

## The Effect of Interpolating the Coefficients in Nonlinear Parabolic Galerkin Procedures\*

By Jim Douglas, Jr. and Todd Dupont

**Abstract.** Error estimates are derived for a class of Galerkin methods for a quasilinear parabolic equation. In these Galerkin methods, both continuous and discrete in time, the nonlinear coefficient in the differential equation is interpolated into a finite-dimensional function space in order to compute the integrals involved. Asymptotic error estimates of optimal order are produced.

**Introduction.** In order to use Galerkin methods for parabolic problems, it is necessary to compute large numbers of integrals involving the coefficients in the differential equation. An efficient and practically successful method of approximating these integrals is to interpolate or project the coefficients and evaluate the integrals by formula. It is possible to show, for a rather general collection of approximation schemes, that the resulting approximate solution is essentially as good as if the integrals had been evaluated exactly.

These procedures are particularly useful on nonlinear parabolic problems in which we have used a Galerkin-type procedure in the space variables and have discretized the time variable. For these procedures, it is necessary to reform the matrices at every time step; the matrix elements are the integrals referred to above. The effect of this is that much of the computation time is spent forming matrices. Hence, economies in the formation of the matrices have a very important effect on the total cost of the computation.

We present here several error estimates for approximations of the solution of a particular nonlinear parabolic problem. In the process of proving these estimates, we develop some approximation theory which may be useful in producing similar estimates for other problems.

In Section 1, we illustrate how to handle a very simple, specific example. In Section 2, we define the principal differential problem and present several error estimates under abstract hypotheses on the approximation scheme to be used for the coefficients. In Section 3, we develop examples of function spaces and interpolation methods which satisfy the abstract hypotheses of Section 2. Finally, in Section 4, we state some specific applications of the results of Sections 2 and 3.

---

Received November 9, 1973.

AMS (MOS) subject classifications (1970). Primary 65N30; Secondary 35K60.

\*This research was partially supported by the National Science Foundation Grant NSF GP-42228.

1. An Example. Consider the problem

$$(1.1) \quad \begin{aligned} \partial u / \partial t - (\partial / \partial x)(a(x, u) \partial u / \partial x) &= 0, \quad (x, t) \in I \times (0, T], \\ (\partial u / \partial x)(0, t) &= (\partial u / \partial x)(1, t) = 0, \quad 0 < t < T, \\ u(x, 0) &= u_0(x), \quad x \in I, \end{aligned}$$

where  $I = (0, 1)$ . Assume that there are positive constants  $\alpha_0$  and  $\alpha_1$  such that for any  $(x, r) \in I \times \mathbb{R}$ ,  $0 < \alpha_0 \leq a(x, r) \leq \alpha_1$ . Let  $\delta = \{x_j\}_{j=0}^N$ ,  $0 = x_0 < x_1 < \dots < x_N = 1$ ,  $h_j = x_j - x_{j-1}$ ,  $I_j = (x_{j-1}, x_j)$ , and  $h = \max_{1 \leq j \leq N} h_j$ . Take  $M = M_1(3, \delta) = \{v \in C^1(\bar{I}) : v \text{ is a cubic polynomial on } I_j, j = 1, \dots, N\}$ ; i.e.,  $M$  is the Hermite piecewise cubic polynomial space over the partition  $\delta$ . Also, take  $\sigma = \{t_k\}_{k=0}^M$ ,  $0 = t_0 < t_1 < \dots < t_M = T$ ,  $\Delta t_k = t_k - t_{k-1}$ ,  $\Delta t = \max_{1 \leq k \leq M} \Delta t_k$ . The sequence  $\{U_n\}_{n=0}^M$  in  $M$  will be taken to approximate  $u$  at times  $t_n$ ,  $n = 0, \dots, M$ . First, choose  $U_0 \in M$  such that  $u_0 - U_0$  is small;  $U_0$  can be, for example, the  $L^2$  or  $H^1$  projection of  $u_0$  or its Hermite cubic interpolant. The successive  $U_n$ 's are defined by the relations

$$(1.2) \quad (\partial_t U_{n+1/2}, V) + \left( \tilde{a}(x, E_{n+1/2}) \frac{\partial U_{n+1/2}}{\partial x}, \frac{d}{dx} V \right) = 0, \quad V \in M,$$

where  $(f, g) = \int_I f g dx$ ,  $U_{n+1/2} = \frac{1}{2}(U_{n+1} + U_n)$  and  $\partial_t U_{n+1/2} = (U_{n+1} - U_n)/\Delta t_{n+1}$ . The function  $E_{n+1/2}$  is a prediction of  $U_{n+1/2}$  and is given by

$$(1.3) \quad E_{n+1/2} = \begin{cases} U_{n-1/2} + \left[ \frac{\Delta t_{n+1} + \Delta t_n}{\Delta t_n + \Delta t_{n-1}} \right] (U_{n-1/2} - U_{n-3/2}), & n \geq 2, \\ U_1 + (\Delta t_2 / 2 \Delta t_1)(U_1 - U_0), & n = 1, \\ \frac{1}{2}(Y_1 + U_0), & n = 0, \end{cases}$$

where  $Y_1 \in M$  satisfies

$$(1.4) \quad ((Y_1 - U_0)/\Delta t_1, V) + \frac{1}{2} \left( \tilde{a}(x, U_0) \frac{\partial}{\partial x} (Y_1 + U_0), \frac{d}{dx} V \right) = 0, \quad V \in M.$$

The function  $\tilde{a}(x, Z)$ , for any  $Z \in M$ , is to be an approximation of  $a(x, Z)$ ; the detailed construction of  $\tilde{a}$  will be discussed later.

A convenient basis for  $M$  is the set  $\{V_i\}_{i=0}^{2N+1}$ , chosen so that for all  $i$  and  $j$  between 0 and  $N$

$$(1.5) \quad V_{2i}(x_j) = V'_{2i+1}(x_j) = \delta_{ij}, \quad V'_{2i}(x_j) = V_{2i+1}(x_j) = 0,$$

where  $\delta_{ij}$  is the Kronecker delta. The functions  $V_{2i}$  and  $V_{2i+1}$  are the "value" and "slope" functions at the knot  $x_i$ . Note that, with this basis, it is easy to construct the Hermite cubic interpolant of a differentiable function; we could, for example, define

$$(1.6) \quad U_0(x) = \sum_{i=0}^N (u_0(x_i) V_{2i}(x) + u'_0(x_i) V_{2i+1}(x)).$$

Equation (1.2) is equivalent to

$$(1.7) \quad (C + \Delta t_{n+1} A_{n+1/2}/2)(\gamma_{n+1} - \gamma_n) = -\Delta t_{n+1} A_{n+1/2} \gamma_n,$$

where

$$(1.8) \quad \gamma_n = (\gamma_{n,0}, \gamma_{n,1}, \dots, \gamma_{n,2N+1})^T,$$

$$U_n = \sum_{j=0}^{2N+1} \gamma_{n,j} V_j, \quad C = (c_{ij}) = ((V_j, V_i)),$$

$$(1.9) \quad A_{n+1/2} = (a_{n+1/2,i,j}) = ((\tilde{a}(x, E_{n+1/2}) V_j', V_i')).$$

The  $(2N+2) \times (2N+2)$  matrix  $C$  can be written as

$$(1.10) \quad C = \begin{pmatrix} D_0 & F_1 & & & 0 \\ & F_1^T & D_1 & & \\ & & \ddots & \ddots & \\ & & & F_N & \\ 0 & & & & F_N^T & D_N \end{pmatrix} = \text{tridiag}\{F_i^T, D_i, F_{i+1}\},$$

where  $D_i = DL_i + DR_i$  and, using  $h_0 = h_{N+1} = 0$ ,

$$(1.11) \quad DL_i = (2520)^{-1} \begin{pmatrix} 936h_i & -132h_i^2 \\ -132h_i^2 & 24h_i^3 \end{pmatrix}, \quad DR_i = (2520)^{-1} \begin{pmatrix} 936h_{i+1} & 132h_{i+1}^2 \\ 132h_{i+1}^2 & 24h_{i+1}^3 \end{pmatrix},$$

$$F_i = (2520)^{-1} \begin{pmatrix} 324h_i & -78h_i^2 \\ 78h_i^2 & -18h_i^3 \end{pmatrix}.$$

The matrix  $A = A_{n+1/2}$  can be written as

$$(1.12) \quad A = \begin{pmatrix} G_0 & H_1 & & & 0 \\ & H_1^T & G_1 & & \\ & & \ddots & \ddots & \\ & & & H_N & \\ 0 & & & & H_N^T & G_N \end{pmatrix} = \text{tridiag}\{H_i^T, G_i, H_{i+1}\},$$

where  $G_i = GL_i + GR_i$  with  $GL_0 = GR_N = 0$ ; the formulas for  $GL_i$ ,  $GR_i$ , and  $H_i$  will be given after the description of  $\tilde{a}$ .

For  $W \in M$ , we shall define  $\tilde{a}(x, W)$  by first defining  $\hat{a}(x, W)$  and then modifying

$\hat{a}(x, W)$ , if necessary, to insure that  $\tilde{a}(x, W)$  is bounded above and below by  $3\alpha_1/2$  and  $\alpha_0/2$ , respectively. The functions  $\hat{a}$  and  $\tilde{a}$  will be cubic polynomials on each  $I_j$ . The function  $\hat{a}(x, W)$  can be defined by taking

$$(1.13) \quad VL_j = a(x_{j-1}, W(x_{j-1})), \quad VR_j = a(x_j, W(x_j)),$$

$$SL_j = (d/dx)(a(x, W))(x_{j-1}), \quad SR_j = (d/dx)(a(x, W))(x_j),$$

where these four numbers give the values and slopes of  $\hat{a}(x, W)$  at the left and right ends of  $I_j$ ; in this case,  $\hat{a}$  is the Hermite cubic interpolant of  $a(x, W)$ . Another reasonable choice for  $\hat{a}$  is to take  $VL_j$  and  $VR_j$  as above and choose  $SL_j$  and  $SR_j$  such that  $\hat{a}(x, W)$  interpolates  $a(x, W)$  at the two points  $\frac{1}{2}(x_{j-1} + x_j) \pm \theta h_j$ , where  $0 < \theta < \frac{1}{2}$ ; in this case,  $\hat{a}(x, W)$  is just the cubic Lagrange interpolant of  $a(x, W)$  using the four points  $x_{j-1}$ ,  $\frac{1}{2}(x_{j-1} + x_j) \pm \theta h_j$ ,  $x_j$ . This second technique gives a better approximation to  $a(x, W)$  and can be more or less work to produce than the Hermite interpolant, depending on the form of the function  $a(x, r)$ . There are many other useful ways in which we could choose  $\hat{a}$ , and some of them will be discussed in Section 3.

The function  $\tilde{a}(x, W)$  can be obtained from  $\hat{a}(x, W)$  as follows. Let  $\alpha = \min\{VL_j, VR_j\} \geq \alpha_0$ . If  $|SL_j| > \gamma = \alpha/h_j$ , then change  $SL_j$  to be such that  $|SL_j| = \gamma$  and its sign is unchanged. Similarly, if  $|SR_j| > \gamma$ , multiply it by  $\gamma/|SR_j|$ . With  $\tilde{a}$  defined in this fashion, we know that

$$\alpha_0/2 \leq \alpha/2 \leq \tilde{a}(x, W) \leq \alpha/2 + \max\{VL_j, VR_j\} \leq 3\alpha_1/2, \quad x \in I_j.$$

Thus, to produce  $\tilde{a}$  from  $\hat{a}$ , we simply check the size of  $SL_j$  and  $SR_j$  and replace them if they are too large.

