

# Discrete Galerkin and Related One-Step Methods for Ordinary Differential Equations\*

By **Bernie L. Hulme**

**Abstract.** New techniques for numerically solving systems of first-order ordinary differential equations are obtained by finding local Galerkin approximations on each subinterval of a given mesh. Different classes of methods correspond to different quadrature rules used to evaluate the innerproducts involved. At each step, a polynomial of degree  $n$  is constructed and the arcs are joined together continuously, but not smoothly, to form a piecewise polynomial of degree  $n$  and class  $C^0$ . If the  $n$ -point quadrature rule used for the innerproducts is of order  $\nu + 1$ ,  $\nu \geq n$ , then the Galerkin method is of order  $\nu$  at the mesh points. In between the mesh points, the  $j$ th derivatives have accuracy of order  $O(h^{\min(\nu, n+1)})$ , for  $j = 0$  and  $O(h^{n-j+1})$  for  $1 \leq j \leq n$ .

**1. Introduction.** This paper extends the concept of discrete Galerkin methods from those based on Gauss-Legendre quadrature [12] to methods based on any interpolatory quadrature formula. Basically, the idea is to approximate each element in the solution of a system of first-order ordinary differential equations by a continuous piecewise polynomial on one subinterval at a time. Two other methods using piecewise polynomial approximation are shown to be equivalent to the discrete Galerkin methods. One is a one-step collocation method which Wright [13] has studied and which is related to some more general methods of Cooper [5], and the other is a quadrature method similar to that studied by Axelsson [1].

Having shown the equivalence of these methods, we easily obtain order of convergence results for all the methods. At the mesh points the errors are of order  $O(h^\nu)$  in the step size  $h$ , where  $\nu + 1$  is the order of the quadrature formula used in the Galerkin approach. Special classes of methods are discussed along with their stability properties, and numerical examples are given.

**2. The Problem and the Approximating Subspaces.** Let us consider solving only a single ordinary differential equation

$$(1) \quad u'(t) = f(t, u(t)), \quad t_0 \leq t,$$

$$(2) \quad u(t_0) = u_0$$

on a finite interval  $[t_0, t_N]$ . It is assumed that  $f(t, x) \in C^r$  in  $[t_0, t_N] \times R$ , where  $R = (-\infty, \infty)$ , so that the exact solution  $u(t) \in C^{r+1}[t_0, t_N]$ , where  $r \geq 1$ , and that  $f$  has a Lipschitz constant in this same region.

---

Received December 20, 1971.

AMS 1969 subject classifications. Primary 6561.

Key words and phrases. Discrete Galerkin methods, initial value problems, ordinary differential equations, piecewise polynomials, collocation, quadrature, implicit Runge-Kutta methods,  $A$ -stable.

\* This work was supported by the United States Atomic Energy Commission.

Copyright © 1972, American Mathematical Society

For the sake of simplicity, we define a *uniform mesh*  $\pi: t_i = t_0 + ih, 0 \leq i \leq N$ , in order that we may approximate  $u(t)$  on the partitioned interval  $[t_0, t_N]$  by a *piecewise  $n$ th degree polynomial*

$$(3) \quad y(t) = \sum_{j=1}^{n+1} b_j^{(i)} \varphi_{i,j}(t), \quad t_i \leq t \leq t_{i+1}, \quad 0 \leq i \leq N - 1,$$

where  $\varphi_{i,j}(t)$  are  $n$ th degree piecewise polynomial basis functions and  $n \geq 1$ . Let us denote by  $S_{n,0}(\pi)$  the class of all such  $y \in C^0[t_0, t_N]$ .

In the next three sections, we present three different methods for obtaining the same approximate solution  $y(t) \in S_{n,0}(\pi)$ . It will be noticed that, in these one-step methods,  $h$  could be changed at each step, and therefore the methods and results would hold for a *variable mesh* also. Moreover, if  $u$  and  $f$  were vector functions,  $y$  could be taken to be a vector of piecewise  $n$ th degree polynomials, and the methods and results would carry over to *systems* of first-order equations.

**3. Discrete Galerkin Methods.** If we require that  $y(t) \in S_{n,0}(\pi)$  provides a local Galerkin approximation to  $u(t)$  on each subinterval of  $\pi$ , then  $y(t)$  must satisfy

$$(4) \quad \begin{aligned} y(t_{i+}) &= u_0, & i &= 0, \\ &= y(t_{i-}), & i &\geq 1, \end{aligned}$$

and

$$(5) \quad \int_{t_i}^{t_{i+1}} [y' - f(t, y)] \varphi_{i,k} dt = 0, \quad 1 \leq k \leq n, \quad 0 \leq i \leq N - 1.$$

To obtain a one-step numerical method, however, we replace the integral in (5) by an interpolatory quadrature formula

$$(6) \quad \int_{t_i}^{t_{i+1}} v(t) dt = h \sum_{m=1}^n w_m v(\sigma_{i,m}) + O(h^{p+1}), \quad v \geq n,$$

$$(7) \quad \sigma_{i,m} = t_i + \theta_m h, \quad 1 \leq m \leq n,$$

where  $0 \leq \theta_1 < \theta_2 < \dots < \theta_n \leq 1$  and  $w_m \neq 0, 1 \leq m \leq n$ . The result is

$$(8) \quad h \sum_{m=1}^n w_m [y'(\sigma_{i,m}) - f(\sigma_{i,m}, y(\sigma_{i,m}))] \varphi_{i,k}(\sigma_{i,m}) = 0, \quad 1 \leq k \leq n.$$

We shall call any  $y(t) \in S_{n,0}(\pi)$  which satisfies (4) and (8) for  $0 \leq i \leq N - 1$  a one-step *discrete Galerkin solution* to (1)–(2).

