

## A High-Order Difference Method for Differential Equations\*

By Robert E. Lynch and John R. Rice

**Abstract.** This paper analyzes a high-accuracy approximation to the  $m$ th-order linear ordinary differential equation  $Mu = f$ . At mesh points,  $U$  is the estimate of  $u$ ; and  $U$  satisfies  $M_n U = I_n f$ , where  $M_n U$  is a linear combination of values of  $U$  at  $m + 1$  stencil points (adjacent mesh points) and  $I_n f$  is a linear combination of values of  $f$  at  $J$  auxiliary points, which are between the first and last stencil points. The coefficients of  $M_n$ ,  $I_n$  are obtained "locally" by solving a small linear system for each group of stencil points in order to make the approximation exact on a linear space  $S$  of dimension  $L + 1$ . For separated two-point boundary value problems,  $U$  is the solution of an  $n$ -by- $n$  linear system with full bandwidth  $m + 1$ . For  $S$  a space of polynomials, existence and uniqueness are established, and the discretization error is  $O(h^{L+1-m})$ ; the first  $m - 1$  divided differences of  $U$  tend to those of  $u$  at this rate. For a general set of auxiliary points one has  $L = J + m$ ; but special auxiliary points, which depend upon  $M$  and the stencil points, allow larger  $L$ , up to  $L = 2J + m$ . Comparison of operation counts for this method and five other common schemes shows that this method is among the most efficient for given convergence rate. A brief selection from extensive experiments is presented which supports the theoretical results and the practicality of the method.

**1. Introduction.** We consider some aspects of a new flexible finite-difference method which gives high-accuracy approximation to solutions  $u$  of linear differential equations  $Mu = f$  subject to rather general initial or boundary conditions. The approximation to  $u$  is taken as  $U$  defined at mesh points as the solution of a system of difference equations  $M_n U = I_n f$  together with appropriate boundary conditions;  $n$  is used to identify a particular partition of the domain of  $u$ .  $M_n$  is a difference operator and  $M_n U$  is a linear combination of values of  $U$  at a small\*\* number of mesh points of a standard stencil; the value of  $I_n f$  is equal to a linear combination of values of  $f$  at several *auxiliary points* close to the stencil points and always inside the domain of  $u$ .

With appropriate normalization of the coefficients of  $M_n$  and  $I_n$ , then  $I_n f$  is  $f + O(h)$ , where  $h$  is a norm of the partition; and thus, the operator  $I_n$  can be regarded as a perturbation, or an expansion, of the identity operator as is commonly done for

---

Received July 1, 1976; revised September 23, 1977 and July 16, 1979.

AMS (MOS) subject classifications (1970). Primary 65L10.

\*Work supported in part by National Science Foundation Grants GP3290X and 7610225.

\*\*For an  $m$ th-order ordinary differential equation, one uses  $m + 1$  adjacent mesh points as stencil points; for a second-order elliptic partial differential equation, one uses a nine-point stencil for two independent variables, a 27-point stencil for three independent variables, and so on.

such operators in Approximation Theory. We have named this method *High-Order Difference approximation with Identity Expansions* which leads to the pronounceable acronym HODIE.

In this paper the application of the HODIE method to *ordinary differential equation problems* is treated. The analysis and results presented here give insight into the more complicated—and more important—application of HODIE to the solution of partial differential equations. Results about the multi-dimensional applications are given by Lynch and Rice [1975], [1977], [1978a], [1978b], by Lynch [1977a], [1977b], and by Boisvert [1978]; more detailed analyses will be presented at a later time. The method was discovered by R. E. Lynch during a study of methods for approximating solutions of elliptic partial differential equations in two independent variables.

Some of the key features of the method include: (a) the small number of stencil points which leads to a matrix with small bandwidth; (b) the coefficients of the operators  $M_n, I_n$  are determined so that the approximation is exact on a linear space of functions and their values are obtained by solving a small local system of linear algebraic equations whose size is fixed independent of the mesh length; (c) high accuracy is obtained by the use of values of  $f$  at the auxiliary points rather than with additional stencil points; (d) a variety of boundary conditions can easily be approximated with high accuracy; (e) the method is computationally efficient.

Although the difference equation is similar to one obtained by the Mehrstellenverfahren (or the “Hermitian” method) of Collatz [1960] after one replaces derivatives of  $f$  with divided differences, the method of obtaining the coefficients of the difference equation is different from that of Mehrstellenverfahren.

For ordinary differential equations, the HODIE method gives the same difference equations as obtained by Osborne [1967] who generalized the Størmer-Numerov scheme. Osborne was pessimistic about its practicality; he did not prove convergence results. More recently and independently, Doedel [1976], [1978] presented an essentially equivalent method for the ordinary differential equation case and he proved some results. Doedel also presents results about difference schemes which use more than the minimal number,  $m + 1$ , of stencil points for an  $m$ th-order ordinary differential operator; we do not consider this case. The results presented below are more complete than those of Doedel for the cases we treat. Both Osborne’s and Doedel’s approaches lead to different and less efficient implementations than the one described below.

This paper is briefly summarized as follows. In Section 2, a description of the HODIE method for ordinary differential equations is presented and some simple examples are given. In Section 3, a bound on the truncation error is given when the HODIE method is exact on  $P_L$ , the space of polynomials of degree at most  $L$ . For an  $m$ th-order operator and when the location of the auxiliary points is independent of  $M$ , the order of the truncation error is  $L - m + 1$ . Higher-order (“superconvergence”) is obtained with special auxiliary points (“Gauss-type points”) whose location depends upon  $M$ . In Section 4, approximation for the simple operator  $M = d^m/dt^m$

is treated in detail. A direct relationship between the truncation error and quadrature error is demonstrated and Gauss-type auxiliary points are introduced and analyzed. These Gauss-type points are the zeros of polynomials orthogonal with respect to an integral inner product with weight function a polynomial  $B$ -spline. In Section 5, we extend the results of Section 4 to the general linear variable coefficient differential operator with leading term  $d^m/dt^m$ . In Section 6, we show that the HODIE method gives a stable difference approximation. In Section 7 we show that the order of the discretization error is equal to the order of the truncation error for the initial value problem; in Section 8 the same result is shown for the separated two-point boundary value problem. In these convergence results, the convergence and rate of convergence applies for the first  $m - 1$  divided differences of  $U$  to those of  $u$ . Section 9 contains a comparison of the computational effort for the HODIE method and five other methods; this suggests that the HODIE method is among the most efficient methods available for solving second-order boundary value problems. Finally, in Section 10, we give a small sample of extensive experimental results which verify that the HODIE method works as the theory predicts and that there are no unforeseen difficulties in its implementation. Numerical examples indicate that the use of the Gauss-type auxiliary point for the operator  $D^m$ , which can be calculated in advance for a uniform mesh and which are then independent of the mesh size, can yield enhanced accuracy for the general  $m$ th-order problem with operator  $M$ .

**2. Approximation of Differential Operators.** We construct and analyze high-accuracy  $(m + 1)$ -point difference approximation to  $m$ th-order differential equations  $M[u, f] = 0$  subject to appropriate initial or two-point boundary conditions  $M^k[u, c_k] = 0, k = 0, \dots, m - 1$ , where

$$(2-1a) \quad M[u, f](t) = Mu(t) - f(t), \quad A < t < B,$$

$$(2-1b) \quad Mu(t) = D^m u(t) + \sum_{i=0}^{m-1} a_i(t)D^i u(t), \quad D = d/dt,$$

$$(2-1c) \quad M^k[u, c_k] = M^k u(A) + M^k u(B) - c_k, \quad k = 0, \dots, m - 1,$$

$$(2-1d) \quad M^k u(t) = \sum_{i=0}^{m-1} a_{k,i}(t)D^i u(t).$$

For the initial value problem,  $a_{k,i}(A) = 0$  if  $i \neq k, a_{k,k}(A) = 1$ , and  $a_{k,i}(B) = 0, i, k = 0, \dots, m - 1$ ; for the separated two-point boundary value problem, either  $M^k u(A)$  or  $M^k u(B)$  is zero,  $k = 0, \dots, m - 1$ .

The interval  $A \leqq t \leqq B$  is partitioned into  $n$  subintervals by  $n + 1$  mesh points  $t_k: A = t_0 < t_1 < \dots < t_n = B$ , with  $m \leqq n$ . The approximation  $U$  to the solution  $u$  is obtained at these mesh points as the solution of a system of  $m$ th-order difference equations subject to appropriate boundary conditions.

The approximation  $M_n[u, f] = M_n u - I_n f$  of the differential operator  $M$  is obtained locally by use of a pair of point sets and a set of basis functions. The  $m + 1$  stencil points are  $m + 1$  adjacent mesh points:  $\bar{t}_k = (t_k, t_{k+1}, \dots, t_{k+m})$ , and we

set  $h_k = (t_{k+m} - t_k)/m$ . The difference operator  $M_n$  with coefficients  $\alpha$  is

$$M_n U(t_k) = (1/h_k^m) \sum_{i=0}^m \alpha_{k,i} U_{k+i}, \quad U_{k+i} \equiv U(t_{k+i}).$$

The second set of points comprise  $J$  distinct auxiliary points  $\bar{\tau}_k = (\tau_{k,1}, \dots, \tau_{k,J})$ , subject to the restrictions  $t_k \leq \tau_{k,1} < \dots < \tau_{k,J} \leq t_{k+m}$ . The identity expansion  $I_n$  with coefficients  $\beta$  is

$$I_n f(t_k) = \sum_{j=1}^J \beta_{k,j} f_{k,j}, \quad f_{k,j} \equiv f(\tau_{k,j}).$$

For a given  $f$ ,  $U$  is the solution of  $M_n[U, f] = M_n U - I_n f = 0$  subject to appropriate boundary conditions.

The coefficients  $\alpha, \beta$  of the operators  $M_n$  and  $I_n$  are determined so that the approximation is exact on an  $(L + 1)$ -dimensional linear space  $S$  of functions. A basis  $s_0, s_1, \dots, s_L$  for  $S$  is chosen, and the coefficients are made to satisfy the HODIE equations  $M_n[s_l, Ms_l] = 0, l = 0, \dots, L$ ; that is,

$$(2-2) \quad (1/h_k^m) \sum_{i=0}^m \alpha_{k,i} s_l(t_{k+i}) - \sum_{j=1}^J \beta_{k,j} Ms_l(\tau_{k,j}) = 0, \quad l = 0, \dots, L.$$

The system (2-2) is homogeneous in the coefficients  $\alpha, \beta$ ; hence, in addition to (2-2), we take some convenient normalization equation such as

$$(2-3) \quad (a) \ \beta_{k,1} = 1; \quad (b) \ \sum_j |\beta_{k,j}| = 1; \quad \text{or} \quad (c) \ \sum_j \beta_{k,j} = 1.$$

The first is used in actual computation since it simplifies the calculation. The second and third are useful at various places in the theoretical treatment. It is a consequence of the analysis in Sections 4 and 5 that the third normalization can be used. Remarks about bases and efficient methods of solving the HODIE equations (2-2) are given in Section 9.

Boundary conditions for  $U$  are obtained in a similar way; they are treated in Sections 7 and 8.

The truncation error is defined with respect to a space of functions  $\Sigma$  in terms of the truncation operator

$$(2-4) \quad T_n[\sigma] = M_n[\sigma, M\sigma] - M[\sigma, M\sigma] = M_n\sigma - I_n(M\sigma), \quad \sigma \in \Sigma.$$

For  $\sigma \in \Sigma$ , the truncation error is the value of the max-norm of  $T_n[\sigma]$ , namely,  $\|T_n[\sigma]\| = \max_k |T_n[\sigma]_k|$ . In Section 3, we obtain a bound on the truncation error for the case that  $S$  is a space of polynomials and  $\Sigma$  is a space of sufficiently smooth functions.

The truncation error is related to the discretization error, defined as the max-norm of the error  $e = u - U$  at mesh points. This is because if  $u \in \Sigma$ , then  $M_n e = M_n u - M_n U = M_n u - I_n(Mu) = T_n u$ ; that is,  $e$  satisfies the equation  $M_n e = T_n u$ . In Sections 7 and 8 we show with natural hypotheses and appropriate boundary condition approximation, a bound on the truncation error yields a similar bound on the discretization error.

*Examples.* We consider a few examples for equal-spaced mesh points with spacing  $h$  and the operator  $Mu = D^2u + a_1Du + a_0u$ . It is sufficient to consider  $t_k = -h, t_{k+1} = 0, t_{k+2} = h$ . For brevity, we use a single subscripted notation for the coefficients  $\alpha, \beta$  and the auxiliary points  $\tau$ . For approximation which is exact on the space  $S = \mathbf{P}_L$  of polynomials of degree at most  $L$ , we use the Lagrange basis for quadratic interpolation together with elements which are zero at the three stencil points, specifically:

$$s_0(t) = t(t - h)/(2h^2), \quad s_1(t) = (h^2 - t^2)/h^2, \quad s_2(t) = t(t + h)/(2h^2),$$

$$s_l(t) = t^{l-2}(t^2 - h^2)/h^l, \quad l = 3, 4, \dots, L.$$

For normalization, we take the sum of the  $\beta$ 's to be equal to unity. After division by appropriate powers of  $h$ , the HODIE equations (2-2) and the normalization equation become

$$l = 0 \quad 0 = \alpha_0 - \sum_{j=1}^J \beta_j \{ 1 + a_1(\tau_j)[\tau_j - h/2] + a_0(\tau_j)[\tau_j^2 - \tau_j h]/2 \},$$

$$l = 1 \quad 0 = \alpha_1 - \sum_{j=1}^J \beta_j \{ -2 + a_1(\tau_j)[-2\tau_j] + a_0(\tau_j)[h^2 - \tau_j^2] \},$$

$$l = 2 \quad 0 = \alpha_2 - \sum_{j=1}^J \beta_j \{ 1 + a_1(\tau_j)[\tau_j + h/2] + a_0(\tau_j)[\tau_j^2 + \tau_j h]/2 \}$$

$$\text{normalization } 1 = \sum_{j=1}^J \beta_j,$$

$$l = 3 \quad 0 = \sum_{j=1}^J \beta_j \{ 6\tau_j + a_1(\tau_j)[3\tau_j^2 - h^2] + a_0(\tau_j)[\tau_j^3 - \tau_j h^2] \},$$

$$l = 4 \quad 0 = \sum_{j=1}^J \beta_j \{ 12\tau_j^2 - 2h^2 + a_1(\tau_j)[4\tau_j^3 - 2\tau_j h^2] + a_0(\tau_j)[\tau_j^4 - \tau_j^2 h^2] \},$$

$$l = 5 \quad 0 = \sum_{j=1}^J \beta_j \{ 20\tau_j^3 - 6\tau_j h^2 + a_1(\tau_j)[5\tau_j^4 - 3\tau_j h^2] + a_0(\tau_j)[\tau_j^5 - \tau_j^3 h^2] \},$$

and so on.

Note that the first three equations give the  $\alpha$ 's in terms of the  $\beta$ 's. Also note that in all the equations above, the terms which involve the coefficients  $a_1$  and  $a_0$  are order  $h$  and  $h^2$ , respectively, compared with the leading term in each of the curly brackets.

For specific examples, we consider  $M = D^2$  in which  $a_1 = a_0 = 0$ . One obtains immediately that  $\alpha_0 = \alpha_2 = 1, \alpha_1 = -2$ , so that the difference operator  $M_n$  is the usual divided difference approximation for the second-derivative operator

$$M_n U(t_k) = (U_k - 2U_{k+1} + U_{k+2})/h^2 = U[t_k, t_k + h, t_k + 2h].$$

However, the operator  $I_n$  changes when  $J$  or the locations of the auxiliary points change. Below  $O(h^p)$  denotes the truncation error with respect to the space of functions  $\Sigma = C^{(p+2)}$  of functions with continuous  $(p + 2)$ nd derivative. Below, we also abbreviate  $I_n f(t_k)$  with just  $I_n f$ .

*Example 2-1.* For  $J = 1$  and  $-h = t_k \leq \tau_1 \leq t_{k+2} = h, \tau_1 \neq 0$ , the equation for  $l = 3$  is not satisfied [ $a_1 = a_0 = 0$ ] and with  $I_n f = f(\tau_1)$  we obtain an  $O(h)$  scheme which is exact on  $P_2$ .

*Example 2-2.* For  $J = 1$  and  $\tau_1 = t_{k+1} = 0$ , the equation for  $l = 3$  is satisfied, but the one for  $l = 4$  is not satisfied, and with  $I_n f = f(\tau_1)$  we obtain an  $O(h^2)$  scheme which is exact on  $P_3$ .

*Example 2-3.* For  $J = 2$  and  $-\tau_1 = \tau_2 = h(1/6)^{1/2}$  we obtain an  $O(h^4)$  scheme which is exact on  $P_5$  with  $I_n f = [f(\tau_1) + f(\tau_2)]/2$ .

*Example 2-4.*  $J = 3$ , exact on  $P_5, O(h^4)$  Størmer-Numerov approximation  $I_n f = [f(-h) + 10f(0) + f(h)]/12$ .

