

Aplikace matematiky

Jitka Segethová

Reducing the bandwidth in solving linear algebraic systems arising in the finite element method

Aplikace matematiky, Vol. 25 (1980), No. 4, 286–304

Persistent URL: <http://dml.cz/dmlcz/103862>

Terms of use:

© Institute of Mathematics AS CR, 1980

Institute of Mathematics of the Czech Academy of Sciences provides access to digitized documents strictly for personal use. Each copy of any part of this document must contain these *Terms of use*.



This document has been digitized, optimized for electronic delivery and stamped with digital signature within the project *DML-CZ: The Czech Digital Mathematics Library* <http://dml.cz>

REDUCING THE BANDWIDTH IN SOLVING LINEAR
ALGEBRAIC SYSTEMS ARISING
IN THE FINITE ELEMENT METHOD

JITKA SEGETHOVÁ

(Received June 1, 1978)

1. INTRODUCTION

The matrix of the system of linear algebraic equations, arising in the application of the finite element method to one-dimensional problems, is a bandmatrix. In approximation of high order, the band is very wide but the elements situated far from the diagonal of the matrix are negligibly small as compared with the diagonal elements. The situation is similar but much more complex for two or three-dimensional problems.

In numerical practice, the outer diagonals of the matrix influence the solution of the system very little. The aim of this paper is to show that it is possible to work with a matrix the bandwidth of which is reduced, i.e., some non-zero diagonals of the original matrix are replaced by zeros. For a sufficiently high order of approximation (i.e., for a sufficiently wide band) the error caused by this replacement may be negligible in comparison with the rounding errors involved in the numerical solution of the system.

We choose the hill functions ω_N [1], [5] for the trial functions in solving a model problem by the finite element method since we can study their asymptotic behavior for $N \rightarrow \infty$ by means similar to those used in [4]. This as well as the formulation of the model problem is the contents of Section 2.

In Section 3 we investigate the behavior of the elements of the matrix of the system as $N \rightarrow \infty$. This study gives us a key to the introduction of a matrix the band of which is reduced. In conclusion, the difference of the original and the reduced matrices is estimated.

Section 4 is devoted to the analysis of the rounding errors, which is based on [6], [7]. In this part of the paper we obtain error bounds for numerical solution of the original and the reduced systems. These bounds are quantities of the same order. Moreover, a simple numerical example illustrates the statements of Section 4.

2. ASYMPTOTIC BEHAVIOR OF THE HILL FUNCTIONS.
THE MODEL PROBLEM

We employ the following particular type of hill functions (called also B-splines) as trial functions.

Definition 2.1. Let \mathbb{R} be a one-dimensional Euclidean space. Let us put

$$(2.1) \quad \begin{aligned} \omega_1(t) &= 1, \quad |t| \leq \frac{1}{2}, \\ &= 0, \quad |t| > \frac{1}{2}, \end{aligned}$$

$$(2.2) \quad \omega_N = \omega_{N-1} * \omega_1$$

where $*$ denotes the convolution.

The above introduced hill function ω_N is a piecewise polynomial function of degree $N - 1$, it is continuous together with its derivatives up to the order $N - 2$ and its support is the interval $\langle -\frac{1}{2}N, \frac{1}{2}N \rangle$ (see [1], [5]). Further,

$$(2.3) \quad \omega_N(-t) = \omega_N(t), \quad \omega'_N(-t) = -\omega'_N(t).$$

The asymptotic behavior of the hill functions ω_N for $N \rightarrow \infty$ is studied in [4]. Let us perform a more detailed analysis necessary for our further purposes.

Theorem 2.1. Let ω_N be given by (2.1), (2.2). Then there exists a positive constant C independent of N such that

$$(2.4) \quad |\omega_N(t\sqrt{N})\sqrt{N} - \sqrt{(6\pi^{-1})e^{-6t^2}}| \leq CN^{-1},$$

$$(2.5) \quad |\omega'_N(t\sqrt{N})N + 12\sqrt{(6\pi^{-1})te^{-6t^2}}| \leq CN^{-1}$$

uniformly with respect to all $t \in \mathbb{R}$.

Proof. The convolution formula (2.2) is analogous to the convolution formula of the probability theory describing the probability density of the sum of N independent random variables with the same density. The bound (2.4) follows immediately from the central limit theorem ([2] Ch. XV, Sec. 5, Theorem 2, and Ch. XVI, Sec. 2, Theorem 2, cf. also [4] Theorem 3.2).

Moreover, we obtain (2.5) in the same way as (2.4) by substituting ω'_N for ω_N in the proofs of the corresponding theorems of [2].

We confine ourselves to a simple model problem in order to be able to present the main idea of our investigation of reducing the bandwidth. The analysis can be done in an analogous way for a very wide class of problems.

Definition 2.2. Let us put

$$(2.6) \quad Lu = -u'' + cu, \quad c > 0,$$

and consider the boundary-value problem

$$(2.7) \quad Lu = f \quad \text{on} \quad \left(-\frac{1}{2}\pi, \frac{1}{2}\pi\right),$$

$$(2.8) \quad u'(-\frac{1}{2}\pi) = u'(\frac{1}{2}\pi) = 0$$

where $f \in L_2(-\frac{1}{2}\pi, \frac{1}{2}\pi)$.

Let N, M be positive integers,

$$(2.9) \quad 2M > N.$$

Let us put

$$(2.10) \quad Q = [\frac{1}{2}(N - 1)]^*,$$

$$(2.11) \quad h = \pi/(2M),$$

$$(2.12) \quad P = 2M + 2Q + 1.$$

The function

$$u_N(t) = \sum_{j=-M-Q}^{M+Q} x_{N,h,j} \omega_{N,h,j}(t)$$

where

$$(2.13) \quad \omega_{N,h,j}(t) = \omega_N(t/h - j)$$

is said to be the approximate Ritz-Galerkin solution of the problem (2.7), (2.8) if the coefficients $x_{N,h,j}$; $j = -M - Q, \dots, M + Q$, are the solution of the system

$$(2.14) \quad A_{N,h} x_{N,h} = s_{N,h}$$

of P linear algebraic equations, where

$$A_{N,h} = \{a_{N,h,j,k}\}_{j,k=-M-Q}^{M+Q},$$

$$x_{N,h} = \{x_{N,h,j}\}_{j=-M-Q}^{M+Q},$$

$$s_{N,h} = \{s_{N,h,k}\}_{k=-M-Q}^{M+Q}.$$

The elements of the matrix $A_{N,h}$ and of the vector $s_{N,h}$ are given by the relations

$$(2.15) \quad a_{N,h,j,k} = I'_{N,h,j,k} + cI_{N,h,j,k} = (\omega'_{N,h,j}, \omega_{N,h,k}^3) + c(\omega_{N,h,j}, \omega_{N,h,k});$$

$$j, k = -M - Q, \dots, M + Q,$$

$$(2.16) \quad s_{N,h,k} = (f, \omega_{N,h,k}); \quad k = -M - Q, \dots, M + Q,$$

where (\cdot, \cdot) is the scalar product in $L_2(-\frac{1}{2}\pi, \frac{1}{2}\pi)$.