We may write (4) and (8) in matrix form as

$$(9) \quad \mathbf{A}^\sigma \mathbf{b}^{(i)} = \mathbf{c}^\sigma(\mathbf{b}^{(i)}), \quad 0 \leq i \leq N - 1,$$

where

$$(10) \quad \mathbf{b}^{(i)} = \{b_1^{(i)}, b_2^{(i)}, \dots, b_{n+1}^{(i)}\}^T,$$

$$(11) \quad \begin{aligned} A_{k,i}^\sigma &= \varphi_{i,j}(t_i), & k &= 1, \\ &= h \sum_{m=1}^n w_m \varphi_{i,k-1}(\sigma_{i,m}) \varphi'_{i,j}(\sigma_{i,m}), & 2 \leq k \leq n+1, \quad 1 \leq j \leq n+1, \end{aligned}$$

$$\begin{aligned}
 c_k^G(\mathbf{b}^{(i)}) &= y_i = y(t_i), & k &= 1, \\
 (12) \quad &= h \sum_{m=1}^n w_m f\left(\sigma_{i,m}, \sum_{j=1}^{n+1} b_j^{(i)} \varphi_{i,j}(\sigma_{i,m})\right) \varphi_{i,k-1}(\sigma_{i,m}), & 2 \leq k \leq n+1.
 \end{aligned}$$

To be assured that  $\mathbf{A}^G$  is nonsingular, we assume that the determinant

$$(13) \quad |\varphi_{i,k}(\sigma_{i,m})|_{1 \leq k, m \leq n} \neq 0.$$

Then  $\mathbf{A}^G \mathbf{b}^{(i)} = \mathbf{0}$  implies that  $y(t_i) = 0$  and, in view of (8) and (13), that  $hw_m y'(\sigma_{i,m}) = 0, 1 \leq m \leq n$ . Since  $hw_m \neq 0$  by assumption, we have that  $y'(\sigma_{i,m}) = 0, 1 \leq m \leq n$ , and consequently  $y \equiv 0$  on  $[t_i, t_{i+1}]$ ,  $\mathbf{b}^{(i)} = \mathbf{0}$ , and  $\mathbf{A}^G$  is nonsingular. Thus, our numerical method depends on the solution of

$$(14) \quad \mathbf{b}^{(i)} = (\mathbf{A}^G)^{-1} \mathbf{c}^G(\mathbf{b}^{(i)}).$$

The existence of a unique solution to (14) is guaranteed for sufficiently small  $h$  as follows. Since

$$\|(\mathbf{A}^G)^{-1} \mathbf{c}^G(\mathbf{b}) - (\mathbf{A}^G)^{-1} \mathbf{c}^G(\mathbf{b}^*)\|_\infty \leq \|(\mathbf{A}^G)^{-1}\|_\infty hLQ_G \|\mathbf{b} - \mathbf{b}^*\|_\infty$$

where  $L$  is the Lipschitz constant for  $f$  on  $[t_0, t_N] \times R$  and

$$(15) \quad Q_G = \max_{2 \leq k \leq n+1} \sum_{m=1}^n |w_m \varphi_{i,k-1}(\sigma_{i,m})| \sum_{j=1}^{n+1} |\varphi_{i,j}(\sigma_{i,m})|,$$

the right side of (14) is a contraction mapping on  $R^{n+1}$  when

$$(16) \quad h < H_G = (LQ_G \|(\mathbf{A}^G)^{-1}\|_\infty)^{-1},$$

and a successive substitution iteration will converge to the unique solution of (14).

Notice that (8) generalizes the discrete Galerkin scheme in [12] where the inner products  $(\varphi_{i,k}, \varphi'_{i,j})_h$  are done exactly. That scheme coincides with the present method in the case of Gauss-Legendre quadrature because this rule gives exact results for the innerproducts just mentioned.

It should be pointed out that there exist basis functions  $\varphi_{i,k}$  and abscissae  $\sigma_{i,m}$  for which the determinant in (13) vanishes. For example, when  $n = 2$  there is  $\varphi_{i,1} = 1, \varphi_{i,2} = [(2/h)(t - t_i) - 1]^2$  and  $\sigma_{i,1}, \sigma_{i,2}$  symmetrically placed about  $(t_i + t_{i+1})/2$ . Therefore, one must choose carefully the basis functions and abscissae in the ‘‘orthogonality’’ equations (8) to ensure that  $\mathbf{A}^G$  is nonsingular. An obvious choice is the basis  $\varphi_{i,k} = ((t - t_i)/h)^{k-1}, 1 \leq k \leq n$ , which is *unisolvant* on  $[t_i, t_{i+1}]$ , i.e., (13) holds for *all* sets of  $n$  distinct points  $\sigma_{i,m} \in [t_i, t_{i+1}]$ .

**4. Collocation Methods.** If we require that  $y(t) \in S_{n,0}(\pi)$  collocate to (1) at the points  $\sigma_{i,m}$  of (7), then  $y(t)$  must satisfy (4) and

$$(17) \quad hy'(\sigma_{i,m}) = hf(\sigma_{i,m}, y(\sigma_{i,m})), \quad 1 \leq m \leq n, \quad 0 \leq i \leq N - 1.$$

We shall call any  $y(t) \in S_{n,0}(\pi)$  which satisfies (4) and (17) a *one-step collocation solution* to (1)–(2). The stability of such solutions has been studied by Wright [13], and Cooper [5] has derived similar methods for equations of arbitrary order.

