)$ operations. By a special elimination we show that A can be transformed to a band matrix. First we consider the matrix B . By subtracting from row 0 half of row 2 and further from row i the row $i+2$ for $i = 1, 2, \dots, N-2$, we get a band matrix TB (T represents the elimination process). TB has the following structure:

$$TB = \begin{pmatrix} 2 & 0 & 0 & & & & & & & & & \\ & 6 & 0 & -2 & & & & & & & & 0 \\ & & 12 & 0 & -6 & & & & & & & \\ & & & \ddots & & \ddots & & & & & & \\ & & & & 0 & & \ddots & & & & & \\ & & & & & (i+1)(i+2) & 0 & 6(i+2) - (i+3)(i+4) & & & & \\ & & & & & & \ddots & & & & & \\ & & & & & & & 0 & & & & \\ & & & & & & & & (N-1)N & & & \\ \mathbf{0} & & & & & & & & & 0 & 6N - (N+1)(N+2) & \\ & & & & & & & & & N(N+1) & 0 & \\ & & & & & & & & & & 0 & (N+1)(N+2) \end{pmatrix}.$$

TABLE 2
Smallest and largest eigenvalues λ_{\min} and λ_{\max} of H

N	λ_{\min}	λ_{\max}	condition number	condition number/ N^2
4	$1.97075 \cdot 10^0$	$3.38728 \cdot 10^1$	$1.71878 \cdot 10^1$	$1.07424 \cdot 10^0$
8	$1.11919 \cdot 10^0$	$2.20707 \cdot 10^2$	$1.97202 \cdot 10^2$	$3.08128 \cdot 10^0$
16	$5.55226 \cdot 10^{-1}$	$8.87355 \cdot 10^2$	$1.59819 \cdot 10^3$	$6.24298 \cdot 10^0$
32	$1.84507 \cdot 10^{-1}$	$3.77197 \cdot 10^3$	$2.04435 \cdot 10^4$	$1.99644 \cdot 10^1$

We remark that by an alternating direction implicit (ADI) method (proposed, e.g., by Haidvogel and Zang [7]) the amount of work can be reduced to $O(N^3)$ operations. But for an efficient determination of the parameters of ADI the eigenvalues have to be calculated. These pre-computations are already quite costly, and hence the total amount of work is comparable to our direct solver. Furthermore, for our method standard routines from programming libraries are available.

We finally mention that this efficient direct solver can also be used to derive defect corrections for variable-coefficient problems. Here we refer to D'yakonov's method [16] where the Laplace operator is used for preconditioning.

5. Inhomogeneous Boundary Conditions. We now consider the Poisson equation with inhomogeneous boundary conditions, given by

$$(5.1) \quad -\Delta u = f \quad \text{on } \Omega, \quad u = g \quad \text{on } \partial\Omega.$$

We describe a collocation method for problem (5.1) where the boundary conditions are collocated, too. First, we determine a smooth function u_N^1 which satisfies

$$u_N^1(x_i, x_j) = g(x_i, x_j) \quad \text{for } (x_i, x_j) \in \partial\Omega.$$

Usually, u_N^1 is also chosen as a suitable polynomial of degree $\leq N$ in x, y . Further, let u_N^0 denote the collocation approximation (as in Section 3) for the solution u^0 of the problem

$$-\Delta u^0 = f + \Delta u_N^1, \quad u^0 = 0 \quad \text{on } \partial\Omega.$$

Then $u_N = u_N^1 + u_N^0$ is a collocation approximation for the solution of problem (5.1). The approximation error of this method can be estimated by the interpolation error of u_N^1 on the boundary $\partial\Omega$ and the collocation error for u_N^0 . A detailed analysis of convergence can be found in [11].

We remark that the above treatment can be extended to variable coefficient problems in a straightforward manner.