*Example 2-5.*  $J = 3$ , exact on  $P_7, O(h^6)$  approximation of Osborne [1967]  $I_n f = [5f(\tau_1) + 14f(0) + 5f(\tau_3)]/24, -\tau_1 = \tau_3 = h(2/5)^{1/2}$ .

*Example 2-6.*  $J = 5$ , exact on  $P_{11}$ , a new  $O(h^{10})$  approximation  $I_n f = \beta_1 f(\tau_1) + \dots + \beta_5 f(\tau_5)$ ,

$$\beta_1 = \beta_5 = 0.0516582578, \quad \beta_2 = \beta_4 = 0.2394732407, \quad \beta_3 = 0.4177370031,$$

$$-\tau_1 = \tau_5 = 0.8214405997h, \quad -\tau_2 = \tau_4 = 0.4499203525h, \quad \tau_3 = 0.$$

*Example 2-7.*  $J = 3$ , exact on the space of cubic splines with joints at the equal spaced mesh points (see, for example, Birkhoff and de Boor [1965, p. 189]),  $O(h^2)$  approximation  $I_n f = [f(-h) + 4f(0) + f(h)]/6$ .

To use these schemes for the Dirichlet problem, one solves the system

$$(2-5a) \quad U_0 = u(A), \quad U_n(B) = u(B),$$

$$(2-5b) \quad (U_{k-1} - 2U_k + U_{k+1})/h^2 = g_k, \quad k = 1, \dots, n - 1,$$

with

$$h = (B - A)/n, \quad t_k = A + kh \quad \text{and} \quad g_k = I_n f(t_{k-1}) = \sum_{j=1}^J \beta_j f(A + kh + \tau_j),$$

where  $g$  differs from example to example.

For the initial value problem, the value of  $Du(A) = du(A)/dt$  is given. One can approximate this with the forward divided difference and by equating its value to a linear combination of values of  $f$ , one can obtain higher accuracy (see Section 7). Use of polynomial spaces  $S$  and simplification leads to the pair of initial values for the solution of the second-order difference equation (2-5b)

$$(2-5a') \quad U_0 = u(A), \quad U_1 = u(A) + hDu(A) + h^2 g^0/2,$$

and the following gives the value of  $g^0$  for accuracy comparable to that for the schemes given above:

*Example 2-1'.*  $O(h), g^0 = f(A)$ .

*Example 2-2'.*  $O(h^2), g^0 = f(A + h/3)$ .

Example 2-4'.  $O(h^4), g^0 = [9f(A) + 25f(A + 2h/5) + 2f(A + h)]/36.$

Example 2-5'.  $O(h^6), g^0 = \beta_1 f(A + \tau_1) + \beta_2 f(A + \tau_2) + \beta_3 f(A + \tau_3),$

$$\begin{aligned} \beta_1 &= 0.4018638275, & \beta_2 &= 0.4584822127, & \beta_3 &= 0.1396539598, \\ \tau_1 &= 0.0885879595h, & \tau_2 &= 0.4094668644h, & \tau_3 &= 0.7876594618h. \end{aligned}$$

Example 2-6'.  $O(h^{10}), g^0 = \beta_1 f(A + \tau_1) + \dots + \beta_5 f(A + \tau_5),$

$$\begin{aligned} \beta_1 &= 0.1935631805, & \beta_2 &= 0.3343492762, & \beta_3 &= 0.2927739742, \\ \beta_4 &= 0.1478177401, & \beta_5 &= 0.0314958290, \\ \tau_1 &= 0.0398098571h, & \tau_2 &= 0.1980134179h, & \tau_3 &= 0.4379748102h, \\ \tau_4 &= 0.6954642734h, & \tau_5 &= 0.9014649142h. \end{aligned}$$

**3. Truncation Error for Polynomial Approximation.** We only consider approximation away from boundaries and approximation which is exact on a polynomial space  $P_L$  for some  $L \geq m$ . Results for approximation of boundary conditions are obtained by an easy modification. Results for other spaces, such as those appropriate for approximation near singular points of differential equations, will be presented elsewhere.

We use  $\xi_{k,j}, j = 0, 1, \dots,$  to denote distinct points such that  $t_k \leq \xi_{k,j} \leq t_{k+m}$  and set  $\bar{\xi}_{k,j} = (\xi_{k,0}, \dots, \xi_{k,j})$ ; we also set

$$(3-1) \quad \Delta \bar{\xi}_{k,j} = \min_{i,q=0,\dots,j,i \neq q} |\xi_{k,i} - \xi_{k,q}|.$$

We use the polynomials

$$(3-2a) \quad w(\bar{\xi}_{k,j}; t) = \prod_{q=0}^j (t - \xi_{k,q}) / (j + 1)!, \quad j = 0, 1, \dots;$$

and, to simplify notation below, we set

$$(3-2b) \quad w(\bar{\xi}_{k,-1}; t) = 1.$$

We also use the Lagrange polynomial interpolation basis with respect to the points in  $\bar{\xi}_{k,j}$

$$(3-3) \quad l_r(\bar{\xi}_{k,j}; t) = w(\bar{\xi}_{k,j}; t) / [(t - \xi_{k,r}) w'(\bar{\xi}_{k,j}; \xi_{k,r})], \quad r = 0, \dots, j.$$

For all  $t$  between  $t_k$  and  $t_{k+m}$ , it follows from (3-2a) that  $|w|$  is bounded above by a constant times  $h_k^{j+1}$ , and  $|w'(\bar{\xi}_{k,j}; \xi_{k,r})|$  is bounded below by a constant times  $(\Delta \bar{\xi}_{k,j})^j$ . Then from (3-3) it follows that  $|l_r|$  is bounded by a constant times  $(h_k / \Delta \bar{\xi}_{k,j})^j$ . Similarly, there is a constant  $K$  which depends on  $j$  but does not depend on  $h_k$  or  $\xi_{k,j}$ , such that for all  $t, t_k \leq t \leq t_{k+m}$

$$(3-4) \quad \begin{aligned} |D^i w(\bar{\xi}_{k,j}; t)| &\leq K h_k^{j-i+1}, & i &= 0, \dots, j + 1, \\ |D^i l_r(\bar{\xi}_{k,j}; t)| &\leq K h_k^{j-i} / \Delta \bar{\xi}_{k,j}^j, & i, r &= 0, \dots, j. \end{aligned}$$

Because  $T_n u = M_n u - I_n [Mu]$  involves derivatives of  $u$  only up to order  $m \leq L$ , it follows (see, for example, Theorem 2.1 of de Boor and Lynch [1966]) that for

$$(3-5) \quad u \in \mathbf{F}^{L+1} [t_k, t_{k+m}] = \{v \mid D^L v \text{ is absolutely continuous,} \\ D^{L+1} v \text{ is square integrable on } t_k \leq t \leq t_{k+m}\}$$

we have

$$(3-6a) \quad T_n u(t_k) = \sum_{i=0}^L T_n [l_i(\bar{\xi}_{k,L}; t)]_k u(\xi_{k,i}) \\ + \int_{t_k}^{t_k+m} T_n(t) [q(\bar{\xi}_{k,L}; t, x)]_k D^{L+1} u(x) dx,$$

where

$$(3-6b) \quad q(\bar{\xi}_{k,L}; t, x) = \left[ (t-x)_+^L - \sum_{i=0}^L l_i(\bar{\xi}_{k,L}; t) (\xi_{k,i} - x)_+^L \right] / L!,$$

$$(3-6c) \quad (t-x)_+^L = \begin{cases} (t-x)^L & \text{for } t-x > 0, \\ 0 & \text{for } t-x \leq 0, \end{cases}$$

and where the subscript  $(t)$ , as in  $T_n(t)$ , denotes that the operator is applied to a function of  $t$ .

Suppose that the HODIE approximation is exact on the space  $\mathbf{P}_L$  of polynomials of degree at most  $L$  so that the coefficients  $\alpha, \beta$  satisfy (2-2) for a set of basis elements for  $\mathbf{P}_L$ . Then because  $l_i(\bar{\xi}_{k,L}; \cdot) \in \mathbf{P}_L$ , the sum in (3-6b) is equal to zero. The sum in the definition (3-6b) of  $q$  is that element of  $\mathbf{P}_L$  which interpolates to  $(t-x)_+^L$  at the points  $t = \xi_{k,j}, j = 0, \dots, L$ . By taking  $m+1$  of the points in  $\bar{\xi}_{k,L}$  to be the stencil points  $t_k$ , one has  $q = 0$  on the stencil points and hence  $M_n(t)q = 0$ ; (3-6a) then reduces to

$$(3-7a) \quad T_n u(t_k) = - \sum_{j=1}^J \beta_{k,j} \int_{t_k}^{t_k+m} M_{(t)} q(\bar{\xi}_{k,L}; t, x) |_{t=\tau_{k,j}} D^{L+1} u(x) dx,$$

where

$$(3-7b) \quad M_{(t)} q(\bar{\xi}_{k,L}; t, x) \\ = \sum_{i=0}^m a_i(t) \left\{ (t-x)_+^{L-i} / (L-i)! - \sum_{j=0}^L (\xi_{k,L} - x)_+^L D^i l_j(\bar{\xi}_{k,L}; t) / L! \right\}.$$

In (3-7a), points  $\tau_{k,j}, x$ , and those in  $\bar{\xi}_{k,L}$  are between  $t_k$  and  $t_{k+m}$ . Therefore, by (3-4) we can bound the quantity in curly brackets in (3-7b) by  $h_k^{L-m} (K_1 + K_2 [h_k / \Delta \bar{\xi}_{k,L}]^L)$ , where  $K_1, K_2$  are constants which do not depend on  $h_k$  or  $\bar{\xi}_{k,L}$ . Consequently, if  $D^{L+1} u$  and the coefficients  $a_i$  are continuous, then

$$|T_n u(t_k)| \leq K_3 (1 + [h_k / \Delta \bar{\xi}_{k,L}]^L) \left\{ \sum_{j=1}^J |\beta_{k,j}| \right\} \|D^{L+1} u\|_\infty h_k^{L-m+1},$$



where  $\|\cdot\|_\infty$  denotes the max-norm. The constant  $K_3$  depends on  $\max_i \|a_i\|_\infty$ , on  $L$ , but not on  $h_k$  or  $\bar{\xi}_{k,L}$ .

We have introduced the restriction that  $m + 1$  of the points in  $\bar{\xi}_{k,L}$  are the stencil points in  $\bar{t}_k$ ; the other  $L - m$  points in  $\bar{\xi}_{k,L}$  are arbitrary, and we can choose them to maximize  $\Delta \bar{\xi}_{k,L}$ . Clearly, this maximum depends only on the stencil points and  $L$ . For  $L \geq m$ , set

$$(3-8) \quad \begin{aligned} R_L(\bar{t}_k) &= h_k / \max \Delta \bar{\xi}_{k,L} \text{ where the maximization is over} \\ \text{all points } \xi_{k,L} \text{ such that } t_k &\leq \xi_{k,l} \leq t_{k+m}, l = 0, \dots, L, \\ \text{and } m + 1 \text{ of the points } \xi_{k,l} &\text{ are equal to } t_{k+j}, j = 0, \dots, m. \end{aligned}$$

Furthermore, set

$$(3-9) \quad H_n = \max_{j=0, \dots, n-m} (t_{j+m} - t_j) / m,$$

and we have the following

**THEOREM 3-1.** *Suppose the coefficients  $a_i$  of  $M$  are continuous. Let  $A = t_0 < t_1 < \dots < t_n = B, n \geq m$ , be a set of mesh points and  $\bar{t}_k, k = 0, \dots, n - m$ , sets of auxiliary points. Suppose that for  $k = 0, \dots, n - m$ , there are coefficients  $\alpha_{k,i}, \beta_{k,j}$  which satisfy (2-2) and (2-3b) for  $s_0, \dots, s_L, L \geq m$ , a basis for  $P_L$ . Then there is a constant  $K$  which depends only on  $B - A$ , the order  $m$  of  $M, L$ , and the coefficients  $a_i$  such that for any  $u$  with continuous  $(L + 1)$ st derivative*

$$|T_n u(t_k)| \leq K \left[ 1 + \max_{j=0, \dots, n-m} R_L(\bar{t}_j)^L \right] \|D^{L+1} u\|_\infty H_n^{L-m+1}, \quad k = 0, \dots, n - m.$$

Note that  $R_L(\bar{t}_k)$  in (3-8) is related to a localized mesh ratio. This is because the stencil points  $t_k$  are included in  $\bar{\xi}_{k,L}$  and thus  $R_L(\bar{t}_k)$  exceeds  $h_k / (t_{k+i} - t_{k+i-1})$ ; it also exceeds  $L/m$ , obtained for the case that the points in  $\bar{\xi}_{k,L}$  are equally spaced. Consequently,

$$(3-10) \quad R_L(\bar{t}_j) \geq \max \left\{ L, \max_{k=0, \dots, n-m} \left[ \max_{i=1, \dots, m} (t_{k+m} - t_k) / (t_{k+i} - t_{k+i-1}) \right] \right\} / m.$$

**4. Analysis of the Special Case  $M = D^m$ .** The main results about the special case  $M = D^m$  carry over to the general case of the variable coefficient operator  $M$  in (2-1b). In this section, we consider in detail the special case. To distinguish between the two cases, we use the superscript 0 for quantities which apply to the special case, in particular, we use  $\alpha^0, \beta^0, M_n^0$ , and  $I_n^0$  for the coefficients and the operators when  $M = D^m$ .

In (2-2) set  $M = D^m$ , replace  $\alpha, \beta$  with  $\alpha^0, \beta^0$ , and use the following basis for  $P_L$  (see (3-2) and (3-3)):

$$(4-1a) \quad s_i(t) = \begin{cases} I_i(\bar{t}_k; t), & i = 0, \dots, m, \\ w(\bar{\xi}_{k,i-1}; t), & i = m + 1, \dots, L, \end{cases}$$

where

$$(4-1b) \quad \bar{\xi}_{k,i-1} = (\xi_{k,0}, \dots, \xi_{k,i-1}), \text{ the points } \xi_{k,l} \text{ are distinct,}$$

$$t_k \leq \xi_{k,l} \leq t_{k+m}, \text{ and } \xi_{k,j} = t_{k+j}, j = 0, \dots, m.$$

In this section we use the normalization (2-3c), and it is a consequence of the analysis below that this is allowed, i.e.  $\sum_j \beta_j \neq 0$ . Note that  $D^m w(\bar{\xi}_{k,m-1}; t) = D^m(t - t_k) \cdots (t - t_{k+m-1})/m! = 1$  so that (2-3c) can be written as

$$\sum_{j=1}^J \beta_{k,j}^0 D^m w(\bar{\xi}_{k,m-1}; \tau_{k,j}) = 1.$$

Since the Lagrange basis element  $l_i(\bar{t}_k; \cdot)$  is in  $\mathbf{P}_m$ , its  $m$ th derivative is a constant.

The HODIE equations for the special case then become

$$(4-2a) \quad \alpha_{k,i}^0/h_k^m - [m!/w'(\bar{t}_k; t_{k+i})] \sum_{j=1}^J \beta_{k,j}^0 = 0, \quad i = 0, \dots, m,$$

$$(4-2b) \quad \sum_{j=1}^J \beta_{k,j}^0 D^m w(\bar{\xi}_{k,m+l-2}; \tau_{k,j}) = \delta_{l,1}, \quad l = 1, \dots, L - m + 1,$$

where  $\delta_{l,i}$  denotes the Kronecker delta function.

Since the sum of the  $\beta^0$ 's is unity, (4-2a) shows that the operator  $M_n^0$  is  $m!$  times the usual divided-difference approximation to  $M = D^m$

$$(4-3) \quad M_n^0 u(t_k) = \sum_{i=0}^m \alpha_{k,i}^0 u(t_{k+i})/h_k^m = m! \sum_{i=0}^m u(t_{k+i})/w'(\bar{t}_k; t_{k+i})$$

$$= m! u[t_k, t_{k+1}, \dots, t_{k+m}],$$

that is,  $M_n^0 u(t_k)$  is the  $m$ th derivative of the unique polynomial in  $\mathbf{P}_m$  which interpolates to the values  $u(t_{k+i})$  at  $t_{k+i}, i = 0, \dots, m$ .