In matrix form, (4) and (17) are

$$(18) \quad \mathbf{A}^C \mathbf{b}^{(i)} = \mathbf{c}^C(\mathbf{b}^{(i)}), \quad 0 \leq i \leq N - 1,$$

where  $\mathbf{b}^{(i)}$  is given by (10), and

$$(19) \quad \begin{aligned} A_{k,i}^C &= \varphi_{i,i}(t_i), & k &= 1, \\ &= h\varphi'_{i,i}(\sigma_{i,k-1}), & 2 \leq k \leq n+1, \quad 1 \leq j \leq n+1, \end{aligned}$$

$$(20) \quad \begin{aligned} c_k^C(\mathbf{b}^{(i)}) &= y_i, & k &= 1, \\ &= hf\left(\sigma_{i,k-1}, \sum_{j=1}^{n+1} b_j^{(i)} \varphi_{i,i}(\sigma_{i,k-1})\right), & 2 \leq k \leq n+1. \end{aligned}$$

Again, it is clear that  $\mathbf{A}^C$  is nonsingular because  $\mathbf{A}^C \mathbf{b}^{(i)} = \mathbf{0}$  implies  $y(t_i) = y'(\sigma_{i,m}) = 0, 1 \leq m \leq n, y \equiv 0$  on  $[t_i, t_{i+1}]$  and  $\mathbf{b}^{(i)} = \mathbf{0}$ . Also, we have

$$\|(\mathbf{A}^C)^{-1} \mathbf{c}^C(\mathbf{b}) - (\mathbf{A}^C)^{-1} \mathbf{c}^C(\mathbf{b}^*)\|_\infty \leq \|(\mathbf{A}^C)^{-1}\|_\infty hLQ_C \|\mathbf{b} - \mathbf{b}^*\|_\infty$$

where

$$(21) \quad Q_C = \max_{2 \leq k \leq n+1} \sum_{j=1}^{n+1} |\varphi_{i,i}(\sigma_{i,k-1})|,$$

so that (18) has a unique solution when

$$(22) \quad h < H_C = (LQ_C \|(\mathbf{A}^C)^{-1}\|_\infty)^{-1}.$$

When  $h$  satisfies both (16) and (22), it is obvious that the collocation solution is identical with the discrete Galerkin solution because the collocation solution also satisfies (8). Moreover, the collocation equations (4) and (17) are simpler than the discrete Galerkin equations (4) and (8) since no quadratures are involved, and no additional assumption such as (13) is required to guarantee the existence and uniqueness of the collocation solution.

Wright [13] has already pointed out that these collocation methods are “equivalent” to a subclass of implicit Runge-Kutta methods in the sense that they produce the same discrete approximations. For some implicit Runge-Kutta methods which are not equivalent to collocation techniques, see Ehle [8, Chapter 4] and Chipman [4, Chapter 3].

**5. Interpolatory Quadrature Methods.** There is yet another class of methods to which the previous two schemes are equivalent in the sense that they all produce the same approximate solution given the same abscissae  $\sigma_{i,m}$  in (7). Let us define the array

$$(23) \quad a_{m,k} = h^{-1} \int_{t_i}^{\sigma_{i,m}} l_k(t) dt, \quad 1 \leq k, m \leq n,$$

where the  $l_k(t)$  are the Lagrange interpolation coefficients

$$(24) \quad l_k(t) = \prod_{j=1; j \neq k}^n \frac{t - \sigma_{i,j}}{\sigma_{i,k} - \sigma_{i,j}}, \quad 1 \leq k \leq n.$$

We shall call any  $y(t) \in S_{n,0}(\pi)$  which satisfies (4) and

$$(25) \quad y(\sigma_{i,m}) = y_i + h \sum_{k=1}^n a_{m,k} f(\sigma_{i,k}, y(\sigma_{i,k})), \quad 1 \leq m \leq n,$$

for  $0 \leq i \leq N - 1$ , a one-step *interpolatory quadrature solution* to (1)–(2). Since (4) is the same as (25) for  $m = 1$  when  $\sigma_{i,1} = t_i$ , we should replace the equation (25) for  $m = 1$  with

$$(25') \quad y'(\sigma_{i,1}) = f(\sigma_{i,1}, y(\sigma_{i,1})), \quad \text{when } \sigma_{i,1} = t_i.$$

Certain methods of this type have been studied by Axelsson [1] and Hammer and Hollingsworth [9].

Equations (4) and (25) can also be put into matrix form

$$(26) \quad \mathbf{A}^Q \mathbf{b}^{(i)} = \mathbf{c}^Q(\mathbf{b}^{(i)}), \quad 0 \leq i \leq N - 1,$$

where

$$(27) \quad \begin{aligned} A_{k,i}^Q &= \varphi_{i,j}(t_i), & k &= 1, \\ &= \varphi_{i,j}(\sigma_{i,k-1}), & 2 \leq k \leq n + 1, & \quad 1 \leq j \leq n + 1, \end{aligned}$$

and

$$(28) \quad \begin{aligned} c_k^Q(\mathbf{b}^{(i)}) &= y_i, & k &= 1, \\ &= y_i + h \sum_{l=1}^n a_{k-1,l} f\left(\sigma_{i,l}, \sum_{j=1}^{n+1} b_j^{(i)} \varphi_{i,j}(\sigma_{i,l})\right), & 2 \leq k \leq n + 1. \end{aligned}$$

$\mathbf{A}^Q$  is obviously nonsingular and (26) has a unique solution when

$$(29) \quad h < H_Q = (LQ_Q \|(\mathbf{A}^Q)^{-1}\|_\infty)^{-1},$$

where

$$(30) \quad Q_Q = \max_{2 \leq k \leq n+1} \sum_{l=1}^n |a_{k-1,l}| \sum_{j=1}^{n+1} |\varphi_{i,j}(\sigma_{i,l})|.$$

When  $\sigma_{i,1} = t_i$ , the second equation in (26) is defined from (25') by

$$(27') \quad A_{2,i}^Q = \varphi'_{i,j}(\sigma_{i,1}), \quad 1 \leq j \leq n + 1,$$

and

$$(28') \quad c_2^Q(\mathbf{b}^{(i)}) = f\left(\sigma_{i,1}, \sum_{j=1}^{n+1} b_j^{(i)} \varphi_{i,j}(\sigma_{i,1})\right)$$

and  $\mathbf{A}^Q$  is still nonsingular.