6. Convergence. We give a proof of convergence in the one-dimensional case. The pseudospectral approximation $u_N \in \mathbf{P}_{N+2}^0$ can be equivalently written as the solution of the discrete problem

$$a_{N,\omega}(u_N, v_N) = (f, v_N)_{N,\omega} \quad \text{for all } v_N \in \mathbf{P}_N,$$

where $\mathbf{P}_{N+2}^0 = \{u_N = (1 - x^2)z_N : z_N \in \mathbf{P}_N\}$ and \mathbf{P}_N denotes the space of polynomials of degree $\leq N$. Here, $a_{N,\omega}$ and $(\cdot, \cdot)_{N,\omega}$ are given by

$$a_{N,\omega}(u, v) = (-u'', v)_{N,\omega} \quad \text{and} \quad (v, w)_{N,\omega} = \sum_{j=0}^N v(x_j)w(x_j)\omega_j,$$

where $x_j = \cos(j\pi/N)$ and $\omega_0 = \omega_N = \pi/2N$, $\omega_j = \pi/N$ ($j = 1, \dots, N - 1$) are the weights of the Chebyshev-Gauss-Lobatto quadrature formula. Convergence will be proved by using variational principles (see the Babuška/Nečas conditions [1], [15]). Here the stability and convergence of the spectral scheme is assured, provided the form $a_{N,\omega}$ fulfills some properties of continuity and coerciveness (see Canuto and Quarteroni [3]). The main problem is to show the coerciveness; we show that for each $u_N \in \mathbf{P}_{N+2}^0$ there exists a $v_N \in \mathbf{P}_N$ such that

$$(6.1) \quad a_{N,\omega}(u_N, v_N) \geq c \|u_N\|_V \|v_N\|_W, \quad c > 0,$$

where $V = H_0^{1,\omega} \cap H^{2,\omega}$ with the norm of $H^{2,\omega}$ and $W = H^{0,\omega} = L^\omega$. ($H_0^{1,\omega}$ denotes the subspace of functions from $H^{1,\omega}$ with compact support; L^ω the space of ω -integrable functions where $\omega(x) = (1 - x^2)^{-1/2}$.) Since $u_N v_N \in \mathbf{P}_{2N}$, and in view of the equivalence of discrete and continuous norms, we get for $v_N = u_N''$:

$$\begin{aligned} |a_{N,\omega}(u_N, v_N)| &= \left| \sum_{j=0}^N u_N''(x_j) u_N''(x_j) \omega_j \right| \geq c \int_{-1}^1 |u_N''(x)|^2 \omega(x) dx \\ &= c \|v_N\|_\omega^2 \geq c' \|v_N\|_{L^\omega} \|u_N\|_{H^{2,\omega}}. \end{aligned}$$

The last inequality follows from Poincaré's theorem. Hence we get the estimate (see [3, Theorem 1.2]):

$$\begin{aligned} \|u - u_N\|_{2,\omega} \leq \inf_{w_N \in \mathbf{P}_{N+2}^0} & \left\{ c_1 \|u - w_N\|_{2,\omega} \right. \\ & \left. + c_2 \sup_{v_N \in \mathbf{P}_N} \left[\frac{|a_{N,\omega}(w_N, v_N) - a_\omega(w_N, v_N)|}{\|v_N\|_\omega} \right. \right. \\ & \left. \left. + \frac{|(f, v_N)_\omega - (f, v_N)_{N,\omega}|}{\|v_N\|_\omega} \right] \right\}. \end{aligned}$$

For $u \in H^{s,\omega}$, $s \geq 2$, the approximation error can be estimated by (see [2], [4])

$$\|u - w_N\|_{2,\omega} \leq cN^{2-s} \|u\|_{s,\omega}.$$

Since $a_{N,\omega}(w_N, v_N) - a_\omega(w_N, v_N) = (-w_N'', v_N)_{N,\omega} - (-w_N'', v_N)_\omega$, and using the estimate of [3, Lemma 2.5], resp. [4, Lemma 3.2], we get

$$\begin{aligned} |(-w_N'', v_N)_{N,\omega} - (-w_N'', v_N)_\omega| &\leq C \|v_N\|_\omega \|w_N'' - P_{N-1} w_N''\|_\omega \\ &\leq CN^{2-s} \|w_N\|_{s,\omega} \|v_N\|_\omega. \end{aligned}$$

Here we have used $P_c w_N'' = w_N''$ (P_c denotes the interpolation operator onto \mathbf{P}_N) and the approximation property of the orthogonal projection P_{N-1} (see the approximation results of Canuto and Quarteroni [2]). By the same arguments we get

$$\frac{|(f, v_N)_\omega - (f, v_N)_{N,\omega}|}{\|v_N\|_\omega} \leq CN^{-\sigma} \|f\|_{\sigma,\omega} \quad \text{for } f \in H^{\sigma,\omega}, \sigma > 1/2.$$

Altogether, we therefore conclude that

$$\|u - u_N\|_{2,\omega} \leq CN^{2-s} \|u\|_{s,\omega} + C' N^{-\sigma} \|f\|_{\sigma,\omega}$$

for $u \in H^{s,\omega}$, $s \geq 2$ and $f \in H^{\sigma,\omega}$, $\sigma > 1/2$. Hence, the proposed method allows stronger convergence estimates than those obtained for the usual collocation method.