We omit the subscripts N, h and write $I_{jk}, I'_{jk}, a_{jk}, x_j, s_k, A, x, s$ wherever it is not ambiguous.

The following theorem shows some basic properties of the elements of the matrix A .

*) $[z]$ denotes the largest integer which fulfils $[z] \leq z$.

Theorem 2.2. Let a_{jk} be given by (2.15) of Definition 2.2. Then

$$(2.17) \quad I_{jk} = I_{kj} = I_{-j,-k} = I_{-k,-j},$$

$$(2.18) \quad I'_{jk} = I'_{kj} = I'_{-j,-k} = I'_{-k,-j},$$

$$(2.19) \quad a_{jk} = a_{kj} = a_{-j,-k} = a_{-k,-j}$$

for $j, k = -M - Q, \dots, M + Q$. Further,

$$(2.20) \quad I_{jk} = I'_{jk} = a_{jk} = 0 \quad \text{for } |j - k| \geq N; \\ j, k = -M - Q, \dots, M + Q,$$

i.e., A is a symmetric bandmatrix of bandwidth $2N - 1$.

Proof. The symmetry relations (2.17), (2.18), (2.19) are an easy consequence of the substitution of the equalities (2.3) into (2.15). Since the support of ω_N is finite we obtain (2.20). ■

The relations (2.17), (2.18), (2.19) allow us to study only one quarter of the elements of the matrix A . In the following, we choose those pairs of indices $j, k = -M - Q, \dots, M + Q$ for which $0 \leq k - j, 0 \leq k + j$. With respect to (2.20) we introduce the following four sets of pairs of indices j, k .

Definition 2.3. Let N, M be positive integers satisfying (2.9), let Q be given by (2.10). We put

$$S^{++} = \{(j, k); |j| \leq M + Q, |k| \leq M + Q, 0 \leq k - j < N, 0 \leq k + j\}, \\ S^{+-} = \{(j, k); |j| \leq M + Q, |k| \leq M + Q, 0 \leq k - j < N, k + j \leq 0\}, \\ S^{-+} = \{(j, k); |j| \leq M + Q, |k| \leq M + Q, -N < k - j \leq 0, 0 \leq k + j\}, \\ S^{--} = \{(j, k); |j| \leq M + Q, |k| \leq M + Q, -N < k - j \leq 0, k + j \leq 0\}.$$

In the notation introduced in Definition 2.3 we confine ourselves to the analysis of only the elements of the set S^{++} in the following. The properties of the elements of the other three sets can be established according to Theorem 2.2 in an apparent way.

Let us mention that all the four sets $S^{++}, S^{+-}, S^{-+}, S^{--}$ depend on the integers N, M . Therefore we sometimes write $S^{++}(N, M)$ or $S^{++}(N)$ instead of S^{++} etc.

The knowledge of the asymptotic behavior of the hill functions enables us to obtain asymptotic formulae for the elements a_{jk}, s_k as $N \rightarrow \infty$.

Theorem 2.3. Let a_{jk}, s_k be given by (2.15), (2.16) of Definition 2.2. Then

$$(2.21) \quad I_{jk} = I_{jk}^*(1 + O(N^{-1})),$$

$$(2.22) \quad I'_{jk} = I_{jk}^{*'}(1 + O(N^{-1})),$$

$$(2.23) \quad a_{jk} = (I_{jk}^* + cI_{jk}^{*'})(1 + O(N^{-1}))$$

for $N \rightarrow \infty$ and $(j, k) \in S^{++}$ where*

$$(2.24) \quad I_{jk}^* = (6N^{-1})^{3/2} \pi^{-1} h^{-1} \frac{1}{6} N h^2 \sqrt{(\frac{1}{2}\pi)} \exp(-3N^{-1}(j-k)^2) \times \\ \times \Phi(\sqrt{(3N^{-1})(N+j-k)}),$$

$$(2.25) \quad I_{jk}^* = (6N^{-1})^{3/2} \pi^{-1} h^{-1} \exp(-3N^{-1}(j-k)^2) \times \\ \times (\sqrt{(\frac{1}{2}\pi)} (1 - 6N^{-1}(j-k)^2) \Phi(\sqrt{(3N^{-1})(N+j-k)}) - \\ - \sqrt{(6N^{-1})(N+j-k)} \times \exp(-3N^{-1}(N+j-k)^2))$$

for $j < M - Q$ and

$$(2.26) \quad I_{jk}^* = \frac{1}{2}(6N^{-1})^{3/2} \pi^{-1} h^{-1} \frac{1}{6} N h^2 \sqrt{(\frac{1}{2}\pi)} \exp(-3N^{-1}(j-k)^2) \times \\ \times (\Phi(\sqrt{(3N^{-1})(2M-j-k)}) - \Phi(\sqrt{(3N^{-1})(-N-j+k)})),$$

$$(2.27) \quad I_{jk}^* = \frac{1}{2}(6N^{-1})^{3/2} \pi^{-1} h^{-1} \exp(-3N^{-1}(j-k)^2) (\sqrt{(\frac{1}{2}\pi)} \times \\ \times (1 - 6N^{-1}(j-k)^2) (\Phi(\sqrt{(3N^{-1})(2M-j-k)}) - \Phi(\sqrt{(3N^{-1})} \times \\ \times (-N-j+k))) + \sqrt{(6N^{-1})((-N-j+k)} \times \\ \times \exp(-3N^{-1}(-N-j+k)^2) - (2M-j-k) \times \\ \times \exp(-3N^{-1}(2M-j-k)^2)))$$

for $M - Q \leq j$.

