

Finite Element Multistep Discretizations of Parabolic Boundary Value Problems

By Miloš Zlámal

Abstract. The initial-boundary value problem for a linear parabolic equation in an infinite cylinder under the Dirichlet boundary condition is solved by applying the finite element discretization in the space dimension and A_0 -stable multistep discretizations in time. No restriction on the ratio of the time and space increments is imposed. The methods are analyzed and bounds for the discretization error in the L_2 -norm are given.

1. Introduction. The problem we are considering is the initial-boundary value problem

$$(1.1) \quad \begin{aligned} \partial u / \partial t &= Lu & \text{for } (x, t) \in \Omega \times (0, \infty), \\ u &= 0 & \text{on } \Gamma \times (0, \infty), \\ u(x, 0) &= g(x) & \text{in } \Omega. \end{aligned}$$

Here,

$$(1.2) \quad \begin{aligned} Lu &= \sum_{i,j=1}^N \frac{\partial}{\partial x_i} \left(a_{ij}(x) \frac{\partial u}{\partial x_j} \right) - a(x)u, \\ a_{ij}(x) &= a_{ji}(x), \quad \sum_{i,j=1}^N a_{ij}(x) \xi_i \xi_j \geq \alpha \sum_{i=1}^N \xi_i^2, \quad \alpha = \text{const} > 0, \quad a(x) \geq 0, \end{aligned}$$

and $x = (x_1, \dots, x_N)$ is a point of a bounded domain Ω in Euclidean N -space R^N with a smooth boundary Γ .

Before formulating (1.1) in the weak variational form, let us introduce some notations. The norm $\|\cdot\|_{L_2(\Omega)}$ of the space $L_2(\Omega)$ and the scalar product are denoted by $\|\cdot\|_0$ and $(\cdot, \cdot)_0$, respectively. $H^m \equiv W_2^{(m)}(\Omega)$, $m = 0, 1, \dots$, denotes the Sobolev space defined by

$$\|v\|_{H^m} = \left(\sum_{|j| \leq m} \|D^j v\|_0^2 \right)^{1/2}.$$

Instead of $\|v\|_{H^m}$, we write $\|v\|_m$. H_0^1 is the closure of $\mathcal{D}(\Omega)$, the set of infinitely differentiable functions with compact support in Ω , in the norm $\|\cdot\|_1$. The energy norm $|v|_1$ is defined by $|v|_1^2 = a(v, v)$, where $a(v, w)$ is the energy bilinear functional

Received December 27, 1973.

AMS (MOS) subject classifications (1970). Primary 65N30.

Copyright © 1975, American Mathematical Society

$$a(v, w) = \int_{\Omega} \left[\sum_{i,j=1}^N a_{ij}(x) \frac{\partial v}{\partial x_i} \frac{\partial w}{\partial x_j} + a(x)vw \right] dx.$$

The weak form of (1.1) is to find for $t > 0$ the function $u \in H_0^1$ such that, besides the initial condition, it satisfies

$$(1.3) \quad (\dot{u}, \varphi)_0 + a(u, \varphi) = 0 \quad \forall \varphi \in H_0^1.$$

A well-known approach for getting an approximate solution of the problem (1.1) consists in first applying the Galerkin principle to (1.3). Let S be a finite-dimensional subspace of H_0^1 . The Galerkin solution is the function $U \in S$ which satisfies

$$(1.4) \quad (\dot{U}, \varphi)_0 + a(U, \varphi) = 0 \quad \forall \varphi \in S.$$

The Galerkin formulation yields a system of ordinary differential equations in time. A suitable discretization in time will give a computable approximate solution of the problem (1.1). The choice of finite element subspaces for S and of Crank-Nicolson and other one step discretizations in time was considered in several papers published in recent years (see references in [9]). In [9], we chose for S finite-dimensional subspaces V_h^p of H_0^1 which have the following approximation property: for any $v \in H^{p+1} \cap H_0^1$, there exists a function $\hat{v} \in V_h^p$ such that

$$(1.5) \quad \|v - \hat{v}\|_j \leq Ch^{p+1-j} \|v\|_{p+1}, \quad j = 0, 1,$$

C being a constant independent of the small positive parameter h and of the function v . Finite element subspaces constructed first for special domains, later for arbitrary curved domains (see [8], [2], [10], [11]), possess this property. The parameter h is, in general, the maximum diameter of all elements.