Since the Galerkin and collocation solutions,  $\bar{y}(t)$ , satisfy (17), we see that

$$(31) \quad \begin{aligned} \bar{y}(t) &= \bar{y}_i + \int_{t_i}^t \bar{y}'(s) ds = \bar{y}_i + \int_{t_i}^t \sum_{k=1}^n f(\sigma_{i,k}, \bar{y}(\sigma_{i,k})) l_k(s) ds \\ &= \bar{y}_i + \sum_{k=1}^n f(\sigma_{i,k}, \bar{y}(\sigma_{i,k})) \int_{t_i}^t l_k(s) ds, \quad t_i \leq t \leq t_{i+1}, \end{aligned}$$

and, in particular, that  $\bar{y}(t)$  satisfies (25) and (25'). Thus, when  $h$  is small enough to satisfy (16), (22) and (29), all three schemes provide the same approximate solution  $\bar{y}(t) = y(t)$ . The collocation method still seems to be the simplest since it does not require the computation of the  $a_{m,k}$ . However, the interpolatory quadrature viewpoint makes it clear that the weights in the Galerkin method are

$$(32) \quad w_m = h^{-1} \int_{t_i}^{t_{i+1}} l_m(t) dt, \quad 1 \leq m \leq n,$$

and, from (31), the solution for all three methods satisfies

$$(33) \quad y_{i+1} = y_i + h \sum_{m=1}^n w_m f(\sigma_{i,m}, y(\sigma_{i,m})), \quad 0 \leq i \leq N - 1.$$

**6. Order of Convergence.** In this section, we use the theory given in Henrici [10, Chapter 2] of discrete one-step methods to derive asymptotic error bounds for the discrete values  $y(t_i) = y_i$  given by the three methods above. Continuous error bounds are then obtained from the discrete ones.

It follows immediately from (33) that all three methods may be written in terms of an increment function  $\Phi$  as

$$(34) \quad y_{i+1} = y_i + h\Phi(t_i, y_i; h), \quad 0 \leq i \leq N - 1,$$

where

$$(35) \quad \Phi(t_i, y_i; h) = \sum_{m=1}^n w_m f(\sigma_{i,m}, y(\sigma_{i,m})).$$

In order for the discrete one-step theory to apply, we must show that  $\Phi$  is Lipschitz continuous with respect to  $y$  in  $\Omega \equiv [t_0, t_N] \times R \times [0, h_0]$ . If, for any  $i, 0 \leq i \leq N - 1$ , and any  $y_i^* \in R$ ,  $y^*(t)$  is the approximate solution to  $u' = f(t, u)$ ,  $u(t_i) = y_i^*$ ,  $t_i \leq t \leq t_{i+1}$ , given by the above methods, then (31) holds for  $y^*$ :

$$(36) \quad y^*(t) = y_i^* + \sum_{k=1}^n f(\sigma_{i,k}, y^*(\sigma_{i,k})) \int_{t_i}^t l_k(s) ds, \quad t_i \leq t \leq t_{i+1}.$$

Subtracting (36) from (31) and letting

$$(37) \quad \sum_{k=1}^n \max_{t_i \leq t \leq t_{i+1}} \left| \int_{t_i}^t l_k(s) ds \right| \leq hB_0, \quad 0 \leq i \leq N - 1,$$

we find that

$$(38) \quad \max_{t_i \leq t \leq t_{i+1}} |y(t) - y^*(t)| \leq \frac{1}{1 - h_0 B_0 L} |y_i - y_i^*|, \quad 0 \leq i \leq N - 1,$$

where  $0 \leq h \leq h_0 < (B_0 L)^{-1}$ . The Lipschitz condition then follows from (35) and (38) since, for  $0 \leq h \leq h_0$  and  $0 \leq i \leq N - 1$ ,

$$(39) \quad |\Phi(t_i, y_i; h) - \Phi(t_i, y_i^*; h)| \leq \frac{LW}{1 - h_0 B_0 L} |y_i - y_i^*|$$

where  $W = \sum_{m=1}^n |w_m|$ .

The discrete error bounds are now derived in

**THEOREM 1.** Assume that  $f(t, x) \in C^r$  in  $[t_0, t_N] \times R$  so that  $u(t) \in C^{r+1}[t_0, t_N]$ , and denote by  $L$  the Lipschitz constant for  $f$  in this region. Given a mesh  $\pi$  of size  $h$ , some basis of piecewise polynomials of degree  $n$  for the space  $S_{n,0}(\pi)$ , distinct abscissae  $\theta_m \in [0, 1]$ ,  $1 \leq m \leq n$ , and the associated interpolatory quadrature formula (6) of order  $\nu + 1$ ,  $n \leq \nu \leq r$ , let the constants  $H_G, H_C, H_Q$  and  $B_0$  be defined as in (16),

(22), (29) and (37), respectively, and let the mesh size  $h$  satisfy  $0 < h \leq h_0 < \min\{H_G, H_C, H_Q, (B_0L)^{-1}\}$ . If  $y(t)$  is the discrete Galerkin, collocation and interpolatory quadrature solution to (1)–(2) defined in Sections 3, 4 and 5, respectively, then there exists a constant  $M$  such that

$$(40) \quad |u_i - y_i| \leq Mh^\nu, \quad 0 \leq i \leq N.$$

*Proof.* The local truncation error  $\tau_i$  is defined from (33) by

$$\begin{aligned} \tau_i &= u_{i+1} - u_i - h \sum_{m=1}^n w_m f(\sigma_{i,m}, u(\sigma_{i,m})) \\ &= \int_{t_i}^{t_{i+1}} f(t, u(t)) dt - h \sum_{m=1}^n w_m f(\sigma_{i,m}, u(\sigma_{i,m})). \end{aligned}$$

Thus, in view of (6),  $|\tau_i| \leq Kh^{\nu+1}$ , where  $K$  depends on the maximum value of  $u^{(\nu+1)}(t)$  on  $[t_0, t_N]$ , and (40) follows immediately from Henrici's Theorem 2.2 [10]. Q.E.D.