Further,

$$(2.28) \quad |s_k|^2 \leq \|f\|_{L_2}^2 I_{kk} \quad \text{for } k = -M - Q, \dots, M + Q.$$

Proof. Considering that $(j, k) \in S^{++}$ and that (2.9) holds, we distinguish two possibilities in the proof, $j < M - Q$ and $j \geq M - Q$, which correspond to different limits of integration in the calculation of (2.15). In both the cases we substitute the expressions

$$\omega_N(t) = \sqrt{(6\pi^{-1}N^{-1})} \exp(-6N^{-1}t^2) + O(N^{-3/2}), \\ \omega'_N(t) = 12N^{-1} \sqrt{(6\pi^{-1}N^{-1})} t \exp(-6N^{-1}t^2) + O(N^{-2}),$$

following from (2.4), (2.5) of Theorem 2.1, into (2.15). Taking into account that the functions $\Phi(t)$ and $\exp(-t^2)$ are bounded in \mathbb{R} and considering that the length w of the interval of integration in (2.15) fulfils the inequality

$$w \leq Nh,$$

we obtain (2.21) and (2.22) as well as (2.23) by direct calculation. The terms of higher degree in N^{-1} are neglected in these formulae. Applying the Schwarz inequality to (2.16), we obtain (2.28).

*) We use the usual notation

$$\Phi(t) = \frac{2}{\sqrt{\pi}} \int_0^t \exp(-u^2) du.$$

3. SYSTEM WITH REDUCED MATRIX

In this section we investigate the asymptotic behavior of the elements a_{jk} , s_k as $N \rightarrow \infty$. We introduce the matrix B of a reduced bandwidth as well as the associate approximate solution of the model problem making use of the knowledge of this asymptotic behavior.

In conclusion, we estimate the difference $B - A$ of the original and the reduced matrices and the difference $t - s$ of the original and the reduced right-hand parts of the system.

Let us start with the study of the elements a_{jk} of the matrix A .

Definition 3.1. *Let N, M be positive integers satisfying (2.9). Moreover, let there exist positive constants C, D independent of N such that*

$$(3.1) \quad M = M(N) \leq CN^D.$$

Let Q be given by (2.10). Let us choose $\frac{1}{2} > \varepsilon > 0$ and $Z > 0$ (independently of N) and denote by K the least odd integer not less than $ZN^{1/2+\varepsilon}$. Further let us put

$$(3.2) \quad R = \frac{1}{2}(K - 1).$$

Let us introduce the following five sets of pairs of indices j, k :

$$S_1^{++}(N) = \{(j, k); j = -Q, \dots, M - Q - 1; k = j + K, \dots, j + N - 1\} \cap S^{++}(N),$$

$$S_2^{++}(N) = \{(j, k); j = M - Q, \dots, M - R - 1; k = j + K, \dots, M + Q\},$$

$$S_3^{++}(N) = \{(j, k); j = M - R, \dots, M; k = M + R + 1, \dots, M + Q\},$$

$$S_4^{++}(N) = \{(j, k); j = M + 1, \dots, M + R - 1; k = M + R + 1, \dots, M + Q\},$$

$$S_5^{++}(N) = \{(j, k); j = M + R, \dots, M + Q - 1; k = j + 1, \dots, M + Q\}.$$

Moreover, we write

$$\tilde{S}^{++}(N) = \bigcup_{i=1}^5 S_i^{++}(N).$$

In an analogous way we can introduce the sets $\tilde{S}^{+-}(N)$, $\tilde{S}^{-+}(N)$, $\tilde{S}^{--}(N)$ (cf. Theorem 2.2 and Definition 2.3) and put

$$\tilde{S}(N) = \tilde{S}^{++}(N) \cup \tilde{S}^{+-}(N) \cup \tilde{S}^{-+}(N) \cup \tilde{S}^{--}(N).$$

Apparently, $\tilde{S}^{++}(N) \subset S^{++}(N)$. The sets $S_i^{++}(N)$ are mutually disjoint for N sufficiently large by definition.

In the following we always suppose that N, M are positive integers satisfying (2.9), (3.1). Moreover, P, Q always denote the integers given by (2.12), (2.10) of Definition

2.2 while R is given by (3.2) of Definition 3.1. With respect to (2.11), (3.1) we also have

$$(3.3) \quad h = h(N) = \pi/(2M(N)).$$

Theorem 3.1. *Let $\{j(N), k(N)\}_{N=1}^{\infty}$ be a sequence of pairs of integers. Let N_0 be a positive integer such that*

$$(3.4) \quad (j(N), k(N)) \in \mathfrak{S}^{++}(N) \quad \text{for } N \geq N_0.$$

Let us put

$$(3.5) \quad \alpha_{j(N), k(N)} = \pi h(\frac{1}{6}N)^{3/2} a_{j(N), k(N)}.$$

Let J be an arbitrary real number. Then there exists a positive constant C_J independent of N such that

$$(3.6) \quad |\alpha_{j(N), k(N)}| \leq C_J N^{-J} \quad \text{for } N \geq N_0.$$

Proof. 1. First we will prove that, choosing a fixed i and considering a sequence $\{j(N), k(N)\}_{N=1}^{\infty}$ such that

$$(3.7) \quad (j(N), k(N)) \in S_i^{++}(N) \quad \text{for } N \geq N_i$$

with some positive N_i , we find a positive constant C_{iJ} independent of N such that

$$(3.8) \quad |\alpha_{j(N), k(N)}| \leq C_{iJ} N^{-J} \quad \text{for } N \geq N_i.$$

Having proved (3.8) for $i = 1, \dots, 5$, we obtain the statement of the theorem by the following argument.

Let us consider the sequence (3.4). Let us split its subsequence $\{j(N), k(N)\}_{N=N_0}^{\infty}$ with some $N'_0 \geq N_0$ into five sets U_i ; $i = 1, \dots, 5$, such that

$$U_i = \{(j(N), k(N)); (j(N), k(N)) \in S_i^{++}(N)\}.$$

Apparently

$$(3.9) \quad \{j(N), k(N)\}_{N=N_0}^{\infty} = \bigcup_{i=1}^5 U_i.$$

For a sufficiently large N'_0 , each of the sets U_i is either empty or infinite and at least one of them is infinite with respect to (3.4). Let us put

$$I = \{i; U_i \neq \emptyset\}.$$

Then $I \neq \emptyset$ and if $i \in I$ we can write

$$U_i = \{j(N_n^i), k(N_n^i)\}_{n=1}^{\infty}, \quad N'_0 \leq N_n^i, \quad N_n^i < N_{n+1}^i.$$

For each $i \in I$ we obtain from (3.8) that

$$|\alpha_{j(N_n^i), k(N_n^i)}| \leq C_{iJ} (N_n^i)^{-J} \quad \text{for } N_n^i \geq N_i.$$

Putting $C'_J = \max_{i \in I} C_{iJ}$, we have

$$|\alpha_{j(N),k(N)}| \leq C'_J N^{-J} \quad \text{for } N \geq N''_0$$

where $N''_0 = \max(N'_0, \max_{i \in I} N_i)$, from which (3.6) follows immediately.