In this paper, we again choose the subspaces V_h^p for S , and we discretize (1.4) by a A_0 -stable linear multistep method. A_0 -stable linear multistep methods were introduced for ordinary differential equations by Cryer [3]. When we apply the multistep method (ρ, σ) , where

$$\rho(\xi) = \sum_{j=0}^{\nu} \alpha_j \xi^j, \quad \alpha_{\nu} > 0, \quad \sigma(\xi) = \sum_{j=0}^{\nu} \beta_j \xi^j,$$

to the scalar equation $\dot{x}(t) = -\lambda x(t)$, $x(0) = 1$, the approximate values x^n of $x(nk)$ (k is the time increment) are determined by $\sum_{j=0}^{\nu} \alpha_j x^{n+j} = -k\lambda \sum_{j=0}^{\nu} \beta_j x^{n+j}$. A_0 -stability requires that $x^n \rightarrow 0$ as $n \rightarrow \infty$ for all positive λ . This is fulfilled iff all roots $\xi_j(\tau)$, $j = 1, \dots, \nu$, of the polynomial

$$(1.6) \quad p(\xi) = \rho(\xi) + \tau\sigma(\xi)$$

satisfy $|\xi_j(\tau)| < 1$ for every $\tau > 0$.

Denote by U^n the approximate values of U at the time level $t = nk$, $n = 0, 1, \dots$, and assume that $U^0, U^1, \dots, U^{\nu-1}$ are given. If we apply the scheme (ρ, σ) to (1.4), we get the recurrence relationship for U^n :

$$(1.7) \quad \left(\sum_{j=0}^{\nu} \alpha_j U^{n+j}, \varphi \right)_0 + ka \left(\sum_{j=0}^{\nu} \beta_j U^{n+j}, \varphi \right) = 0 \quad \forall \varphi \in V_h^p, \quad n = 0, 1, \dots.$$

Besides A_0 -stability, we require that the method (ρ, σ) be stable in the sense of Dahlquist and of the order $q \geq 1$ and that the roots of the polynomial $\sigma(\xi)$ with modulus equal to one be simple. Under the assumption that the solution of (1.1) is smooth enough, we prove the following bound which is uniform for $\nu \leq n < \infty$ and which holds without any restriction on the ratio kh^{-1} :

$$\sup_{\nu \leq n < \infty} \|u^n - U^n\|_0 \leq C \left[\sum_{j=0}^{\nu-1} \|u^j - U^j\|_0 + (h^{p+1} + k^q) \log \frac{1}{k} \|g\|_m \right];$$

here u^n are the exact values $u(x, nk)$.

To see what computations the relationship (1.7) represents, let us choose a basis v_1, \dots, v_l of V_h^p (of course, in finite element subspaces, we do not choose an arbitrary basis). Let M be the so-called mass matrix, $M = \{(v_i, v_j)_0\}_{i,j=1}^l$, and K the stiffness matrix, $K = \{a(v_i, v_j)\}_{i,j=1}^l$. Both these matrices are positive definite. If $\mathbf{v} = (v_1, \dots, v_l)^T$ (the superscript T denotes transposition) and $U^n = (\alpha^n)^T \mathbf{v}$, where $\alpha^n = (\alpha_1^n, \dots, \alpha_l^n)^T$, then setting in (1.7) the basis functions v_1, \dots, v_l for φ , we get

$$(1.8) \quad \sum_{j=0}^{\nu} (\alpha_j M + \beta_j k K) \alpha^{n+j} = 0, \quad n = 0, 1, \dots.$$

Evidently, at every time step we have to solve a linear system with the same matrix $B = \alpha_\nu M + \beta_\nu k K$. This matrix is positive definite (from A_0 -stability it follows $\beta_\nu > 0$; see [3, Theorem 3.1]) sparse, banded, and its condition number does not grow too fast. In the case of finite element subspaces, it follows that

$$(1.9) \quad \text{cond}(B) = O(kh^{-2}).$$

Multistep methods require the determination of starting values $U^0, \dots, U^{\nu-1}$, and it is desirable that these values be calculated to an accuracy as high as the local accuracy of the method. This disadvantage of multistep methods can be overcome, at least for ν small ($\nu \leq 4$), by computing U^0, \dots, U^ν by the Crank-Nicolson (i.e., trapezoidal) method or by the Calahan method (a third-order one-step method, see [9]) and using a step sufficiently small with respect to the step k of the main method.