By Taylor's Theorem, any  $u$  in  $\mathbf{F}^m$  can be represented as

$$u(t) = \sum_{i=0}^{m-1} D^i u(t_k)(t - t_k)^i/i! + \int_{t_k}^t (t - x)^{m-1} D^m u(x) dx/(m - 1)!.$$

Because the  $m$ th divided difference of an element of  $\mathbf{P}_{m-1}$  is zero, we have

$$(4-4) \quad M_n^0 u(t_k) = m! u[t_k, \dots, t_{k+m}] = \int_{t_k}^{t_k+m} B_m(\bar{t}_k; x) D^m u(x) dx,$$

where  $B_m(\bar{t}_k; x)$  is the  $m$ th divided difference  $g_m[t_k, \dots, t_{k+m}; x]$  with respect to  $t$  of

$$g_m(t; x) = (t - x)_+^{m-1} = \begin{cases} (t - x)^{m-1}/(m - 1)! & \text{if } t \geq x, \\ 0 & \text{if } t < x, \end{cases}$$

so that  $B_m(\bar{t}_k; \cdot)$  is the  $(m - 1)$ st degree polynomial  $B$ -spline with joints at the stencil points in  $\bar{t}_k$ . This  $B$ -spline satisfies (Curry and Schoenberg [1966])

$$(4.5a) \quad B_m(\bar{t}_k; x) = \begin{cases} >0, & t_k < x < t_{k+m}, \\ =0, & x \leq t_k \text{ or } t_{k+m} \leq x, \end{cases}$$

$$(4-5b) \quad \int_{t_k}^{t_k+m} B_m(\bar{t}_k; x) dx = 1.$$

Therefore, we have

$$\begin{aligned} T_n^0 u(t_k) &= M_n^0(t_k) - I_n^0 [D^m u]_k \\ &= \int_{t_k}^{t_k+m} B_m(\bar{t}_k; x) D^m u(x) dx - \sum_{j=1}^J \beta_{k,j}^0 D^m u(\tau_{k,j}) \\ &\equiv E_n^0 [D^m u]_k; \end{aligned}$$

and in this we have defined the operator  $E_n^0$ . Clearly,  $E_n^0 u(t_k)$  is the quadrature error in using  $I_n^0 u(t_k)$  as an approximation to the integral of  $B_m(\bar{t}_k; x)u(x)$ . This quadrature error is zero for any  $v$  in  $\mathbf{P}_{J-1}$  if and only if

$$(4-6) \quad \beta_{k,j}^0 = \int_{t_k}^{t_k+m} B_m(\bar{t}_k; x) l_{j-1}(\bar{\tau}_k; x) dx, \quad j = 1, \dots, J.$$

Since the sum over  $j$  of  $l_{j-1}(\bar{\tau}_k; x)$  is unity, it follows from (4-5b) and (4-6) that the sum of  $\beta^0$ 's is unity. But then with the  $\beta$ 's in (4-6), for any  $u$  in  $\mathbf{P}_{m+J-1}$ ,  $T_n^0 u = E_n^0 [D^m u] = 0$ .

Consequently, for any stencil points  $\bar{t}_k$  and any  $J$  auxiliary points  $\bar{\tau}_k$ , there is a unique HODIE scheme with normalization (2-3c) which is exact on  $\mathbf{P}_{m+J-1}$ . One obtains a family of HODIE schemes which are exact on  $\mathbf{P}_L$  for any  $L$  such that  $0 \leq L < m + J - 1$ . We now show that there exist special sets of auxiliary points which make the approximation exact on  $\mathbf{P}_L$  for  $L$  up to  $m + 2J - 1$ .

Since  $B_m(\bar{t}_k; \cdot)$  is positive on the range of integration, we can define the following inner product

$$(u, v) = \int_{t_k}^{t_k+m} B_m(\bar{t}_k; x) u(x) v(x) dx.$$

For fixed  $m, k$ , and  $B_m(\bar{t}_k; \cdot)$ , let  $b_0, b_1, \dots$  with  $b_i$  in  $\mathbf{P}_i$  denote the normalized orthogonal polynomials with respect to this inner product; we call these the *B-spline orthogonal polynomials*. Based on the well-known theory of orthogonal polynomials,  $b_i$  has  $i$  distinct real zeros in  $t_k < t < t_{k+m}$ , and, for fixed  $i$ , we call these the *B-spline Gauss points*.

When the  $J$  auxiliary points in  $\bar{\tau}_k$  are the *B-spline Gauss points* for  $b_J$ , then the unique HODIE approximation which is exact on  $\mathbf{P}_{J+m-1}$  is also exact on  $\mathbf{P}_{2J+m-1}$ . In this case, the  $\beta^0$ 's are the coefficients of the  $J$ -point Gauss quadrature formula with weight functions  $B_m(\bar{t}_k; \cdot)$  and each  $\beta_{k,j}^0, j = 1, \dots, J$ , is positive. Since  $(b_J, b_J)$  is positive and  $I_n^0 [b_J^2]_k = 0$ , this HODIE approximation is not exact on  $\mathbf{P}_{2J+m}$ .

The *B-spline Gauss points* and the quadrature coefficients have been tabulated by Phillips and Hanson [1974] for a number of degrees and for a normalized interval and equally spaced joints.

The preceding results are summarized below.

**THEOREM 4-1.** *Let  $M = D^m$ , and let the normalization for HODIE approximation be (2-3c). For any set of  $m + 1$  stencil and  $J > 0$  auxiliary points  $\bar{t}_k, \bar{\tau}_k$ , there is a HODIE approximation with coefficients  $\alpha_{k,j} = \alpha_{k,j}^0, \beta_{k,j} = \beta_{k,j}^0$ , which is exact on*

$\mathbf{P}_L$  for any  $L$  with  $0 \leq L - m \leq J - 1$ . The operator  $M_n = M_n^0$  is unique; it is  $m!$  times the divided difference operator with respect to the stencil points. There are sets of  $J$  auxiliary points for which a HODIE approximation is exact for  $L$  with  $J \leq L - m \leq 2J - 1$ . If  $L - m \geq J - 1$ , then the coefficients  $\beta^0$  of  $I_n^0$  are unique and are given by (4-6). The  $J$  auxiliary points which give exactness on  $\mathbf{P}_{2J+m-1}$  are the zeros of the  $J$ th degree  $B$ -spline orthogonal polynomial  $b_J$  associated with the  $B$ -spline  $B_m(\bar{t}_k; \cdot)$  with joints at the stencil points.

The examples for  $M = D^2$  in Section 2 illustrate various special cases of the results stated in Theorem 4-1. Examples 2-2, 2-3, 2-5, and 2-6 use  $J$   $B$ -spline Gauss points for  $J = 1, 2, 3$ , and  $5$ , respectively.

Examples 2-4 (Störmer-Numerov) and 2-7 (exact on cubic splines) both use the same set of three auxiliary points. Both are exact on  $\mathbf{P}_3$  and for this  $L - m = 3 - 2 = 1 < J - 1 = 2$ ; their different sets of  $\beta$ 's illustrate the nonuniqueness for  $L - m < J - 1$ . Since the Störmer-Numerov scheme is also exact on  $\mathbf{P}_4$  and since  $L - m = J - 1$  for this case, the scheme is the unique HODIE scheme with those three auxiliary points which is exact on  $\mathbf{P}_4$ . One of the auxiliary points,  $\tau_{k,2} = t_{k+1}$  is a  $B$ -spline Gauss point for all odd degree  $B$ -spline orthogonal polynomials associated with  $B_{2,k}$  with equally spaced joints; because of this (or, alternatively, symmetry), the scheme is exact on  $\mathbf{P}_5$ . Another set of three auxiliary points (Example 2-5) yields an approximation exact on  $\mathbf{P}_7$ .

We now derive bounds on the elements of the inverse of the coefficient matrix of the system in (4-2b) with  $L - m + 1 = J$ ; these are used in the next section. For fixed  $i = 1, \dots, J$  consider the systems

$$(4-7a) \quad \sum_{j=1}^J x_{j,i} D^m w(\bar{\xi}_{k,m+l-2}; \tau_{k,j}) = \delta_{l,i}, \quad l = 1, \dots, J.$$

For fixed  $r = 1, \dots, J$ , multiply the  $l$ th equation by the number  $\pi_{r-1,m+l-2}$  (determined below) and sum with respect to  $l$  to obtain

$$(4-7b) \quad \sum_{j=1}^J x_{j,i} D^m \sum_{l=1}^J \pi_{r-1,m+l-2} w(\bar{\xi}_{k,m+l-2}; \tau_{k,j}) = \pi_{r-1,m+i-2}.$$

Define the polynomial  $p_{r-1} \in \mathbf{P}_{J+m-1}$  by

$$\begin{aligned} p_{r-1}(t) &= \sum_{l=1}^J \pi_{r-1,m+l-2} w(\bar{\xi}_{k,m+l-2}; t) \\ &= \sum_{l=1}^J \pi_{r-1,m+l-2} (t - \xi_{k,0}) \cdots (t - \xi_{k,m+l-2}) / (m+l-1)!. \end{aligned}$$

Then the  $\pi$ 's can be expressed in terms of divided differences of  $p_{r-1}$ :

$$\pi_{r-1,m+l-2} = (m+l-1)! p_{r-1}[\bar{\xi}_{k,0}, \dots, \xi_{k,m+l-1}], \quad l = 1, \dots, J.$$

Choose these constants so that  $D^m p_{r-1}(t) = l_{r-1}(\bar{t}_k; t)$ , where  $l_{r-1}$  is the  $(r-1)$ st element of the Lagrange basis (3-3) associated with  $\bar{t}_k$ . The  $J$ -by- $J$  matrix  $\Pi$  with

elements  $\pi_{r-1,m+l-2}$  is the nonsingular matrix associated with change in basis for  $\mathbf{P}_J$  from  $D^m w(\xi_{k,m+j-1}; \cdot)$  to the Lagrange basis.

Since (4-7b) is the product of  $\Pi$  with the system (4-7a), the  $x_j$ 's also satisfy (4-7a). That is, the left side of (4-7b) can be written as

$$\sum_{j=1}^J x_j i l_{r-1}(\bar{\tau}_k, \tau_{k,j}).$$

Consequently, existence and uniqueness of polynomial interpolation shows that the solution of (4-7b) is given by  $x_{j,i} = \pi_{j-1,m+i-2}, j = 1, \dots, J$ . The points  $\xi_{k,l}$  are distinct, are between  $t_k$  and  $t_{k+m}$  and  $\xi_{k,l} = t_{k+l}, l = 0, \dots, m$ . Hence it follows from (4-4) that

$$\begin{aligned} x_{r,i} &= \int_{t_k}^{t_{k+m}} B_{m+i-1}(\bar{\xi}_{k,m+i-1}; x) D^{m+i-1} p_{r-1}(x) dx \\ &= \int_{t_k}^{t_{k+m}} B_{m+i-1}(\bar{\xi}_{k,m+i-1}; x) D^{i-1} l_{r-1}(\bar{\tau}_k; x) dx, \end{aligned}$$

where  $B_{m+i-1}(\bar{\xi}_{m+i-1}; \cdot)$  denotes the polynomial  $B$ -spline of degree  $m + i - 2$  with joints at  $\xi_{k,l}, l = 0, \dots, m + i - 1$ . For the case  $i = 1$ , this reduces to  $x_{j,1} = \beta_{k,j}$  with  $\beta_{k,j}$  given in (4-6). By (4-5) and (3-4), we have, therefore, the following result:

LEMMA 4-1. Let  $\bar{\xi}_k = (\xi_{k,0}, \dots, \xi_{k,m+J-2})$ , where the points  $\xi_{k,l}$  are distinct and between  $t_k$  and  $t_{k+m}$  and  $\xi_{k,l} = t_{k+l}, l = 0, \dots, m$ . Let  $B^0$  denote the matrix with elements

$$(B^0)_{l,i} = D^m w(\xi_{k,m+l-2}; \tau_{k,i}), \quad l, i = 1, \dots, J,$$

where  $w$  is as in (3-2a) and  $\bar{\tau}_k = (\tau_{k,1}, \dots, \tau_{k,J})$  is a set of auxiliary points between  $t_k$  and  $t_{k+m}$ . There exist constants  $K_{l,i}$  which are independent of  $h_k$  and  $\bar{\tau}_k$  (see (3-1)) such that

$$|([B^0]^{-1})_{l,i}| \leq K_{l,i} (h_k / \Delta \bar{\tau}_k)^{J-1} / h_k^{i-1}.$$

**5. Analysis of the Variable-Coefficient Case.** Let  $\lambda$  and  $\psi$  denote the functions obtained by applying  $M$  to the basis elements  $l$  and  $w$  in (4-1a)

$$(5-1a) \quad \lambda_i(t) = Ms_i(t) = Ml_i(\bar{t}_k; t); \quad i = 0, \dots, m,$$

$$(5-1b) \quad \psi_l(t) = Ms_{m+l}(t)/(h_k)^l = Mw(\bar{\xi}_{k,m+l-1}; t)/(h_k)^l, \quad l = 1, \dots, L - m,$$

and set

$$(5-1c) \quad \psi_0(t) \equiv 1.$$

We use  $\lambda_i^0, \psi_l^0$  to denote these functions in the special case  $M = D^m$ .

The HODIE equations are then

$$(5-2a) \quad \alpha_{k,i} / h_k^m - \sum_{j=1}^J \beta_{k,j} \lambda_i(\tau_{k,j}) = 0, \quad i = 0, \dots, m,$$

$$(5-2b) \quad \sum_{j=1}^J \beta_{k,j} \psi_{l-1}(\tau_{k,j}) = \delta_{l,1}, \quad l = 1, \dots, L - m + 1.$$

To see that  $\lambda, \psi$  differ from  $\lambda^0, \psi^0$  by  $O(h_k)$ , express the variables in terms of nondimensional parameters  $\gamma, \gamma_{k,j}, \rho_{k,j}$

$$t = t_k + \gamma h_k; \quad \xi_{k,i} = t_k + \gamma_{k,i} h_k, \quad i = 0, \dots, L; \quad \tau_{k,j} = t_k + \rho_{k,j} h_k, \\ j = 1, \dots, J,$$

$$\bar{\gamma}_{k,l} = (\gamma_{k,0}, \dots, \gamma_{k,l}), \quad \bar{\rho}_k = (\rho_{k,1}, \dots, \rho_{k,J});$$

and then  $\gamma_{k,j}$  and  $\rho_{k,j}$  are between 0 and  $m$ . From (3-2a) we have

$$w(\bar{\xi}_{k,l}; t) = h_k^{l+1} w(\bar{\gamma}_{k,l}; \gamma), \quad D^r w(\bar{\xi}_{k,l}; t) = h_k^{l+1-r} D^r w(\bar{\gamma}_{k,l}; \gamma).$$

Since

$$\lambda_i^0(\tau_{k,j}) = m! / w'(t_k; t_{k+i}) = m! / [h_k^m w'(\bar{\gamma}_{k,m}; \gamma_{k,i})],$$

for  $i = 0, \dots, m$  we have

$$(5-3a) \quad \lambda_i(\tau_{k,j}) = \lambda_i^0(\tau_{k,j}) \left[ 1 + \left\{ h_k a_{m-1}(\tau_{k,j}) \sum_{q=0, q \neq i}^m (\rho_{k,j} - \gamma_{k,q}) \right. \right. \\ \left. \left. + \dots + h_k^m a_0(\tau_{k,j}) \prod_{q=0, q \neq i}^m (\rho_{k,j} - \gamma_{k,q}) \right\} / m! \right].$$

For  $l = 1, \dots, L - m$  we have

$$(5-3b) \quad \psi_l(\tau_{k,j}) = \psi_l^0(\tau_{k,j}) + \sum_{p=0}^{m-1} a_p(\tau_{k,j}) D^p w(\bar{\xi}_{k,m+l-1}; \tau_{k,j}) \\ = D^m w(\bar{\gamma}_{k,m+l-1}; \rho_{k,j}) + \sum_{p=0}^{m-1} h_k^{m-p} a_p(\tau_{k,j}) D^p w(\bar{\gamma}_{k,m+l-1}; \rho_{k,j}).$$

The following establishes existence and uniqueness of HODIE schemes for  $L - m = J - 1$  and  $h_k$  sufficiently small:

**THEOREM 5-1.** *Let the normalization of the HODIE approximation be (2-3c). Suppose that the coefficients  $a_i$  of  $M$  are continuous. There is a positive  $H$  such that if the stencil points  $\bar{t}_k$  satisfy  $0 < h_k = (t_{k+m} - t_k) / m \leq H$ , then for any set of  $J$  auxiliary points  $\bar{\tau}_k$ , there is a unique HODIE approximation which is exact on  $\mathbf{P}_{J+m-1}$ . Its coefficients  $\alpha, \beta$  are the solution of (5-2) with  $L - m = J - 1$ .*

*Proof.* For details, see the end of this section; the main line of the argument is as follows: By hypothesis, the coefficients  $a_i$  are bounded; hence, so are the functions  $\lambda_i, \psi_i$  in (5-1). By (5-3) the values of these functions differ by  $O(h_k)$  at the auxiliary points from  $\lambda_i^0, \psi_i^0$  for the special case  $M = D^m$ . Because of the uniqueness of the coefficients  $\alpha_{k,i}^0, \beta_{k,j}^0$ , there is positive  $H$  so that the coefficient matrix of the system in (5-2) is a nonsingular matrix (which is essentially independent of  $h_k$ ) plus  $O(h_k)$  and, thus, is itself nonsingular for  $h_k < H$ .  $\square$

To show that HODIE approximations exist for  $L - m = 2J - 1$  with special auxiliary points, we need some preliminary results. After changing to nondimensional

parameters, the functions  $\psi_l$  in (5-1) have the same form as the functions in the next theorem. This theorem shows that the set of functions  $\psi_l, l = 0, \dots, L - m$ , is a Chebyshev set.