With  $\tilde{a}(x, W)$  defined on  $I_j$  by  $VL, VR, SL, SR$ , we can give the formulas for  $H_j$ ,  $GR_{j-1}$ , and  $GL_j$ . Let

$$(1.14) \quad H_j = (2520)^{-1} \begin{pmatrix} p_1 & p_2 \\ p_3 & p_4 \end{pmatrix}, \quad GL_i = (2520)^{-1} \begin{pmatrix} q_1 & q_2 \\ q_3 & q_4 \end{pmatrix},$$

$$GR_{j-1} = (2520)^{-1} \begin{pmatrix} g_1 & g_2 \\ g_3 & g_4 \end{pmatrix}.$$

Then

$$(1.15) \quad \begin{aligned} p_1 &= -(1512(VL + VR) + 324(SL - SR))h_j^{-1}, \\ p_2 &= 288VL - 36VR + 54SL - 18SR, \\ p_3 &= 36VL - 288VR - 18SL + 54SR, \\ p_4 &= -(42(VL + VR) + 3(SL - SR))h_j, \\ q_1 &= -p_1, \quad q_2 = q_3 = -p_2, \\ q_4 &= (78VL + 258VR + 15SL - 21SR)h_j, \\ g_1 &= -p_1, \quad g_2 = g_3 = -p_3, \\ g_4 &= (258VL + 78VR + 21SL - 15SR)h_j. \end{aligned}$$

It is possible to show that if  $a(x, u)$  and  $u$  are sufficiently smooth and  $\Delta t_{k+1}/\Delta t_k$  is bounded, then there is a  $C$ , independent of the  $h_i$ 's and  $\Delta t_k$ 's such that

$$(1.16) \quad \max_{0 \leq n \leq M} \|U_n - u(\cdot, t_n)\|_{L^2(I)} \leq C(h^4 + (\Delta t)^2).$$

This estimate is of the same form as we would expect if the integrals had been evaluated exactly.

**2. Procedures and Estimates.** In this section, we shall demonstrate several error estimates for approximate solutions of the following parabolic problem:

$$\frac{\partial u}{\partial t} - \nabla \cdot (a(x, u) \nabla u) = \frac{\partial u}{\partial t} - \sum_{j=1}^p \frac{\partial}{\partial x_j} \left( a(x, u) \frac{\partial u}{\partial x_j} \right) = 0, \quad (x, t) \in \Omega \times (0, T],$$

$$(2.1) \quad \frac{\partial u}{\partial \nu}(x, t) = 0, \quad (x, t) \in \partial\Omega \times (0, T],$$

$$u(x, 0) = u_0(x), \quad x \in \Omega,$$

where  $\Omega$  is a bounded domain in  $\mathbf{R}^p$  with boundary  $\partial\Omega$ ,  $p \leq 3$ , and  $\partial/\partial\nu$  is outward normal differentiation on  $\partial\Omega$ , and the function  $a(x, r)$  is such that there are positive constants  $\alpha_0$  and  $\alpha_1$  such that for all  $(x, r) \in \Omega \times \mathbf{R}$ ,  $0 < \alpha_0 \leq a(x, r) \leq \alpha_1$ . The error estimates presented here are abstractions and slight improvements of those of the authors [6] and of Wheeler [11].

We shall assume that  $u(x, t)$  is a solution of the weak form of (2.1) [8] in the sense that for each time

$$(2.1') \quad (\partial u / \partial t, V) + (a(u) \nabla u, \nabla V) = 0,$$

for all  $v$  in the Sobolev space  $H^1(\Omega)$ , where  $(f, g)$  is the  $L^2(\Omega)$  inner product, and we have suppressed writing the  $x$  argument of  $a(x, u)$ . If  $u$  and  $\partial\Omega$  are smooth, then  $u$  is a solution of (2.1) if and only if it is a solution of (2.1'), provided of course that the initial values coincide.

In each of the procedures that we consider for approximating the solution of (2.1), we shall use two spaces  $M$  and  $N$  of functions defined on  $\Omega$ . The space  $M$  will be a finite-dimensional subspace of  $H^1(\Omega)$ , and the approximate solution will be an element of  $M$  for each time. The space  $N$  will be a subspace of  $L^\infty(\Omega)$ , and an element of  $N$  will be used to approximate the coefficient  $a$  at each time. Assume that there is a map  $\tilde{a}: M \rightarrow N$ , such that for all  $W \in M$

$$(2.2) \quad \alpha_0/2 \leq \tilde{a}(W)(x) \leq 3\alpha_1/2.$$

Take  $U_0 \in M$  to be an approximation of  $u_0$ . The continuous-time Galerkin approximation of the solution of (2.1) is a differentiable map  $U: [0, T] \rightarrow M$  such that

$$(2.3) \quad (\partial U / \partial t, V) + (\tilde{a}(U) \nabla U, \nabla V) = 0, \quad V \in M, 0 < t \leq T,$$

$$U(0) = U_0.$$

A discrete-time Galerkin approximation of the solution of (2.1) is a sequence  $\{U_n\}_{n=0}^M$  in  $M$ , where  $U_n$  is to approximate  $u(\cdot, t_n)$  and  $0 = t_0 < t_1 < \dots < t_M = T$ . The sequence  $\{U_n\}_{n=0}^M$  will then be required to satisfy an approximation to (2.3) of the form

$$(2.4) \quad (\partial_t U_{n+1/2}, V) + (\tilde{a}(E_{n+1/2}(U)) \nabla U_{n+1/2}, \nabla V) = 0, \quad V \in M, 0 \leq n < M,$$

where  $\Delta t_n = t_n - t_{n-1}$ ,  $\partial_t r_{n+1/2} = (r_{n+1} - r_n)(\Delta t_{n+1})^{-1}$ ,  $r_{n+1/2} = \frac{1}{2}(r_{n+1} + r_n)$  and  $E_{n+1/2}(U)$  is an approximation of  $U_{n+1/2}$ . The function  $E_{n+1/2}(U)$  will be taken to depend on a certain number of  $U_k$ 's with  $k \leq n+1$ . Note that if we take

$$(2.5) \quad E_{n+1/2}(U) = U_{n+1/2},$$

then (2.4) is the Crank-Nicolson approximation to (2.3) [6]. If we chose

$$(2.6) \quad E_{1/2}(U) = U_{1/2}, \quad E_{n+1/2}(U) = U_n + \frac{\Delta t_{n-1}}{2\Delta t_n}(U_n - U_{n-1}), \quad n \geq 1,$$

or use (2.6) for  $n = 0$  and 1 and

$$(2.7) \quad E_{n+1/2} = U_{n-1/2} + \frac{\Delta t_{n+1} + \Delta t_n}{\Delta t_n + \Delta t_{n-1}}(U_{n-1/2} - U_{n-3/2}), \quad n \geq 2,$$

then for  $n \geq 1$ , (2.4) defines  $\{U_n\}$  in terms of a sequence of *linear* algebraic equations that are second-order correct (in  $\Delta t_n$ ) approximations of (2.3). In practice, we might replace the first step of (2.4) by a predictor-corrector procedure and employ (2.6) or (2.7) thereafter, as was indicated by the example of Section 1. This additional complication can be treated by arguments similar to those in [6], [11], but will not be discussed here.

In both the continuous and discrete time cases, we shall present bounds for the error in the "natural" or "energy" norm and in the  $L^2(\Omega)$ -norm.

For integer  $s \geq 0$ , use the norm on the Sobolev space  $H^s(\Omega)$  given by

$$(2.8) \quad \|\phi\|_s^2 = \sum_{|\alpha| \leq s} \|D^\alpha \phi\|^2,$$

where  $\|\phi\|^2 = (\phi, \phi)$ ,  $D^\alpha = (\partial/\partial x_1)^{\alpha_1} \dots (\partial/\partial x_p)^{\alpha_p}$  for  $p$ -tuples of nonnegative integers  $\alpha = (\alpha_1, \dots, \alpha_p)$ , and  $|\alpha| = \alpha_1 + \dots + \alpha_p$ . Extend  $(\cdot, \cdot)$  to give the duality between  $H^1(\Omega)$  and  $(H^1(\Omega))'$ . Let  $H^{-1}(\Omega) = (H^1(\Omega))'$ ; some authors use  $H^{-1}(\Omega) = (H_0^1(\Omega))'$ , but this is not convenient here. Equip  $H^{-1}(\Omega)$  with the operator norm; i.e., take

$$(2.9) \quad \|\phi\|_{-1} = \sup\{(\phi, \Psi) : \Psi \in H^1(\Omega), \|\Psi\|_1 = 1\}.$$

If  $\phi: [0, T] \rightarrow X$ , where  $X$  is a normed space with norm  $\|\cdot\|_X$ , then we shall use the

following notations:

$$(2.10) \quad \|\phi\|_{L^2(X)}^2 = \int_0^T \|\phi(t)\|_X^2 dt, \quad \|\phi\|_{L^\infty(X)} = \sup_{0 \leq t \leq T} \|\phi(t)\|_X.$$

In the special case,  $X = \mathbf{R}$ , we shall use the usual notation  $\|\phi\|_{L^2(0,T)}$ , and in the case,  $X = H^s(\Omega)$  or  $L^s(\Omega)$ , we shall write  $\|\phi\|_{L^2(H^s)}$ ,  $\|\phi\|_{L^2(L^s)}$ , etc.

We shall assume that the solution  $u$  of (2.1') has a uniformly bounded gradient. Also, we shall assume that the solution  $u$  and the mapping  $\tilde{a}$  are such that there exist a constant  $L$  and a nonnegative function  $\theta(t)$  such that for all  $W \in M$  and all  $t \in [0, T]$

$$(2.11) \quad \|\tilde{a}(W) - a(\cdot, u)\| \leq L\|W - u\| + \theta(t).$$

In the examples which we shall consider in the next section,  $L$  is approximately the size of the Lipschitz constant for  $a(x, r)$  in the variable  $r$ , and  $\theta(t)$  is approximately the size of

$$\inf\{\|a(\cdot, u) - \Psi\| + \|u - \chi\| : \Psi \in N, \chi \in M\}.$$

The relations (2.2) and (2.11) are all that we assume about the approximation process for  $a(x, u)$ ; these two assumptions allow estimates to be derived by methods that are very close to those of [6] and [11]. First, we shall estimate the error of the continuous-time Galerkin approximation in the "natural norm" for this problem.