For simplicity, we restricted ourselves to the homogeneous problem (1.1). The generalization of (1.7) for the nonhomogeneous equation $\partial u / \partial t = Lu + F(x, t)$ is immediate:

$$\left(\sum_{j=0}^{\nu} \alpha_j U^{n+j}, \varphi \right)_0 + ka \left(\sum_{j=0}^{\nu} \beta_j U^{n+j}, \varphi \right) = k \left(\sum_{j=0}^{\nu} \beta_j F^{n+j}(x), \varphi \right)_0$$

$$\forall \varphi \in V_h^p, \quad F^n(x) = F(x, nk).$$

The same bound for the error can be proved if t runs through a finite interval $\langle 0, T \rangle$. For the infinite interval $\langle 0, \infty \rangle$, such a bound cannot, of course, be proved unless some assumption on the growth of $F_t(x, t)$ is imposed.

The exact solution of the problem (1.1) has the property that

$$(1.10) \quad \|u(x, t)\|_0 \leq e^{-\lambda_1 t} \|g\|_0$$

for any $g \in L_2(\Omega)$. Here λ_1 is the smallest (positive) eigenvalue of the operator $-Lu$. Under the additional assumptions that the roots of the polynomial $\rho(\zeta)$ with modulus equal to one are real and the modulus of all roots of the polynomial $\sigma(\zeta)$ is less than one, we prove that scheme (1.7) preserves the asymptotic behavior characteristic of (1.1), again without placing any restriction on the ratio kh^{-1} : if $U^j \in L_2(\Omega)$, $j=0, \dots, \nu-1$, then $\|U^n\|_0$ decreases exponentially,

$$(1.11) \quad \|U^n\|_0 \leq Ce^{-\alpha_0 n k} \max_{0 \leq j \leq \nu-1} \|U^j\|_0, \quad \alpha_0 = \text{const} > 0, \quad n \geq \nu.$$

The backward differentiation multistep methods (see [7, p. 242]) with the step number $\nu \leq 6$ possess all the above properties.

2. Preliminaries. For simplicity, we assume that

$$(2.1) \quad a_{ij}(x), a(x), g(x) \in C^\infty(\bar{\Omega}), \quad \Gamma \in C^\infty.$$

We state some facts about the solution $u(x, t)$ of (1.1). It is of the form $\sum_{i=1}^\infty g_i e^{-\lambda_i t} \psi_i(x)$ where λ_i and $\psi_i(x)$ are (positive) eigenvalues and (orthonormal) eigenfunctions, respectively, of the problem

$$(2.2) \quad -L\psi = \lambda\psi, \quad \psi|_\Gamma = 0,$$

and g_i are the Fourier coefficients of the initial value $g(x)$. Ladyženskaja [6, Chapter III, Section 17] showed that if $g \in H^m$ and

$$(2.3) \quad g|_\Gamma = Lg|_\Gamma = \dots = L^{[(m-1)/2]}g|_\Gamma = 0,$$

then $u(x, t) \in H^m$ for $t \geq 0$. Conversely, if $u(x, t) \in H^m$ for $t \geq 0$ then $g \in H^m$ and (2.3) is satisfied. The proof is based on two inequalities. The first holds for any series $\sum_{i=1}^\infty g_i \psi_i(x)$:

$$(2.4) \quad \left\| \sum_{i=1}^\infty g_i \psi_i(x) \right\|_m^2 \leq C \sum_{i=1}^\infty \lambda_i^m g_i^2.$$

(In the sequel, C is a generic constant, not necessarily the same in any two places, which does not depend on h, k, n, l, τ, g .) Concerning the other, we need only the following consequence: if $g \in H^m$ and (2.3) is satisfied, then it holds that

$$(2.5) \quad \sum_{i=1}^\infty \lambda_i^m g_i^2 \leq C \|g\|_m^2.$$

3. Convergence. The main results of the paper are contained in the following

THEOREM. *Let the linear multistep method (ρ, σ) be stable in the sense of Dahlquist, A_0 -stable and of the order $q \geq 1$, and let the roots of the polynomial $\sigma(\xi)$ with modulus equal to one be simple. Let (1.2) and (2.1) hold and g satisfy (2.3) with $m = \max(p + 1, 2q)$ (this requirement is equivalent to the assumption $u(x, t) \in H^m$ for $t \geq 0$). Then, for arbitrary h, k , the discretization error is bounded by*

$$(3.1) \quad \sup_{\nu \leq n < \infty} \|u^n - U^n\|_0 \leq C \left[\sum_{j=0}^{\nu-1} \|u^j - U^j\|_0 + (h^{p+1} + k^q) \log \frac{1}{k} \|g\|_m \right].$$

If, in addition, the roots of the polynomial $\rho(\xi)$ with modulus equal to one are real and the modulus of all roots of $\sigma(\xi)$ is less than one, then (1.11) holds for any $U^j \in L_2(\Omega)$, $j = 0, \dots, \nu - 1$.