**THEOREM 5-2.** *Let  $K$  and  $m$  denote positive integers. Let  $\gamma_k, k = 0, \dots, K + m - 1$ , denote distinct points in the unit interval. Let the functions  $\Phi_l$  have the form*

$$\Phi_0(h; \gamma) \equiv 1, \quad \Phi_l(h; \gamma) = D^m \prod_{k=0}^{m+l} (\gamma - \gamma_k) + \phi_l(h; \gamma), \quad l = 1, \dots, K - 1,$$

where  $\phi_l$  is continuous and  $O(h)$  on  $0 \leq \gamma \leq 1$ . Let  $\bar{\rho} = (\rho_1, \dots, \rho_K)$  have distinct components such that  $0 \leq \rho_k \leq 1$ . There is a positive  $H$  such that for any  $h, 0 \leq h \leq H$ ,

$$\sum_{l=0}^{K-1} c_l \Phi_l(h; \rho_k) = 0, \quad k = 1, \dots, K, \text{ implies } c_l = 0, l = 0, \dots, K - 1.$$

The result in Theorem 5-2 is easy to prove for fixed  $\bar{\rho}$ . But, in addition, we must show that  $H$  can be chosen independent of  $\bar{\rho}$ .

*Proof.* First, let  $\bar{\rho}$  be fixed and consider the  $K$ -by- $K$  matrix  $V(h)$  with elements  $V(h)_{k,l} = \Phi_{l+1}(h; \rho_k)$ . The product of  $V(0)$  and a diagonal matrix is equal to a Vandermonde matrix; therefore,  $V(0)$  is nonsingular. By continuity of the elements of  $V(h)$ , there is an  $H(\bar{\rho})$  such that  $V(h)$  is nonsingular for all  $h, 0 \leq h \leq H(\bar{\rho})$ .

Second, suppose that there is no positive  $H_0$  independent of  $\bar{\rho}$  such that  $V(h)$  is nonsingular for all  $h, 0 \leq h \leq H_0$ . Then there are sequences with index  $i = 1, 2, \dots, \rightarrow \infty$ ,

$$H_i \downarrow 0, \quad c_l(H_i), \quad l = 0, \dots, K - 1, \text{ with } \max_l |c_l(H_i)| = 1,$$

$$\bar{\rho}(H_i) = (\rho_1(H_i), \dots, \rho_K(H_i)), \quad P_i(\rho) = \sum_{l=0}^{K-1} c_l(H_i) \Phi_l(H_i; \rho),$$

where  $P_i$  has zeros at  $\rho = \rho_j(H_i), j = 1, \dots, K$ . There exist, therefore, convergent subsequences (whose elements we also denote as above) such that

$$c_l(H_i) \rightarrow c_l^*, \quad \rho_j(H_i) \rightarrow \rho_j^*, \quad \text{and } P_i \rightarrow P^*.$$

By continuity and the form of the functions  $\Phi_l$ , the limiting function  $P^*$  is a polynomial of degree at most  $K - 1$ . Again by continuity,  $P^*(\rho_j^*) = 0, j = 1, \dots, K$ . Since  $\max_l |c_l^*| = 1, P^*$  is not identically equal to zero; consequently, the  $K$  points  $\rho_j^*$  are not distinct. Suppose that there are exactly  $N > 1$  zeros of  $P^*$  which are equal to  $\rho_k^*$  and  $\rho_k^* = \rho_{k+1}^* = \dots = \rho_{k+N-1}^*$ . Then we can write

$$P_i(\rho) = p(H_i; \rho) \prod_{q=0}^{N-1} [\rho - \rho_{k+q}(H_i)] \rightarrow p(0; \rho) [\rho - \rho_k^*]^N,$$

which shows that the  $(K - 1)$ st degree polynomial  $P^*$  has  $K$  zeros counting multiplicities. This contradiction establishes the theorem.  $\square$

The application of Theorem 5-2 to HODIE approximation follows from representations of moments of a Chebyshev set. Such moments are discussed in detail in Chapter 2 of Karlin and Studden [1966]; see, especially, pages 38-46. We summarize the pertinent information in the next paragraph.

Let  $\mu$  denote any nondecreasing right continuous function of bounded variation on  $t_k \leqq t \leqq t_{k+m}$ . Let  $\psi_l, l = 0, \dots, L$ , denote functions of a Chebyshev set on this interval. The  $l$ th moment  $q_l$  of the set with respect to the measure  $d\mu$  is

$$q_l = \int_{t_k}^{t_{k+m}} \psi_l(x) d\mu(x), \quad l = 0, \dots, L.$$

For each measure, one gets a set of moments  $\bar{q} = (q_0, \dots, q_L)$  and the set of all such  $\bar{q}$  is a subset  $Q$  of Euclidean  $(L + 1)$ -space which is called the *moment space* of the Chebyshev set. This moment space is the smallest cone with vertex at the origin which contains the curve  $\Psi(t) = (\psi_0(t), \dots, \psi_L(t)), t_k \leqq t \leqq t_{k+m}$ ; this curve is not in Euclidean  $L$ -space. If  $L = 2J - 1, J \geqq 1$ , and  $\bar{q} \in Q$  is an interior point of  $Q$ , then there is a unique principal representation of  $\bar{q}$  which involves  $J$  points  $\tau_{k,1} < \dots < \tau_{k,J}$  in the open interval  $t_k < t < t_{k+m}$ ; that is, there are positive values  $\beta_{k,j}$  such that

$$q_l = \sum_{j=1}^J \beta_{k,j} \psi_l(\tau_{k,j}), \quad l = 0, \dots, L = 2J - 1.$$

Clearly, the principal representation gives an abstract setting for Gauss quadrature.

Let  $Q_0$  denote the moment space for the Chebyshev set  $1, D^m s_{m+l}, l = 1, \dots, L$ , where  $s_{m+l}$  are the basis elements in (4-1). The results in Section 4 for the case  $M = D^m$  show that with  $d\mu(x) = B_m(\bar{t}_k; x) dx$ , then  $\bar{q}_0 = (q_{0,0}, q_{0,1}, \dots, q_{0,L}), q_{0,l-1} = \delta_{l,1}$ , is in  $Q_0$ . Thus, the principal representation is given with  $\tau_{k,j}$ , the zeros of the  $J$ th degree  $B$ -spline orthogonal polynomial and with  $\beta_{k,j}$  equal to  $\beta_{k,j}^0$  in (4-6). By uniqueness,  $\bar{q}_0$  is an interior point of the moment space  $Q_0$  and so there is a closed sphere  $S_0$  with center  $\bar{q}_0$  in the interior of  $Q_0$ .

It follows from Theorem 5-2 that if the coefficients  $a_i$  of  $M$  are continuous, then the functions  $\psi_l$  in (5-1b) form a Chebyshev set for all  $h_k$  sufficiently small. Let  $Q$  denote the moment space for this Chebyshev set. The curve  $\Psi(t) = (\psi_0, \dots, \psi_L)$  converges uniformly to the curve  $\Psi_0(t) = (1, D^m s_{m+1}(t), \dots, D^m s_{m+L}(t))$  on  $t_k \leqq t \leqq t_{k+m}$ ; hence, for sufficiently small  $h_k$ , the sphere  $S_0$  is in the interior of the moment space  $Q$ . This establishes the next theorem.

**THEOREM 5-3.** *Suppose the coefficients  $a_i$  of  $M$  are continuous. For a HODIE approximation with  $J$  auxiliary points, there is a positive  $H$  such that for any set of  $m + 1$  stencil points  $\bar{t}_k$  with  $0 < h_k = (t_{k+m} - t_k)/m \leqq H$ , there is a set of  $\beta_{k,j}$ 's and a unique set of  $J$  auxiliary points  $\bar{\tau}_k$  with  $t_k < \tau_{k,1} < \dots < \tau_{k,J} < t_{k+m}$  such that the HODIE approximation is exact on  $P_{2J+m-1}$ . The  $\beta_{k,j}$ 's are nonzero, all have the same sign and are unique for a given normalization.*

We call the special set of  $J$  auxiliary points which makes the HODIE approximation exact on  $P_{2J+m-1}$  the *generalized B-spline Gauss points*.



We now obtain a specific uniform bound on the  $\beta$ 's for the variable coefficient case.

The system (5-2b) with  $L - m + 1 = J$  can be written in matrix form as

$$(B^0 + B^1)\bar{\beta} = \bar{e}_1, \quad \bar{e}_1^t = (1, 0, \dots, 0),$$

where  $B^0$  is the unperturbed matrix in Lemma 4-1. Note that its elements are, essentially, independent of  $h_k$ . With  $\bar{\beta} = \bar{\beta}^0 + \delta\bar{\beta}$ , where  $\bar{\beta}^0$  has components  $\beta_{k,j}^0$  from the special case  $M = D^m$ , we have

$$(I + [B^0]^{-1}B^1)\delta\bar{\beta} = - [B^0]^{-1}B^1\bar{\beta}^0.$$

From (5-3b) it follows that elements of  $B^1$  are given by

$$(B^1)_{1,j} = 0, \quad (B^1)_{i,j} = \sum_{p=0}^{m-1} h_k^{m-p} a_p(\tau_{k,j}) D^p w(\bar{\gamma}_{k,m+l-1}; \rho_{k,j}), \quad i = 2, \dots, J.$$

Thus, there are constants  $k_{i,j}$  which depend only on max-norm bounds on  $a_0, a_1, \dots, a_{m-1}$  but not on  $h_k$  such that for all sufficiently small  $h_k$ ,

$$|(B^1)_{i,j}| \leq h_k k_{i,j}, \quad i, j = 1, \dots, J.$$

Hence, from Lemma 4-1 we obtain bounds on the elements of the product  $[B^0]^{-1}B^1$

$$|([B^0]^{-1}B^1)_{r,j}| \leq (h_k/\Delta\bar{\tau}_k)^{J-1} h_k \sum_i K_{r,i} k_{i,j}.$$

Consequently, for all sufficiently small  $h_k$ ,  $I + [B^0]^{-1}B^1$  is invertible; and we have the bound

$$\|\delta\bar{\beta}\|_\infty \leq \| [B^0]^{-1}B^1 \|_\infty \|\bar{\beta}^0\|_\infty / (1 - \| [B^0]^{-1}B^1 \|_\infty),$$

where the norms are the vector max-norm and the matrix row-sum-norm. For all sufficiently small  $h_k$  there is, therefore, a constant  $K_0$  such that

$$(5-4) \quad \|\delta\bar{\beta}\|_\infty \leq h_k (h_k/\Delta\bar{\tau}_k)^{J-1} K_0 \max_j |\beta_{k,j}^0|.$$

Lemma 4-1 with  $i = 1$  gives a bound on  $\beta_{k,j}^0$  which yields

$$\max_j |\beta_{k,j}^0| \leq (h_k/\Delta\bar{\tau}_k)^{J-1} \max_r (K_{r,1}) [1 + h_k (h_k/\Delta\bar{\tau}_k)^{J-1} K_0].$$

This gives the following result:

LEMMA 5-1. *Under the same hypotheses as Theorem 5-3, there is a constant  $K$  which is independent of  $h_k/\Delta\bar{\tau}_k$  such that for all sufficiently small  $h_k$ ,*

$$\max_j |\beta_{k,j}| \leq K (h_k/\Delta\bar{\tau}_k)^{J-1}.$$

Equations (5-2a), (5-3a), and (5-4) yield the following result:

COROLLARY 5-1. *Under the same hypotheses as Theorem 5-3,*

$$\alpha_{k,j} = \alpha_{k,j}^0 + O(h_k [h_k/\Delta\bar{\tau}_k]^{J-1}), \quad \beta_{k,j} = \beta_{k,j}^0 + O(h_k [h_k/\Delta\bar{\tau}_k]^{J-1}).$$

For  $R = h_k/\Delta\bar{\tau}_k$  bounded,  $|\alpha_{k,m}|$  converges to a positive value as  $h_k$  goes to zero.

**6. Stability of the HODIE Difference Operator.** In this section we show that the initial value problem for the HODIE difference equation is stable.

We first obtain bounds on the coefficients of  $M_n V$  when this is expressed in terms of divided differences of  $V$ . Let  $V$  denote a function defined at mesh points, and let  $M_n$  denote a difference operator obtained from a HODIE approximation exact on  $\mathbf{P}_L$  with  $L \geq m$ . For fixed  $k$ , let  $p$  denote that unique element in  $\mathbf{P}_m$  which interpolates to  $V(t_j), j = k, \dots, k + m$ . Writing  $p$  in the Newton form of the interpolation polynomial, we have

$$p(t) = V[t_k]s_{k,0}(t) + V[t_k, t_{k+1}]s_{k,1}(t) + \dots + V[t_k, \dots, t_{k+m}]s_{k,m}(t),$$

where

$$s_{k,0}(t) = 1, \quad s_{k,l+1}(t) = s_{k,l}(t)(t - t_{k+l}), \quad l = 0, \dots, m - 1.$$

Hence

$$(6-1) \quad M_n V(t_k) = M_n p(t_k) = V[t_k]C_{k,0} + \dots + V[t_k, \dots, t_{k+m}]C_{k,m},$$

where the  $C$ 's are independent of  $V$  and are given, for  $l = 0, 1, \dots, m$ , by

$$C_{k,l} = M_n s_{k,l}(t_k) = I_n [Ms_{k,l}](t_k) = \sum_{j=1}^J \beta_{k,j} \sum_{i=0}^{m-1} a_i(\tau_{k,j}) D^i s_{k,l}(\tau_{k,j});$$

the last equality holds because the HODIE approximation is exact on  $\mathbf{P}_L, L \geq m$ .

Using the normalization (2-3)(c),  $D^m s_{m,m}(t) = m!$ , and  $a_m(t) \equiv 1$ , we have

$$C_{k,m} = m! + \sum_{j=1}^J \beta_{k,j} \sum_{i=0}^{m-1} a_i(\tau_{k,j}) D^i s_{k,m}(\tau_{k,j}).$$

Set  $H_n = \max_k \{h_k\}$ . Because the auxiliary points are between  $t_k$  and  $t_{k+m}$ , there are constants  $K_l$  which depend on  $\max_i \|a_i\|_\infty$ , but not on the mesh points nor on the auxiliary points nor on  $H_n$  such that for  $H_n < 1$

$$|C_{k,l}| \leq \max_j |\beta_{k,j}| \sum_{i=0}^{m-1} \|a_i\|_\infty |D^i s_{k,l}(\tau_{k,j})| \leq K_l \max_j |\beta_{k,j}|,$$

$$|C_{k,m} - m!| \leq H_n K_m \max_j |\beta_{k,j}|.$$

By Lemma 5.1,  $\max_j |\beta_{k,j}| \leq KR^{J-1}, R = h_k/\Delta\tau_k$ . Consequently, if the ratio  $R$  is uniformly bounded as  $H_n \downarrow 0$ , then the coefficients  $C_{k,l}$  are uniformly bounded (as  $H_n \downarrow 0$ ) independent of  $k$  and  $H_n$ . Furthermore, for all sufficiently small  $H_n, C_{k,m} \approx m!$ , so that  $C_{k,m}$  is uniformly bounded above zero

$$(6-2) \quad C_{k,m} \geq \delta > 0.$$

Next we show that, for sufficiently small mesh spacing, there is a unique solution of the initial value problem

$$(6-3a) \quad M_n V(t_k) = F(t_k), \quad k = 0, 1, \dots,$$

$$(6-3b) \quad V[t_0], V[t_0, t_1], \dots, V[t_0, \dots, t_{m-1}] \text{ are given.}$$

Because of (6-2), we can divide the difference equation (6-3a) by  $C_{k,m}$ . By using (6-1),  $V[t_k, \dots, t_{k+m}]$  can be expressed explicitly in terms of  $V[t_k], \dots, V[t_k, \dots, t_{k+m-1}]$ ; and so, there is a unique solution of (6-3).

We now obtain a bound on the solution of (6-3) for the homogeneous case  $F \equiv 0$ . Using the definition of the  $m$ th divided difference and the difference equation, we obtain an expression for the  $(m - 1)$ st divided difference at  $t_{k+1}$

$$(6.4a) \quad \begin{aligned} V[t_{k+1}, \dots, t_{k+m}] &= (1 - mh_k C_{k,m-1}/C_{k,m})V[t_k, \dots, t_{k+m-1}] \\ &\quad - (mh_k C_{k,m-2}/C_{k,m})V[t_k, \dots, t_{k+m-2}] \\ &\quad - \dots - (mh_k C_{k,0}/C_{k,m})V[t_k]. \end{aligned}$$

We also have

$$(6.4b) \quad \begin{aligned} V[t_{k+1}, \dots, t_{k+i}] &= V[t_k, \dots, t_{k+i-1}] \\ &\quad + (t_{k+i} - t_k)V[t_k, \dots, t_{k+i}], \quad i = 1, \dots, k + m - 1. \end{aligned}$$

Let  $\|V(t_k)\|_{m-1}$  be defined by

$$(6.5) \quad \begin{aligned} \|V(t_k)\|_{m-1} &= |V[t_k]| + |V[t_k, t_{k+1}]| \\ &\quad + \dots + |V[t_k, \dots, t_{k+m-1}]|. \end{aligned}$$

From (6-4) we obtain

$$\|V(t_{k+1})\|_{m-1} \leq (1 + H_n K) \|V(t_k)\|_{m-1},$$

where

$$K = \max_{k,i} \{ 1 + mC_{k,m-i}/C_{k,m} + (t_{k+i} - t_k)/H_n \};$$

and, with the assumption introduced above,  $K$  can be taken independent of  $H_n$  for all sufficiently small  $H_n$ . Consequently, we have (for  $H_n K < 1$  and for  $M_n V = 0$ )

$$\begin{aligned} \|V(t_k)\|_{m-1} &\leq (1 + H_n K)^k \|V(t_0)\|_{m-1} \\ &\leq \|V(t_0)\|_{m-1} \exp(H_n K k), \quad k = 0, \dots, n - m. \end{aligned}$$

From this inequality we can obtain a bound on the Green's function for the initial value difference equation problem.