Therefore it remains to show (3.8) for $i = 1, \dots, 5$. We will write (j, k) instead of $(j(N), k(N))$ wherever it is not ambiguous.

2. Let $\{j(N), k(N)\}_{N=1}^\infty$ be a sequence satisfying (3.7) for $i = 1$. From Definition 3.1 and Theorem 2.3 we find that a_{jk} is given by (2.23), (2.24), (2.25) in this case, i.e.,

$$(3.10) \quad \alpha_{jk} = \exp(-3N^{-1}(j-k)^2) (\sqrt{(\frac{1}{2}\pi)} (\frac{1}{6}ch^2N + 1 - 6N^{-1}(j-k)^2) \times \\ \times \Phi(\sqrt{(3N^{-1})(N+j-k)} - \sqrt{(6N^{-1})(N+j-k)}) \times \\ \times \exp(-3N^{-1}(N+j-k)^2)) (1 + O(N^{-1})).$$

Since $j + K \leq k \leq j + N - 1$ in $S_1^{++}(N)$ we have

$$k - j \geq K \geq ZN^{1/2+\epsilon},$$

i.e.,

$$(3.11) \quad (k-j)^2 N^{-1} \geq Z^2 N^{2\epsilon} \rightarrow \infty \quad \text{for } N \rightarrow \infty.$$

Therefore there exists a positive constant C'_{1J} independent of N such that

$$(3.12) \quad \exp(-3N^{-1}(j-k)^2) < C'_{1J} N^{-J}$$

for an arbitrary real J . The functions $\Phi(t)$ and $t \exp(-t^2)$ are bounded in \mathbb{R} . For any sequence the elements of which lie in $S^{++}(N)$, we have $0 \leq k - j < N$ from Definition 2.3 and, moreover,

$$(3.13) \quad (k-j)^2 N^{-1} < N.$$

Therefore we finally obtain

$$(3.14) \quad |\alpha_{jk}| \exp(3N^{-1}(j-k)^2) \leq D_1 N$$

with some positive constant D_1 from (2.9) and (3.3). Considering (3.12) and (3.14), we find a positive constant C_{1J} independent of N such that (3.8) holds with $i = 1$.

3. Let $\{j(N), k(N)\}_{N=1}^\infty$ be a sequence satisfying (3.7) for $i = 2$. From Definition 3.1 and Theorem 2.3 we find that a_{jk} is given by (2.23), (2.26), (2.27) in this case, i.e.,

$$(3.15) \quad \alpha_{jk} = \frac{1}{2} \exp(-3N^{-1}(j-k)^2) (\sqrt{(\frac{1}{2}\pi)} (\frac{1}{6}ch^2N + 1 - 6N^{-1}(j-k)^2) \times \\ \times (\Phi(\sqrt{(3N^{-1})(2M-j-k)} - \Phi(\sqrt{(3N^{-1})(-N-j+k)})) + \\ + \sqrt{(6N^{-1})((-N-j+k) \exp(-3N^{-1}(-N-j+k)^2) - \\ - (2M-j-k) \exp(-3N^{-1}(2M-j-k)^2))}) (1 + O(N^{-1})).$$

Since $j + K \leq k$ in $S_2^{++}(N)$ we obtain (3.11) also in this case. Obviously (3.13) also holds because $S_2^{++}(N) \subset S^{++}(N)$. Therefore we can repeat the argument of part 2 of the proof in order to find a positive constant C_{2J} independent of N such that (3.8) holds with $i = 2$.

4. Let $\{j(N), k(N)\}_{N=1}^{\infty}$ be a sequence satisfying (3.7) for $i = 3$. From Definition 3.1 and Theorem 2.3 we find that α_{jk} is given by (3.15). Further we obtain

$$k - j \geq R + 1 > \frac{1}{2}ZN^{1/2+\varepsilon}$$

with respect to (3.2) and the inequality

$$(3.16) \quad K \geq ZN^{1/2+\varepsilon}$$

Therefore

$$(3.17) \quad (k - j)^2 N^{-1} \geq \frac{1}{4}Z^2N^{2\varepsilon} \rightarrow \infty \quad \text{for } N \rightarrow \infty .$$

Now we can repeat the argument of part 2 in order to find a positive constant C_{3J} independent of N such that (3.8) holds with $i = 3$.

5. Let $\{j(N), k(N)\}_{N=1}^{\infty}$ be a sequence satisfying (3.7) for $i = 4$. We find that α_{jk} is given by (3.15). Further we obtain

$$k - j \leq Q - 1 < \frac{1}{2}N$$

with respect to (2.10). Hence we finally have

$$(-N + k - j)N^{-1/2} < -\frac{1}{2}N^{1/2} \rightarrow -\infty \quad \text{for } N \rightarrow \infty .$$

Using an asymptotic expansion for $\Phi(t)$ at $+\infty$ we obtain (see e.g. [3])

$$(3.18) \quad \Phi(t) = 1 - \exp(-t^2)((t\sqrt{\pi})^{-1} + O(t^{-3})), \quad t \rightarrow \infty .$$

Therefore there exists a positive constant \tilde{C}'_{4J} independent of N such that

$$(3.19) \quad 0 < 1 + \Phi(\sqrt{(3N^{-1})(-N - j + k)}) \leq \tilde{C}'_{4J}N^{-J}$$

for an arbitrary real J . Analogously there exists a constant \tilde{C}''_{4J} such that

$$(3.20) \quad 0 < (N + j - k) \exp(-3N^{-1}(-N - j + k)^2) \leq \tilde{C}''_{4J}N^{-J} .$$

From Definition 3.1 we further obtain that

$$k + j - 2M \geq R + 2 > \frac{1}{2}ZN^{1/2+\varepsilon}$$

with respect to (3.2), (3.16). Then

$$(3.21) \quad (k + j - 2M)N^{-1/2} \geq \frac{1}{2}ZN^{\varepsilon} \rightarrow \infty \quad \text{for } N \rightarrow \infty$$

and we find positive constants $\hat{C}'_{4J}, \hat{C}''_{4J}$ independent of N such that

$$(3.22) \quad 0 < \Phi(\sqrt{(3N^{-1})(2M - j - k)}) + 1 \leq \hat{C}'_{4J}N^{-J} ,$$

$$(3.23) \quad 0 < (-2M + k + j) \exp(-3N^{-1}(2M - j - k)^2) \leq \hat{C}'_{4j} N^{-j}$$

by the same argument as above. Since the function $\exp(-t^2)$ is bounded in \mathbb{R} and

$$\frac{1}{6}ch^2N + 1 - 6N^{-1}(j - k)^2 \leq C'_{4j}N$$

with some constant C'_{4j} we can find a positive constant C_{4j} independent of N such that (3.8) holds with $i = 4$ if we take into account the inequalities (3.19), (3.20), (3.22), (3.23).