**THEOREM 2.1.** *There is a constant  $C$ , depending only on  $\alpha_0, \alpha_1, L, \|\nabla u\|_{L^\infty(L^\infty(\Omega))}$ , and  $T$ , such that if  $U$  is the solution of (2.3) and  $u$  is the solution of (2.1'), then*

$$(2.12) \quad \|U - u\|_{L^\infty(L^2(\Omega))} + \|U - u\|_{L^2(H^1(\Omega))} \leq C [\|(U - u)(0)\| + E + \|\theta\|_{L^2(0,T)}],$$

where

$$(2.13) \quad E = \inf \left\{ \|u - Z\|_{L^\infty(L^2)} + \|u - Z\|_{L^2(H^1)} + \left\| \frac{\partial}{\partial t} (u - Z) \right\|_{L^2(H^{-1})} : Z : [0, T] \rightarrow M, \right. \\ \left. Z \text{ continuously differentiable on } [0, T] \right\}.$$

*Proof.* Let  $Z$  be an arbitrary differentiable map of  $[0, T]$  into  $M$ . Then, with  $\eta = u - Z$  and  $\vartheta = U - Z$ , we see from (2.1') and (2.3) that for  $V \in M$

$$(2.14) \quad \left( \frac{\partial}{\partial t} \vartheta, V \right) + (\tilde{a}(U) \nabla \vartheta, \nabla V) \\ = \left( \frac{\partial \eta}{\partial t}, V \right) + (\tilde{a}(U) \nabla \eta + [a(u) - \tilde{a}(U)] \nabla u, \nabla V),$$

where we have suppressed the writing of the  $x$  variable in  $a(x, u)$ . Using  $V = \vartheta$  at each time, we see that

$$\begin{aligned}
(2.15) \quad & \frac{1}{2} \frac{d}{dt} (\|\vartheta(t)\|^2) + \frac{1}{2} \alpha_0 \|\nabla \vartheta\|^2 \\
& \leq \left\| \frac{\partial \eta}{\partial t} \right\|_{-1} \|\vartheta\|_1 + \frac{3}{2} \alpha_1 \|\nabla \eta\| \|\nabla \vartheta\| \\
& \quad + \|\nabla u\|_{L^\infty(\Omega)} [L(\|\vartheta\| + \|\eta\|) + \theta] \|\nabla \vartheta\| \\
& \leq \frac{1}{4} \alpha_0 \|\nabla \vartheta\|^2 + C \left[ \left\| \frac{\partial \eta}{\partial t} \right\|_{-1}^2 + \|\eta\|_1^2 + \|\vartheta\|^2 + \theta^2 \right],
\end{aligned}$$

where  $C$  depends only on  $\alpha_0, \alpha_1, L$  and  $\|\nabla u\|_{L^\infty(\Omega)}$ ; in deriving (2.15), we used the inequality  $cd \leq \epsilon c^2 + d^2/4\epsilon$ , valid for all  $c, d$  and all positive  $\epsilon$ . The estimate (2.15) and Gronwall's Lemma imply that there is a  $C$  depending on the permitted quantities such that

$$\begin{aligned}
(2.16) \quad & \|\vartheta\|_{L^\infty(L^2)} + \|\vartheta\|_{L^2(H^1)} \\
& \leq C \left[ \|\vartheta(0)\| + \|\eta\|_{L^2(H^1)} + \left\| \frac{\partial \eta}{\partial t} \right\|_{L^2(H^{-1})} + \|\theta\|_{L^2(0,T)} \right].
\end{aligned}$$

The triangle inequality and taking the infimum on  $Z$  give the estimate (2.12).

Now assume that if  $V \in M$  then  $\|\nabla V\|_{L^\infty(\Omega)}$  is finite; this is, of course, valid for the piecewise polynomial spaces frequently employed. For each  $t$ , define  $W(t) \in M$  by

$$(2.17) \quad (a(u)\nabla(u - W), \nabla V) + (u - W, V) = 0, \quad V \in M.$$

Thus,  $W: [0, T] \rightarrow M$  is a weighted  $H^1(\Omega)$  projection of the solution for each  $t$ . The next theorem, which says that  $\|u - U\|_{L^\infty(L^2)}$  is about the size of  $\|u - W\|_{L^\infty(L^2)}$ , will be used to derive  $L^2$  error bounds.

**THEOREM 2.2.** *There is a constant  $C$ , depending only on  $\alpha_0, \alpha_1, L, \|\nabla W\|_{L^\infty(L^\infty)}$ , and  $T$ , where  $W$  is defined by (2.17), such that if  $u$  and  $U$  are the solutions of (2.1') and (2.3), respectively, then, with  $\xi = u - U$  and  $\eta = u - W$ ,*

$$(2.18) \quad \|\xi\|_{L^\infty(L^2)} \leq C \left[ \|\xi(0)\| + \|\eta\|_{L^\infty(L^2)} + \left\| \frac{\partial \eta}{\partial t} \right\|_{L^2(H^{-1})} + \|\theta\|_{L^2(0,T)} \right].$$

*Proof.* Note that, from the definition of  $W$  and  $\eta$ ,

$$(2.19) \quad (a(u)\nabla \eta, \nabla V) = -(\eta, V), \quad V \in M.$$

Letting  $\vartheta = U - W$  and using (2.19), we see that for  $V \in M$ ,

$$\begin{aligned}
(2.20) \quad & (\partial \vartheta / \partial t, V) + (\tilde{a}(U)\nabla \vartheta, \nabla V) = (\partial \eta / \partial t, V) + (a(u)\nabla u - \tilde{a}(U)\nabla W, \nabla V) \\
& = (\partial \eta / \partial t, V) + ([a(u) - \tilde{a}(U)]\nabla W, \nabla V) - (\eta, V).
\end{aligned}$$

With  $V = \vartheta$ , (2.20) implies that



$$\begin{aligned}
(2.21) \quad & \frac{1}{2} \frac{d}{dt} (\|\vartheta\|^2) + \frac{1}{2} \alpha_0 \|\nabla \vartheta\|^2 \\
& \leq \left\| \frac{\partial \eta}{\partial t} \right\|_{-1} \|\vartheta\|_1 + \|\nabla W\|_{L^\infty(\Omega)} [L(\|\eta\| + \|\vartheta\|) + \theta] \|\nabla \vartheta\| + \|\eta\| \|\vartheta\| \\
& \leq \frac{1}{4} \alpha_0 \|\nabla \vartheta\|^2 + C \left[ \left\| \frac{\partial \eta}{\partial t} \right\|_{-1}^2 + \|\eta\|^2 + \|\vartheta\|^2 + \theta^2 \right],
\end{aligned}$$

where  $C$  depends only on  $\alpha_0$ ,  $L$  and  $\|\nabla W\|_{L^\infty(\Omega)}$ . Gronwall's Lemma and (2.21) imply that

$$(2.22) \quad \|\vartheta\|_{L^\infty(L^2)} \leq C \left[ \|\vartheta(0)\| + \|\eta\|_{L^2(L^2)} + \left\| \frac{\partial \eta}{\partial t} \right\|_{L^2(H^{-1})} + \|\theta\|_{L^2(0,T)} \right].$$

The triangle inequality then implies the conclusion (2.18).

Asymptotic error estimates are easily obtained from Theorem 2.2 provided one can demonstrate uniform boundedness of  $\nabla W$ ; one way this can be done is by using so-called inverse assumptions [11].

In order to derive estimates for the discrete time procedures, we need some assumptions on the functions  $E_{n+1/2}(U)$ . Each  $E_{n+1/2}(U)$  will be assumed to be defined by a function, which we shall also call  $E_{n+1/2}$ , of  $U_{n+1}$ ,  $U_n$ ,  $U_{n-1}$ ,  $U_{n-2}$ ; in the cases  $E_{1/2}$  and  $E_{3/2}$ , we of course assume there is no dependence on  $U_{-1}$  and  $U_{-2}$ . It will be assumed that the rules which define the  $E_{n+1/2}$ 's are such that there is a constant  $K_1$  such that, for any permitted partition  $\{t_n\}_{n=0}^M$  of  $[0, T]$  and any  $V_1, \dots, V_4, Z_1, \dots, Z_4 \in L^2(\Omega)$ ,

$$(2.23) \quad \|E_{n+1/2}(V_1, \dots, V_4) - E_{n+1/2}(Z_1, \dots, Z_4)\|^2 \leq K_1 \sum_{l=1}^4 \|V_l - Z_l\|^2.$$

The functions  $E_{n+1/2}$  will also be assumed to be second-order correct in the sense that there is a constant  $K_2$  independent of the partition  $\{t_n\}_{n=0}^M$  such that, if  $\|\partial^2 w / \partial t^2\|_{L^2(L^2(\Omega))} < \infty$ ,

$$\begin{aligned}
(2.24) \quad & \|w(t_{n+1/2}) - E_{n+1/2}(w_{n+1}, w_n, w_{n-1}, w_{n-2})\|^2 \\
& \leq K_2 (\Delta t)^3 \int_{t_{n-2}^*}^{t_{n+1}} \left\| \frac{\partial^2 w}{\partial t^2} \right\|^2 dt,
\end{aligned}$$

where  $t_{n-2}^* = \max\{0, t_{n-2}\}$  and  $\Delta t = \max\{\Delta t_{n+1}, \Delta t_n, \Delta t_{n-1}\}$ . Note that (2.23) is always satisfied by  $E_{n+1/2}$  defined by (2.5) with  $K_1 = 1/2$ , but that we need  $\Delta t_{n+1}/\Delta t_n$  bounded to get (2.23) for  $E_{n+1/2}$  defined by (2.6) or (2.7). The relation (2.24) is satisfied by each of the examples of  $E_{n+1/2}$ .

In order to simplify the analysis, we shall consider only the case of uniform time steps; i.e., take  $t_n = n\Delta t$ , where  $\Delta t = T/M$ . The function  $a(x, r)$  will be assumed to

have uniformly bounded first and second derivatives with respect to  $r$ . Take  $L$  of (2.11) large enough to bound  $|\partial a(x, r)/\partial r|$ , and let  $L_1$  bound  $|\partial^2 a(x, r)/\partial r^2|$ . In addition, we shall assume that the solution  $u$  of the differential problem is such that

$$(2.25) \quad \left\| \frac{\partial^3 u}{\partial t^3} \right\|_{L^2(H^{-1})} + \left\| \frac{\partial^2 u}{\partial t^2} \right\|_{L^2(H^1)} + \left\| \frac{\partial u}{\partial t} \right\|_{L^\infty(L^\infty)} + \left\| \frac{\partial u}{\partial t} \right\|_{L^\infty(H^1)} < \infty.$$

Adopt the following discrete analogues of the notations defined in (2.10); if  $\phi: \{t_n\}_{n=0}^M \rightarrow X$ ,  $X$  a normed space with norm  $\|\cdot\|_X$ , then

$$(2.26) \quad \begin{aligned} \|\phi\|_{L_{\Delta t}^\infty(X)} &= \max_{0 \leq n \leq M} \|\phi_n\|_X, \\ \|\phi\|_{L_{\Delta t}^2(X)}^2 &= \sum_{n=0}^{M-1} \|\phi_{n+1/2}\|_X^2 \Delta t. \end{aligned}$$