For each  $l = 1, \dots, n - m$ , let  $G_l$  denote the solution of

$$(6.6a) \quad G_l(t_k) = 0, \quad k = 0, 1, \dots, l + m - 2,$$

$$(6.6b) \quad M_n G_l(t_{l-1}) = 1,$$

$$(6.6c) \quad M_n G_l(t_k) = 0, \quad k = l, l + 1, \dots, n - m.$$

Because  $G_l[t_{l-1}] = G_l[t_{l-1}, t_l] = \dots = G_l[t_{l-1}, \dots, t_{l+m-2}] = 0$ , we have

$$M_n G_l(t_{l-1}) = C_{l-1,m} G_l[t_l, \dots, t_{l+m-1}] / mh_{l-1} = 1.$$

Hence, it follows from (6-6c) that for  $k = l, l + 1, \dots, n - m$ ,

$$\begin{aligned} \|G_l(t_k)\|_{m-1} &\leq \|G_l(t_l)\|_{m-1} \exp(H_n K [k - l]) \\ &= |G_l[t_l, \dots, t_{l+m-1}]| \exp(H_n K [k - l]) \\ &= mh_{l-1} \exp(H_n K [k - l]) / C_{l,m} \\ &\leq H_n \exp(H_n K [k - l]) / \delta, \end{aligned}$$

where  $\delta$  is as in (6-2).

The solution  $V$  of the initial value problem

$$\begin{aligned} M_n V(t_k) &= F(t_k), \quad k = 0, 1, \dots, n - m, \\ V[t_0] &= V[t_0, t_1] = \dots = V[t_0, \dots, t_{m-1}] = 0 \end{aligned}$$

is, therefore,

$$V(t_k) = \sum_{l=0}^{k-m} G_l(t_k) F(t_l), \quad k = m, \dots, n,$$

and is bounded by

$$\|V(t_k)\|_{m-1} \leq K_0 \max_l |F(t_l)| \exp(H_n K k),$$

where  $K_0 \leq (H_n + 1/K) / \delta$ . Thus, we have the following:

**THEOREM 6-1.** *Suppose the coefficients of  $M$  in (2.1b) are continuous. For the partition and set of auxiliary points*

$$A = t_0 < \dots < t_n = B, \quad t_k \leq \tau_{k,1} < \dots < \tau_{k,J} \leq t_{k+m},$$

let  $h_k = (t_{k+m} - t_k) / m$  and  $H_n = \max_k h_k$ . Let  $H_n$  be sufficiently small so that there are HODIE coefficients  $\alpha, \beta$  such that

$$(1/h_k)^m \sum_{i=0}^m \alpha_{k,i} s(t_{k+i}) = \sum_{j=1}^J \beta_{k,j} M s_l(\tau_{k,j}), \quad l = 0, \dots, L \geq J + m - 1,$$

where  $s_0, \dots, s_L$  is a basis for  $\mathbf{P}_L$ . If  $H_n$  is sufficiently small, then for given values  $F(t_k), k = 0, 1, \dots, n - m$ , and  $c_q, q = 0, \dots, m - 1$ , the unique solution of

$$(1/h_k)^m \sum_{i=0}^m \alpha_{k,i} V(t_{k+i}) = F(t_k), \quad k = 0, \dots, n - m,$$

$$V[t_0, \dots, t_q] = c_q, \quad q = 0, \dots, m - 1,$$

is bounded by

$$\|V(t_k)\|_{m-1} \leq \left( \|V(t_0)\|_{m-1} + K_0 \max_j |F(t_j)| \right) \exp(H_n K k),$$

where the constants  $K_0$  and  $K$  depend only on

$$H_n / \min_k [t_k - t_{k-1}] \quad \text{and} \quad H_n / \min_k \left[ \min_j [\tau_{k,j} - \tau_{k,j-1}] \right],$$

and bounds on the coefficients  $a_i$  of  $M$ .

**7. Discretization Error for the Initial Value Problem.** In this section we show that the first  $m - 1$  divided differences of the HODIE approximation  $U$ , subject to appropriate initial conditions, converge as  $O(h_n^{L-m+1})$  to the first  $m - 1$  divided differences of the solution  $u$  of the differential equation

$$(7-1a) \quad Mu = f, \quad A < t \leq B,$$

subject to the initial conditions at  $t = A = t_0$

$$(7-1b) \quad D^q u(A) = D^q u(t_0) = c_q, \quad q = 0, \dots, m - 1.$$

The initial conditions for the HODIE difference equation

$$(7-2) \quad M_n U(t_k) = I_n f(t_k), \quad k = 0, 1, \dots,$$

are taken as

$$(7-3a) \quad U_0 = u(t_0).$$

$$(7-3b) \quad U_p = \sum_{q=0}^{m-1} \gamma_{p,q} D^q u(t_0) (t_p - t_0)^q / q! + h^m \sum_{j=1}^J \beta_{p,j} f(\tau_j), \quad p = 1, \dots, m - 1,$$

where  $h = (t_{m-1} - t_0) / (m - 1)$  and the auxiliary points satisfy

$$(7-3c) \quad A = t_0 \leq \tau_1 < \tau_2 < \dots < \tau_J \leq t_{m-1}.$$

The coefficients  $\gamma_{p,q}, \beta_{p,j}$  are chosen to make the boundary conditions (7-3b) exact on  $P_L$  and we now use the basis

$$(7-4) \quad s_l(t) = (t - t_0)^l / h^l, \quad l = 0, \dots, L.$$

The  $\beta$ 's are determined *first* by using the basis elements (7-4) with  $l = m, m + 1, \dots, L$ . The coefficients of the  $\gamma$ 's in (7-3b) are then *zero*, and this gives the  $m - 1$  systems for  $p = 1, \dots, m - 1$

$$(7-5) \quad \frac{(t_p - t_0)^l}{h^l} = \sum_{j=1}^J \beta_{p,j} \left[ \frac{l!}{(l - m)!} \frac{(\tau_j - t_0)^{l-m}}{h^{l-m}} + h \frac{l!}{(l - m + 1)!} a_{m-1}(\tau_j) \frac{(\tau_j - t_0)^{l-m+1}}{h^{l-m+1}} + \dots + h^m a_0(\tau_j) \frac{(\tau_j - t_0)^l}{h^l} \right], \quad l = m, \dots, L.$$

For  $J = L - m + 1$ , this system consists of  $J$  equations in  $J$  unknowns. There is no normalization equation for the  $\beta$ 's because of the inherent normalization imposed by the left side of (7-5).

For the special case that  $M = D^m$ , the functions  $a_{m-1}, \dots, a_0$ , are identically equal to zero, and the coefficient matrix of the  $\beta$ 's in (7-5) has elements

$$(7-5a) \quad l! (\tau_j - t_0)^{l-m} / [(l - m)! h^{l-m}].$$

This matrix is equal to the product of a Vandermonde matrix and a nonsingular

diagonal matrix. Hence, in this case ( $M = D^m$  and  $J = L - m + 1$ ) there are unique sets  $\beta_{p,1}, \dots, \beta_{p,J}, p = 1, \dots, m - 1$ , which solve the systems. Moreover, the quantities

$$(t_p - t_0)^l/h^l \quad \text{and} \quad (\tau_j - t_0)^{l-m}/h^{l-m}$$

are order 1, hence so are the  $\beta$ 's; in fact, for fixed spacings

$$t_p = t_0 + \delta_p h, \quad \tau_j = t_0 + \epsilon_j h, \quad \text{with } \delta_p, \epsilon_j \text{ constants,}$$

the  $\beta$ 's are independent of  $h$ . For the case that the relative spacings change as  $t_{m-1} \rightarrow t_0$ , one obtains bounded  $\beta$ 's provided there is a positive constant  $R_{\max}$  such that

$$(7-6a) \quad \max_{p=0, \dots, m-2} \{(t_{m-1} - t_0)/(t_{p+1} - t_p)\} \leq R_{\max}$$

and

$$(7-6b) \quad \max_{j=1, \dots, J-1} \{(t_{m-1} - t_0)/(\tau_{j+1} - \tau_j)\} \leq R_{\max}.$$

We call this the *relative spacing condition* and throughout we assume it holds with a fixed value of  $R_{\max}$ .

If the  $a$ 's in (7-5) are bounded, then the facts above remain true for a general operator  $M \neq D^m$ , provided  $h$  is sufficiently small, because the terms in (7-5) involving  $a_{m-1}, \dots, a_0$  merely perturb the coefficients (7-5a) of the  $\beta$ 's by at most  $O(h)$ .

After computing the  $\beta$ 's, one computes the  $\gamma$ 's in (7-3b) by using basis elements (7-4) with  $l = 0, 1, \dots, m - 1$ . For each  $p = 1, \dots, m - 1$ , one solves

$$(t_p - t_0)^l/h^l = \gamma_{p,l}(t_p - t_0)^l/h^l + \sum_{j=1}^J \beta_{p,j} [h^{m-l}! a_l(\tau_j) + \dots + h^m a_0(\tau_j)(\tau_j - t_0)^l/h^l],$$

$$l = 0, 1, \dots, m - 1.$$

For the special case that  $M = D^m$ , all the  $a$ 's are zero; hence, all the  $\gamma$ 's are equal to unity. For general  $M \neq D^m$ , the coefficients of the  $\beta$ 's might be as small as  $O(h^m)$  or as large as  $O(h^{m-l})$ , so if the relative spacing conditions hold as  $h \downarrow 0$ , then

$$\gamma_{p,l} = 1 + O(h^{m-l}), \quad l = 0, 1, \dots, m - 1, p = 1, \dots, m - 1.$$

Consequently, for any  $u \in C^{L+1}$  we have by construction

$$(7-7a) \quad u(t_p) = \sum_{q=0}^{m-1} \gamma_{p,q} D^q(t_0)(t_p - t_0)^q/q! + h^m \sum_{j=1}^J \beta_{p,j} M[u(\tau_j)] + T_{n,p}[u],$$

where the truncation error of the boundary conditions approximation is

$$(7-7b) \quad T_{n,p}[u] = O(h^{L+1})$$

provided the relative spacing condition holds. (We omit a formal proof of (7-7b); one can be obtained easily by modifying the discussion in Section 3.) When one writes a specific bound, such as

$$|T_{n,p}[u]| \leq K_{n,p} h^{L+1}, \quad h \text{ sufficiently small,}$$

the constant  $K_{n,p}$  depends on  $u$ , the coefficients  $a_0, \dots, a_{m-1}$ , and  $R_{\max}$  of the relative spacing condition.

With the relative spacing condition (7-6) and the conditions in Theorem 6-1, then for sufficiently small  $H_n$ , as in Theorem 6-1, there is one and only one solution of (7-2)–(7-3).

To prove convergence, we consider a solution  $v$  of the differential equation  $Mv = f$  which takes on the same initial values as  $U$ . Thus,  $v$  satisfies

$$(7-8) \quad Mv = f, \quad v(t_p) = U_p, \quad p = 0, 1, \dots, m - 1.$$

The function  $v$  depends on  $h$  and the locations of the mesh points  $t_0, \dots, t_{m-1}$ .

Since  $u$  and  $v$  satisfy the same differential equation, we can write  $v$  as  $u$  plus an element of the null space of  $M$ . We choose the basis  $u^{(l)}$ ,  $l = 0, \dots, m - 1$ , of the null space which satisfies

$$(7-9) \quad Mu^{(l)} = 0, \quad D^p u^{(l)}(t_0) = \delta_{p,l}, \quad l, p = 0, \dots, m - 1.$$

where  $\delta_{p,l}$  is the Kronecker delta function. Thus, for some coefficients  $b_k$ , we can write (since  $u(t_0) = v(t_0) = U_0$ )

$$(7-10) \quad v = u + b_1 u^{(1)} + \dots + b_{m-1} u^{(m-1)}.$$

The  $b$ 's depend on  $h$ ; we obtain bounds on them below.

We assume that  $u \in C^{L+1}$  and  $u^{(l)} \in C^{L+1}$ ,  $l = 0, \dots, m - 1$ ; these assumptions hold, for example, if the coefficients of the differential operator and the right side of the differential equation are in  $C^{L-m+1}$ . It then follows for  $p = 0, \dots, m - 1$ , that  $u(t_p) - v(t_p) = O(h^{L+1})$ ; this is because by (7-7), (7-3), and (7-8), we have

$$(7-11) \quad \begin{aligned} u(t_p) &= \sum_{q=0}^{m-1} \gamma_{p,q} D^q u(t_0) (t_p - t_0)^q / q! + h^m \sum_{j=1}^J \beta_{p,j} Mu(\tau_j) + O(h^{L+1}) \\ &= U_p + O(h^{L+1}) = v(t_p) + O(h^{L+1}), \quad p = 1, \dots, m. \end{aligned}$$

Consequently, because of this and (7-10), we have

$$(7-12) \quad b_1 u^{(1)}(t_p) + \dots + b_{m-1} u^{(m-1)}(t_p) = O(h^{L+1}), \quad p = 1, \dots, m - 1.$$

The functions  $u^{(l)}$  are *fixed* elements in the null space of  $M$  which satisfy (7-9). From Taylor's Theorem we obtain

$$u^{(l)}(t_p) = (t_p - t_0)^l / l! + O(h^m).$$

Therefore, since (7-12) holds,  $h^i b_i / i!$  satisfies a linear system whose right side is  $O(h^{L-1})$  and whose matrix is an  $O(h)$  perturbation of a Vandermonde matrix involving points in  $[0, 1]$  which are well separated. Hence

$$h^i b_i = O(h^{L+1}) \quad \text{or} \quad b_i = O(h^{L+1-i}), \quad i = 1, \dots, m - 1.$$

We now obtain bounds on the differences between divided differences of  $u$  and  $v$ .

We have

$$\begin{aligned} & u[t_k, \dots, t_{k+p}] - v[t_k, \dots, t_{k+p}] \\ &= \sum_{l=1}^{m-1} b_l u^{(l)}[t_k, \dots, t_{k+p}] \\ &= \sum_{l=1}^{m-1} b_l D^p u^{(l)}(t_k + \xi_{k,l,p} [t_{k+p} - t_k])/p!, \end{aligned}$$

where the  $\xi_{k,l,p}$ 's are some values between 0 and 1. Hence

$$(7-13) \quad \begin{aligned} & |u[t_k, \dots, t_{k+p}] - v[t_k, \dots, t_{k+p}]| \\ & \leq \|D^p u^{(l)}(t)\|_\infty \max_l \{|b_l|\} = O(h^{L-m+2}), \end{aligned}$$

where  $\|\cdot\|_\infty$  denotes the max norm on  $[A, B]$ .