Proof. We first write u^n in the form $u^n = \xi^n + \eta^n$ with $\eta^n \in V_h^p$ being the Ritz approximation of u^n , i.e., the orthogonal projection of u^n onto V_h^p with the energy norm $[a(\cdot, \cdot)]^{1/2}$ (several authors have used this decomposition; we learned it from Bramble, Thomée [1]). Hence,

$$(3.2) \quad a(\eta^n, \varphi) = a(u^n, \varphi) = (-Lu^n, \varphi)_0 = (-\dot{u}^n, \varphi)_0, \quad \forall \varphi \in V_h^p,$$

and with respect to (1.5), we find (see, e.g., [10]) that

$$\|\xi^n\|_0 \leq Ch^{p+1} \|Lu^n\|_{p-1} \leq Ch^{p+1} \|u^n\|_{p+1}.$$

By means of (2.4) and (2.5), we immediately obtain

$$(3.3) \quad \|\xi^n\|_0 \leq Ch^{p+1} \|g\|_{p+1}.$$

Therefore it is sufficient to prove for $\epsilon^n = \eta^n - U^n$

$$(3.4) \quad \max_{\nu \leq n < \infty} \|\epsilon^n\|_0 \leq C \left[\sum_{j=0}^{\nu-1} \|\epsilon^j\|_0 + (h^{p+1} + k^q) \log \frac{1}{k} \|g\|_m \right].$$

If we use (3.2), we see that

$$(3.5) \quad \left(\sum_{j=0}^{\nu} \alpha_j \eta^{n+j}, \varphi \right)_0 + ka \left(\sum_{j=0}^{\nu} \beta_j \eta^{n+j}, \varphi \right) = (\pi^n - \omega^n, \varphi)_0, \quad \forall \varphi \in V_h^p,$$

where

$$\pi^n = \sum_{j=0}^{\nu} (\alpha_j u^{n+j} - k\beta_j \dot{u}^{n+j}), \quad \omega^n = \sum_{j=0}^{\nu} \alpha_j \xi^{n+j}.$$

Subtracting (1.7) from (3.5), we get

$$(3.6) \quad \left(\sum_{j=0}^{\nu} \alpha_j \epsilon^{n+j}, \varphi \right)_0 + ka \left(\sum_{j=0}^{\nu} \beta_j \epsilon^{n+j}, \varphi \right) = (\pi^n - \omega^n, \varphi)_0, \quad \forall \varphi \in V_h^p.$$

We write (3.6) in a matrix form. For this purpose, let \mathbf{w} be the vector $\mathbf{w} = M^{-1/2}\mathbf{v}$ (\mathbf{v} is the basis vector; see Introduction) and let us set $\epsilon^n = (\epsilon^n)^T \mathbf{w}$ (notice that $\epsilon^n = \eta^n - U^n \in V_h^p$). Since $(\mathbf{v}, \mathbf{v}^T)_0 = M$ and $a(\mathbf{v}, \mathbf{v}^T) = K$, we have $(\mathbf{w}, \mathbf{w}^T)_0 = I$ and $a(\mathbf{w}, \mathbf{w}^T) = M^{-1/2}KM^{-1/2}$. The matrix $S = M^{-1/2}KM^{-1/2}$ is symmetric and positive definite. Putting the components w_i ($i = 1, \dots, l$) of the vector \mathbf{w} for φ in (3.6), we get

$$\sum_{j=0}^{\nu} (\alpha_j I + \beta_j kS) \epsilon^{n+j} = \mathbf{c}^n,$$

where

$$(3.7) \quad \mathbf{c}^n = (\pi^n - \omega^n, \mathbf{w})_0.$$

Denote

$$(3.8) \quad \delta_j(\tau) = \frac{\alpha_j + \beta_j \tau}{\alpha_\nu + \beta_\nu \tau}, \quad j = 0, \dots, \nu \quad (\delta_\nu(\tau) \equiv 1), \quad \mathbf{d}^n = (\alpha_\nu I + \beta_\nu kS)^{-1} \mathbf{c}^n$$

(the matrix $\alpha_\nu I + \beta_\nu kS$ is positive definite since $\alpha_\nu > 0$, $\beta_\nu > 0$). Then

$$(3.9) \quad \sum_{j=0}^{\nu} \delta_j(kS) \epsilon^{n+j} = \mathbf{d}^n$$

and this difference equation will be solved in the way described by Henrici [5, pp. 242–244].