Continuous error bounds are derived in

**THEOREM 2.** *Let the hypotheses of Theorem 1 hold. Then there exist constants  $E_j$ ,  $0 \leq j \leq n$ , such that*

$$(41) \quad \max_{t_0 \leq t \leq t_N} |u(t) - y(t)| \leq E_0 h^{\min(\nu, n+1)}, \quad n \leq \nu \leq r,$$

and

$$(42) \quad \max_{t_i \leq t \leq t_{i+1}} |u^{(j)}(t) - y^{(j)}(t)| \leq E_j h^{n-j+1}, \quad 1 \leq j \leq n, \quad 0 \leq i \leq N-1.$$

*Proof.* We write  $u(t)$  as follows, using the  $n$ -point Lagrange interpolatory quadrature formula:

$$(43) \quad u(t) = u_i + \sum_{k=1}^n f(\sigma_{i,k}, u(\sigma_{i,k})) \int_{t_i}^t l_k(s) ds + R_n(t), \quad t_i \leq t \leq t_{i+1},$$

where  $R_n(t) = O(h^{n+1})$ . Subtracting (43) from (31) and using (37), we discover that

$$\max_{t_i \leq t \leq t_{i+1}} |u(t) - y(t)| \leq \frac{1}{1 - h_0 B_0 L} |u_i - y_i| + O(h^{n+1}), \quad 0 \leq i \leq N-1,$$

and (41) follows from (40). If we differentiate (43) and (31)  $j$  times, using  $R_n^{(j)}(t) = O(h^{n-j+1})$ , and subtract, we have

$$\max_{t_i \leq t \leq t_{i+1}} |u^{(j)}(t) - y^{(j)}(t)| \leq LB_j h^{1-j} \max_{1 \leq k \leq n} |u(\sigma_{i,k}) - y(\sigma_{i,k})| + O(h^{n-j+1})$$

for  $1 \leq j \leq n, 0 \leq i \leq N-1$ , where

$$\sum_{k=1}^n \max_{t_i \leq t \leq t_{i+1}} |l_k^{(j-1)}(t)| \leq h^{1-j} B_j.$$

Then (42) follows from (41). Q.E.D.

Theorems 1 and 2 also hold for a *variable mesh*  $\pi$  since  $h$  can be changed at each step, and the methods and the theorems can be carried over to *systems* of first-order equations by applying the single equation techniques to each equation in the system.

**7. Some Specific Methods and Their Stability Properties.** We shall use the following definitions of  $A$ -stability due to Dahlquist [6] and of strong  $A$ -stability due to Chipman [4].

*Definition 1.* A  $k$ -step method is called  $A$ -stable, if all its solutions tend to zero, as  $i \rightarrow \infty$ , when the method is applied with fixed positive  $h$  to any differential equation of the form  $u' = \lambda u$ , where  $\lambda$  is a complex constant with negative real part.

For a one-step method applied to  $u' = \lambda u$ , the approximate solution may be written as  $y_{i+1} = E(\lambda h)y_i$ , where  $E(\lambda h)$  is a rational approximation to  $\exp(\lambda h)$ . Therefore, a one-step method is  $A$ -stable if and only if  $|E(\lambda h)| < 1$  for  $\text{Re}(\lambda h) < 0$ . It is possible, however, that  $|E(\lambda h)| \rightarrow 1$  as  $|\lambda h| \rightarrow \infty$ . In order to define a special kind of  $A$ -stability, which guarantees that  $y_i$  tends to zero rapidly when  $|\lambda h| \gg 1$ , we use

*Definition 2.* A one-step method is *strongly A-stable* if  $y_{i+1} = E(\lambda h)y_i$ ,  $|E(\lambda h)| < 1$  for  $\text{Re}(\lambda h) < 0$ , and  $|E(\lambda h)| \rightarrow 0$  as  $\text{Re}(\lambda h) \rightarrow -\infty$ .

We remark that strongly  $A$ -stable methods should be especially effective on *stiff* systems of equations since rapidly decaying components of the exact solution can be approximated by rapidly decaying components of the approximate solution for any step size  $h$ .

Several classes of collocation and quadrature methods already have been investigated for stability, and these results therefore hold for the equivalent Galerkin methods. Since the order of accuracy of a method is one less than the order of the associated quadrature (Theorem 1), it is convenient to classify the methods according to the abscissae used and thus deduce the order of the method.

The most accurate methods are those using the  $n$  Gauss-Legendre points. The order of these methods is  $O(h^{2n})$ , and Ehle [7] has shown that they are all  $A$ -stable by proving that  $y_{i+1} = P_{n,n}(\lambda h)y_i$ , where  $P_{n,n}(\lambda h)$  is the  $n$ th diagonal Padé rational approximation to  $\exp(\lambda h)$  with the properties  $|P_{n,n}(\lambda h)| < 1$  for  $\text{Re}(\lambda h) < 0$  and  $|P_{n,n}(\lambda h)| \rightarrow 1$  as  $\text{Re}(\lambda h) \rightarrow -\infty$ . Ehle was studying the implicit Runge-Kutta methods of Butcher [2], but in the Gauss-Legendre case these are equivalent at the mesh points to our Galerkin methods [12]. An alternate proof of  $A$ -stability is given by Wright [13].

The next most accurate methods are those based on the  $n$  Radau points, with either the left or the right endpoint fixed, i.e., either  $\sigma_{i,1} = t_i$  or  $\sigma_{i,n} = t_{i+1}$ . Although they are both of order  $O(h^{2n-1})$ , Wright [13] has shown that the left endpoint methods are not  $A$ -stable, while Axelsson [1] has shown that the right endpoint methods are strongly  $A$ -stable because  $y_{i+1} = P_{n,n-1}(\lambda h)y_i$  and this subdiagonal Padé approximation is such that  $|P_{n,n-1}(\lambda h)| < 1$  for  $|\lambda h| < 0$  and  $|P_{n,n-1}(\lambda h)| \rightarrow 0$  as  $\text{Re}(\lambda h) \rightarrow -\infty$ .