In the case  $\phi: [0, T] \rightarrow \mathbf{R}$ , use the notation

$$\|\phi\|_{L_{\Delta t}^2}^2 = \sum_{n=0}^{M-1} |\phi(t_{n+1/2})|^2 \Delta t.$$

**THEOREM 2.3.** *There are positive constants  $C$  and  $\tau_0$ , depending only on  $\alpha_0, \alpha_1, L, \|\nabla u\|_{L^\infty(L^\infty)}, K_1, K_2$ , and  $T$ , such that if  $u$  is the solution of (2.1'),  $\{U_n\}_{n=0}^M$  is the solution of (2.4), and  $\xi_n = u_n - U_n$ , then, for  $0 < \Delta t < \tau_0$ ,*

$$(2.27) \quad \|\xi\|_{L_{\Delta t}^\infty(L^2)} + \|\xi\|_{L_{\Delta t}^2(H^1)} \leq C\{\|\xi_0\| + E + (\Delta t)^2 \gamma + \|\theta\|_{L_{\Delta t}^2}\},$$

where

$$(2.28) \quad \begin{aligned} E &= \inf \{ \|\eta\|_{L_{\Delta t}^\infty(L^2)} + \|\eta\|_{L_{\Delta t}^2(H^1)} + \|\partial \eta / \partial t\|_{L^2(H^{-1})} : \eta = u - Z, \\ Z: [0, T] &\rightarrow \mathcal{M} \text{ a continuously differentiable map} \}, \\ \gamma &= \|\partial^2 u / \partial t^2\|_{L^2(L^2)} + \|\partial^3 u / \partial t^3\|_{L^2(H^{-1})}. \end{aligned}$$

*Proof.* Let  $Z: [0, T] \rightarrow \mathcal{M}$  be an arbitrary continuously differentiable map of  $[0, T]$  into  $\mathcal{M}$ . Using  $\eta = u - Z$  and  $\vartheta = U - Z$ , we can see that for  $V \in \mathcal{M}$

$$(2.29) \quad \begin{aligned} (\partial_t \vartheta_{n+1/2}, V) + (\tilde{a}(E_{n+1/2}) \nabla \vartheta_{n+1/2}, \nabla V) &= (\partial_t \eta_{n+1/2} + \rho_n, V) \\ &+ (g_n, \nabla V) + (\tilde{a}(E_{n+1/2}) - a(u(t_{n+1/2}))) \nabla u_{n+1/2}, \nabla V) \\ &+ (\tilde{a}(E_{n+1/2}) \nabla \eta_{n+1/2}, \nabla V), \end{aligned}$$

where

$$\begin{aligned} E_{n+1/2} &= E_{n+1/2}(U), \\ \rho_n &= \frac{\partial u}{\partial t}(t_{n+1/2}) - \partial_t u_{n+1/2} \\ &= \frac{-1}{8\Delta t} \int_{t_n}^{t_{n+1}} (\Delta t - 2|\tau - t_{n+1/2}|)^2 \frac{\partial^3 u}{\partial t^3}(\tau) d\tau, \end{aligned}$$

$$\begin{aligned}
(2.30) \quad g_n &= a(u(t_{n+1/2})) \nabla [u(t_{n+1/2}) - u_{n+1/2}] \\
&= \frac{-1}{4} a(u(t_{n+1/2})) \nabla \int_{t_n}^{t_{n+1}} (\Delta t - 2|\tau - t_{n+1/2}|) \frac{\partial^2 u}{\partial t^2}(\tau) d\tau.
\end{aligned}$$

Taking  $V = \vartheta_{n+1/2}$  in (2.29) gives

$$\begin{aligned}
(2.31) \quad & \frac{1}{2\Delta t} [\|\vartheta_{n+1}\|^2 - \|\vartheta_n\|^2] + \frac{1}{2} \alpha_0 \|\nabla \vartheta_{n+1/2}\|^2 \\
& \leq (\|\partial_t \eta_{n+1/2}\|_{-1} + \|\rho_n\|_{-1}) \|\vartheta_{n+1/2}\|_1 + \|g_n\| \|\nabla \vartheta_{n+1/2}\| \\
& \quad + (L\|E_{n+1/2} - u(t_{n+1/2})\| + \theta(t_{n+1/2})) \|\nabla u_{n+1/2}\|_{L^\infty(\Omega)} \|\nabla \vartheta_{n+1/2}\| \\
& \quad + \frac{3}{2} \alpha_1 \|\nabla \eta_{n+1/2}\| \|\nabla \vartheta_{n+1/2}\| \\
& \leq \frac{1}{4} \alpha_0 \|\nabla \vartheta_{n+1/2}\|^2 + C \left[ \|\partial_t \eta_{n+1/2}\|_{-1}^2 + \|\rho_n\|_{-1}^2 + \|g_n\|^2 \right. \\
& \quad \left. + \|E_{n+1/2}(u_{n+1}, u_n, u_{n-1}, u_{n-2}) - u(t_{n+1/2})\|^2 \right. \\
& \quad \left. + \theta^2(t_{n+1/2}) + \sum_{l=n-2}^{n+1} (\|\eta_l\|^2 + \|\vartheta_l\|^2) + \|\nabla \eta_{n+1/2}\|^2 \right],
\end{aligned}$$

where  $C$  depends only on  $\alpha_0, \alpha_1, L, \|\nabla u\|_{L^\infty(L^\infty)}$ , and  $K_1$ ; in the cases  $n = 0$  or  $1$ , the  $\vartheta$  and  $\eta$  terms with negative indexes are omitted. It is easily seen that

$$(2.32) \quad \|\rho_n\|_{-1}^2 \leq (\Delta t)^3 (320)^{-1} \int_{t_n}^{t_{n+1}} \left\| \frac{\partial^3 u}{\partial t^3} \right\|_{-1}^2 d\tau,$$

$$\|g_n\|^2 \leq (\Delta t)^3 (\alpha_1)^2 (48)^{-1} \int_{t_n}^{t_{n+1}} \left\| \frac{\partial^2 u}{\partial t^2} \right\|_1^2 d\tau.$$

The discrete analogue of Gronwall's inequality implies that there exist  $\tau_0 > 0$  and  $C$ , depending only on the permitted quantities, such that for  $0 < \Delta t < \tau_0$ ,

$$\begin{aligned}
(2.33) \quad & \|\vartheta\|_{L_{\Delta t}^\infty(L^2)} + \|\vartheta\|_{L_{\Delta t}^2(H^1)} \\
& \leq C [\|\vartheta_0\| + \|\eta\|_{L^\infty(L^2)} + \|\eta\|_{L_{\Delta t}^2(H^1)} \\
& \quad + \|\partial_t \eta\|_{L_{\Delta t}^2(H^{-1})} + (\Delta t)^2 \gamma + \|\theta\|_{L_{\Delta t}^2(0,T)}].
\end{aligned}$$

Note that since  $\partial_t \eta_{n+1/2}$  is the average of  $\partial \eta / \partial t$  on  $[t_n, t_{n+1}]$ ,

$$(2.34) \quad \|\partial_t \eta\|_{L_{\Delta t}^2(H^{-1})} \leq \|\partial \eta / \partial t\|_{L^2(H^{-1})}.$$

Using (2.34) in (2.33), applying the triangle inequality, and taking the infimum over all possible  $Z$ 's gives the conclusion.

In order to get results that will be used to give  $L^2$ -norm estimates, we shall again compare the parabolic approximation with the solution of the elliptic problem (2.17).

**THEOREM 2.4.** *There are positive constants  $\tau_0$  and  $C$ , depending only on  $\alpha_0, \alpha_1, L, \|W\|_{L^\infty(L^\infty)}, K_1$  and  $T$  such that if  $u$  is the solution of (2.1'),  $\{U_n\}_{n=0}^M$  is the solution of (2.4),  $W: [0, T] \rightarrow M$  is given by (2.17),  $\eta = u - W$ , and  $\xi_n = u_n - U_n$ , then, for  $0 < \Delta t < \tau_0$ ,*

$$(2.35) \quad \|\xi\|_{L_{\Delta t}^\infty(L^2)} \leq C \left[ \|\xi_0\| + \|\eta\|_{L^\infty(L^2)} + \left\| \frac{\partial \eta}{\partial t} \right\|_{L^2(H^{-1})} + \|\theta\|_{L_{\Delta t}^2} + (\Delta t)^2 \gamma \right],$$

where  $\gamma$  depends on  $L_1, K_2$ , the norms of  $u$  in (2.25), and the parameters listed above.

*Proof.* Let  $\vartheta = U - W$ . Then, from (2.4) and the average of (2.1') at  $t_n$  and  $t_{n+1}$ , we see that for  $V \in M$

$$(2.36) \quad \begin{aligned} & (\partial_t \vartheta_{n+1/2}, V) + (\tilde{a}(E_{n+1/2}) \nabla \vartheta_{n+1/2}, \nabla V) \\ &= (\partial_t \eta_{n+1/2} + \rho_n, V) + ((a(u) \nabla \eta)_{n+1/2}, \nabla V) \\ & \quad + ((a(u) \nabla W)_{n+1/2} - \tilde{a}(E_{n+1/2}) \nabla W_{n+1/2}, \nabla V), \end{aligned}$$

where  $E_{n+1/2} = E_{n+1/2}(U)$  and

$$(2.37) \quad \rho_n = \left( \frac{\partial u}{\partial t} \right)_{n+1/2} - \partial_t u_{n+1/2} = \frac{1}{2\Delta t} \int_{t_n}^{t_{n+1}} (t_{n+1} - \tau)(\tau - t_n) \frac{\partial^3 u}{\partial t^3}(\tau) d\tau.$$

Note that we can replace the term containing  $\nabla \eta$  by using (2.19) just as in the proof of Theorem 2.2. Next note that

$$(2.38) \quad \begin{aligned} & (a(u) \nabla W)_{n+1/2} - \tilde{a}(E_{n+1/2}) \nabla W_{n+1/2} \\ &= (a(u(t_{n+1/2})) - \tilde{a}(E_{n+1/2})) \nabla W_{n+1/2} \\ & \quad + (a(u)_{n+1/2} - a(u(t_{n+1/2}))) \nabla W_{n+1/2} \\ & \quad + \frac{1}{4}(a(u)_{n+1} - a(u)_n) \nabla (W_{n+1} - W_n). \end{aligned}$$

The first term in the right-hand side of (2.38) is treated as in the proof of Theorem 2.3. The second term is bounded as

$$(2.39) \quad \begin{aligned} & \|(a(u)_{n+1/2} - a(u(t_{n+1/2}))) \nabla W_{n+1/2}\|^2 \\ & \leq \|\nabla W\|_{L^\infty(L^\infty)}^2 (\Delta t)^3 (48)^{-1} \int_{t_n}^{t_{n+1}} \left\| \frac{\partial^2}{\partial t^2} (a(u)) \right\|^2 d\tau; \end{aligned}$$

the integral is bounded in terms of  $L, L_1$ , and the norms in (2.25). The last term in (2.38) is bounded as

$$(2.40) \quad \begin{aligned} & \|(a(u)_{n+1} - a(u)_n) \nabla (W_{n+1} - W_n)\|^2 \\ & \leq \frac{(\Delta t)^3}{2} L^2 \left\| \frac{\partial u}{\partial t} \right\|_{L^\infty(L^\infty)}^2 \int_{t_n}^{t_{n+1}} \left\| \nabla \frac{\partial W}{\partial t} \right\|^2 d\tau. \end{aligned}$$

The integral in this term is easily bounded by using the defining relation for  $W$  to see that

$$(2.41) \quad \left\| \frac{\partial \eta}{\partial t} \right\|_1 \leq C \left[ \left\| \frac{\partial u}{\partial t} \right\|_1 + \|\eta\|_1 \right] \leq C \left[ \left\| \frac{\partial u}{\partial t} \right\|_1 + \|u\|_1 \right],$$

where the constant involves  $\alpha_0, \alpha_1, L$  and  $\|\partial u / \partial t\|_{L^\infty(L^\infty)}$ .