We define the coefficients $\gamma_l(\tau)$ ($l = 0, 1, \dots$) by

$$(3.10) \quad \frac{1}{\hat{p}(\zeta, \tau)} \equiv [\delta_\nu(\tau) + \delta_{\nu-1}(\tau)\zeta + \dots + \delta_0(\tau)\zeta^\nu]^{-1} = \gamma_0(\tau) + \gamma_1(\tau)\zeta + \dots$$

Similarly as in [5, p. 242], we can prove the estimate

$$(3.11) \quad |\gamma_l(\tau)| \leq C, \quad \tau \geq 0, \quad l = 0, 1, \dots$$

(we leave out the proof even when it is not a trivial matter). We also get the identity (see [5, (5–160), p. 243])

$$(3.12) \quad \delta_\nu(\tau)\gamma_l(\tau) + \delta_{\nu-1}(\tau)\gamma_{l-1}(\tau) + \dots + \delta_0(\tau)\gamma_{l-\nu}(\tau) \equiv \begin{cases} 1, & l = 0, \\ 0, & l > 0. \end{cases}$$

Now we write (3.9) with $n - \nu - l$ instead of n , multiply by $\gamma_l(kS)$ and sum for $l = 0, 1, \dots, n - \nu$. After some rearranging and using (3.12), we obtain (see [5, p. 243])

$$(3.13) \quad \begin{aligned} \epsilon^n = & - [\delta_{\nu-1}(kS)\gamma_{n-\nu}(kS) + \dots + \delta_0(kS)\gamma_{n-2\nu+1}(kS)] \epsilon^{\nu-1} - \dots \\ & - \delta_0(kS)\gamma_{n-\nu}(kS) \epsilon^0 + \sum_{l=0}^{n-\nu} \gamma_l(kS) \mathbf{d}^{n-\nu-l}. \end{aligned}$$

The coefficients $\delta_j(\tau)$ are bounded functions in the interval $(0, \infty)$. Therefore, $\|\delta_j(kS)\| = \max_\Lambda |\delta_j(k\Lambda)|$ (by $\|\cdot\|$ we denote the Euclidean norm of a vector or a matrix) where Λ

are the eigenvalues of S . These are positive, consequently $\|\delta_j(kS)\| \leq \sup_{0 < \tau < h} |\delta_j(\tau)| = O(1)$. The coefficients $\gamma_l(\tau)$ are bounded by (3.11), hence $\|\gamma_l(kS)\| = O(1)$, $l = 0, 1, \dots$. Furthermore, the starting errors are bounded by

$$\|\epsilon^j\|_0 = \|u^j - U^j - \xi^j\|_0 \leq \|u^j - U^j\|_0 + \|\xi^j\|_0 \leq \|u^j - U^j\|_0 + Ch^{p+1} \|g\|_{p+1}.$$

Hence

$$\sum_{j=0}^{\nu-1} \|\epsilon^j\|_0 \leq \sum_{j=0}^{\nu-1} \|u^j - U^j\|_0 + Ch^{p+1} \|g\|_m.$$

Also $\|d^n\| \leq \|(\alpha_\nu I + \beta_\nu kS)^{-1}\| \|c^n\| \leq \alpha_\nu^{-1} \|c^n\|$. Thus, we see that from (3.13) it follows

$$\|\epsilon^n\| \leq C \left(\sum_{j=0}^{\nu-1} \|\epsilon^j\| + \sum_{r=0}^{n-\nu} \|c^r\| \right).$$

Since

$$\|\epsilon^n\|_0^2 = (\epsilon^n)^T (\mathbf{w}, \mathbf{w}^T)_0 \epsilon^n = \|\epsilon^n\|^2,$$

we have

$$(3.14) \quad \|\epsilon^n\|_0 \leq C \left(\sum_{j=0}^{\nu-1} \|u^j - U^j\|_0 + h^{p+1} \|g\|_m + \sum_{r=0}^{n-\nu} \|c^r\| \right).$$