Further convergence estimates in strong C -norms are given in [12]. There, we deduce estimates such as

$$\|u - u_N\|_{C^2} \leq C \ln NN^{2-s} \|u\|_{C^s}, \quad s \geq 2,$$

where $\|\cdot\|_{C^s}$ denotes the strong norm of $C^s[-1, 1]$.

Convergence estimates for the corresponding Galerkin (-Petrov) methods are given by Krasnosel'skii et al. [13]. Theorem 16.6 of [13] yields asymptotically optimal estimates in the norms of $H^{2,2}$ and C^1 .

7. Examples. We examine numerical examples introduced in [7], [8], [9] and compare the collocation method with the proposed spectral method. For this purpose we calculate the absolute discretization error, measured in the pointwise maximum norm ($\|u_N - u\|_{\text{MAX}}$). We implemented both methods on a Siemens 7.570-P computer and used double-precision arithmetic with an accuracy of 14 digits. The first example is given by

$$(7.1) \quad -\Delta u = \frac{\pi^2}{2} \cos\left(\frac{\pi}{2}x\right) \cos\left(\frac{\pi}{2}y\right) \quad (\text{on } \Omega), \quad u = 0 \quad (\text{on } \partial\Omega)$$

with the exact solution $u = \cos(\frac{\pi}{2}x) \cos(\frac{\pi}{2}y)$.

TABLE 3

$\ u - u_N\ _{\text{MAX}}$ for example (7.1)		
N	Collocation	Our Method
16	$8.08 \cdot 10^{-13}$	$6.11 \cdot 10^{-16}$
32	$1.59 \cdot 10^{-12}$	$1.12 \cdot 10^{-15}$

Inspection of the results in Table 3 shows that the collocation method strongly propagates rounding errors and its inherently high accuracy is somewhat disturbed. Here our method yields an increase in accuracy of about 2-3 digits, which is due to the smaller condition number. The second example is given by

$$(7.2) \quad -\Delta u = 1 \quad (\text{on } \Omega), \quad u = 0 \quad (\text{on } \partial\Omega),$$

where the exact solution can be written as an infinite series [7]. Now collocation is more accurate, which may be due to the fact that we neglected the highest modes in our method. Here, the phenomena of the first example cannot be observed since the discretization error is too large to be influenced by rounding errors. The numerical results for this example are presented in Table 4. For comparison, the corresponding results of [7] are also given.

For reasons mentioned in Section 3, our method yields a higher accuracy.

TABLE 4

$\ u - u_N\ _{\text{MAX}}$ for example (7.2)			
N	Collocation	Our Method	Tau Method [7]
16	$7.47 \cdot 10^{-7}$	$3.99 \cdot 10^{-6}$	$3.52 \cdot 10^{-5}$
32	$5.51 \cdot 10^{-8}$	$2.07 \cdot 10^{-7}$	$2.23 \cdot 10^{-6}$

8. Variable-Coefficient Problems. For nonlinear or nonconstant coefficient problems the presented tau method is no longer efficiently implementable (see also Gottlieb and Orszag [6, Section 10]). The corresponding spectral matrices are dense and the spectral systems cannot be solved efficiently by using FFT's. In this case, collocation (or pseudospectral) methods are recommended. The resulting spectral systems can be efficiently solved by means of iterative methods [7]–[10], [16], [17], [19]. For an iterative method the reduced condition number further leads to an improved convergence property (see, e.g., [14]).

However, a straightforward adaptation of the preceding ideas to collocation with the usual Chebyshev nodes (extrema of Chebyshev polynomials) is not possible. This can easily be seen by considering the expression (3.3). In the four edges, the expression $-\Delta((1-x^2)(1-y^2)T_k T_l)$ is always zero, and hence collocation at these points is not possible for elliptic problems where the right-hand side is not compatible with the boundary conditions. For instance, we refer to the example (7.2) of Section 7. This means that the Gauss-Lobatto nodes cannot be used, and the advantage of FFT is lost. However, Gauss nodes can still be used. We recommend the use of the Chebyshev polynomials for which fast transforms are still available (see [12]). But in comparison with FFT they still need twice the amount of work and are hardly competitive computationally.