From (7-10) it follows that all derivatives of  $v$  are uniformly bounded in terms of derivatives of  $u, u^{(1)}, \dots, u^{(m-1)}$ . From Theorem 3-1, the truncation error for  $v$  satisfies

$$\|T_n[v]\| = O(H_n^{L-m+1}).$$

Since  $v(t_p) = U_p, p = 0, \dots, m - 1$ , and since we are assuming that the mesh points satisfy the restrictions of Theorem 6-1, we can apply this theorem with  $c_q = 0$  and  $F = T_n[v]$  to obtain

$$|v[t_k, \dots, t_{k+p}] - U[t_k, \dots, t_{k+p}]| \leq K_2 H_n^{L-m+1} \exp(K_1 n H_n),$$

where  $K_i$  denote constants. Thus, from (7-13) we obtain

$$|u[t_k, \dots, t_{k+p}] - U[t_k, \dots, t_{k+p}]| \leq K_3 H_n^{L-m+1} \exp(K_1 n H_n).$$

The preceding analysis applies to HODIE solutions on each of a sequence of mesh points:

$$(7-14a) \quad A = t_{0,n} < t_{1,n} < \dots < t_{n,n} = B,$$

and sets of auxiliary points:

$$(7-14b) \quad t_{k,n} \leq \tau_{k,1,n} < \dots < \tau_{k,J,n} \leq t_{k+m,n}, \quad k = 0, \dots, n - m,$$

and sets of boundary auxiliary points:

$$(7-14c) \quad t_{0,n} \leq \tau_{A,1,n} < \dots < \tau_{A,J,n} \leq t_{m-1,n}.$$

Set

$$(7-14d) \quad \begin{aligned} & h_{k,n} = (t_{k+m,n} - t_{k,n})/(m + 1), \quad h_{A,n} = (t_{m-1,n} - t_{0,n})/m, \\ & H_n = \max_{k=0, \dots, n-m} \{h_{k,n}\} \end{aligned}$$



and

$$(7-14e) \quad R_n = \max \left[ \begin{aligned} &\max_{p=0, \dots, m-2} \{h_{A,n}/(t_{p+1,n} - t_{p,n})\}, \\ &\max_{j=1, \dots, J-1} \{h_{A,n}/(\tau_{A,j+1,n} - \tau_{A,j,n})\}, \\ &\max_{k=0, \dots, n-1} \{H_n/(t_{k+1,n} - t_{k,n})\}, \\ &\max_{\substack{k=0, \dots, m-n \\ j=1, \dots, J}} \{H_n/(\tau_{k,j+q,n} - \tau_{k,j,n})\} \end{aligned} \right].$$

The HODIE initial value difference equations for one of the sets of mesh points are then

$$(7-15a) \quad (1/h_{k,n})^m \sum_{i=0}^m \alpha_{k,i,n} U_{k+i} = \sum_{j=1}^J \beta_{k,j,n} f(\tau_{k,j,n}), \quad k = 0, \dots, m-1,$$

$$(7-15b) \quad U_0 = u(A),$$

$$(7-15c) \quad U_p = \sum_{q=0}^{m-1} \gamma_{p,q,n} c_p (t_p - A)^q / q! + (h_{A,n})^m \sum_{j=1}^J \beta_{A,p,j,n} f(\tau_{A,j,n}),$$

$$p = 1, \dots, m-1.$$

Then, we have the following result:

**THEOREM 7-1.** *Suppose  $a_0, \dots, a_{m-1}, f \in C^{L-m+1}$ . Let  $u$  denote the solution of (7-1). For mesh points, auxiliary points, and parameters in (7-14), assume that*

$$R_n \leq R_{\max} \quad \text{and} \quad nH_n \leq \text{const} \quad \text{as } n \rightarrow \infty.$$

*Then for sufficiently small  $H_n$ , there are coefficients  $\alpha, \beta, \gamma$  such that the HODIE approximation (7-15) is exact on  $P_L$  and, for each  $n$ , defines a unique solution  $U = U^{(n)}$ . Its first  $(m-1)$ st divided differences are  $O(H_n^{L-m+1})$  approximations to the divided differences of  $u$*

$$U^{(n)}[t_{k,n}, \dots, t_{k+p,n}] = u[t_{k,n}, \dots, t_{k+p,n}] + O(H_n^{L-m+1}),$$

for  $k = 0, \dots, n-p, p = 0, \dots, m-1$ .

**COROLLARY 7-1.** *Let  $u^{(l)}, l = 0, \dots, m$ , denote the solutions of the problems*

$$Mu^{(l)} = 0, \quad l = 0, \dots, m-1, \quad Mu^{(m)} = f, \quad A < t \leq B,$$

$$D^p u^{(l)}(t_0) = \delta_{p,l}, \quad l, p = 0, 1, \dots, m-1,$$

$$D^p u^{(m)}(t_0) = 0, \quad p = 0, 1, \dots, m-1.$$

Let  $U^{(n,l)}, l = 0, \dots, m$ , denote corresponding HODIE approximations. Under the

same hypotheses of Theorem 7-1, for all sufficiently small  $H_n$ ,

$$U^{(n,l)}[t_k, \dots, t_{k+p}] = u^{(l)}[t_k, \dots, t_{k+p}] + O(H_n^{L-m+1}),$$

$$k = 0, \dots, n - p, p = 0, \dots, m - 1, l = 0, \dots, m.$$

**8. Discretization Error for the Separated Two-Point Boundary Value Problem.**

We now treat the separated two-point boundary value problem:

(8-1) 
$$Mu(t) = f(t), \quad A < t < B,$$

(8-2a) 
$$M_{Aq}u = \sum_{p=0}^{m-1} a_{Aqp} D^p u(A) = c_{Aq}, \quad q = 0, \dots, q_A,$$

(8-2b) 
$$M_{Bq}u = \sum_{p=0}^{m-1} a_{Bqp} D^p u(B) = c_{Bq}, \quad q = q_A + 1, \dots, m - 1.$$

We assume that this problem is well-posed. We also assume that  $0 \leq q_A \leq m - 2$  and then neither (8-2a) nor (8-2b) is vacuous, otherwise this reduces to the initial value problem treated in Section 7.

The general solution of (8-1) can be written in terms of a set of constants  $b_0, \dots, b_{m-1}$  and the functions  $u^{(l)}, l = 0, \dots, m$ , given in Corollary 7.1

(8-3a) 
$$u = b_0 u^{(0)} + \dots + b_{m-1} u^{(m-1)} + u^{(m)}.$$

The solution of (8-1)–(8-2) is then (8-3a) where the  $b$ 's satisfy the algebraic system

$$\sum_{p=0}^{m-1} a_{Aqp} b_p = c_{Aq}, \quad q = 0, \dots, q_A,$$

(8-3b) 
$$\sum_{p=0}^{m-1} a_{Bqp} \sum_{l=0}^{m-1} b_l D^p u^{(l)}(B)$$

$$= c_{Bq} - \sum_{p=0}^{m-1} a_{Bqp} D^p u^{(m)}(B), \quad q = q_A + 1, \dots, m - 1.$$

The assumption that the boundary value problem (8-1)–(8-2) is well-posed implies that the coefficients  $a_{Aq}, a_{Bq}$  are such that this system is satisfied by one and only one set of  $b$ 's.

In the following, we assume that the mesh points are as in Section 7 and as  $n \rightarrow \infty$  the conditions on the mesh points given in Theorem 7.1 are satisfied. We set

$$h_A = (t_{m-1} - t_0)/(m - 1), \quad h_B = (t_n - t_{n-m-1})/(m - 1),$$

suppressing the dependence on  $n$  for brevity, and take  $H_n$  as in Section 7.

For the HODIE estimate, we obtain approximations to the boundary conditions at  $t = A = t_0$  in (8-2a) from the initial conditions (7-3). The system (7-3) can be solved for  $u(A), Du(A), \dots, D^{m-1}u(A)/(m - 1)!$ . This is because the  $\gamma_{p,q}$ 's are  $1 + O(h_A)$  and thus the matrix multiplying the vector of these derivatives, after multiplication by a diagonal matrix with elements  $d_{ii} = h_A^{1-i}$ , is an  $O(h)$  perturbation of a Vandermonde matrix with well separated points on the interval  $[0, 1]$ . We can write the result as

$$(8-4a) \quad \sum_{i=0}^{m-1} \epsilon_{A p i} U_i = D^p u(A) + \sum_{j=1}^J \phi_{A p j} f(\tau_{A j}), \quad p = 0, \dots, m-1,$$

where  $\tau_{A j}, j = 1, \dots, J$ , denote the auxiliary points in (7-3c). By construction, (8-4a) is exact for  $u \in \mathbf{P}_L, L = J + m - 1$ , if  $f = Mu$ . Furthermore, the  $\beta$ 's in (7-3) are uniformly bounded; hence, we conclude that the resulting  $\phi$ 's in (8-4a) are  $O(h^{m-p})$ . Finally, if the values of  $u \in C^{l+1}$  and  $Mu$  at the stencil and auxiliary points are substituted in (7-3) and if the resulting system (7-7) is similarly solved for the derivatives, we find

$$(8-4b) \quad \sum_{i=1}^{m-1} \epsilon_{A p i} u(t_i) = D^p u(A) + \sum_{j=1}^J \phi_{A p j} Mu(\tau_{A j}) + O(h_A^{L+1-p}),$$

$p = 0, \dots, m-1.$

Multiplication of (8-4) by  $a_{A q p}$  and summation with respect to  $p$  gives

$$(8-5a) \quad \sum_{p=0}^{m-1} a_{A q p} \left[ \sum_{i=0}^{m-1} \epsilon_{A p i} U_i - \sum_{p=0}^{m-1} D^p u(A) - \sum_{j=1}^J \phi_{A p j} f(\tau_{A j}) \right] = 0,$$

$$(8-5b) \quad \sum_{p=0}^{m-1} a_{A q p} \left[ \sum_{i=0}^{m-1} \epsilon_{A p i} u(t_i) - \sum_{p=0}^{m-1} D^p u(A) - \sum_{j=1}^J \phi_{A p j} Mu(\tau_{A j}) \right] = O(h_A^{L-m+1}),$$

where, as above,  $u \in C^{L+1}$ .

Similar equations in terms of values at the right end of the interval are obtained for distinct auxiliary points  $\tau_{B j}$  such that  $t_{n-m+1} \leq \tau_{B j} \leq t_n = B$ . The analogues of (8-5) are

$$(8-6a) \quad \sum_{p=0}^{m-1} a_{B q p} \left[ \sum_{i=0}^{m-1} \epsilon_{B p i} U_{n-i} - \sum_{p=0}^{m-1} D^p u(B) - \sum_{j=1}^J \phi_{B p j} f(\tau_{B j}) \right] = 0,$$

$$(8-6b) \quad \sum_{p=0}^{m-1} a_{B q p} \left[ \sum_{i=0}^{m-1} \epsilon_{B p i} u(t_{n-i}) - \sum_{p=0}^{m-1} D^p u(B) - \sum_{j=1}^J \phi_{B p j} Mu(\tau_{B j}) \right]$$

$= O(h_B^{L-m+1}).$

HODIE boundary conditions are obtained from (8-5a) and (8-6a)

$$(8-7a) \quad \sum_{p=0}^{m-1} \sum_{i=0}^{m-1} a_{A q p} \epsilon_{A p i} U_i = c_{A q} + \sum_{p=0}^{m-1} \sum_{j=1}^J a_{A q p} \phi_{A p j} f(\tau_{A j}),$$

$$q = 0, \dots, q_A,$$

$$(8-7b) \quad \sum_{p=0}^{m-1} \sum_{i=0}^{m-1} a_{B q p} \epsilon_{B p i} U_{n-i} = c_{B q} + \sum_{p=0}^{m-1} \sum_{j=1}^J a_{B q p} \phi_{B p j} f(\tau_{B j}),$$

$$q = q_A + 1, \dots, m-1.$$

The HODIE approximation  $U$  of the solution  $u$  of (8-1)–(8-2) is obtained by solving  $M_n U_j = I_n f_j$  subject to (8-7). It is a consequence of the proof of convergence below that this system has a unique solution for all sufficiently small  $h = \max\{h_A, h_B, H_n\}$ .

The sums over  $i$  above can be expressed in terms of divided differences, for example, for any function  $g$ ,

$$\sum_{i=0}^{m-1} \epsilon_{A p i} g(t_i) = \sum_{i=0}^{m-1} \nu_{A p i} g[t_0, \dots, t_i].$$

We now exhibit the structure of these  $\nu$ 's; (8-4b) becomes

$$\sum_{i=0}^{m-1} \nu_{A p i} u[t_0, \dots, t_i] = D^p u(A) + \sum_{j=1}^J \phi_{A p j} M u(\tau_{A j}) + O(h_A^{L+1-p}).$$

By construction, the  $O(h_A^{L+1-p})$  term vanishes if  $u \in \mathbf{P}_L$  and for  $u(t) = s_l(t) = (t - A)^l/l!$ , we have

$$s_l[t_0, \dots, t_i] = \begin{cases} O(h_A^{l-i}), & i = 0, \dots, l-1, \\ 1/l!, & i = l, \\ 0, & i > l. \end{cases}$$

Thus, for each fixed  $p = 0, 1, \dots, m-1$ ,

$$\begin{aligned} \nu_{A p 0} &= D^p s_0(A) + \sum_{j=1}^J \phi_{A p j} M s_0(\tau_{A j}), \\ \nu_{A p l}/p! &= D^p s_l(A) + \sum_{j=1}^J \phi_{A p j} M s_l(\tau_{A j}) \\ &\quad - \sum_{k=0}^{l-1} O(\nu_{A p k} h_A^{l-k}), \quad l = 1, \dots, m-1; \end{aligned}$$

and

$$D^p s_l(A) = \begin{cases} 1, & p = l, \\ 0, & p \neq l. \end{cases}$$

Therefore, since  $\phi_{A p j}$  is  $O(h_A^{m-p})$ , we have

$$(8-8) \quad \nu_{A p i} = \begin{cases} O(h_A^{m-p}), & i = 0, \dots, p-1, \\ p! + O(h_A^{m-p}), & i = p, \\ O(h_A), & i = p+1, \dots, m-1. \end{cases}$$

Subtract (8-5b) from (8-5a) and express the sums over  $i$  in terms of divided differences to obtain

$$\begin{aligned} &\sum_{p=0}^{m-1} \sum_{i=0}^{m-1} a_{A q p} \nu_{A p i} U[t_0, \dots, t_i] \\ &= \sum_{p=0}^{m-1} \sum_{i=0}^{m-1} a_{A q p} \nu_{A p i} u[t_0, \dots, t_i] + O(h_A^{L-m+1}). \end{aligned}$$

The general solution of  $M_n U_j = I_n f_j$  can be written in terms of constants  $B_j$  and functions  $U^{(n,l)}$  given in Corollary 7-1 as

$$(8-9) \quad U = B_0 U^{(n,0)} + \dots + B_{m-1} U^{(n,m-1)} + U^{(n,m)}.$$

Substitute from (8-3) and (8-9) for  $u$  and  $U$  to obtain

$$\begin{aligned} & \sum_{p=0}^{m-1} \sum_{i=0}^{m-1} a_{Aqp} \nu_{Api} \left[ U^{(n,m)} [t_0, \dots, t_i] + \sum_{l=0}^{m-1} B_l U^{(n,l)} [t_0, \dots, t_i] \right] \\ &= \sum_{p=0}^{m-1} \sum_{i=0}^{m-1} a_{Aqp} \nu_{Api} \left[ u^{(m)} [t_0, \dots, t_i] + \sum_{l=0}^{m-1} b_l u^{(l)} [t_0, \dots, t_i] \right] \\ &+ O(h_A^{L-m+1}). \end{aligned}$$

We assume  $u^{(l)} \in C^{L+1}$ ,  $l = 0, \dots, m$ , and thus by Corollary 7-1;

$$(8-10) \quad u^{(l)} [t_0, \dots, t_i] = U^{(n,l)} [t_0, \dots, t_i] + O(H_n^{L-m+1}).$$

Therefore, because of (8-8),

$$(8-11) \quad \begin{aligned} & \sum_{p=0}^{m-1} \sum_{i=0}^{m-1} a_{Aqp} \nu_{Api} \sum_{l=0}^{m-1} (B_l - b_l) U^{(n,l)} [t_0, \dots, t_i] \\ &= O(H_n^{L-m+1}) + O(h_A^{L-m+1}), \quad q = 0, \dots, q_A. \end{aligned}$$

The left side can be expressed in terms of the vectors

$$\begin{aligned} \mathbf{a}_{Aq}^T &= (a_{Aq0}, \dots, a_{Aq,m-1}), \quad q = 0, \dots, q_A, \\ \mathbf{D}^T &= (B_0 - b_0, \dots, B_{m-1} - b_{m-1}), \end{aligned}$$

and the matrices  $(\nu_A)$ ,  $(\Delta U_A)$  with elements

$$\begin{aligned} (\nu_A)_{pi} &= \nu_{Api}, \quad p, i = 0, \dots, m-1, \\ (\Delta U_A)_{il} &= U^{(l)} [t_0, \dots, t_i], \quad i, l = 0, \dots, m-1, \end{aligned}$$

as  $\mathbf{a}_{Aq}^T (\nu_A) (\Delta U_A) \mathbf{D}$ ,  $q = 0, \dots, q_A$ . The elements of  $(\Delta U_A)$  satisfy

$$\begin{aligned} (\Delta U_A)_{il} &= u^{(l)} [t_0, \dots, t_i] + O(H_n^{L-m+1}) \\ &= D^i u^{(l)}(A) + O(h_A) + O(H_n^{L-m+1}). \end{aligned}$$

Thus, because of (8-8), the elements of the product  $(\nu_A) (\Delta U_A)$  differ by at most  $O(h_A)$  from elements of the Wronskian matrix  $(Du_A)$ ,

$$(Du_A)_{il} = D^i u^{(l)}(A), \quad i, l = 0, \dots, m-1,$$

of the functions  $u^{(l)}$ ,  $l = 0, \dots, m-1$ , evaluated at  $A$  (which, in this instance, makes  $Du_A$  the  $m$ -by- $m$  identity matrix).

Hence the  $(q_A + 1)$ -by- $m$  matrix  $\mathbf{a}_{Aq}^T (\nu_A) (\Delta U_A)$  is an  $O(H_n)$  perturbation of the first  $q_A + 1$  rows of the linear system (8-3b).