6. Let $\{j(N), k(N)\}_{N=1}^{\infty}$ be a sequence satisfying (3.7) for $i = 5$. From Definition 3.1 we obtain

$$\begin{aligned} k - j &\leq Q - R < \frac{1}{2}N - \frac{1}{2}ZN^{1/2+\varepsilon}, \\ k + j - 2M &\geq 2R + 1 \geq ZN^{1/2+\varepsilon} \end{aligned}$$

with respect to (2.10), (3.2), (3.16). Repeating the argument of part 5 we finally find a positive constant C_{5j} independent of N such that (3.8) holds with $i = 5$.

The proof of the theorem is complete.

Corollary. Let $\{j(N), k(N)\}_{N=1}^{\infty}$ be a sequence of pairs of integers. Let N_0 be a positive integer such that

$$(3.24) \quad (j(N), k(N)) \in \tilde{S}(N) \quad \text{for } N \geq N_0.$$

Let $\alpha_{j(N), k(N)}$ be given by (3.5), let J be an arbitrary real number. Then there exists a positive constant \tilde{C}_J independent of N such that

$$|\alpha_{j(N), k(N)}| \leq \tilde{C}_J N^{-J} \quad \text{for } N \geq N_0.$$

Proof. Theorem 3.1 treats the case when (3.4) holds for the sequence $\{j(N), k(N)\}$. From Theorem 2.2 we obtain the statement of the corollary for any sequence such that $(j(N), k(N)) \in \tilde{S}^{+-}(N)$ (or $\tilde{S}^{-+}(N)$ or $\tilde{S}^{--}(N)$). Proceeding in the way analogous to that of part 1 of the proof of Theorem 3.1, we obtain the statement for an arbitrary sequence satisfying (3.24). ■

A similar statement holds for the elements s_k of the right-hand part s .

Theorem 3.2. Let $\{k(N)\}_{N=1}^{\infty}$ be a sequence of integers. Let N_0 be a positive integer such that

$$M + R + 1 \leq |k(N)| \leq M + Q \quad \text{for } N \geq N_0.$$

Let us put

$$\sigma_{k(N)} = \pi h \left(\frac{1}{6}N\right)^{3/2} s_{k(N)}.$$

Let J be a real number. Then there exists a positive constant C_J^* independent of N such that

$$|\sigma_{k(N)}| \leq C_J^* \|f\|_{L_2} N^{-J} \quad \text{for } N \geq N_0.$$

Proof. According to (2.28) of Theorem 2.3, it is sufficient to show that

$$(3.25) \quad |\pi^2 h^2 (\frac{1}{6}N)^3 I_{kk}| \leq C_J^* N^{-J} \quad \text{for } N \geq N_0.$$

We have $I_{kk} = I_{-k, -k}$ from (2.17) of Theorem 2.2, which allows us to confine ourselves only to the study of sequences $\{k(N)\}_{N=1}^\infty$ such that

$$(3.26) \quad M + R + 1 \leq k(N) \leq M + Q \quad \text{for } N \geq N_0.$$

The result obtained for these sequences can be generalized in the way indicated in the proof of Corollary.

Thus let $\{k(N)\}_{N=1}^\infty$ be a sequence satisfying (3.26). From Theorem 2.3 we find that I_{kk} is given by (2.21), (2.26) in this case, i.e.,

$$\pi^2 h^2 (\frac{1}{6}N)^3 I_{kk}^* = (\frac{1}{2}\pi)^{3/2} h^3 (\frac{1}{6}N)^{5/2} (\Phi(2\sqrt{(3N^{-1})(M-k)}) - \Phi(-\sqrt{(3N)})).$$

Moreover, we obtain that

$$k - M \geq R + 1 > \frac{1}{2}ZN^{1/2+\varepsilon}$$

with respect to (3.2), (3.16), (3.26). Therefore

$$(3.27) \quad (k - M)N^{-1/2} \geq \frac{1}{2}ZN^\varepsilon \rightarrow \infty \quad \text{for } N \rightarrow \infty.$$

Using the expansion (3.18) for $\Phi(t)$ and proceeding in the way analogous to part 5 of the proof of Theorem 3.1, we find a positive constant C_J^* independent of N such that (3.25) is fulfilled for $k(N)$ satisfying (3.26). This implies the statement of the theorem immediately. ■

We showed in Theorem 3.1, Corollary, and Theorem 3.2 that some of the elements a_{jk} of the matrix A and some of the elements s_k of the vector s (premultiplied by the factor $\pi h (\frac{1}{6}N)^{3/2}$) converge rapidly to 0 as $N \rightarrow \infty$.

Working with the quantities α_{jk}, σ_k we obtain that $\alpha_{00} = O(1)$ as $N \rightarrow \infty$. We will employ this fact in the proof of Theorem 3.3. Apparently, $\alpha_{jk} \rightarrow 0$ implies $a_{jk} \rightarrow 0$ and $\sigma_k \rightarrow 0$ implies $s_k \rightarrow 0$. This is the behavior of e.g. all the elements a_{jk} lying outside the band of width $2K - 1$.

Let us note that the statement is based on the condition (3.16). In general, if we put $K \geq ZN^\beta$ with $\beta \leq \frac{1}{2}$ then none of the relations (3.11), (3.17), (3.21), (3.27) can be fulfilled. In this sense, the condition (3.16) determines the narrowest band with the property that the elements a_{jk} situated outside the band tend to 0.

We will study a matrix and a vector formed from the original matrix A and vector s by substituting zeros for the elements tending to zero. First we will introduce the necessary notation.