Hence we prefer a somewhat different treatment. Let us consider a variable coefficient problem given by

$$(8.1) \quad Lu = -au_{xx} - bu_{yy} + cu_x + du_y + eu = f \quad \text{on } \Omega$$

with homogeneous boundary conditions on $\partial\Omega$. Here a, b, c, d, e and f denote given functions defined on Ω , where $a > 0$ and $b > 0$ on $\bar{\Omega}$. In place of (8.1) we numerically solve the modified singular problem

$$(8.2) \quad L_s u = g \quad \text{on } \Omega,$$

where

$$L_s u = (1-x^2)(1-y^2)Lu \quad \text{and} \quad g = (1-x^2)(1-y^2)f.$$

This can also be written as

$$(8.3) \quad \begin{aligned} & - (1-y^2)a((1-x^2)u_{xx}) - (1-x^2)b((1-y^2)u_{yy}) \\ & + (1-y^2)c((1-x^2)u_x) + (1-x^2)d((1-y^2)u_y) \\ & + (1-x^2)(1-y^2)eu = g. \end{aligned}$$

Now, for an iterative solution of spectral systems the preconditioned Richardson relaxation [9], [19] is recommended. For a more detailed description, denote

$$\bar{\Omega}_N = \left\{ (x_i, y_j) = \left(\cos \frac{i\pi}{N}, \cos \frac{j\pi}{N} \right) : i, j = 0, \dots, N \right\},$$

$\Omega_N = \Omega \cap \bar{\Omega}_N$, $\partial\Omega_N = \partial\Omega \cap \bar{\Omega}_N$ for $\Omega = (-1, 1)^2$, and let $G(\bar{\Omega}_N)$ be the space of grid functions defined on $\bar{\Omega}_N$. Let further $g_N \in G(\bar{\Omega}_N)$ with components $g_N^{i,j} = g(x_i, y_j)$

and $u_N^j \in G(\bar{\Omega}_N)$, $u_N^j = 0$ on $\partial\Omega_N$, be given. The iteration then proceeds as follows:

$$\begin{aligned} u_N^{j+1} &= u_N^j + \omega_N^j P_N(L_s^N u_N - g_N) \quad \text{on } \Omega_N, \\ u_N^{j+1} &= 0 \quad \text{on } \partial\Omega_N \text{ for } j = 0, 1, 2, \dots \end{aligned}$$

Here we denote by

- L_s^N the spectral discretization operator of L_s ,
- P_N the preconditioning operator,
- ω_N^j the relaxation parameter.

For an efficient relaxation these components have to be chosen in a suitable manner. A stable evaluation of $L_s^N u_N^j$ can be done using the splitting of formula (8.3). The evaluation of $(1-x^2)u_{xx}$, $(1-y^2)u_{yy}$ and $(1-x^2)u_x$, $(1-y^2)u_y$ can be accomplished in a stable way with a rounding error propagation of at most $O(N^2)$. This is an easy consequence of the following lemma:

LEMMA 8.1. *Let $v_N(x) = \sum_{p=0}^N a_p T_p(x)$; then*

$$(8.4) \quad -(1-x^2)v_N(x)'' = \sum_{p=0}^N \sigma_p T_p(x)$$

and

$$(8.5) \quad (1-x^2)v_N(x)' = \sum_{p=0}^{N+1} \tau_p T_p(x),$$

where

$$\sigma_p = p(p-1)a_p - \alpha_p^{-1} \sum_{\substack{k=p+2 \\ k+p \text{ even}}}^N (2k)a_k$$

and

$$\tau_p = \frac{p-1}{2}a_{p-1} + \frac{p+1}{2}a_{p+1} \quad \text{for } p = 0, 1, \dots, N+1.$$

Here, α_p is defined as in Lemma 2.1 and we set $a_{-1} = a_{N+1} = a_{N+2} = 0$.