Similarly, we obtain

$$(8-12) \quad \sum_{p=0}^{m-1} \sum_{i=0}^{m-1} a_{Bqp} \nu_{Bpi} \sum_{l=0}^{m-1} (B_l - b_l) U^{(n,l)} [t_0, \dots, t_i] \\ = O(H_n^{L-m+1}) + O(h_B^{L-m+1}), \quad q = q_A - 1, \dots, m - 1,$$

and the left sides can be written as

$$a_{Bq}^T (\nu_B) (\Delta U_B) \mathbf{D}, \quad q = q_A + 1, \dots, m - 1,$$

where  $(\nu_B) (\Delta U_B)$  is an  $O(H_n)$  perturbation of the Wronskian of the functions  $u^{(l)}$ ,  $l = 0, \dots, m - 1$ , evaluated at  $B$ . Hence, for all sufficiently small  $h_A, h_B, H_n$  the matrix

$$(8-13) \quad a_{Bq}^T (\nu_B) (\Delta U_B), \quad q = q_A + 1, \dots, m - 1,$$

consists of an  $O(H_n)$  perturbation of the last  $m - q_A - 1$  rows of the linear system (8-3b).

It then follows from (8-3b), (8-11), and (8-12), that

$$B_l - b_l = O(h^{L-m+1}), \quad h = \max\{h_A, h_B, H_n\};$$

and thus, because of (8-10)

$$U[t_k, \dots, t_{k+l}] = u[t_k, \dots, t_{k+l}] + O(h^{L-m+1}), \quad l = 0, \dots, m - 1.$$

We augment (7-14) with boundary auxiliary points:

$$(8-14a) \quad t_{n-m+1,n} \leq \tau_{B,1,n} < \dots < \tau_{B,J,n} \leq t_{n,n}$$

and set

$$(8-14b) \quad h_{B,n} = (t_{n,n} - t_{n-m+1,n}) / (m - 1),$$

$$(8-14c) \quad R'_n = \max \left[ R_n, \max_{j=1, \dots, J-1} \{h_{B,n} / (\tau_{B,j+1,n} - \tau_{B,j,n})\} \right].$$

We have then proved the following

**THEOREM 8-1.** *Assume that the problem (8-1)–(8-2) is well-posed. In addition to the hypotheses and notation of Theorem 7-1, include boundary auxiliary points and parameters of (8-14). Suppose that*

$$R'_n \leq R_{\max} \quad \text{and} \quad nH_n \leq \text{const} \quad \text{as } n \rightarrow \infty.$$

*Then for all sufficiently small  $h = \max\{h_A, h_B, H_n\}$ , there are coefficients  $\alpha, \beta, \epsilon, \phi$  such that the HODIE approximation given by  $M_n U = I_n f$  and (8-7) is exact on  $\mathbf{P}_L$ . The system has a unique solution  $U = U^{(n)}$  whose first  $(m - 1)$ st divided differences are  $O(h^{L-m+1})$  approximations to the corresponding divided differences of  $u$ .*

**9. Computation Analysis.** In this section, we consider the computational aspects of the HODIE method. We discuss specific features of our implementation, and we compare the amount of work with other available methods. The discussion

is restricted to the case of second-order equations subject to Dirichlet boundary conditions for four reasons: it is simple, it is the most important case, it is readily generalized, and there are detailed analyses of other methods available for comparison.

The differential equation problem is

$$Mu(t) = a_2(t)u''(t) + a_1(t)u'(t) + a_0(t)u(t) = f(t), \quad A \leq t \leq B,$$

$u(A)$  and  $u(B)$  given,

where, for generality, we have taken the coefficient of  $u''$  in  $M$  to be a positive function  $a_2$  rather than unity. Estimates  $U_k = U(t_k)$  of  $u(t_k)$  at mesh points  $A = t_0 < t_1 < \dots < t_n = B$  are obtained by solving the HODIE difference equation problem for  $k = 0, \dots, n - 2$

$$M_n U_k \equiv [\alpha_{k,0} U_k + \alpha_{k,1} U_{k+1} + \alpha_{k,2} U_{k+2}] / h_k^2 = \sum_{j=1}^J \beta_{k,j} f(\tau_{k,j}) \equiv I_n f_k,$$

$$U_0 = u(A), \quad U_n = u(B), \quad h_k = (t_{k+2} - t_k) / 2,$$

where the coefficients  $\alpha, \beta$  satisfy  $M_n [s_l]_k = I_n [Ms_l]_k$  for  $s_l, l = 0, \dots, L$ , a basis for  $\mathbf{P}_L$ . We consider two choices of the auxiliary points  $\tau_{k,j}, j = 1, \dots, J$ :

*Regular auxiliary points:*  $\tau_{k,j} = t_k + (j - 1)h_k / J$ ,

*Gauss-type auxiliary points:* the generalized  $B$ -spline Gauss points.

There are two distinct parts in an implementation of a specific HODIE approximation. The first part consists in the determination of the values of the coefficients  $\alpha_{k,i}, i = 0, 1, 2$ , and  $\beta_{k,j}, j = 1, \dots, J$ , for each  $k = 0, \dots, n - 2$ , and then the determination of the values  $I_n f_k, k = 0, \dots, n - 2$ . The second part is the determination of the values  $U_k, k = 1, \dots, n - 1$ , of the solution of the resulting  $(n - 1)$ -by- $(n - 1)$  tridiagonal system of difference equations.

In the first part, the system of algebraic equations for the  $\alpha$ 's and  $\beta$ 's is reducible: one solves a  $J$ -by- $J$  system for the  $\beta$ 's and then a 3-by-3 system for the  $\alpha$ 's; this is done for each  $k = 0, \dots, n - 2$ . This reducibility results in significant savings of work for the special second-order case,  $m = 2$ , as well as in the general case. Although the Lagrange basis is convenient for theoretical analysis, we have found that it is computationally more efficient to use a different basis:

$$s_0(t) = 1, \quad s_1(t) = t - t_{k+1}, \quad s_2(t) = (t - t_k)(t - t_{k+2}),$$

$$s_{3+l}(t) = (t - t_k)(t - t_{k+1})(t - t_{k+2})p_{l-3}(t),$$

where

$$p_0(t) = 1, \quad p_1(t) = (t - t_{k+1}), \quad p_2(t) = (t - t_{k+1})^2,$$

$$p_3(t) = (t - t_k)p_2(t), \quad p_4(t) = (t - t_k)^2 p_2(t),$$

$$p_5(t) = (t - t_{k+2})p_3(t), \quad p_6(t) = (t - t_{k+2})^2 p_3(t), \quad \text{and so on.}$$

With  $\delta_1 = t_{k+1} - t_k, \delta_2 = t_{k+2} - t_{k+1}$ , this choice leads to the following system

for  $\alpha_{k,i}/h_k^2 = \eta_{k,i}$ ,

$$\eta_{k,0} + \eta_{k,1} + \eta_{k,2} = \sum_j \beta_{k,j} a_0(\tau_{k,j}),$$

$$\delta_1 \eta_{k,0} + \delta_2 \eta_{k,2} = \sum_j \beta_{k,j} [a_1(\tau_{k,j}) + (\tau_{k,j} - t_{k+1}) a_0(\tau_{k,j})],$$

$$\begin{aligned} \delta_1 \delta_2 \eta_{k,1} = \sum_j \beta_{k,j} [2a_2(\tau_{k,j}) + 2(\tau_{k,j} - t_{k+1}) a_1(\tau_{k,j}) \\ + (\tau_{k,j} - t_k)(\tau_{k,j} - t_{k+2}) a_0(\tau_{k,j})]. \end{aligned}$$

Use of the normalization  $\beta_{k,1} = 1$  eliminates one of the  $J$  HODIE equations. The final equations for the  $\beta$ 's are

$$\sum_{j=2}^J \nu_{l,j} \beta_{k,j} = -\nu_{l,1}, \quad l = 1, \dots, J - 1,$$

$$\nu_{l,j} = s_{2+l}''(\tau_{k,j}) a_2(\tau_{k,j}) + s_{2+l}'(\tau_{k,j}) a_1(\tau_{k,j}) + s_{2+l}(\tau_{k,j}) a_0(\tau_{k,j}).$$

The choice of the basis elements makes the evaluation of the coefficients simple and it also gives a structure to the system which allows it to be solved easily. Specifically, for the regular case, three of the auxiliary points are at mesh points. Arranging the system so that its first three columns correspond to  $t_{k+1}, t_k, t_{k+2}$ , one finds that these columns have the special form

$-\delta_1 \delta_2 a_1(t_{k+1})$	$-6\delta_1 a_2(t_k) + 2\delta_1 \delta_2 a_1(t_k)$	$6\delta_2 a_2(t_{k+2}) + 2\delta_1 \delta_2 a_1(t_{k+2})$
$2\delta_1 \delta_2 a_2(t_{k+1})$	X	X
0	X	X
0	X	X
0	0	X
0	0	X
0	0	0
$\vdots$	$\vdots$	$\vdots$
0	0	0

where the X's indicate nonzero elements. This, of course, is very advantageous for solving for the  $\beta$ 's in the regular case.

We consider the computational effort required first for a uniform partition:  $t_k = kh, k = 0, \dots, n$ . We measure the effort in terms of the number  $F$  of function evaluations ( $a_2, a_1, a_0$ , or  $f$ ) and the number  $M$  of multiplications required. In regard to the nonfunction evaluation work, we assume: *the total computational effort is proportional to the number of multiplications*. Table 9-1 lists the effort required for various parts of an implementation of the HODIE scheme.



TABLE 9-1

*Number of multiplications and function evaluations required for each interior mesh point for HODIE approximations of orders 4, 6, 8, 10 for the Regular and the Gauss-type Cases for auxiliary points.*

Computation step	Regular Case				Gauss-type Case			
	$J = 3$	5	7	9	2	3	4	5
Compute the $\beta$ -matrix elements (multiplies)	8	39	89	137	6	14	36	50
Solve for the $\beta$ 's	3	17	47	111	1	7	38	47
Evaluate right sides of the $\alpha$ -equations	13	21	33	43	12	18	24	30
Solve for the $\alpha$ 's	3	3	3	3	3	3	3	3
Solve the tridiagonal system for the $U$ 's	7	9	11	13	6	7	8	9
Total number of multiplications	34M	89M	183M	307M	28M	49M	109M	139M
Total number of function evaluations	4F	12F	20F	28F	8F	12F	16F	20F

The  $\beta$ -matrix elements are found from a simple examination and assuming that the values of  $s''$ ,  $s'$ , and  $s$  have been previously computed and stored (these values are independent of  $k$  since a uniform partition is assumed). The special structure of this matrix for the Regular Case is assumed for estimating the work to solve this matrix equation. For the Gauss-type Case, we have a general  $(J - 1)$ -by- $(J - 1)$  system to solve. Note that we assume that the Gauss-type auxiliary points have been previously computed or are otherwise known. The right sides of the  $\alpha$ -equations are of a special form and the computation is carried out by forming  $\beta_{k,j} a_j(\tau_{k,j})$  and then combining these appropriately. The solution of the  $\alpha$ -equations is trivial and the final multiplications occur in solving the large tridiagonal system plus the evaluation of its right side. In the Regular Case the function values at the mesh points, and the auxiliary points are used more than once without recomputation.

We now use these work estimates to compare, roughly, the work of the HODIE method with other methods. The comparison is presented in Table 9-2 for seven methods, three different orders of accuracy (4, 6, and 8), and for both uniform and non-uniform partitions. The data for collocation by  $C^1$  piecewise polynomials at Gauss points, least squares by splines, and discrete-Ritz are derived from Russell and Varah [1975], where they are described in detail. We have had to modify the multiplication counts in order to account for the slightly different differential equations used here and to rationalize the effect of the  $E_L$  term used by Russell and Varah. Note that the discrete-Ritz method is *limited to selfadjoint problems* and is, therefore, not strictly comparable to the other methods included in Table 9-2. Collocation by splines and extrapolation of the trapezoid rule are analyzed in detail by Russell [1977], and we have adapted his results for our particular equation. Russell also considers collocation with Hermite cubics and quintics in detail.

TABLE 9-2

*Summary of number of multiplications (M) and function evaluations (F) for seven different methods. The counts are given per interior mesh point or interval, and one would hope that methods with the same order give comparable accuracy.*

Method	Order of the method and mesh type					
	Fourth		Sixth		Eighth	
	Uniform	General	Uniform	General	Uniform	General
HODIE, Regular Case	34M + 4F	40M + 4F	89M + 12F	113M + 12F	183M + 20F	241M + 20F
HODIE, Gauss-type Case	28M + 8F	32M + 8F	49M + 12F	57M + 12F	109M + 16F	140M + 16F
Collocation, $C^1$ piecewise polynomials	38M + 8F	42M + 8F	62M + 12F	72M + 12F	145M + 16F	159M + 16F
Collocation, splines	24M + 4F	56M + 4F	37M + 4F	99M + 4F	52M + 4F	152M + 4F
Extrapolation of the trapezoid rule	32M + 8F	32M + 8F	70M + 16F	70M + 16F	165M + 32F	165M + 32F
Least squares, splines	66M + 8F	90M + 8F	198M + 16F	270M + 16F	440M + 24F	580M + 24F
Discrete Ritz, splines or piecewise Hermite	133M + 9F	157M + 9F	465M + 15F	525M + 15F	1200M + 21F	1300M + 21F

We emphasize that the exact values of these operation counts depend on small details of the implementation of a particular algorithm and one can trade multiplications for additions, and so on, in some instances.

The changes for collocation, least squares, and discrete-Ritz from the equal to nonequal spaced meshes come from the need to evaluate the basis functions at each point. The changes for the HODIE method come from the need to evaluate the derivatives of the basis functions in each interval, and we have assumed that two more multiplications are needed for each element of the  $\beta$ -matrix. A minor increase also occurs in the computation of the right side of the  $\alpha$ -matrix equation. There are only insignificant changes in the extrapolation method's work, but it is not clear how effective extrapolation is for nonuniform spacing (consider extrapolation, even for uniform spacing, for a problem for which the error behavior is as in Figure 10-3).

*Considerable caution should be taken in attaching importance to the specific numbers in Table 9-2.* These give only *rough* comparisons and various *other considerations* can completely override the difference between, say, 28 and 35 multiplications per point. *We can only conclude that the first five methods are generally comparable in work and the last two seem unlikely to be competitive.* Collocation with splines seems to gain a work advantage as the order increases, but it is simultaneously increasingly complicated near the boundaries, which may well negate this advantage somewhat.

To obtain a realistic evaluation of these methods, one needs not only actual execution times for the different methods for a range of problems and accuracies, but one also needs to consider other factors such as numerical reliability and stability, ease of programming, and memory requirements.

The operation counts for the HODIE method for *ordinary* differential equations given here indicate that the work is close to the work involved in a number of other

available methods. But, the comparisons for *partial* differential equations indicate that the work for the HODIE method is *significantly less* than for other available methods; see Lynch and Rice [1975, 1978a, 1978b].

**10. Experimental Results.** For two-point, second-order problems we present support for the following points: (1) The HODIE method converges as predicted by theory; there are no unforeseen numerical complications. (2) There are no unforeseen difficulties or complexities in implementation. (3) There is a definite pattern in the relationship among the accuracy actually achieved, the actual computation time, and the order of the method. Specifically, the higher the desired accuracy, the higher should the order of the method be to minimize computation time. (4) The use of the Gauss-type auxiliary points associated with a particular operator  $M$  gives the rate of convergence predicted by theory. (5) The use of Gauss-type auxiliary points for the operator  $D^2$  for a *general* second-order operator  $M$  improves the rate of convergence over that expected for a general set of auxiliary points.

The first two points must be verified for any new method; the third point applies to collections of methods with varying orders; and the last two points apply to the HODIE method and to certain other schemes, such as collocation and Galerkin which have “superconvergence” characteristics.

We note that most of the content of these five points is supported by the theory presented explicitly or implicitly in the preceding sections, or is part of the general folklore about numerical computations. Nevertheless, experience shows that points such as these must be verified experimentally for a new method and, for the rates of convergence, they must be verified in the sense of establishing that asymptotic results are valid in the range of ordinary application.

Accordingly, we have run hundreds of cases for numerous second-order ordinary differential Dirichlet boundary value problems. The results of these experiments support the points listed above, and we have acquired confidence in the reliability of the HODIE method.

The Fortran program we wrote seemed to be as easy to write and to debug as a program for any other method of solving this class of problems. However, we found that in order to verify the rates of convergence for very high-order HODIE schemes, we had to use a very high precision because the accuracies obtained were so high. In the remainder of this section, we discuss only a small subset of the experiments which we performed.

All computation was done on the Purdue University CDC 6500 with double-precision arithmetic, which uses values with about 28 decimal digits. In each experiment, the domain of the problem was partitioned by an equal spaced mesh with  $N$  sub-intervals, so the mesh spacing  $h$  was proportional to  $1/N$ .