Definition 3.2. Let

$$B_{N,h} = B = \{b_{jk}\}_{j,k=-M-Q}^{M+Q}$$

be a square matrix with the elements

$$b_{jk} = 0 \quad \text{for } (j, k) \in \bar{S}(N),$$

$$= a_{jk} \text{ otherwise .}$$

Moreover, let

$$t_{N,h} = t = \{t_k\}_{k=-M-Q}^{M+Q}$$

be the vector the elements of which are given by the relations

$$\begin{aligned} t_k &= 0 \text{ for } M + R + 1 < |k|, \\ &= s_k \text{ otherwise .} \end{aligned}$$

The function

$$(3.28) \quad \hat{u}_N(t) = \sum_{j=-M-Q}^{M+Q} \hat{x}_{N,h,j} \omega_{N,h,j}(t)$$

where $\omega_{N,h,j}$ is given by (2.13) is said to be the associate approximate Ritz-Galerkin solution of the problem (2.7), (2.8) if the coefficients $\hat{x}_{N,h,j}$; $j = -M - Q, \dots, M + Q$, are the solution of the system

$$(3.29) \quad B_{N,h} \hat{x}_{N,h} = t_{N,h}$$

of P linear algebraic equations where

$$\hat{x}_{N,h} = \hat{x} = \{\hat{x}_{N,h,j}\}_{j=-M-Q}^{M+Q} = \{\hat{x}_j\}_{j=-M-Q}^{M+Q} .$$

In addition, we write

$$(3.30) \quad C = B - A, \quad u = t - s .$$

The matrix B is introduced in Definition 3.2 in such a way that

$$(3.31) \quad B = \begin{bmatrix} B' & O & O \\ O & B^* & O \\ O & O & B'' \end{bmatrix}$$

where B^* is a square matrix of order

$$(3.32) \quad P^* = 2M + 2R + 1 ,$$

B', B'' are diagonal matrices of order $Q - R$ and the O 's denote zero matrices of appropriate types.

Similarly

$$(3.33) \quad t^T = (O, t^{*T}, O)$$

where $t^* = \{t_k\}_{k=-M-R}^{M+R}$ and the superscript T denotes the transpose. Let us write

$$(3.34) \quad \hat{x}^T = (\hat{x}'^T, \hat{x}^{*T}, \hat{x}''^T)$$

where $\hat{x}^* = \{\hat{x}_j\}_{j=-M-R}^{M+R}$. In this notation we obtain

$$(3.35) \quad \hat{x}' = \hat{x}'' = O$$

from (3.29). Therefore (3.28) is equivalent to

$$\hat{u}_N(t) = \sum_{j=-M-R}^{M+R} \hat{x}_j \omega_{N,h,j}(t)$$

and \hat{x}_j are determined from the system

$$(3.36) \quad B^* \hat{x}^* = t^*$$

of P^* linear algebraic equations. Moreover, the matrix B^* is a symmetric bandmatrix of bandwidth $2K - 1$ according to Definition 3.2 and

$$b_{jk} = b_{kj} = b_{-j,-k} = b_{-k,-j}; \quad j, k = -M - Q, \dots, M + Q.$$

We will study the relation between the approximate and the associate approximate solution in the next section. First we will estimate the differences $A - B$, $s - t$ making use of the matrix and vector norms.

Definition 3.3. Let $D = \{d_{jk}\}_{j,k=1}^n$ be a square matrix, $r = \{r_k\}_{k=1}^n$ a vector. We write

$$(3.37) \quad \|D\| = \max_{1 \leq j \leq n} \sum_{k=1}^n |d_{jk}|,$$

$$(3.38) \quad \|r\| = \max_{1 \leq k \leq n} |r_k|.$$

The expressions (3.37), (3.38) are the well-known matrix and vector norms (cf. e.g. [7]).

Theorem 3.3. Let A, B, C, s, t, u be the matrices and vectors introduced in Definitions 2.2 and 3.2, let ε_0 be an arbitrary positive number independent of N . Then there exists an integer $N_0 > 0$ such that

$$(3.39) \quad \frac{\|C\|}{\|B\|} < \varepsilon_0, \quad \frac{\|u\|}{\|B\|} < \varepsilon_0$$

for all $N \geq N_0$ (and any M, h satisfying (2.9), (3.1), (3.3)).

Proof. 1. Let us estimate the ratio

$$(3.40) \quad \frac{\|C\|}{\|B\|} = \frac{\pi h (\frac{1}{6}N)^{3/2} \|B - A\|}{\pi h (\frac{1}{6}N)^{3/2} \|B\|}.$$

We have $\pi h (\frac{1}{6}N)^{3/2} \|B\| \geq |\alpha_{00}|$ from (3.5). Further, we obtain that

$$\alpha_{00} = (\sqrt{(\frac{1}{2}\pi)} (\frac{1}{6}ch^2N + 1) \Phi(\sqrt{(3N)}) - \sqrt{(6N) \exp(-3N)}) (1 + O(N^{-1}))$$

from Definition 3.1 and Theorem 2.3. Since there exist positive constants $\bar{C}_{1j}, \bar{C}_{2j}$ independent of N such that

$$1 - \Phi(\sqrt{3N}) \leq \bar{C}_{1J} N^{-J}$$

(cf. (3.18)) and

$$\sqrt{6N} \exp(-3N) \leq \bar{C}_{2J} N^{-J}$$

for any real J we conclude that there exists a positive constant C_3 independent of N such that

$$(3.41) \quad |\alpha_{00}| \geq C_3$$

for a sufficiently large N .

2. With respect to Theorem 3.1 and Definition 3.2, we have

$$|\gamma_{jk}| \leq C_j N^{-J}$$

for all the elements γ_{jk} of the matrix $\pi h(\frac{1}{6}N)^{3/2}C$, any real J , and a sufficiently large N . From (3.40), (3.41) we finally obtain the first inequality in (3.39).

The second inequality in (3.39) follows in an analogous way from Theorem 3.2 and Definition 3.2.

4. A NUMERICAL EXAMPLE. DISCUSSION

In this section we are concerned with the study of the difference of the approximate and the associate approximate Ritz-Galerkin solution of the problem (2.7), (2.8). Keeping the notation of the previous sections, it is sufficient to estimate the difference $x - \hat{x}$ for this purpose.

However, we will consider the situation from the practical, i.e., numerical point of view. Following the ideas of Wilkinson [6], [7], let us distinguish between the true solution x_t of the system (2.14) and the computed solution x_c of the same system; let us introduce the vectors \hat{x}_t, \hat{x}_c with an analogous meaning with respect to the system (3.29).

Let us solve the system of linear algebraic equations by the Gaussian elimination (based on the triangular decomposition of the matrix of the system) with partial pivoting. Let us consider the floating-point arithmetic with the mantissa of t binary digits. (The analysis is analogous and follows also [6], [7] in the case of decimal computation.)