The Lobatto methods, which have both endpoints fixed  $\sigma_{i,1} = t_i$  and  $\sigma_{i,n} = t_{i+1}$ , are of order  $O(h^{2n-2})$ . They are  $A$ -stable [1], [13] since  $y_{i+1} = P_{n-1,n-1}(\lambda h)y_i$ . It should be remarked that our Lobatto and Radau (right endpoint) methods are not equivalent to Butcher's corresponding implicit Runge-Kutta methods II and III [3] since Ehle [8, Chapter 4] has shown the latter not to be  $A$ -stable.

Methods based on *equal weight Chebyshev* formulae [11, Section 8.13] are of order  $O(h^{n+1})$  for  $n$  odd and  $O(h^{n+2})$  for  $n$  even. Also, methods associated with *Newton-Cotes* formulae [11, Section 3.5] are of order  $O(h^n)$  for  $n$  even and  $O(h^{n+1})$



for  $n$  odd. And, finally, arbitrarily chosen abscissae would be expected to yield only  $O(h^n)$  accuracy. Wright [13] has proved that any two, three or four symmetrically placed abscissae will produce  $A$ -stable methods.

**8. Sample Calculations.** In the following tables, we present results for the problem

$$(44) \quad u'(t) = u - (2t/u), \quad u(0) = 1, \quad u(t) = (2t + 1)^{1/2}, \quad 0 \leq t \leq 1,$$

computed from the collocation equations (4) and (17) using the basis functions  $\varphi_{i,j}(t) = ((t - t_i)/h)^{j-1}$ ,  $1 \leq j \leq n + 1$ ,  $0 \leq i \leq N - 1$ . Each table concerns one of the classes of methods discussed in Section 7 and illustrates the order of convergence results of Theorem 1 for each class by showing the *discrete error norms*

$$(45) \quad \|e(t; h)\|' = \max_{0 \leq i \leq N} |e(t_i; h)|,$$

for  $h = 1/N$ ,  $1 \leq N \leq 6$ , where  $e = u - y$ , as well as showing in parentheses the *computed orders of convergence*

$$(46) \quad \omega = \frac{\log[\|e(t; h_1)\|'/\|e(t; h_2)\|']}{\log(h_1/h_2)} \approx \nu$$

based on successive mesh sizes  $h_1$  and  $h_2$ .

The nonlinear equations (18) were iterated at each step  $[t_i, t_{i+1}]$  until  $y_{i+1} = \sum_{j=1}^{n+1} b_j^{(i)}$  satisfied a relative error tolerance of  $10^{-11}$ . The reader will notice irregularities in the computed orders of convergence for the more accurate methods. This occurs when the theoretical maximum relative errors are significantly smaller than the relative error tolerance of  $10^{-11}$ . Also, the two-point Radau method failed to satisfy the relative error tolerance test for  $h = 1$ .

TABLE 1. Error Norms for  $n$ -Point Gauss-Legendre

| $h$ | $n = 2$                       | $n = 3$                       | $n = 4$                        | $n = 5$                         | $n = 6$                         |
|-----|-------------------------------|-------------------------------|--------------------------------|---------------------------------|---------------------------------|
| 1   | 1.47(10) <sup>-2</sup>        | 7.08(10) <sup>-4</sup>        | 2.95(10) <sup>-5</sup>         | 1.17(10) <sup>-6</sup>          | 4.71(10) <sup>-8</sup>          |
| 1/2 | 1.39(10) <sup>-3</sup> (3.40) | 2.22(10) <sup>-5</sup> (4.99) | 3.26(10) <sup>-7</sup> (6.50)  | 4.87(10) <sup>-9</sup> (7.90)   | 7.85(10) <sup>-11</sup> (9.23)  |
| 1/3 | 3.07(10) <sup>-4</sup> (3.72) | 2.40(10) <sup>-6</sup> (5.49) | 1.78(10) <sup>-8</sup> (7.17)  | 1.41(10) <sup>-10</sup> (8.74)  | 2.66(10) <sup>-13</sup> (8.34)  |
| 1/4 | 1.01(10) <sup>-4</sup> (3.84) | 4.67(10) <sup>-7</sup> (5.69) | 2.08(10) <sup>-9</sup> (7.47)  | 9.06(10) <sup>-12</sup> (9.54)  | 2.98(10) <sup>-13</sup> (7.61)  |
| 1/5 | 4.25(10) <sup>-5</sup> (3.90) | 1.28(10) <sup>-7</sup> (5.80) | 3.79(10) <sup>-10</sup> (7.63) | 3.84(10) <sup>-13</sup> (14.17) | 1.08(10) <sup>-12</sup> (-5.76) |
| 1/6 | 2.08(10) <sup>-5</sup> (3.93) | 4.40(10) <sup>-8</sup> (5.86) | 9.38(10) <sup>-11</sup> (7.66) | 4.55(10) <sup>-13</sup> (-0.93) | 1.61(10) <sup>-12</sup> (-2.18) |