We need a bound for $\|c^r\|$. c^r is of the form (3.7). If $f \in L_2(\Omega)$ and $\hat{f} = \hat{f}^T \mathbf{w} \in V_h^p$ is the orthogonal projection of f onto V_h^p with the norm $\|\cdot\|_0$, we easily find that $\hat{f} = (f, \mathbf{w})_0$ and that $\|\hat{f}\|_0 \leq \|f\|_0$. Since $\|\hat{f}\|_0 = \|\hat{f}\|$, we have $\|\hat{f}\| \leq \|f\|_0$. Therefore,

$$(3.15) \quad \|c^r\| \leq \|\pi^r\|_0 + \|\omega^r\|_0.$$

To estimate $\|\pi^r\|_0$, we use the assumption that the scheme (ρ, σ) is of the order q . It means that for any function $y(t) \in C^{(s)}$, $s \leq q + 1$, it holds that

$$(3.16) \quad \sum_{j=0}^{\nu} \alpha_j y(t + jk) - k \sum_{j=0}^{\nu} \beta_j \dot{y}(t + jk) = O \left(k^s \max_{0 < \tau < \nu k} |y^{(s)}(t + \tau)| \right)$$

(it follows from the formula (5.178) in [5, p. 248]). Set $y(t) = e^{-\lambda t}$, $s = q + 1$. After dividing by $e^{-\lambda i t}$, we get

$$\sum_{j=0}^{\nu} (\alpha_j + \beta_j k \lambda_i) e^{-jk \lambda_i} = O(k^{q+1} \lambda_i^{q+1}).$$

The Fourier coefficients π_i^r of π^r are evidently equal to

$$e^{-rk \lambda_i} g_i \sum_{j=0}^{\nu} (\alpha_j + \beta_j k \lambda_i) e^{-jk \lambda_i}.$$

Therefore,

$$(\pi_i^r)^2 \leq C e^{-2rk\lambda_i} k^{2(q+1)} \lambda_i^{2(q+1)} g_i^2.$$

Since $e^{-2rk\lambda_i} k^{2\lambda_i^2} \leq e^{-rk\lambda_1} (e^{-\frac{1}{2}rk\lambda_i} k\lambda_i)^2 \leq r^{-2} e^{-rk\lambda_1}$ if $r \geq 1$ (due to $x e^{-\alpha x} \leq (\alpha)^{-1} < (2\alpha)^{-1}$), we have $(\pi_i^r)^2 \leq C r^{-2} e^{-rk\lambda_1} k^{2q} \lambda_i^{2q} g_i^2$ and

$$(3.17) \quad \|\pi^r\|_0^2 = \sum_{i=1}^{\infty} (\pi_i^r)^2 \leq C r^{-2} e^{-rk\lambda_1} k^{2q} \|g\|_m^2, \quad r \geq 1.$$

Concerning π^0 , we use (3.16) with $s = q$ and we get

$$(3.18) \quad \|\pi^0\|_0^2 \leq C k^{2q} \|g\|_m^2.$$

The estimates for $\|\omega^r\|_0$ can be obtained in a similar way. Set $z^r = \sum_{j=0}^{\nu} \alpha_j u^{r+j}$ and write z^r as the sum $x^r + y^r$, where y^r is the Ritz approximation of z^r . Then $\|x^r\|_0 \leq C h^{p+1} \|z^r\|_{p+1}$. Since η^r is the Ritz approximation of u^r , we have $y^r = \sum_{j=0}^{\nu} \alpha_j \eta^{r+j}$ and $x^r = \sum_{j=0}^{\nu} \alpha_j \xi^{r+j} = \omega^r$, hence

$$(3.19) \quad \|\omega^r\|_0 \leq C h^{p+1} \|z^r\|_{p+1}.$$

The Fourier coefficients z_i^r are equal to $\sum_{j=0}^{\nu} \alpha_j e^{-(r+j)k\lambda_i} g_i = e^{-rk\lambda_i} g_i \sum_{j=0}^{\nu} \alpha_j e^{-jk\lambda_i}$. Because of the consistency of the scheme (ρ, σ) , it follows that $\sum_{j=0}^{\nu} \alpha_j = 0$. Therefore, $\sum_{j=0}^{\nu} \alpha_j e^{-jk\lambda_i} = O(k\lambda_i)$, consequently $(z_i^r)^2 \leq C e^{-2rk\lambda_i} (k\lambda_i)^2 g_i^2$; by means of (3.19), we find in the same way as before that

$$(3.20) \quad \begin{aligned} \|\omega^r\|_0^2 &\leq C r^{-2} e^{-rk\lambda_1} h^{2(p+1)} \|g\|_m^2, \quad r \geq 1, \\ \|\omega^0\|_0^2 &\leq C h^{2(p+1)} \|g\|_m^2. \end{aligned}$$