The detailed analysis of the rounding errors made in the arithmetic operations on digital computers in the floating-point computation is a very complex and deep problem. There are phenomena observed widely in practice but rigorously demonstrated only in certain special cases. In view of these rather unusual conditions we use several statements of [6], [7] the assumptions of which cannot be rigorously verified. Similarly we make some assumptions of this nature in the following. These assumptions are justified by a numerical example presented in the conclusion of this section.

Let us return to the error analysis. The computed solution x_c of the system (2.14) is the true solution of the perturbed system

$$(A + E)x_c = s$$

where

$$(4.1) \quad \|E\| \leq 2^{-t} g q(P),$$

$$q(P) = 2 \cdot 14 \left(\frac{1}{2}P + 1\right) (P - 1),$$

P is the order of the matrices A , E , and

$$g = \max_{1 \leq r \leq P} \max_{-M-Q \leq j, k \leq M+Q} |a_{jk}^{(r)}|$$

with $A^{(r)} = \{a_{jk}^{(r)}\}_{j,k=-M-Q}^{M+Q}$ defined as the matrix appearing in the r th step of the elimination (i.e., the original system is $A^{(1)}x = s$). The error expressed by E originates in the triangular decomposition and represents the prevailing part of the total error. The perturbations resulting from the rounding errors made in the backsubstitution are negligible as compared with E . Further, the bound (4.1) for E may be attained in practice. All these statements are taken from [7] (Ch. 3, Secs. 14–33).

It is

$$\|(T + V)^{-1}\| \leq \frac{\|T^{-1}\|}{1 - \|T^{-1}V\|} \leq \frac{\|T^{-1}\|}{1 - \|T^{-1}\| \|V\|}$$

provided that $\|T^{-1}\| \|V\| < 1$ (cf. e.g. [7]). Considering that $Ax_t = s$ and supposing that $\|A^{-1}\| \|E\| < 1$ (this assumption is quite natural; if it is not true the rounding errors are so large that their analysis is impossible), we arrive at the inequality

$$(4.2) \quad \frac{\|x_c - x_t\|}{\|x_t\|} \leq \frac{\kappa(A) \|E\| \|A\|}{1 - \kappa(A) \|E\| \|A\|}$$

where $\kappa(A) = \|A\| \|A^{-1}\|$ is the condition number of the matrix A .

Taking into account the relations (3.31), (3.33), (3.34), (3.35), we concluded in Section 3 that it is sufficient to solve the system (3.36) for the vector \hat{x}^* in order to obtain the solution \hat{x} of the system (3.29). Let \hat{x}_c^* be the computed solution of the system (3.36), let

$$\hat{x}_c^T = (O, \hat{x}_c^{*T}, O).$$

By the same argument as above we see that \hat{x}_c is the true solution of the perturbed system

$$(B + \hat{E}) \hat{x}_c = t,$$

i.e., with respect to (3.30),

$$(4.3) \quad (A + C + \hat{E}) \hat{x}_c = s + u$$

where

$$\hat{E} = \begin{bmatrix} O & O & O \\ O & \hat{E}^* & O \\ O & O & O \end{bmatrix},$$

\hat{E}^* is a square matrix of order P^* and the O 's denote zero matrices of appropriate types. According to [7] we have

$$(4.4) \quad \|\hat{E}\| = \|\hat{E}^*\| \leq 2^{-t} \hat{g}^* q(P^*)$$

where

$$\hat{g}^* = \max_{1 \leq r \leq P^*} \max_{-M-R \leq j, k \leq M+R} |b_{jk}^{(r)}|$$

and $B^{*(r)} = \{b_{jk}^{(r)}\}_{j,k=-M-R}^{M+R}$ is the matrix appearing in the r th step of the elimination.

Considering the matrix $C + \hat{E}$ and the vector u in (4.3) to be the perturbations of the system (2.14), we obtain

$$(4.5) \quad \|\hat{x}_c - x_t\| \leq \frac{\alpha(B) (\|C\| \|x_t\| \|B\| + \|\hat{E}\| \|x_t\| \|B\| + \|u\| \|B\|)}{1 - \alpha(B) \|\hat{E}\| \|B\|}$$

provided that $\|B^{-1}\| \|\hat{E}\| < 1$.

Let us suppose that there exists a positive constant G independent of P^* such that

$$(4.6) \quad \hat{g}^* \leq G \|B^*\|.$$

Numerical results confirm this assumption in the whole practically investigated range of values of P^* . Further, let G be chosen in such a way that the bound (4.6) for \hat{g}^* may be attained.

From (4.4), (4.6) we thus obtain

$$\|\hat{E}\| \leq 2^{-t} G \|B^*\| q(P^*) \leq 2^{-t} G \|B\| q(P^*).$$

In particular,

$$(4.7) \quad \|\hat{E}\| \|B\| \leq 2^{-t} G q(P^*)$$

where the bound is attainable and tends to infinity as $N \rightarrow \infty$ (and $P \rightarrow \infty$, $P^* \rightarrow \infty$ as well). From Theorem 3.3 we have $\|u\| \|B\| \rightarrow 0$ as $N \rightarrow \infty$. Let us now choose a sufficiently large N and consider such x_t that

$$(4.8) \quad \|\hat{E}\| \|x_t\| \|B\| \gg \|u\| \|B\|.$$

Numerical results in the following correspond apparently to this situation. Neglecting the term $\|u\| \|B\|$ in (4.5) with respect to (4.8), we arrive at

$$(4.9) \quad \frac{\|\hat{x}_c - x_t\|}{\|x_t\|} \leq \frac{\alpha(B) (\|C\| \|B\| + \|\hat{E}\| \|B\|)}{1 - \alpha(B) \|\hat{E}\| \|B\|}.$$

Moreover, we have $\|C\| \|B\| \rightarrow 0$ as $N \rightarrow \infty$ from Theorem 3.2. Neglecting now also the term $\|C\| \|B\|$ in (4.9) with respect to (4.7), we finally obtain

$$(4.10) \quad \frac{\|\hat{x}_c - x_t\|}{\|x_t\|} \leq \frac{\kappa(B) \|\hat{E}\|/\|B\|}{1 - \kappa(B) \|\hat{E}\|/\|B\|}.$$

Comparing the estimates (4.2), (4.10), we conclude that both the bounds are of the same nature as long as we suppose $\kappa(B) \sim \kappa(A)$. In fact, numerical results show that rather $\kappa(B) < \kappa(A)$. Furthermore, it is $P^* < P$ and therefore we can expect that \hat{x}_c (obtained from \hat{x}_c^*) gives better results than x_c for a sufficiently large P .