*Example 10-1.*

$$u''(t) - 4u(t) = 2 \cosh(1), \quad 0 < t < 1,$$

$$u(0) = u(1) = 0,$$

$$\text{solution: } u(t) = \cosh(2t - 1) - \cosh(1).$$

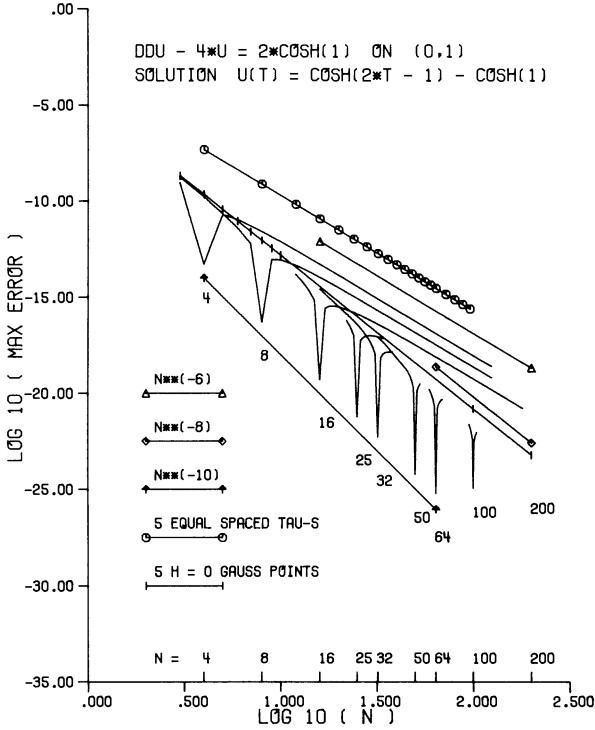


FIGURE 10-1

*Behavior of the error as a function of number N of subintervals for eleven different 5  $\tau$ -point HODIE schemes for Example 10-1.*

This problem has been used by Russell and Shampine [1972], de Boor and Swartz [1973], and others.

Figure 10-1 summarizes one set of experimental results. The logarithm of the maximum error is plotted versus the logarithm of the number of subintervals for eleven different sets of  $J = 5$  auxiliary points. We describe the various curves in this figure and give our interpretation of the results.

(a) The topmost curve gives the results when five regular (equal-spaced) auxiliary points were used. One expects at least  $O(h^5)$  rate of convergence with a set of five auxiliary points because the approximation is locally exact on at least  $P_7$ . The curve shows a very consistent  $O(h^6)$  rate of convergence. The central auxiliary point is the central mesh point of the three-point difference operator  $M_N$  and it is clear from the symmetry of the differential operator that this auxiliary point is a zero of every odd-degree generalized  $B$ -spline orthogonal polynomial. This (or symmetry) shows that one expects  $O(h^6)$  rather than  $O(h^5)$  convergence.

(b) There is a set of nine curves in Figure 10-1 which have sharp downward spikes at  $N = 4, 8, 16, 25, 32, 50, 64, 100,$  and  $200$ , respectively. The set of five auxiliary points used for each one of these curves is the set of five Gauss-type points for that value of  $N$  at which the spike occurs. One has nine different sets of these Gauss-type

points because their locations depend on  $h = 1/N$ . The curve with spike at  $N = 8$  is typical, and we describe some of its features. First, the spike is very abrupt, for the curve also shows the error for the cases of  $N = 7$  and  $N = 9$ . Second, for  $N$  different from 8, the auxiliary points are not the Gauss-type points, hence one expects only  $O(h^6)$ —one of the points is the central mesh point of the operator  $M_N$ —and this behavior can be seen for large values of  $N$ , say  $N$  greater than about 16 for the curve with spike at  $N = 8$ .

(c) Consider the tips of the spikes from the collection of nine curves discussed in (b). If one joins the tips, one sees a very consistent  $O(h^{10})$  rate of convergence for  $N$  up to 64. This is what one expects, since this new curve gives the behavior of the error when five Gauss-type points are used for each  $N$ . The maximum error at  $N = 64$  is about  $10^{-25}$  and the  $O(h^{10})$  rate of convergence breaks down beyond  $N = 64$  because of roundoff error; the values of the Gauss-type points were accurate only to about one part in  $10^{15}$  (single precision on the CDC 6500).

(d) The last curve (with vertical ties above the spikes) is the one for five Gauss-type points for the operator  $M = D^2$ . Except for the central auxiliary point, these are not the Gauss-type points for the operator  $D^2 - 4$ . One expects at least  $O(h^6)$  rate of convergence; however, a very consistent  $O(h^8)$  rate of convergence is observed. As  $h$  tends to zero, the Gauss-type auxiliary points tend to those of the operator  $D^2$ , hence one expects improvement over an arbitrary set of auxiliary points, even a set which contains the central mesh point of the operator  $M_N$ .

*Example 10-2.* Typical of a fairly difficult problem is one taken from Rachford and Wheeler [1974]:

$$\begin{aligned}
 t_0 &= 0.36388, \\
 \frac{d}{dt} \left[ (.01 + 100(t - t_0) \frac{d}{dt} u(t)) \right] \\
 &= -2\{1 + 100(t - t_0)(\tan^{-1} [100(t - t_0)] - \tan^{-1} [100t_0])\}, \\
 u(0) &= u(1) = 0, \\
 \text{solution: } u(t) &= (1 - t)\{\tan^{-1} [100(t - t_0)] + \tan^{-1} [100t_0]\}.
 \end{aligned}$$

The solution has a very sharp rise near  $t = 0.36$ ; it increases from about 0.1 at  $t = 0.3$  to about 1.7 at  $t = 0.4$ , and then it decreases nearly linearly to 0 at  $t = 1$ . See Rachford and Wheeler for a graph of the solution.

Results for two sets of auxiliary points are shown in Figure 10-2, three Regular auxiliary points—which is an extension of the  $O(h^4)$  Störmer-Numerov scheme equivalent to that obtained by Swartz [1974, pp. 304–304]—and the seven Gauss-type auxiliary points for the operator  $D^2$ . One sees that there is a considerable irregularity for  $N$  up to about 100; and then, for large  $N$  the error decreases smoothly at rates of  $O(h^4)$  and  $O(h^{10})$ , respectively. For a general set of seven auxiliary points, one expects  $O(h^7)$  rate of convergence; the use of the Gauss-type points for the operator  $D^2$  apparently improves the rate of convergence to  $O(h^{10})$ .

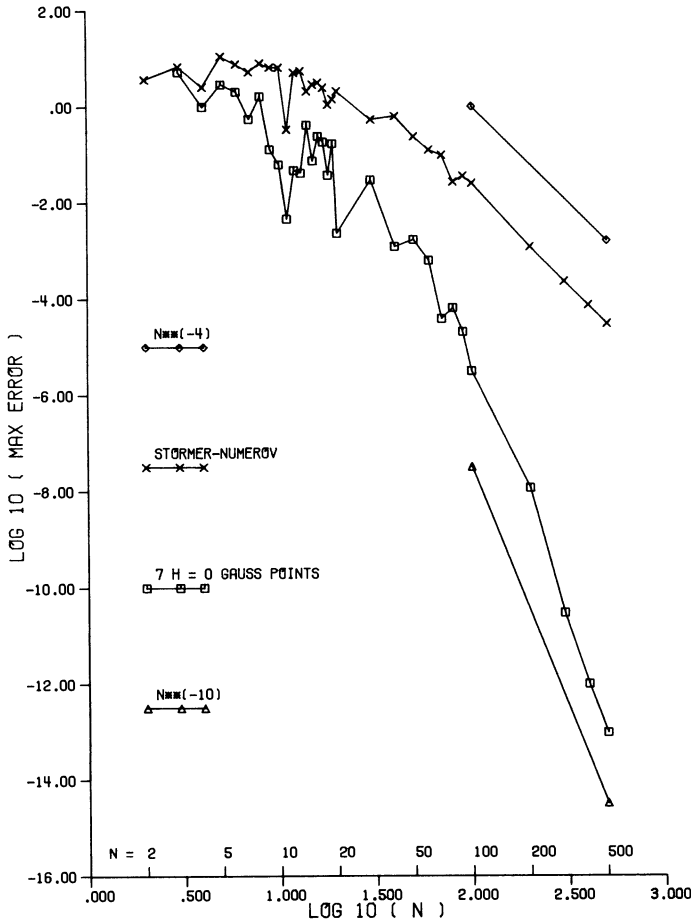


FIGURE 10-2

*Behavior of the error as a function of the number  $N$  of subintervals for two HODIE schemes, one the Størmer-Numerov scheme, and one with seven Gauss-type  $\tau$ -points for the operator  $D^2$  for Example 10-2.*

To compare efficiency, we note that the Størmer-Numerov scheme with  $N = 300$  required almost exactly the same amount of computation time as the seven-point scheme with 100 points. The Størmer-Numerov scheme achieved a maximum error of .00026, which is almost exactly 100 times greater than the error for the higher-order scheme.

Finally, we note that the usefulness of extrapolation techniques is doubtful for either of these schemes for  $N$  less than about 100.

*Example 10-3.* The final example we discuss is:

$$u''(t) + \sin(t)u'(t) + 4t^2u(t) = 2[1 + t \sin(t)] \cos(t^2), \quad 0 \leq t \leq 5,$$

$$u(0) = u(5) = 0,$$

$$\text{solution: } u(t) = \sin(t^2).$$

The solution has several oscillations as  $t$  ranges from 0 to 5.

We solved this problem with a wide variety of HODIE schemes and Figure 10-3 summarizes the results for a selection of them. This figure shows the relationship between work, order, and accuracy. The logarithm of the execution time is plotted versus the logarithm of the maximum error. Since the error is, asymptotically, proportional to  $N^{-p}$  and the time is proportional to  $N$ , one expects straight-line graphs for large  $N$ ; the slope gives  $p$ .

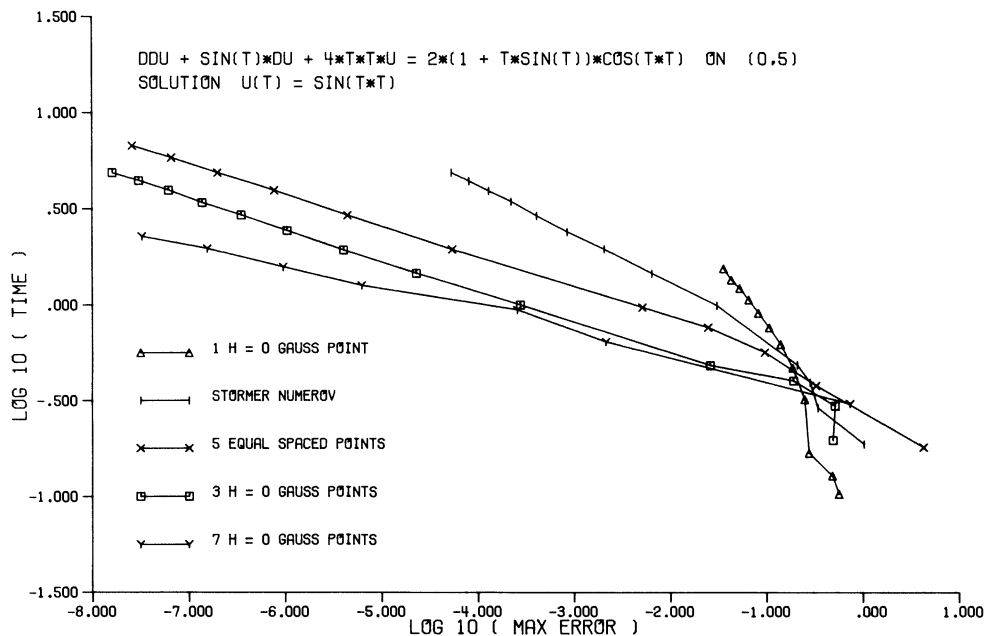


FIGURE 10-3

*Illustration of the relationship between work (execution time), accuracy achieved, and order of the HODIE method for Example 10-3*

One sees the advantage that comes from using a higher-order method for higher accuracy. All of the methods require a fairly large value of  $N$  to achieve any significant accuracy. The low-order methods are competitive only for very low accuracy requirements. The 5-point regular method and the 3-point  $D^2$  Gauss-type method both are  $O(h^6)$  methods, but the maximum error of the regular method is about 10 times larger than the Gauss-type method for the same execution time.

**11. Acknowledgements.** We thank the National Science Foundation for partial support of this research. We thank Carl de Boor for helpful suggestions and comments on early drafts of this paper. We especially thank Blair Swartz for the many hours he spent making incisive comments and suggesting clarifications; his helpful and encouraging attitude greatly increased the clarity of this presentation.

- G. BIRKHOFF & C. R. DE BOOR [1965], "Piecewise polynomial interpolation and approximation," *Approximation of Functions* (H. L. Garabedian, Ed.), Elsevier, Amsterdam, pp. 164–190.
- G. BIRKHOFF & G.-C. ROTA [1969], *Ordinary Differential Equations*, 2nd ed., Blaisdell, New York.
- R. F. BOISVERT [1978], *The Effect on Accuracy of the Placement of Auxiliary Points in the HODIE Method for the Helmholtz Problem*, Dept. of Comput. Sci. Report CSD-TR 266, Purdue University, June.
- L. COLLATZ [1960], *The Numerical Treatment of Differential Equations*, 3rd ed., Springer-Verlag, Berlin.
- H. B. CURRY & I. J. SCHOENBERG [1966], "On Polya frequency functions. IV. The spline functions and their limits," *J. Analyse Math.*, v. 17, pp. 71–107.
- C. DE BOOR & R. E. LYNCH [1966], "On splines and their minimum properties," *J. Math. Mech.*, v. 15, pp. 953–970.
- C. DE BOOR & B. SWARTZ [1973], "Collocation at Gaussian points," *SIAM J. Numer. Anal.*, v. 10, pp. 582–606.
- E. J. DOEDEL [1976], *The Construction of Finite Difference Approximations to Ordinary Differential Equations*, Report, Dept. of Appl. Math., California Institute of Technology, Pasadena, California.
- E. J. DOEDEL [1978], "The construction of finite difference approximations to ordinary differential equations," *SIAM J. Numer. Anal.*, v. 15, pp. 450–465.
- S. J. KARLIN & W. J. STUDDEN [1966], *Tchebycheff Systems: With Applications in Analysis and Statistics*, Interscience, New York.
- R. E. LYNCH & J. R. RICE [1975], *The HODIE Method: A Brief Introduction with Summary of Computational Properties*, Dept. of Comput. Sci. Report 170, Purdue University, Nov. 18.
- R. E. LYNCH & J. R. RICE [1977], *High Accuracy Finite Difference Approximation to Solutions of Elliptic Partial Differential Equations*, Dept. of Comput. Sci. Report CSD-TR 223, Purdue University, Feb. 21.
- R. E. LYNCH [1977a],  $O(h^6)$  Accurate Finite Difference Approximation to Solutions of the Poisson Equation in Three Variables, Dept. of Comput. Sci. Report CSD-TR 221, Purdue University, Feb. 15.
- R. E. LYNCH [1977b],  $O(h^6)$  Discretization Error Finite Difference Approximation to Solutions of the Poisson Equation in Three Variables, Dept. of Comput. Sci. Report CSD-TR 230, Purdue University, April 19.
- R. E. LYNCH & J. R. RICE [1978a], "High accuracy finite difference approximation to solutions of elliptic partial differential equations," *Proc. Nat. Acad. Sci. U.S.A.*, v. 75, pp. 2541–2544.
- R. E. LYNCH & J. R. RICE [1978b], "The HODIE method and its performance for solving elliptic partial differential equations," *Recent Advances in Numerical Analysis* (C. de Boor and G. H. Golub, Eds.), Academic Press, New York, pp. 143–175.
- M. R. OSBORNE [1967], "Minimizing truncation error in finite difference approximation to ordinary differential equations," *Math. Comp.*, v. 21, pp. 133–145.
- J. L. PHILLIPS & R. J. HANSON [1974], "Gauss quadrature rules with B-spline weight functions," *Math. Comp.*, v. 28, p. 666 and microfiche supplement.
- H. H. RACHFORD & M. F. WHEELER [1974], "An  $H^{-1}$  Galerkin procedure for the two-point boundary value problem," *Mathematical Aspects of Finite Elements in Partial Differential Equations* (C. de Boor, Ed.), Academic Press, New York, pp. 253–382.
- R. D. RUSSELL & L. F. SHAMPINE [1972], "A collocation method for boundary value problems," *Numer. Math.*, v. 19, pp. 1–28.
- R. D. RUSSELL & J. M. VARAH [1975], "A comparison of global methods for linear two-point boundary value problems," *Math. Comp.*, v. 29, pp. 1007–1019.
- R. D. RUSSELL [1977], "A comparison of collocation and finite differences for two-point boundary value problems," *SIAM J. Numer. Anal.*, v. 14, pp. 19–39.
- B. K. SWARTZ [1974], "The construction and comparison of finite difference analogs of some finite element schemes," *Mathematical Aspects of Finite Elements in Partial Differential Equations* (C. de Boor, Ed.), Academic Press, New York, pp. 297–312.