(3.17), (3.18), (3.20) and (3.15) give

$$\begin{aligned} \|c^r\| &\leq C r^{-1} e^{-\frac{1}{2}rk\lambda_1} (h^{p+1} + k^q) \|g\|_m, \quad r \geq 1, \\ \|c^0\| &\leq C (h^{p+1} + k^q) \|g\|_m. \end{aligned}$$

We come back to (3.14). If we find out that

$$\sum_{r=1}^{\infty} \|c^r\| \leq C (h^{p+1} + k^q) \log \frac{1}{k} \|g\|_m,$$

the bound (3.1) is proved. For this purpose, it is sufficient to realize that

$$\sum_{r=1}^{\infty} r^{-1} e^{-\frac{1}{2}rk\lambda_1} = \sum_{r=1}^{\infty} r^{-1} (e^{-\frac{1}{2}k\lambda_1})^r = -\log(1 - e^{-\frac{1}{2}k\lambda_1}) = O\left(\log \frac{1}{k}\right).$$

To prove the second part of the Theorem, we need a better estimate of the coefficients $\gamma_l(\tau)$. According to our assumptions, the roots λ of $\rho(\zeta)$ with modulus equal to one are simple and real (i.e., $\lambda = \pm 1$). The associated roots $\zeta(\tau)$ of the polynomial $p(\zeta)$ defined by (1.6) have the expansion $\zeta(\tau) = \lambda + a_1\tau + O(\tau^2)$, with $a_1 = -\sigma(\lambda)/\rho'(\lambda)$

$\neq 0$ (since $\sigma(\lambda) = 0$ means that λ would be a root of $p(\zeta)$). Therefore, $|\zeta(\tau)| = |1 + \lambda^{-1}a_1\tau + O(\tau^2)|$. The growth parameter $a = \lambda^{-1}a_1$ is different from zero and real, hence $|\zeta(\tau)| = [1 + 2a\tau + O(\tau^2)]^{1/2}$. As a must be negative (otherwise $|\zeta(\tau)| > 1$), we see that, for τ sufficiently small, $|\zeta(\tau)| \leq 1 - c\tau$, $c = \frac{1}{2}|a| > 0$. The other roots of $p(\zeta)$ have modulus less than one for $\tau \geq 0$. On the basis of these facts, we can prove that $|\gamma_i(\tau)| \leq C(1 - \frac{1}{2}c\tau)^i$ for $0 < \tau \leq \tau_1$ if τ_1 is a sufficiently small number. In the interval (τ_1, ∞) , it follows easily (from the assumption that the modulus of all roots of $\sigma(\zeta)$ is less than one) that $|\gamma_i(\tau)| \leq C(1 - \vartheta)^i$, $0 < \vartheta < 1$. Hence

$$|\gamma_i(\tau)| \leq C[\max(1 - \frac{1}{2}c\tau, 1 - \vartheta)]^i, \quad \tau > 0,$$

and, because the eigenvalues of the matrix kS are bounded from below by $k\lambda_1$, it holds for k sufficiently small that

$$(3.21) \quad \|\gamma_i(kS)\| \leq C(1 - \frac{1}{2}c\lambda_1 k)^i \leq Ce^{-\alpha_0 ik}, \quad \alpha_0 = \frac{1}{4}c\lambda_1 > 0.$$

To complete the proof, we set $U^n = (a^n)^T w$ and get (in the same way as we got (3.13))

$$a^n = -[\delta_{\nu-1}(kS)\gamma_{n-\nu}(kS) + \dots + \delta_0(kS)\gamma_{n-2\nu+1}(kS)]a^{\nu-1} \\ - \dots - \delta_0(kS)\gamma_{n-\nu}(kS)a^0.$$

The estimate (1.12) follows immediately from (3.21).