TABLE 2. Error Norms for  $n$ -Point Radau (right endpoint)

| $h$ | $n = 2$                       | $n = 3$                       | $n = 4$                       | $n = 5$                        | $n = 6$                         |
|-----|-------------------------------|-------------------------------|-------------------------------|--------------------------------|---------------------------------|
| 1   | nonconv.                      | 3.60(10) <sup>-3</sup>        | 1.75(10) <sup>-4</sup>        | 7.49(10) <sup>-6</sup>         | 3.08(10) <sup>-7</sup>          |
| 1/2 | 1.01(10) <sup>-2</sup>        | 2.14(10) <sup>-4</sup> (4.07) | 3.63(10) <sup>-6</sup> (5.59) | 5.68(10) <sup>-8</sup> (7.04)  | 9.08(10) <sup>-10</sup> (8.41)  |
| 1/3 | 3.33(10) <sup>-3</sup> (2.74) | 3.45(10) <sup>-5</sup> (4.50) | 2.93(10) <sup>-7</sup> (6.20) | 2.38(10) <sup>-9</sup> (7.82)  | 2.09(10) <sup>-11</sup> (9.30)  |
| 1/4 | 1.47(10) <sup>-3</sup> (2.85) | 8.95(10) <sup>-6</sup> (4.69) | 4.54(10) <sup>-8</sup> (6.48) | 2.23(10) <sup>-10</sup> (8.23) | 2.09(10) <sup>-12</sup> (8.00)  |
| 1/5 | 7.70(10) <sup>-4</sup> (2.90) | 3.08(10) <sup>-6</sup> (4.78) | 1.03(10) <sup>-8</sup> (6.64) | 3.35(10) <sup>-11</sup> (8.50) | 1.02(10) <sup>-12</sup> (3.23)  |
| 1/6 | 4.52(10) <sup>-4</sup> (2.92) | 1.27(10) <sup>-6</sup> (4.84) | 3.03(10) <sup>-9</sup> (6.73) | 6.03(10) <sup>-12</sup> (9.40) | 1.47(10) <sup>-12</sup> (-2.03) |

TABLE 3. Error Norms for  $n$ -Point Lobatto

| $h$ | $n = 2$               | $n = 3$               | $n = 4$               | $n = 5$                | $n = 6$                 |
|-----|-----------------------|-----------------------|-----------------------|------------------------|-------------------------|
| 1   | $2.68(10)^{-1}$       | $1.59(10)^{-2}$       | $1.11(10)^{-3}$       | $5.78(10)^{-5}$        | $2.62(10)^{-6}$         |
| 1/2 | $5.24(10)^{-2}(2.36)$ | $1.73(10)^{-3}(3.20)$ | $4.00(10)^{-5}(4.79)$ | $7.30(10)^{-7}(6.31)$  | $1.22(10)^{-8}(7.75)$   |
| 1/3 | $2.32(10)^{-2}(2.01)$ | $4.00(10)^{-4}(3.62)$ | $4.55(10)^{-6}(5.36)$ | $4.21(10)^{-8}(7.04)$  | $3.68(10)^{-10}(8.63)$  |
| 1/4 | $1.31(10)^{-2}(1.99)$ | $1.35(10)^{-4}(3.78)$ | $9.08(10)^{-7}(5.61)$ | $5.04(10)^{-9}(7.38)$  | $2.63(10)^{-11}(9.16)$  |
| 1/5 | $8.37(10)^{-3}(1.99)$ | $5.71(10)^{-5}(3.86)$ | $2.52(10)^{-7}(5.74)$ | $9.31(10)^{-10}(7.57)$ | $2.84(10)^{-12}(9.98)$  |
| 1/6 | $5.82(10)^{-3}(1.99)$ | $2.80(10)^{-5}(3.90)$ | $8.75(10)^{-8}(5.81)$ | $2.30(10)^{-10}(7.67)$ | $3.27(10)^{-13}(11.86)$ |

TABLE 4. Error Norms for  $n$ -Point Equal Weight Chebyshev

| $h$ | $n = 2$               | $n = 3$               | $n = 4$               | $n = 5$               | $n = 6$                |
|-----|-----------------------|-----------------------|-----------------------|-----------------------|------------------------|
| 1   | $1.47(10)^{-2}$       | $4.46(10)^{-3}$       | $5.78(10)^{-4}$       | $2.47(10)^{-4}$       | $3.97(10)^{-5}$        |
| 1/2 | $1.39(10)^{-3}(3.40)$ | $4.40(10)^{-4}(3.34)$ | $2.15(10)^{-5}(4.75)$ | $9.81(10)^{-6}(4.66)$ | $6.10(10)^{-7}(6.02)$  |
| 1/3 | $3.07(10)^{-4}(3.72)$ | $1.00(10)^{-4}(3.65)$ | $2.48(10)^{-6}(5.32)$ | $1.18(10)^{-6}(5.23)$ | $3.80(10)^{-8}(6.85)$  |
| 1/4 | $1.01(10)^{-4}(3.84)$ | $3.38(10)^{-5}(3.79)$ | $4.98(10)^{-7}(5.58)$ | $2.42(10)^{-7}(5.51)$ | $4.71(10)^{-9}(7.25)$  |
| 1/5 | $4.25(10)^{-5}(3.90)$ | $1.43(10)^{-5}(3.86)$ | $1.39(10)^{-7}(5.72)$ | $6.83(10)^{-8}(5.66)$ | $8.87(10)^{-10}(7.48)$ |
| 1/6 | $2.08(10)^{-5}(3.93)$ | $7.01(10)^{-6}(3.90)$ | $4.83(10)^{-8}(5.80)$ | $2.39(10)^{-8}(5.76)$ | $2.21(10)^{-10}(7.63)$ |