Proof. The proof can easily be accomplished using the standard formulas noted in the proof of Lemma 2.1. \square

Since FFT's should be applicable, we use in (8.5) only the N -series and set $\tau_{N+1} = 0$. After an FFT into physical space, the multiplication by the coefficient functions $(1-y^2)a$, $(1-x^2)b$, $(1-y^2)c$, $(1-x^2)d$ and $(1-x^2)(1-y^2)e$ can be accomplished pointwise. In algebraic notation, this can be represented by a matrix multiplication. This is stable and only takes $O(N^2)$ arithmetic operations. Hence we get an efficient computation of $L_s^N u_N^j$ implementable with FFT's, which causes rounding error propagation of at most $O(N^2)$.

For an anisotropic test problem with coefficients

$$\begin{aligned} a(x, y) &= 1 + \varepsilon \exp(x), \\ b(x, y) &= 1 + \varepsilon \exp(y), \end{aligned} \quad \varepsilon = 0., 0.1, 1$$

and $c = a_x, d = b_y, e = 0$, we have numerically calculated the smallest and largest eigenvalues of L_s^N (see Table 5). The calculations were done by using certain variants of the power method. The results once more indicate an $O(N^2)$ condition number of L_s^N . For the Poisson equation, the smallest eigenvalues approximate $\lambda_{\min} = 3.138$, which should also be a good prediction for the smallest eigenvalue of the singular operator $L_s = -(1 - x^2)(1 - y^2)\Delta$.

TABLE 5
Smallest and largest eigenvalues of L_s^N

ε	N	λ_{\min}	λ_{\max}
0.	16	3.138	402.82
	32	3.138	1752.79
0.1	16	3.541	457.02
	32	3.542	1976.13
1.0	16	6.293	1002.90
	32	6.294	4247.37

From these considerations it can be deduced that pure Richardson relaxation without defect correction yields a convergence factor of $1 - O(N^{-2})$. Clearly, better factors can be obtained by using a suitable preconditioning. Here we recommend a correction based on the five-point finite difference discretization (see [9], [19]) of (8.2) at the Chebyshev nodes. Here the term $(1 - x^2)(1 - y^2)$ has to be incorporated into the difference formula. In the selfadjoint case with $c = a_x, d = b_y$ and $e = 0$, we make use of the formula in [9]. In $(x_i, y_j) \in \Omega_N$ we obtain

$$L_{FD}^{i,j} = -\beta \begin{bmatrix} 0 & \beta_{0,1}^{i,j} & 0 \\ \beta_{-1,0}^{i,j} & \beta_{0,0}^{i,j} & \beta_{1,0}^{i,j} \\ 0 & \beta_{0,-1}^{i,j} & 0 \end{bmatrix},$$

where $\beta = 1/(2s_{1/2}s_1)$ and

$$\begin{aligned} \beta_{0,1}^{i,j} &= -b \left(x_i, \frac{1}{2}(y_j + y_{j+1}) \right) \frac{s_i^2 s_j}{s_{j+1/2}}, \\ \beta_{0,-1}^{i,j} &= -b \left(x_i, \frac{1}{2}(y_j + y_{j-1}) \right) \frac{s_i^2 s_j}{s_{j-1/2}}, \\ \beta_{-1,0}^{i,j} &= -a \left(\frac{1}{2}(x_i + x_{i-1}), y_j \right) \frac{s_j^2 s_i}{s_{i-1/2}}, \\ \beta_{1,0}^{i,j} &= -a \left(\frac{1}{2}(x_i + x_{i+1}), y_j \right) \frac{s_j^2 s_i}{s_{i+1/2}}, \\ \beta_{0,0}^{i,j} &= -(\beta_{-1,0}^{i,j} + \beta_{1,0}^{i,j} + \beta_{0,-1}^{i,j} + \beta_{0,1}^{i,j}). \end{aligned}$$

Here we denote $s_i = \sin(i\pi/N)$, $s_{i+1/2} = \sin((i+1/2)\pi/N)$ for $i = 0, \dots, N$. We are now in a position to prove some nice properties of the corresponding discretization matrix L_{FD}^N .