The above relations, used with (2.36) with  $V = \mathcal{V}_{n+1/2}$ , give the conclusion just as in the proof of Theorem 2.3.

**3. Construction of  $\tilde{a}(U)$ .** In this section, we shall present several examples of triples  $(M, N, \tilde{a})$  where the map  $\tilde{a}: M \rightarrow N$  satisfies (2.2) and (2.11). We shall be particularly interested in sequences of such triples for which  $L$  of (2.11) is bounded and the norm of  $\theta$  tends to zero. The map  $\tilde{a}$  will be constructed in two steps. We shall first study a map  $\hat{a}: M \rightarrow N$  which satisfies (2.11) but which may fail to satisfy (2.2). The map  $\tilde{a}$  will then be constructed using  $\hat{a}$ . There are two reasons for studying  $\hat{a}$  independently of  $\tilde{a}$ . The map  $\hat{a}$  is useful by itself in situations where (2.2) is not needed; this would be the case for the numerical quadratures of a term  $f(x, t, u)$  added to the right-hand side of (2.1). Also, if certain "inverse assumptions" hold on  $M$  and  $N$ , then  $\tilde{a}(W) = \hat{a}(W)$  for the functions  $W$  in which we are interested, though not necessarily for all  $W \in M$ .

The function  $\hat{a}(W)$  will be the projection into  $N$  of  $a(x, W)$  using an inner product  $(\cdot, \cdot)_p$ . The choice of reasonable inner products is quite large and, by varying  $N$  and  $(\cdot, \cdot)_p$ , we can, for example, make  $\hat{a}(W)$  a pointwise least-squares fit, a Lagrange or Hermite interpolant or a smooth-spline-type interpolant. The interpolation schemes so constructed can be local or global in character, as is convenient.

Suppose that an inner product  $(\cdot, \cdot)_p$  is defined on a class of functions which contains  $M \cup N \cup \{a(\cdot, u): 0 \leq t \leq T\} \cup a(\cdot, M)$ , where  $u(x, t)$  is the solution of (2.1) and assume that  $(\cdot, \cdot)_p$  defines a norm  $\|\cdot\|_p$  on  $N$ . Then define  $\hat{a}(W)$  for  $W \in M$  to be such that  $\hat{a}(W) \in N$  and

$$(3.1) \quad (\hat{a}(W) - a(x, W), V)_p = 0, \quad V \in N;$$

i.e., take  $\hat{a}(W)$  to be the unique element of  $N$  such that  $\|\hat{a}(W) - a(x, W)\|_p$  is minimized on  $N$ . A useful example of such an inner product is

$$(f, g)_p = \sum_{k=1}^J f(Z_k)g(Z_k),$$

where the  $Z_k$ 's are in  $\Omega$ ; all that is required for  $(\cdot, \cdot)_p$  to be defined is that the functions in  $M \cup N \cup \{a(\cdot, u): 0 \leq t \leq T\} \cup a(\cdot, M)$  have pointwise values at the points  $Z_k$ .

The following lemma will be used in finding  $L$  and  $\theta$  such that (2.11) is valid.

**LEMMA 3.1.** *Suppose that  $(\cdot, \cdot)_p$  is also defined on a vector space  $\tilde{M}$  containing  $M$ . Let  $u(x, t)$  be the solution of (2.1). Suppose that there exist constants  $\gamma_0, \gamma_1, \gamma_2$  such*

that

$$\begin{aligned}
 \|V\|_p &\leq \gamma_0 \|V\|, & V \in \mathring{M}, \\
 \|V\| &\leq \gamma_1 \|V\|_p, & V \in N, \\
 \|a(\cdot, W) - a(\cdot, u(\cdot, t))\|_p &\leq \gamma_2 \|W - u\|_p, & W \in M, 0 \leq t \leq T.
 \end{aligned}
 \tag{3.2}$$

Then

$$\| \hat{a}(W) - a(\cdot, u(\cdot, t)) \| \leq \gamma_0 \gamma_1 \gamma_2 \|W - u\| + \theta(t),
 \tag{3.3}$$

where

$$\begin{aligned}
 \theta(t) = \inf \{ &\|a^* - a(\cdot, u(\cdot, t))\| + \gamma_1 \|a(\cdot, u(\cdot, t)) - a^*\|_p \\
 &+ \gamma_1 \gamma_2 \|u^+ - u(\cdot, t)\|_p + \gamma_0 \gamma_1 \gamma_2 \|u - u^+\| : a^* \in N, u^+ \in \mathring{M} \}.
 \end{aligned}
 \tag{3.4}$$

*Proof.* For any  $a^* \in N$ ,

$$\begin{aligned}
 \| \hat{a}(W) - a(u) \| &\leq \| \hat{a}(W) - a^* \| + \| a^* - a(u) \| \\
 &\leq \gamma_1 \| \hat{a}(W) - a^* \|_p + \| a^* - a(u) \| \\
 &\leq \gamma_1 \| a(W) - a^* \|_p + \| a^* - a(u) \|,
 \end{aligned}
 \tag{3.5}$$

where  $a(u) = a(\cdot, u(\cdot, t))$  and  $a(W) = a(\cdot, W(\cdot))$ . In the last step, we used the fact that

$$\|a(W) - a^*\|_p^2 = \|a(W) - \hat{a}(W)\|_p^2 + \|\hat{a}(W) - a^*\|_p^2 \geq \|\hat{a}(W) - a^*\|_p^2,$$

which follows from (3.1). Next note that

$$\|a(W) - a^*\|_p \leq \|a(W) - a(u)\|_p + \|a(u) - a^*\|_p \leq \gamma_2 \|W - u\|_p + \|a(u) - a^*\|_p.
 \tag{3.6}$$

For any  $u^+ \in \mathring{M}$ , we have

$$\begin{aligned}
 \|W - u\|_p &\leq \|W - u^+\|_p + \|u^+ - u\|_p \leq \gamma_0 \|W - u^+\| + \|u^+ - u\|_p \\
 &\leq \gamma_0 \|W - u\| + \gamma_0 \|u - u^+\| + \|u^+ - u\|_p.
 \end{aligned}
 \tag{3.7}$$

Use (3.7) in (3.6) and the result in (3.5) to obtain the conclusion.

In most of the applications of this lemma,  $M \subset N = \mathring{M}$ . In each example in which we consider a sequence of pairs of spaces  $M$  and  $N$ , we shall be able to choose the inner product  $(\cdot, \cdot)_p$  so that  $\gamma_0, \gamma_1, \gamma_2$  are independent of particular element of the sequence.

Before we discuss the construction of  $\tilde{a}$  from  $\hat{a}$ , we shall present several examples of spaces  $M$  and  $N$  and inner products  $(\cdot, \cdot)_p$ . In the first two examples,  $\hat{a}(W)$  is a weighted least-squares fit of  $a(W) = a(\cdot, W(\cdot))$  at a finite set of points; included in each example is the important special case in which  $\hat{a}(W)$  is obtained by Lagrange interpolation of  $a(W)$ .

It is important to note that if  $(f, g)_p = \sum c_i f(x_i)g(x_i)$  for positive  $c_i$ 's, then the  $\gamma_2$  of Lemma 3.1 can be taken equal to the Lipschitz constant for  $a(x, r)$  as a function

of  $r$ . In the examples that follow we shall assume that

$$(3.8) \quad |a(x, r_1) - a(x, r_2)| \leq L|r_1 - r_2|,$$

for all  $x \in \bar{\Omega}$  and  $r_1, r_2 \in \mathbf{R}$ .

*Example 1: Local, One-Dimensional, Polynomial Least-Squares Fits.* In this example,  $\Omega$  is to be the interval  $(0, 1)$  and  $\hat{a}(W)$  is to be a polynomial on certain subintervals that is a discrete least-squares approximation of  $a(W)$  on each subinterval. Take a partition  $\delta = \{x_j\}_{j=0}^J$ ,  $0 = x_0 < x_1 < \cdots < x_J = 1$ . Let  $I_j = (x_{j-1}, x_j)$ ,  $h_j = x_j - x_{j-1}$ , and  $h = \max h_j$ . For some positive integer  $N$ , take  $N$  to be the set of all functions on  $\Omega$  which are in  $P_N(I_j)$  for each  $j$ , where  $P_N(S)$  is the set of all polynomials of degree less than  $N + 1$  on  $S$ . Let  $M = N \cap H^1(\Omega)$ ; i.e.,  $M$  is the set of all continuous functions on  $\Omega$  which are in  $P_N(I_j)$  for  $j = 1, \dots, J$ .

For some  $K \geq N$ , let  $0 \leq e_0 < e_1 < \cdots < e_K \leq 1$  and  $0 < b_k, k = 0, \dots, K$ .

Define

$$(3.9) \quad (f, g)_p = \sum_{j=1}^J \sum_{k=0}^K b_k f(x_{j-1} + h_j e_k) g(x_{j-1} + h_j e_k);$$

where  $f(x_{j-1} + h_j e_0)$  and  $f(x_{j-1} + h_j e_K)$  are taken to be limits from the right and left, respectively, in case  $e_0 = 0$  or  $e_K = 1$ . Note that, with this  $N$  and  $(\cdot, \cdot)_p$ , the construction of  $\hat{a}(W)$  is done locally, that is, subinterval by subinterval. If we set  $\overset{\circ}{M} = N$ , then a simple homogeneity argument shows that, for  $l = 0$  and  $1$ ,

$$(3.10) \quad \gamma_l = \sup \left\{ \left[ \|f\|^2 / \sum_{k=0}^K b_k f^2(e_k) \right]^{l-1/2} : 0 \neq f \in P_N((0, 1)) \right\}.$$

Thus, the  $\gamma_0, \gamma_1$  and  $\gamma_2$  of Lemma 3.1 are independent of the partition  $\delta$ ; they are completely determined by  $L, N, \{e_k\}_{k=0}^K$ , and  $\{b_k\}_{k=0}^K$ .

In order to estimate the infimum in (3.4), we shall use the following lemma.

**LEMMA 3.2.** *There is a constant  $C$ , independent of the partition  $\delta$ , but depending on  $N, \{e_k\}_{k=0}^K$  and  $\{b_k\}_{k=0}^K$ , such that if  $1 \leq S \leq N + 1$ , then, for all  $g \in H^S(I_j)$ ,  $j = 1, \dots, J$ ,*

$$(3.11) \quad \inf \{ \|g - g^+\| + \|g - g^+\|_p : g^+ \in N \} \leq Ch^S \left[ \sum_{j=1}^J \int_{I_j} (g^{(S)}(x))^2 dx \right]^{1/2}.$$

This lemma is easy to prove by taking  $g^+$  as a Lagrange interpolant of  $g$  on each subinterval; we shall not prove it here since it follows from a more general lemma in the next example.