4. Some Remarks. 1. To get the estimate (1.9) for the condition number of the matrix $B = \alpha_\nu M + \beta_\nu kK$, we assume the following additional properties of the basis $\{v_1, \dots, v_l\}$ of the space V_h^p : if $\varphi = x^T v \in V_h^p$, then

$$(a) \quad ch^{-N} \|\varphi\|_0^2 \leq \|x\|^2 \leq Ch^{-N} \|\varphi\|_0^2, \quad c = \text{const} > 0,$$

$$(b) \quad a(\varphi, \varphi) \leq Ch^{-2} \|\varphi\|_0^2.$$

The finite element subspaces used in applications possess these properties.

Let Λ be an eigenvalue of the matrix B and x the corresponding eigenvector. We have $\alpha_\nu Mx + \beta_\nu kKx = \Lambda x$; multiplying this equation by x^T and setting $\varphi = x^T v$, we get

$$\Lambda = \alpha_\nu \frac{\|\varphi\|_0^2}{\|x\|^2} + \beta_\nu k \frac{a(\varphi, \varphi)}{\|x\|^2}.$$

Hence

$$\Lambda_{\max} \leq C(1 + kh^{-2})h^N, \quad \Lambda_{\min} \geq c_1 h^N, \quad c_1 = \text{const} > 0,$$

and, if we exclude the uninteresting case $kh^{-2} \rightarrow 0$, we have $\text{cond}(B) = \Lambda_{\max}/\Lambda_{\min} = O(kh^{-2})$.

2. The backward differentiation methods with $\nu = q \leq 6$ are stable in the sense

of Dahlquist, the only root of $\rho(\xi)$ with modulus equal to one is the principal root $\xi = 1$ (see Cryer [4]), and they are A_0 -stable (actually, they are $A(\alpha)$ -stable with $90^\circ \geq \alpha \geq 18^\circ$, see [7, p. 242]). Further, the only root of $\sigma(\xi)$ is zero. Hence, these methods fulfill all assumptions of the Theorem.

3. It is known (see [7, p. 243]) that the implicit R -stage Runge-Kutta methods of order $2R$ are A -stable. It will be shown elsewhere that if we discretize (1.4) by means of such a method, the error is bounded by

$$\sup_{1 \leq n < \infty} \|u^n - U^n\|_0 \leq \|u^0 - U^0\|_0 + C(h^{p+1} + k^{2R}) \log \frac{1}{k} \|g\|_m,$$

under the assumption that $g(x)$ satisfies (2.3) with $m = \max(p + 1, 4R)$.

Computing Center of the Technical University
Obránců míru 21
602 00 Brno, Czechoslovakia

1. J. H. BRAMBLE & V. THOMÉE, "Discrete time Galerkin methods for a parabolic boundary value problem." (To appear.)
2. P. G. CIARLET & P. A. RAVIART, "Interpolation theory over curved elements, with applications to finite element methods," *Comp. Meth. Appl. Mech. Eng.*, v. 1, 1972, pp. 217–249.
3. C. W. CRYER, "A new class of highly-stable methods: A_0 -stable methods," *BIT*, v. 13, 1973, pp. 153–159.
4. C. W. CRYER, "On the instability of high order backward-difference multistep methods," *BIT*, v. 12, 1972, pp. 17–25.
5. P. HENRICI, *Discrete Variable Methods in Ordinary Differential Equations*, Wiley, New York-London, 1962. MR 24 #B1772.
6. O. A. LADYŽENSKAJA, V. A. SOLONNIKOV & N. N. URAL'CEVA, *Linear and Quasi-Linear Equations of Parabolic Type*, Transl. Math. Monographs, vol. 23, Amer. Math. Soc., Providence, R. I., 1968. MR 39 #3159b.
7. J. D. LAMBERT, *Computational Methods in Ordinary Differential Equations*, Wiley, London, 1973.
8. O. C. ZIENKIEWICZ, *The Finite Element Method in Engineering Science*, McGraw-Hill, London, 1971. MR 47 #4518.
9. M. ZLÁMAL, "Finite element methods for parabolic equations," *Math. Comp.*, v. 28, 1973, pp. 393–404.
10. M. ZLÁMAL, "Curved elements in the finite element method. I," *SIAM J. Numer. Anal.*, v. 10, 1973, pp. 229–240.
11. M. ZLÁMAL, "Curved elements in the finite element method. II," *SIAM J. Numer. Anal.*, v. 11, 1974, pp. 347–362.