LEMMA 8.2. *The matrix L_{FD}^N has real eigenvalues, the largest one growing as $O(N^2)$. In particular, for the Poisson equation, the smallest eigenvalues have a lower bound of $1/2$. Hence the condition number of L_{FD}^N grows only as $O(N^2)$.*

Proof. We observe that $L_{FD}^{i,j}$ can be written as $L_{FD}^{i,j} = d_{i,j} \hat{L}_{FD}^{i,j}$, $d_{i,j} = s_i s_j$, where $\hat{L}_{FD}^{i,j}$ is a symmetric star. Hence, in matrix notation, we obtain $L_{FD}^N = D_N \hat{L}_{FD}^N$, where $D_N = \text{diag}(d_{i,j})_{i,j=1,\dots,N-1}$. Now the original eigenvalue problem is equivalent to the symmetric problem

$$D_N^{1/2} \hat{L}_{FD}^N D_N^{1/2} e_N = \lambda e_N.$$

This proves that all eigenvalues are real. The fact that the largest eigenvalue grows as $O(N^2)$ easily follows from Gerschgorin type estimates. Since $L_{FD}^{i,j}$ has a Z -structure (i.e., the off-diagonal elements are nonpositive), we only have to find a vector $z_N = (z_N^{i,j})$ with $z_N^{i,j} > 0$ ($i, j = 1, \dots, N-1$) such that $(L_{FD}^N z_N)^{i,j} \geq \frac{1}{2} z_N^{i,j}$ ($i, j = 1, \dots, N-1$). Then it follows that the matrix L_{FD}^N is an M -matrix, and its eigenvalues have a lower bound of $1/2$ (see Schröder [18, Chapter II, Proposition 1.4]). Now we prove that the special vector z_N with $z_N^{i,j} = s_i^{1/2} s_j^{1/2}$ leads to the desired result. First, we consider points of Ω far away from the boundary. For large N we obtain (in obvious notation for stars)

$$L_{FD}^{i,j} = \frac{N^2}{\pi^2} \left(s_j^2 [-1 \ 2 \ -1] + s_i^2 \begin{bmatrix} -1 \\ 2 \\ -1 \end{bmatrix} \right).$$

If we show that

$$(8.6) \quad \frac{N^2}{\pi^2} ([-1 \ 2 \ -1] s^{1/2})_i \geq \frac{1}{4} s_i^{-3/2} \quad (i = 1, \dots, N-1) \text{ for } s^{1/2} = \sin^{1/2} x,$$

then it also follows that

$$(L_{FD}^N z_N)^{i,j} \geq \frac{1}{4} \left(\frac{s_j^2}{s_i^2} + \frac{s_i^2}{s_j^2} \right) s_i^{1/2} s_j^{1/2} \geq \frac{1}{2} s_i^{1/2} s_j^{1/2} = \frac{1}{2} z_N^{i,j}.$$

The last inequality easily follows from $x^2 + y^2 \geq 2xy$ for real x, y . (8.6) is a consequence of the fact that the finite difference operator is identical to the second-order difference star. Hence, the left-hand side in (8.6) approximates the second derivative of the function $-\sin^{1/2} x$ in $x_i = i\pi/N$. Since further

$$(-\sin^{1/2} x)' = \frac{1}{4} \sin^{-3/2} x + \frac{1}{4} \sin^{1/2} x \geq \frac{1}{4} \sin^{-3/2} x \quad \text{for } x \in [0, \pi],$$

we obtain the desired result. A special treatment of points near the boundary confirms the above result. \square

Numerical calculations show that the smallest eigenvalue also approximates $\lambda_{\min} = 3.138$ for increasing N .

The convergence factors after this preconditioning are equal to those obtained without multiplication by $(1 - x^2)(1 - y^2)$. Here we refer to results given in [9], where we have applied this relaxation in a multigrid context. It is still of great interest to find even better or less expensive defect corrections which utilize the improved condition number.

9. First-Order Problems. We consider the first-order problem

$$(9.1) \quad \begin{aligned} u' &= f \quad \text{on } \Omega = (-1, 1), \\ u(1) &= 0 \quad \text{or} \quad u(-1) = 0. \end{aligned}$$

Depending on the boundary conditions, we use a spectral approximation given by

$$u_N^\pm(x) = (1 \pm x) \sum_{p=0}^N a_p T_p(x).$$

The first-order derivative of u_N can be represented in the following way.