We can conclude from Lemmas 3.1 and 3.2 that, provided  $u$  and  $a(u)$  are sufficiently smooth,

$$(3.12) \quad \|\hat{a}(W) - a(u)\| \leq L\|W - u\| + \theta(t),$$

where

$$(3.13) \quad L = \gamma_0 \gamma_1 L, \quad \theta = Ch^S [\|(\partial/\partial x)^S u\| + \|(\partial/\partial x)^S a(x, u(x))\|],$$

for  $1 \leq S \leq N+1$ ,  $\gamma_0$  and  $\gamma_1$  are defined by (3.10), and  $L$  is such that (3.8) holds.

It is of interest to note that if we choose  $\{e_k\}$  and  $\{b_k\}$  such that

$$(3.14) \quad \sum_{k=0}^K b_k f(e_k) = \int_0^1 f(x) dx, \quad f \in P_{2N}((0, 1)),$$

then  $\gamma_0 = \gamma_1 = 1$ . The relation can be achieved by taking  $K \geq N$  and  $\{e_k, b_k\}$  to be the Gaussian quadrature points and weights, respectively. In the case  $K = N$ ,  $\hat{a}(W)$  is just the Lagrange interpolant of  $a(W)$  with respect to the points  $x_{j-1} + e_k h_j$ ; thus the  $b_k$ 's need not be used in computing  $\hat{a}(W)$ . One further remark about the case in which  $K = N \leq 3$  and the  $e_k$ 's are the Gaussian quadrature points is that we can use  $\hat{a}(W)$  directly instead of  $\tilde{a}(W)$  even though  $\hat{a}(W)$  may fail to satisfy (2.2). To see this, note that, for any  $W \in M$ , we have the estimate

$$(3.15) \quad \alpha_0 \|V'\|^2 \leq (\hat{a}(W)V', V') \leq \alpha_1 \|V'\|^2, \quad V \in M;$$

this uses the facts that  $\hat{a}(W) = a(W)$  at the points  $x_{j-1} + e_k h_j$ , that Gaussian quadrature on  $N+1$  points is exact on polynomials of degree  $2N+1$ , and that for  $N \leq 3$ ,  $2N+1 \geq N+2(N-1)$ . The estimate (3.15) and a crude upper bound on  $|\hat{a}(W)|$  in terms of  $\alpha_0, \alpha_1$  and  $N$  allow us to dispense with (2.2) in the proofs of Theorems 2.1 through 2.4 in this case.

*Example 2: General Construction of Local Least-Squares Fits.* Take a finite collection of pairs of sets  $S_k \subset B_k$ ,  $k = 1, \dots, K$ , where  $S_k$  and  $B_k$  are the closures in  $\mathbf{R}^p$  of their nonvoid interiors  $S_k^0$  and  $B_k^0$ . For each  $k = 1, \dots, K$ , let  $N_k$  be a finite-dimensional subspace of the space  $C(B_k)$  of continuous functions on  $B_k$ , and let  $p_k$  be a finite set of pairs  $(e, b)$  where  $e \in S_k$  and  $b > 0$ . Assume that for each  $k = 1, \dots, K$ ,

$$\|\phi\|_{L^2(S_k)} \quad \text{and} \quad \|\phi\|_{p_k} = \left( \sum_{(e,b) \in p_k} b \phi^2(e) \right)^{1/2}$$

define norms on  $N_k$ .

Let  $\Omega$  be a bounded domain on  $\mathbf{R}^p$  that is "triangulated" in the following fashion. Assume that  $\bar{\Omega} = \bigcup_{j=1}^J \sigma_j$ , where, for each  $j = 1, \dots, J$ ,  $\sigma_j$  is the image of a closed set  $\Gamma_j$  under a nonsingular affine map  $T_j$  on  $\mathbf{R}^p$ , where for some positive integer  $k_j \leq K$ ,  $S_{k_j} \subset \Gamma_j \subset B_{k_j}$ . Further assume that the boundary  $\partial\sigma_j$  has zero  $p$ -dimensional measure for each  $j$  and that, for  $l \neq j$ ,  $\sigma_l^0$  is disjoint from  $\sigma_j^0$ .

Let  $N$  be the space of functions  $\phi$  on  $\Omega$  such that the restriction of  $\phi$  to each  $\sigma_j^0$  lies in the set  $\{_j T^* V : V \in N_{k_j}\}$ , where  $_j T$  is the inverse of  $T_j$  and  $(_j T^* V)(x) = V(_j T(x))$ , for  $x \in \sigma_j$ . For functions  $f$  and  $g$  on  $\Omega$  such that their restrictions to  $\sigma_j^0$  have continuous extensions to  $\sigma_j^0 \cup T_j(S_{k_j})$ , define

$$(3.16) \quad (f, g)_p = \sum_{j=1}^J \sum_{(e,b) \in p_{k_j}} b f(T_j e) g(T_j e) |T_j'|,$$



where  $T'_j$  is the linear part of  $T_j$  and  $|T'_j|$  is the absolute value of  $\det(T'_j)$ . In (3.16), if  $(e, b) \in \mathfrak{p}_{k_j}$  is such that  $e \in \partial S_k$ , use the continuous extensions of  $f$  and  $g$  to  $T_j(S_{k_j})$  to evaluate  $f(T_j e)$  and  $g(T_j e)$ . Note, in particular, that  $(\cdot, \cdot)_{\mathfrak{p}}$  is defined on  $N$ . Henceforth, we shall ignore the technicality of how  $f(T_j e)$  is evaluated on  $\partial(T_j(S_{k_j}))$ .

We shall take  $M \subset N \cap H^1(\Omega)$  and  $\overset{\circ}{M} = N$ . For general  $N$ ,  $M$  may not be a good space with which to approximate functions in  $H^1(\Omega)$ , and it is not essential that the choice  $M \subset N$  be made. However, for many special cases of importance in practice, such as the  $N_k$ 's being certain classes of polynomials, the spaces  $M$  can have nice approximation properties.

**LEMMA 3.3.** *There are constants  $\gamma_0$  and  $\gamma_1$ , depending only on  $S_k$ ,  $B_k$ ,  $\mathfrak{p}_k$ , and  $N_k$  for  $k = 1, \dots, K$ , such that*

$$(3.17) \quad \|V\|_{\mathfrak{p}} \leq \gamma_0 \|V\| \quad \text{and} \quad \|V\| \leq \gamma_1 \|V\|_{\mathfrak{p}}, \quad V \in N.$$

In fact, we may take  $\gamma_l = \sup\{\gamma_{lk}, k = 1, \dots, K\}$  for  $l = 0, 1$ , where

$$(3.18) \quad \gamma_{0,k} = \sup\{\|f\|_{\mathfrak{p}_k} / \|f\|_{L^2(S_k)} : 0 \neq f \in N_k\},$$

$$\gamma_{1,k} = \sup\{\|f\|_{L^2(B_k)} / \|f\|_{\mathfrak{p}_k} : 0 \neq f \in N_k\}.$$

Notice in particular that  $\gamma_0$  and  $\gamma_1$  are independent of the sets  $\sigma_j$  and the maps  $T_j$ .

*Proof.* We know that  $\gamma_{0,k}$  and  $\gamma_{1,k}$ , as defined in (3.18), exist since, on the finite-dimensional vector space  $N_k$ , the norms  $\|f\|_{\mathfrak{p}_k}$ ,  $\|f\|_{L^2(S_k)}$ ,  $\|f\|_{L^2(B_k)}$  are equivalent. For any  $f \in L^2(\sigma_j)$ , let  $r(y) = f(T_j y)$  for  $y \in \Gamma_j$ . Then

$$(3.19) \quad \begin{aligned} \int_{\sigma_j} f^2(x) dx &= \int_{\sigma_j} r^2({}_j T x) dx = |T'_j| \int_{\sigma_j} r^2({}_j T x) |{}_j T'| dx \\ &= |T'_j| \int_{\Gamma_j} r^2(y) dy. \end{aligned}$$

Thus, if  $f \in N$ , and therefore  $r \in N_{k_j}$ , we see that

$$(3.20) \quad \|f\|_{L^2(\sigma_j)}^2 / \sum_{(e,b) \in \mathfrak{p}_{k_j}} b f^2(T_j a) |T'_j| = \|r\|_{L^2(\Gamma_j)}^2 / \|r^2\|_{\mathfrak{p}_{k_j}}^2$$

is bounded above and below by  $\gamma_{1,k_j}$  and  $1/\gamma_{0,k_j}$ , respectively. The conclusion follows.

We now need to consider the infimum in (3.4). In order to produce bounds for this infimum, we shall make some approximation assumptions on the spaces  $N_k$  and a weak smoothness assumption on  $\Omega$ . The spaces  $N_k$  will be assumed to include all polynomials of degree less than  $m \geq 3$  and the sets  $B_k$  will be assumed to be the closures of domains having the restricted cone property. It then follows from a result of Bramble and Hilbert [3] that there is a constant  $C_{BH}$  such that if  $V \in H^m(B_k)$

$$(3.21) \quad \inf\{\|V - \Psi\|_{L^2(B_k)}^2 + \|V - \Psi\|_{\mathfrak{p}_k}^2 : \Psi \in N_k\} \leq C_{BH} \sum_{|\alpha|=m} \|D^\alpha V\|_{L^2(B_k)}^2.$$

The domain  $\Omega$  will be assumed to have the restricted cone property; hence the Calderon

extension theorem [1] implies that, for each  $s \geq 0$ , there is a continuous linear map  $E_s: H^s(\Omega) \rightarrow H^s(\mathbf{R}^p)$  such that  $E_s g$  restricted to  $\Omega$  is  $g$  for each  $g \in H^s(\Omega)$ . Let  $T$  denote the triangulation given by the  $\sigma_j$ 's and  $T_j$ 's, and let

$$(3.22) \quad N_T = \left\| \sum_{j=1}^J \chi_{T_j(B_{k_j})} \right\|_{L^\infty(\mathbf{R}^p)},$$

where  $\chi_E$  is the characteristic function of the set  $E$ . We shall assume that  $u$  and  $a(x, u)$  are such that for each  $t, u$  and  $a(x, u)$  belong to  $H^m(\Omega)$ .

For  $\phi \in H^m(\Omega)$ , extend  $\phi$  to  $\tilde{\phi} = E_m \phi \in H^m(\mathbf{R}^p)$ . Note that if  $T_j^* \tilde{\phi} = \tilde{\phi} \circ T_j$ , a change of variables and (3.21) shows that

$$(3.23) \quad \begin{aligned} & \inf \{ \|\phi - \chi\|^2 + \|\phi - \chi\|_{\mathfrak{p}}^2 : \chi \in N \} \\ & \leq \sum_j |T_j'| \inf \{ \|\tilde{\phi} - \chi\|_{L^2(B_{k_j})}^2 + \|\tilde{\phi} - \chi\|_{\mathfrak{p}_{k_j}}^2 : \chi \in N_{k_j} \} \\ & \leq C_{BH} \sum_j |T_j'| \sum_{|\alpha|=m} \|D^\alpha T_j^* \tilde{\phi}\|_{L^2(B_{k_j})}^2. \end{aligned}$$

It is easily seen that there is a constant  $\tilde{C}$ , depending only on  $p$  and  $m$ , such that

$$(3.24) \quad |T_j'| \sum_{|\alpha|=m} \|D^\alpha T_j^* \tilde{\phi}\|_{L^2(B_{k_j})}^2 \leq \tilde{C} \|T_j'\|^{2m} \sum_{|\beta|=m} \|D^\beta \tilde{\phi}\|_{L^2(T_j B_{k_j})}^2,$$

where  $\|T_j'\|$  denotes the norm of the linear map  $T_j'$  with respect to the Euclidean norm on  $\mathbf{R}^p$ . The  $\tilde{C}$  in (3.24) does not involve any properties of the maps  $T_j$ . Thus we see that

$$(3.25) \quad \inf \{ \|\phi - \chi\|^2 + \|\phi - \chi\|_{\mathfrak{p}}^2 : \chi \in N \} \leq C_{BH} \tilde{C} N_T \left( \max_j \|T_j'\| \right)^{2m} \|E_m\|^2 \|\phi\|_m^2,$$

where  $\|E_m\|$  is the norm of  $E_m$  as a map of  $H^m(\Omega)$  to  $H^m(\mathbf{R}^p)$ .