TABLE 5. Error Norms for Newton-Cotes,  $\theta_k = (k - 1)/(n - 1), 1 \leq k \leq n$

| $h$ | $n = 2$               | $n = 3$               | $n = 4$               | $n = 5$               | $n = 6$               |
|-----|-----------------------|-----------------------|-----------------------|-----------------------|-----------------------|
| 1   | $2.68(10)^{-1}$       | $1.59(10)^{-2}$       | $3.58(10)^{-3}$       | $6.48(10)^{-4}$       | $1.90(10)^{-4}$       |
| 1/2 | $5.24(10)^{-2}(2.36)$ | $1.73(10)^{-3}(3.20)$ | $3.57(10)^{-4}(3.33)$ | $2.69(10)^{-5}(4.59)$ | $7.60(10)^{-6}(4.65)$ |
| 1/3 | $2.32(10)^{-2}(2.01)$ | $4.00(10)^{-4}(3.62)$ | $8.26(10)^{-5}(3.61)$ | $3.26(10)^{-6}(5.21)$ | $9.34(10)^{-7}(5.17)$ |
| 1/4 | $1.31(10)^{-2}(1.99)$ | $1.35(10)^{-4}(3.78)$ | $2.81(10)^{-5}(3.75)$ | $6.69(10)^{-7}(5.50)$ | $1.94(10)^{-7}(5.45)$ |
| 1/5 | $8.37(10)^{-3}(1.99)$ | $5.71(10)^{-5}(3.86)$ | $1.19(10)^{-5}(3.83)$ | $1.89(10)^{-7}(5.66)$ | $5.55(10)^{-8}(5.62)$ |
| 1/6 | $5.82(10)^{-3}(1.99)$ | $2.80(10)^{-5}(3.90)$ | $5.89(10)^{-6}(3.88)$ | $6.62(10)^{-8}(5.76)$ | $1.96(10)^{-8}(5.72)$ |

TABLE 6. Error Norms for  $\theta_k = (2k - 1)/2n, 1 \leq k \leq n$

| $h$ | $n = 2$               | $n = 3$               | $n = 4$               | $n = 5$               | $n = 6$               |
|-----|-----------------------|-----------------------|-----------------------|-----------------------|-----------------------|
| 1   | $3.05(10)^{-2}$       | $6.49(10)^{-3}$       | $1.79(10)^{-3}$       | $4.46(10)^{-4}$       | $1.49(10)^{-4}$       |
| 1/2 | $6.99(10)^{-3}(2.13)$ | $6.69(10)^{-4}(3.28)$ | $1.67(10)^{-4}(3.42)$ | $1.88(10)^{-5}(4.57)$ | $5.82(10)^{-6}(4.68)$ |
| 1/3 | $3.00(10)^{-3}(2.08)$ | $1.54(10)^{-4}(3.62)$ | $3.80(10)^{-5}(3.65)$ | $2.29(10)^{-6}(5.19)$ | $7.12(10)^{-7}(5.18)$ |
| 1/4 | $1.67(10)^{-3}(2.04)$ | $5.22(10)^{-5}(3.77)$ | $1.28(10)^{-5}(3.78)$ | $4.73(10)^{-7}(5.49)$ | $1.48(10)^{-7}(5.46)$ |
| 1/5 | $1.06(10)^{-3}(2.03)$ | $2.21(10)^{-5}(3.85)$ | $5.43(10)^{-6}(3.85)$ | $1.34(10)^{-7}(5.65)$ | $4.22(10)^{-8}(5.62)$ |
| 1/6 | $7.35(10)^{-4}(2.02)$ | $1.09(10)^{-5}(3.89)$ | $2.67(10)^{-6}(3.89)$ | $4.71(10)^{-8}(5.75)$ | $1.49(10)^{-8}(5.72)$ |

**Acknowledgment.** The author wishes to thank Mrs. Sharon L. Daniel for her implementation of the methods of this paper in codes for the CDC 6600.

Applied Mathematics Division, 1722  
 Sandia Laboratories  
 Albuquerque, New Mexico 87115

1. O. AXELSSON, "A class of  $A$ -stable methods," *Nordisk Tidskr. Informationsbehandling (BIT)*, v. 9, 1969, pp. 185–199. MR 40 #8266.
2. J. C. BUTCHER, "Implicit Runge-Kutta processes," *Math. Comp.*, v. 18, 1964, pp. 50–64. MR 28 #2641.
3. J. C. BUTCHER, "Integration processes based on Radau quadrature formulas," *Math. Comp.*, v. 18, 1964, pp. 233–244. MR 29 #2973.
4. F. H. CHIPMAN, *Numerical Solution of Initial Value Problems Using A-stable Runge-Kutta Processes*, Ph.D. Thesis, Univ. of Waterloo, Waterloo, Ontario, 1971.
5. G. J. COOPER, "Interpolation and quadrature methods for ordinary differential equations," *Math. Comp.*, v. 22, 1968, pp. 69–76. MR 36 #7333.
6. G. G. DAHLQUIST, "A special stability problem for linear multistep methods," *Nordisk Tidskr. Informationsbehandling (BIT)*, v. 3, 1963, pp. 27–43. MR 30 #715.
7. B. L. EHLE, "High order  $A$ -stable methods for the numerical solution of systems of D. E.'s," *Nordisk Tidskr. Informationsbehandling (BIT)*, v. 8, 1968, pp. 276–278. MR 39 #1119.
8. B. L. EHLE, *On Padé Approximations to the Exponential Function and A-Stable Methods for the Numerical Solution of Initial Value Problems*, Ph.D. Thesis, Univ. of Waterloo, Waterloo, Ontario, 1969.
9. P. C. HAMMER & J. W. HOLLINGSWORTH, "Trapezoidal methods of approximating solutions of differential equations," *MTAC*, v. 9, 1955, pp. 92–96. MR 17, 302.
10. P. HENRICI, *Discrete Variable Methods in Ordinary Differential Equations*, Wiley, New York, 1962. MR 24 #B1772.
11. F. B. HILDEBRAND, *Introduction to Numerical Analysis*, McGraw-Hill, New York, 1956, MR 17, 788.
12. B. L. HULME, "One-step piecewise polynomial Galerkin methods for initial value problems," *Math. Comp.*, v. 26, 1972, pp. 415–426.
13. K. WRIGHT, "Some relationships between implicit Runge-Kutta, collocation and Lanczos  $\tau$  methods, and their stability properties," *Nordisk Tidskr. Informationsbehandling (BIT)*, v. 10, 1970, pp. 217–227.