LEMMA 9.1. *There holds*

$$u_N^+(x)' = \sum_{p=0}^N r_p T_p(x) \quad \text{and} \quad u_N^-(x)' = \sum_{p=0}^N s_p T_p(x),$$

where

$$(9.2) \quad r_p = (p + 1)a_p + \alpha_p^{-1} \sum_{k=p+1}^N (2k)a_k$$

and

$$(9.3) \quad s_p = -(p + 1)a_p - \alpha_p^{-1} \sum_{k=p+1}^N (-1)^{k+p} (2k)a_k.$$

α_p is defined as in Lemma 2.1.

Proof. We have

$$u_N^\pm(x)' = \pm \sum_{p=0}^N a_p T_p(x) + (1 \pm x) \sum_{p=0}^{N-1} h_p T_p(x),$$

where $h_p = \alpha_p^{-1} \sum_{k=p+1; k+p \text{ odd}}^N (2k)a_k$. By further making use of the formulas given in the proof of Lemma 2.1, we easily obtain (9.2) and (9.3). \square

The matrices representing the relations (9.2) and (9.3) are called R and S . Obviously, R and S are upper triangular with real positive and negative eigenvalues $\pm(p + 1)$, $p = 0, \dots, N$. Again, we obtain largest eigenvalues increasing only as $O(N)$ instead of $O(N^2)$.

Next, we consider constant coefficient systems. We shall analyze the case of two equations. After diagonalization, such a system takes the form (see also Funaro [5])

$$(9.4) \quad \begin{aligned} u_x &= f \\ v_x &= g \end{aligned} \quad \text{on } \Omega = (-1, 1)$$

with

$$u(1) = \alpha v(1), \quad u(-1) = \beta v(-1),$$

where α, β are two real numbers.

Problem (9.4) admits a unique solution, provided α and β satisfy the condition

$$(9.5) \quad \alpha\beta \neq 1.$$

Here we use approximations given by

$$u_N = (1-x)w_N + \gamma_N \quad \text{and} \quad v_N = (1+x)z_N + \delta_N,$$

where $w_N, z_N \in \mathbf{P}_N$ and $\gamma_N, \delta_N \in \mathbf{R}$. The w_N and z_N are determined as the spectral tau approximation of the problems

$$((1-x)w)_x = f \quad \text{and} \quad ((1+x)z)_x = g.$$

Subsequently, the real numbers γ_N, δ_N are chosen in such a way that the boundary conditions are fulfilled. By making use of (9.5), we obtain

$$\gamma_N = \frac{2\alpha}{1-\alpha\beta}(z_N(1) + \beta w_N(-1)) \quad \text{and} \quad \delta_N = \frac{2\alpha}{1-\alpha\beta}(\alpha z_N(1) + w_N(-1)).$$

Hence, the approximations u_N, v_N are fully determined and the improvements mentioned above are realized.

Consider now variable coefficient systems given by

$$\begin{aligned} a_1 u_x + a_2 v_x &= f \\ a_3 u_x + a_4 v_x &= g \end{aligned} \quad \text{on } \Omega = (-1, 1)$$

with boundary conditions

$$\begin{aligned} \alpha_1 u(1) + \alpha_2 v(1) &= 0, \\ \alpha_3 u(-1) + \alpha_4 v(-1) &= 0, \end{aligned} \quad \alpha_i \in \mathbf{R}, \quad i = 1, 2, 3, 4.$$

a_1, \dots, a_4 denote continuous functions on $\bar{\Omega}$, prescribed in such a way that the quantity $\det \begin{pmatrix} a_1 & a_2 \\ a_3 & a_4 \end{pmatrix}$ never vanishes on $[-1, 1]$.

As already seen for second-order problems, direct solvers based on the tau method are no longer available; iterative methods based on collocation have to be used. To achieve a reduction of the condition number, the system should be multiplied by $(1-x)$ and $(1+x)$. Hence we get

$$a_1(1-x)u_x + a_2(1-x)v_x = (1-x)f, \quad a_3(1+x)u_x + a_4(1+x)v_x = (1+x)g.$$

The trivial equations at $x = 1$, resp. $x = -1$, are replaced by the corresponding boundary equations. Now pseudospectral discretization can be applied in a straightforward manner. Iterative methods can be efficiently implemented using FFT and, if there is no preconditioning, convergence factors of $1 - O(N^{-1})$ can be expected. Efficient preconditioners for such problems are proposed by Funaro [5]. The corresponding eigenvalues are real positive and lie between 1 and $\pi/2$.