LEMMA 3.4. *There is a constant  $C$ , depending on  $C_{BH}$ ,  $\tilde{C}$  and  $\|E_m\|$  but independent of the triangulation  $T$ , such that with  $\theta(t)$  defined by (3.4)*

$$(3.26) \quad \|\theta\|_{L^2(0,T)} \leq C N_T^{1/2} \left[ \max_j \|T_j'\| \right]^m [\|u\|_{L^2(H^m)} + \|a(u)\|_{L^2(H^m)}].$$

*Example 3: Local Fits of Values and First Derivatives.* Let  $S_k \subset B_k$  be as in Example 2; suppose that  $N_k$  is a finite-dimensional subspace of the space  $C^1(B_k)$  of continuously differentiable functions on  $B_k$ . (A function is in  $C^1(B_k)$  if it can be extended to be in  $C^1(\mathbf{R}^p)$ .) Let  $\mathfrak{p}_k$  be as in Example 2 and take  $\mathfrak{p}'_k$  to be a finite collection of triples  $(e, b, c)$  such that  $e \in S_k$ ,  $b > 0$ ,  $0 \neq c \in \mathbf{R}^p$ . Assume, for each  $k = 1, \dots, K$ , that  $\|\cdot\|_{L^2(S_k)}$  and  $\|\cdot\|_{\mathfrak{p}+\mathfrak{p}'_k}$  define norms on  $N_k$ , where

$$(3.27) \quad \begin{aligned} \|\phi\|_{\mathfrak{p}_k+\mathfrak{p}'_k}^2 &= \|\phi\|_{\mathfrak{p}_k}^2 + \|\phi\|_{\mathfrak{p}'_k}^2 \\ \|\phi\|_{\mathfrak{p}_k}^2 &= \sum_{(e,b) \in \mathfrak{p}_k} b \phi^2(e), \quad \|\phi\|_{\mathfrak{p}'_k}^2 = \sum_{(e,b,c) \in \mathfrak{p}'_k} b \left( \frac{\partial}{\partial c} \phi(e) \right)^2. \end{aligned}$$

Now assume that  $\Omega$  is triangulated as in Example 2; also adopt the definition of  $N$  given there. Take

$$(3.28) \quad \begin{aligned} (f, g)_p = & \sum_{j=1}^J |T'_j| \left\{ \sum_{(e,b) \in \mathfrak{r}_k} b f(T_j e) g(T_j e) \right. \\ & \left. + \sum_{(e,b,c) \in \mathfrak{p}'_k} \frac{\partial f}{\partial(T'_j c)}(T_j e) \frac{\partial g}{\partial(T'_j c)}(T_j e) \right\}. \end{aligned}$$

This inner product gives an analogue of Lemma 3.3 by change of variables.

LEMMA 3.5. *There are constants  $\gamma_0$  and  $\gamma_1$  depending only on  $S_k, B_k, \mathfrak{p}_k, \mathfrak{p}'_k$ , and  $N_k$  for  $k = 1, \dots, K$  such that*

$$(3.29) \quad \|V\|_p \leq \gamma_0 \|V\| \quad \text{and} \quad \|V\| \leq \gamma_1 \|V\|_p, \quad V \in N.$$

It is also the case that the argument leading to Lemma 3.4 can be used almost unchanged to prove an analogue in this case; since the statement of this lemma is exactly the same as Lemma 3.4, it will not be repeated.

The  $\gamma_2$  of Lemma 3.1 cannot, in general, be taken to be the Lipschitz constant for  $a$ .

LEMMA 3.6. *Suppose that  $a(x, r)$ ,  $(\nabla_x a)(x, r)$ , and  $\partial a(x, r)/\partial r$  are uniformly Lipschitz as functions of  $r$  with Lipschitz constants  $L_a, L_{a_x}, L_{a_r}$ , respectively. Also, suppose that  $\|\nabla u\|_{L^\infty(\Omega \times (0, T))}$  is finite. We then have that*

$$(3.30) \quad \|a(x, u) - a(x, W)\|_p \leq \gamma_3 \|u - W\|_p,$$

where (using  $\|\cdot\|$  for Euclidean norm on  $\mathbb{R}^p$ )

$$(3.31) \quad \begin{aligned} \gamma_3 &= 2[L_a + \gamma_4(L_{a_x} + L_{a_r}\|\nabla u\|_{L^\infty(\Omega \times (0, T))})], \\ \gamma_4 &= \max\{\|T'_j c\|: 1 \leq j \leq J, (e, b, c) \in \mathfrak{p}'_{k_j}\}. \end{aligned}$$

*Proof.* The terms in the sums over the  $\mathfrak{p}_{k_j}$ 's are estimated just as before. Note that for  $s \in \mathbb{R}^p$

$$(3.32) \quad \begin{aligned} & \left| \frac{\partial}{\partial s} (a(x, u) - a(x, W)) \right| \\ &= \left| s \cdot [(\nabla_x a)(x, u) - (\nabla_x a)(x, W)] \right. \\ & \quad \left. + \left[ \frac{\partial a}{\partial r}(x, u) - \frac{\partial a}{\partial r}(x, W) \right] \frac{\partial u}{\partial s} + \frac{\partial a}{\partial r}(x, W) \frac{\partial}{\partial s} (u - W) \right| \\ & \leq \|s\| L_{a_x} |(u - W)(x)| + L_{a_r} \|\nabla u\|_{L^\infty(\Omega \times (0, T))} \|s\| + L_a \left| \frac{\partial}{\partial s} (u - W)(x) \right|. \end{aligned}$$

This relation, when summed, gives the conclusion.

Note that this example of  $\hat{a}$  includes the case given by (1.13). In that case,  $K = 1$ ,  $S_1 = B_1 = [0, 1]$ ,  $m = 4$ , and  $N_1$  is the space of cubic polynomials on  $[0, 1]$ . The sets  $p_1$  and  $p'_1$  can be taken as

$$p_1 = \{(0, b), (1, b)\}, \quad p'_1 = \{(0, b, 1), (1, b, 1)\},$$

for any  $b > 0$ .

*Example 4: One-Dimensional Smooth Cubic Spline Interpolation.* It seems likely that nonlocal interpolation processes will be useful mostly in one-dimensional cases (or in situations where a tensor product structure is available) because of the necessity of solving a large linear system to produce local representations of the approximations of the coefficient  $a(W)$ . To illustrate a nonlocal interpolation process, we shall consider the special case of cubic splines on a uniform mesh.

As in Example 1, take  $\Omega = (0, 1)$ ,  $J > 1$ ,  $\delta = \{x_j\}_{j=0}^J$ ,  $I_j = (x_{j-1}, x_j)$ . Let  $h = 1/J = x_j - x_{j-1}$ , for  $j = 1, \dots, J$ . The spaces  $M$ ,  $N$ ,  $\dot{M}$  are all the smooth cubic splines over this partition, i.e.,

$$(3.33) \quad M = N = \dot{M} = \{V \in C^2(\bar{\Omega}) : V \in P_3(I_j), j = 1, \dots, J\}.$$

Define the discrete inner product by

$$(f, g)_p = \left[ f(x_{1/2})g(x_{1/2}) + f(x_{J-1/2})g(x_{J-1/2}) + \sum_{j=0}^J f(x_j)g(x_j) \right] h,$$

where  $x_\sigma = \sigma h$ . That this inner product induces a norm on  $N$  is easily seen; the function  $\hat{a}(W)$  is just the cubic spline that interpolates  $a(x, W(x))$  at knots, with the "end conditions" that it also interpolates at  $x_{1/2}$  and  $x_{J-1/2}$ .

We shall see that the  $\gamma_0, \gamma_1, \gamma_2$  of Lemma 3.1 can all be taken independent of  $h$  and that if  $a$  and  $u$  are sufficiently nice, then  $\theta(t) = O(h^4)$ . Just as in Examples 1 and 2, we may take  $\gamma_2 = L$ . It is clear that we can take

$$\gamma_0^2 = 3 \sup \left\{ \max_{x \in \bar{\Omega}} |V(x)|^2 : V \in P_3(\bar{\Omega}), \|V\| = 1 \right\}.$$

To see that  $\gamma_1$  is independent of  $h$ , we need to do some computation. If we are given  $F_0, F_{1/2}, F_1, F_2, \dots, F_{J-1}, F_{J-1/2}, F_J$  as the values of  $\phi \in N$ , then, using the notation  $S_j = h\phi'(x_j)$ , we see that

$$(3.34) \quad \begin{aligned} S_{j-1} + 4S_j + S_{j+1} &= 3(F_{j+1} - F_{j-1}), \quad j = 1, 2, \dots, J-1, \\ S_0 &= S_1 + 8F_{1/2} - 4(F_0 + F_1), \\ S_J &= S_{J-1} - 8F_{J-1/2} + 4(F_J + F_{J-1}). \end{aligned}$$

These relations are obtained from the facts that  $\phi''$  is continuous at  $x_j$ ,  $j = 1, \dots, J-1$ , that  $\phi(x_{1/2}) = F_{1/2}$  and that  $\phi(x_{J-1/2}) = F_{J-1/2}$ . The last two equations in

(3.34) can be substituted into the first and  $(J-1)$ st equations to give a strictly-diagonally-dominant set of equations for  $S_1, \dots, S_{J-1}$ . From these equations, we see that

$$(3.35) \quad \sum_{j=0}^J S_j^2 \leq C \left[ F_{1/2}^2 + F_{J-1/2}^2 + \sum_{j=0}^J F_j^2 \right],$$

where  $C$  is independent of  $h$ . It follows from homogeneity in  $h$  that, for any cubic  $q(x)$ ,

$$(3.36) \quad \int_0^h q^2(x) dx \leq \gamma h [q^2(0) + q^2(h) + h^2(q'^2(0) + q'^2(h))],$$

where

$$\gamma = \sup \{ \|q\|^2 : q \in P_3(\bar{\Omega}), q^2(0) + q^2(1) + q'^2(0) + q'^2(1) = 1 \}.$$

Thus, from (3.35) and (3.36), we see that there is a  $\gamma_1$ , independent of  $h$ , such that

$$(3.37) \quad \|\phi\| \leq \gamma_1 \|\phi\|_p, \quad \phi \in N.$$

LEMMA 3.7. *There is a constant  $C$ , independent of  $h$ , such that if  $g \in H^4(\Omega)$ ,*

$$\inf \{ \|g - \chi\| + \|g - \chi\|_p : \chi \in M \} \leq Ch^4 \|g\|_4.$$

Hence, if  $a(\cdot, u(\cdot, t))$  and  $u(\cdot, t)$  belong to  $H^4(\Omega)$ , there is a constant  $C$  such that

$$(3.38) \quad \theta(t) \leq Ch^4 [\|u(\cdot, t)\|_4 + \|a(\cdot, u(\cdot, t))\|_4].$$

*Proof.* All that is needed is a local interpolation process which reproduces cubic polynomials; if we have such a process, the proof is an easy application of the Bramble-Hilbert lemma [3]. Such a process can be defined as follows. Let  $J: C^2(\bar{\Omega}) \rightarrow M$  be such that

$$(3.39) \quad \begin{aligned} & \text{(i) for } 0 \leq k, 3k+3 \leq J \\ & \quad (V - JV)^{(l)}(x_j) = 0, \quad 0 \leq l \leq 2, j = 3k, 3k+3, \\ & \text{(ii) for } 0 \leq 3k < J, 3k+3 > J, \\ & \quad (V - JV)^{(l)}(x_{3k}) = 0, \quad 0 \leq l \leq 2, \\ & \quad (V - JV)^{(l)}(x_J) = 0, \quad 0 \leq l \leq J-3k-1. \end{aligned}$$

It is easily checked that (3.39) uniquely defines a  $C^2$  piecewise cubic on the intervals  $(x_0, x_3)$ ,  $(x_3, x_6)$ , etc.; it is clear that these fit together in a  $C^2$  fashion to give an element of  $M$ .