The following simple numerical example illustrates the above statements. Let us seek the approximate Ritz-Galerkin solution of the problem (2.7), (2.8) with the operator L given in (2.6). We put

$$f(x) = -\sin(dx), \quad d > 0.$$

The exact solution of the problem is

$$u(x) = -\frac{1}{d^2 + c} \sin(dx) + \frac{d}{(d^2 + c)\sqrt{c}} \frac{\cos(\frac{1}{2}\pi d)}{\operatorname{ch}(\frac{1}{2}\pi\sqrt{c})} \operatorname{sh}(x\sqrt{c}).$$

The approach of the whole paper is asymptotic; the statements always assume that N is sufficiently large. Under such conditions, Z in Definition 3.1 may be chosen arbitrarily. In numerical computation we are forced to work only with some limited range of values of N . Our experiments included $N \leq 15$. With regard to this fact, we computed the associate approximate Ritz-Galerkin solution (Definition 3.2) for all the values R given (instead of (3.2) of Definition 3.1) by $R = [\frac{1}{2}(K - 1)]$ with $K = N - 1, N - 2, \dots, 2$. The solution corresponding to $K = N$ is the approximate solution (cf. Definition 2.2). In this way, we obtained a series of results by solving systems of linear algebraic equations with matrices of bandwidth $2N - 1, 2N - 3, \dots, 3$.

Table 1.

N	$\eta(N, N)$	$\eta(N, N - 1)$	$\eta(N, N - 2)$	$\eta(N, N - 3)$	$\eta(N, N - 4)$
2	0.587 -3				
3	0.139 -4	0.185 -2			
4	0.451 -5	0.287 -3	0.201 -2		
5	0.348 -6	0.151 -4	0.907 -3	0.164 -2	
6	0.240 -6	0.401 -6	0.108 -3	0.191 -2	0.119 -2
7	0.231 -6	0.232 -6	0.735 -5	0.303 -3	0.246 -2
8	0.231 -6	0.231 -6	0.353 -6	0.844 -4	0.590 -3
9	0.231 -6	0.231 -6	0.356 -6	0.188 -5	0.475 -3
10	0.231 -6	0.231 -6	0.231 -6	0.307 -5	0.898 -5
11	0.233 -6	0.231 -6	0.231 -6	0.239 -6	0.325 -4
12	0.231 -6	0.231 -6	0.231 -6	0.249 -6	0.615 -6
13	0.324 -5	0.232 -6	0.232 -6	0.231 -6	0.970 -6

The results given in Table 1 correspond to $c = 0.25$, $d = 7$, $M = 16$ (i.e., $h \doteq \doteq 0.098$). The results obtained with various other values of these parameters have very similar character.

The computation was carried out on an IBM computer in double precision. The resulting system of linear algebraic equations was solved by the Gaussian elimination. In Table 1 we present the quantity

$$\eta(N, K) = ((4M + 1)^{-1} \sum_{m=-2M}^{2M} |\hat{u}_N(\frac{1}{2}mh) - u(\frac{1}{2}mh)|^2)^{1/2}$$

where the associate approximate solution \hat{u}_N depends on K in the sense of the above remarks and Definitions 3.1 and 3.2. The format of the table entries is a mantissa and a decimal exponent. Let us note that

$$\max_{x \in \langle -\pi/2, \pi/2 \rangle} |u(x)| \doteq 0.020$$

for $c = 0.25$, $d = 7$.

The values of $\eta(N, N)$ decrease (for fixed $M = 16$) with increasing N up to $N = 5$. In double precision (the mantissa of $t = 56$ binary digits), the rounding errors prevail in the total error for $N > 5$ and $\eta(N, N)$ is almost constant up to $N = 12$; then it increases. This implies that it is necessary to use more precision in order to obtain results of higher accuracy.

On the other hand, it enables us to present a true illustration for our statements. The rounding errors are so large that the result is not affected by narrowing the band of the matrix of the system up to some width. This "critical bandwidth" depends on N . For $N = 7$ the critical bandwidth is 11 (i.e., 1 + 1 diagonals may be deleted), for $N = 12$ the critical bandwidth is 17 (i.e., 3 + 3 diagonals may be deleted).

In general, the presented numerical results confirm the statements of this section concerning the bounds (4.2) and (4.10).

References

- [1] *I. Babuška*: Approximation by hill functions. Comment. Math. Univ. Carolinae 11 (1970), 787–811.
- [2] *W. Feller*: An introduction to probability theory and its applications. Vol. 2. Wiley, New York 1966.
- [3] *I. S. Gradštejn, I. M. Ryzhik*: Tables of integrals, sums, series, and products (Russian). 5th edition. Nauka, Moskva 1971.
- [4] *K. Segeth*: Universal approximation by hill functions. Czechoslovak Math. J. 22 (1972), 612–640.
- [5] *J. Segethová*: Numerical construction of the hill functions. SIAM J. Numer. Anal. 9 (1972), 199–204.
- [6] *J. H. Wilkinson*: Error analysis of direct methods of matrix inversion. J. Assoc. Comput. Mach. 8 (1961), 281–330.
- [7] *J. H. Wilkinson*: Rounding errors in algebraic processes. HMSO, London 1963.

Souhrn

ZUŽOVÁNÍ PÁSU PŘI ŘEŠENÍ SOUSTAV LINEÁRNÍCH ALGEBRAICKÝCH ROVNIC VZNIKAJÍCÍCH Z METODY KONEČNÝCH PRVKŮ

JITKA SEGETHOVÁ

Matice soustavy lineárních algebraických rovnic, vznikající při řešení jedno-rozměrných úloh metodou konečných prvků, je pásová. Při aproximaci vysokého řádu je její pás velmi široký, avšak prvky na okraji pásu (daleko od diagonály) jsou zanedbatelně malé vzhledem k diagonálním prvkům.

V práci se ukazuje, že je prakticky možné pracovat s maticí soustavy, jejíž pás je zúžen, tj. některé nenulové diagonály původní matice jsou nahrazeny nulami. Nepřesnost pramenící z tohoto zúžení pásu může být při dostatečně vysokém řádu aproximace (tj. při dostatečně širokém pásu původní matice) zanedbatelná ve srovnání se zaokrouhlovacími chybami vznikajícími při numerickém řešení soustavy.

Celý problém se studuje na modelové okrajové úloze a je doplněn numerickým příkladem.

Author's address: RNDr. Jitka Segethová, CSc., Matematicko-fyzikální fakulta Karlovy univerzity, Malostranské nám. 25, 118 00 Praha 1.