Acknowledgments. I wish to thank Prof. Dr. K. Witsch for helpful discussions. Further, I am indebted to Dipl.-Math. H. Eisen for the implementation of the methods.

Mathematisches Institut der Universität Düsseldorf
Universitätsstr. 1
4000 Düsseldorf 1
Federal Republic of Germany

1. I. BABUŠKA & A. K. AZIZ, "Survey lectures on the mathematical foundations of the finite element method," in *The Mathematical Foundations of the Finite Element Method with Applications to Partial Differential Equations* (A. K. Aziz, ed.), Academic Press, New York and London, 1972, pp. 3–360.
2. C. CANUTO & A. QUARTERONI, "Approximation results for orthogonal polynomials in Sobolev spaces," *Math. Comp.*, v. 38, 1982, pp. 67–86.
3. C. CANUTO & A. QUARTERONI, "Variational methods in the theoretical analysis of spectral methods," in *Spectral Methods for Partial Differential Equations* (R. G. Voigt, D. Gottlieb and M. Y. Hussaini, eds.), SIAM, Philadelphia, Pa., 1984, pp. 55–78.
4. C. CANUTO & A. QUARTERONI, "Spectral and pseudospectral methods for parabolic problems with nonperiodic boundary conditions," *Calcolo*, v. 18, 1981, pp. 197–218.
5. D. FUNARO, "A preconditioning matrix for the Chebyshev differencing operator," *SIAM J. Numer. Anal.*, v. 24, 1987, pp. 1024–1031.
6. D. GOTTLIEB & S. A. ORSZAG, *Numerical Analysis of Spectral Methods: Theory and Applications*, CBMS-NSF Regional Conference Series in Applied Mathematics No. 26, SIAM, Philadelphia, Pa., 1977.
7. D. B. HAIDVOGEL & T. ZANG, "An accurate solution of Poisson's equation by expansion in Chebyshev polynomials," *J. Comput. Phys.*, v. 30, 1979, pp. 167–180.
8. W. HEINRICHS, *Kollokationsverfahren und Mehrgittermethoden bei elliptischen Randwertaufgaben*, GMD-Bericht Nr. 168, R. Oldenbourg Verlag, München/Wien, 1987.
9. W. HEINRICHS, "Line relaxation for spectral multigrid methods," *J. Comput. Phys.*, v. 77, 1988, pp. 166–182.
10. W. HEINRICHS, "Collocation and full multigrid methods," *Appl. Math. Comput.*, v. 26, 1988, pp. 35–45.
11. W. HEINRICHS, "Konvergenzaussagen für Kollokationsverfahren bei elliptischen Randwertaufgaben," *Numer. Math.*, v. 54, 1989, pp. 619–637.
12. W. HEINRICHS, "Strong convergence estimates for pseudospectral methods," Submitted to *SIAM J. Numer. Anal.*
13. M. A. KRASNOSEL'SKII, G. M. VAINIKKO, P. P. ZABREIKO, YA. B. RUTITSKII & Y. YA. STETSENKO, *Approximate Solution of Operator Equations*, Wolters-Noordhoff, Groningen, 1972.
14. TH. MEIS & U. MARCOWITZ, *Numerische Behandlung partieller Differentialgleichungen*, Springer-Verlag, Berlin and New York, 1978.
15. J. NEČAS, "Sur une méthode pour résoudre les équations aux dérivées partielles du type elliptique, voisine de la variationnelle," *Ann. Sci. Norm. Sup. Pisa*, v. 16, 1962, pp. 305–326.
16. S. A. ORSZAG, "Spectral methods for problems in complex geometries," *J. Comput. Phys.*, v. 37, 1980, pp. 70–92.
17. T. N. PHILLIPS, T. A. ZANG & M. Y. HUSSAINI, "Spectral multigrid methods for Dirichlet problems," in *Multigrid Methods for Integral and Differential Equations* (D. J. Paddon and H. Holstein, eds.), Clarendon Press, Oxford, 1985, pp. 231–252.
18. J. SCHRÖDER, *Operator Inequalities*, Academic Press, New York, 1980.
19. T. A. ZANG, Y. S. WONG & M. Y. HUSSAINI, "Spectral multigrid methods for elliptic equations. II," *J. Comput. Phys.*, v. 54, 1984, pp. 489–507.