The lemma now follows by an argument that is very similar to the proof of Lemma 3.4.

*Construction of  $\tilde{a}$  from  $\hat{a}$ .* We shall now consider techniques which we can use to modify the  $\hat{a}$ 's produced in Examples 1–4 to obtain  $\tilde{a}$ 's which satisfy (2.2) as well as (2.11). Two ways of constructing  $\tilde{a}$  will be discussed in some detail. The first is a

theoretical construction that would be difficult to implement computationally but which is easily described and analyzed. The purpose of this construction is to point out that, in many realistic situations, the map  $\tilde{a}$  can be taken to be  $\hat{a}$  on the functions with which we deal. The second construction is a crude modification of  $\hat{a}$  on those regions in which the function  $a(W)$  is so rough that we cannot be assured that  $\hat{a}(W)$  is an approximation satisfying (2.2). This procedure can be easily implemented computationally; however, it seems unlikely that these "corrections" will be needed in practice if we are computing reasonable approximations to the smooth solution  $u$  of the differential problem.

The most straightforward, conceptually at least, construction of  $\tilde{a}$  is to define

$$(3.40) \quad \tilde{a}(W)(x) = \begin{cases} 3\alpha_1/2, & 3\alpha_1/2 \leq \hat{a}(W)(x), \\ \hat{a}(W)(x), & \alpha_0/2 \leq \hat{a}(W)(x) \leq 3\alpha_1/2, \\ \alpha_0/2, & \hat{a}(W)(x) \leq \alpha_0/2. \end{cases}$$

Because of the hypothesis that  $\alpha_0 \leq a(x, r) \leq \alpha_1$  for all  $r$ , we know that for any  $W \in M$  and  $(x, t) \in \bar{\Omega} \times [0, T]$

$$|\tilde{a}(W)(x) - a(x, u(x, t))| \leq |\hat{a}(W)(x) - a(x, u(x, t))|.$$

Thus, if (2.11) has been proved for  $\hat{a}$ , we know (2.11) holds for  $\tilde{a}$  with the same  $L$  and  $\theta$ . The  $\tilde{a}$  may not now be in the space  $N$  used for  $\hat{a}$ ; however, it is in  $L^\infty(\Omega)$  and satisfies (2.2) and (2.11). Since this map  $\tilde{a}$  satisfies the hypothesis of Section 2, we can use the results there to yield error estimates. These error estimates can be used in turn, with a detailed consideration of  $\hat{a}$ , to show that if  $u$  and  $a(x, u)$  are sufficiently smooth, then we may use  $\tilde{a} = \hat{a}$  in (2.3) and (2.4). This will be done in Section 4 for certain special choices of  $M, N, \hat{a}$ . In particular, we shall show in Example 4 of this section, that, for  $h$  and  $\Delta t$  sufficiently small, we may use  $\tilde{a} = \hat{a}$  in (2.4).

We shall describe the second technique for constructing  $\tilde{a}$  from  $\hat{a}$  in the context of Example 2. We shall then indicate the applicability of this technique in other settings and indicate some variants that may be easier to use in certain cases.

After constructing  $\hat{a}(W)$ , we examine it on each set  $\sigma_j$  and either accept it as  $\tilde{a}(W)$  or replace it by a constant approximation to  $a(W)$ . The most natural test to make would be to see if  $\hat{a}(W)$  satisfies (2.2), but this would be difficult in many cases because finding the maximum and minimum of  $\hat{a}$  is a nontrivial problem in all but the simplest examples.

We can, however, easily measure the  $\mathfrak{p}_k$ -norm of the difference between  ${}_jT^*\hat{a}(W)$  and the best constant approximation to it. If this norm is sufficiently small, then we know that (2.2) is satisfied; otherwise, the constant we compared with is an approximation on  $\sigma_j$  to  $a(x, u(x, t))$  that is about as good as  $\hat{a}(W)$  and has the advantage of satisfying (2.2). The detailed construction of  $\tilde{a}$  is as follows.

For each space  $N_k$ , there is a positive number  $d_k$  such that if  $V \in N_k$  satisfies  $(V, 1)_{\mathfrak{p}_k} = 0$  and  $\|V\|_{\mathfrak{p}_k} \leq d_k$ , then  $\sup\{|V(x)|: x \in B_k\} \leq 1$ . Note that since there

are usually a very small number of  $N_k$ 's and their dimensions are not large, the computation of  $d_k$  should not be a difficult problem; of course, it need be estimated only once, not once for each problem.

Fix  $W \in M$  and  $\sigma_j$  in the triangulation. Let  $k = k_j$  and let  $a^*(W)$  on  $\sigma_j$  be the constant

$$(3.41) \quad \begin{aligned} a_j^*(W) &= \left[ \sum_{(e,b) \in \mathfrak{p}_k} b a(T_j e, W(T_j e)) \right] / \sum_{(e,b) \in \mathfrak{p}_k} b \\ &= (T_j^* a(W), 1)_{\mathfrak{p}_k} / \|1\|_{\mathfrak{p}_k}^2. \end{aligned}$$

Note that  $\|1\|_{\mathfrak{p}_k} \neq 0$  since  $\|\cdot\|_{\mathfrak{p}_k}$  is a norm on  $N_k$  and  $1 \in N_k$ , by the assumption that  $N_k$  contains all polynomials of degree  $< m$  for some  $m > 2$ . If

$$(3.42) \quad \|T_j^*(\hat{a}(W)) - a^*(W)\|_{\mathfrak{p}_k} \leq \frac{1}{2} d_k a_j^*(W),$$

then on  $\sigma_j$

$$(3.43) \quad \frac{1}{2} \alpha_0 \leq \frac{1}{2} a_j^*(W) \leq \hat{a}(W) \leq \frac{3}{2} a_j^*(W) \leq \frac{3}{2} \alpha_1.$$

Thus (2.2) holds on  $\sigma_j$ , and we use  $\tilde{a}(W) = \hat{a}(W)$  there. If, on the other hand, (3.42) fails, let  $\tilde{a}(W) = a^*(W)$  on  $\sigma_j$ . In this case, we know that, for any constant  $c$ ,

$$\|T_j^*(\hat{a}(W)) - c\|_{\mathfrak{p}_k} \geq \frac{1}{2} d_k \alpha_0,$$

since the choice  $c = a_j^*(W)$  minimizes the left-hand side. Thus, for any  $c$ ,

$$\|\hat{a}(W) - c\|_{L^2(\sigma_j)}^2 \geq (\frac{1}{2} d_k \alpha_0 / \gamma_0)^2 |T_j'|,$$

where  $\gamma_0$  is as in Lemma 3. With the choice  $c$  equal the average of  $a(u) = a(x, u(x, t))$  on  $\sigma_j$ , we see that

$$\begin{aligned} \|a(u) - \hat{a}(W)\|_{L^2(\sigma_j)} &\geq \|\hat{a}(W) - c\|_{L^2(\sigma_j)} - \|a(u) - c\|_{L^2(\sigma_j)} \\ &\geq \frac{1}{2} (d_k \alpha_0 / \gamma_0) |T_j'|^{1/2} - \frac{1}{2} \|\nabla a(u)\|_{L^\infty} (\text{diam } \sigma_j) (\text{meas } \sigma_j)^{1/2} \\ &\geq \frac{1}{2} |T_j'|^{1/2} [(d_k \alpha_0 / \gamma_0) - \|\nabla a(u)\|_{L^\infty} (\text{diam } B_k) (\text{meas } B_k)^{1/2} |T_j'|]. \end{aligned}$$

Thus, for  $\|T_j'\|$  sufficiently small, we see that

$$(3.44) \quad \|a(u) - \hat{a}(W)\|_{L^2(\sigma_j)}^2 \geq \frac{1}{2} (\frac{1}{2} d_k \alpha_0 / \gamma_0)^2 |T_j'|.$$

But, since  $a^*(W)$  is between  $\alpha_0$  and  $\alpha_1$ , we see that

$$(3.45) \quad \|a(u) - a^*(W)\|_{L^2(\sigma_j)}^2 \leq (\alpha_1 - \alpha_0)^2 |T_j'| \text{meas } B_k.$$

Hence, from (3.44) and (3.45), there is a constant  $C$  such that, for  $\max_j \|T_j'\|$  sufficiently small,

$$(3.46) \quad \|a(u) - a^*(W)\|_{L^2(\sigma_j)}^2 \leq C \|a(u) - \hat{a}(W)\|_{L^2(\sigma_j)}^2.$$

Thus (2.11) holds for  $\tilde{a}$ ; the restriction on the size of  $\|T_j'\|$  can be viewed as adding a term to  $\theta$  where the additional term is zero for  $\max \|T_j'\|$  sufficiently small. This completes the construction of  $\tilde{a}$  for Example 2.

The above goes through almost unchanged in the case of Example 3. The only change is that  $\|\cdot\|_{p_k}$  and  $(\cdot, \cdot)_{p_k}$  are modified to include the derivative terms.

In the above construction, the choice of the  $\|\cdot\|_{p_k}$ -norm is not essential. We could have used the  $L^2(\Gamma_j)$ -norm, for example. Nor is the choice of  $a^*$  as a piecewise constant essential; what is needed for  $a^*$  is something in  $T^*N_k$  that lies between  $\alpha_0$  and  $\alpha_1$ . Both of these points are illustrated in Section 1. In that case, we took  $a^*(W)$  to be the function which interpolated the values of  $\hat{a}(W)$  and had zero slope at the knots, and we used the norm on  $\hat{a}(W) - a^*(W) = g$  to be the maximum of  $g'$  at each end of the subinterval. In that case, we also chose a different replacement for  $\hat{a}$  than  $a^*$ . Finally, note that even if the  $\hat{a}$  is obtained from a global fit of the coefficients, as in Example 4, we can still use the local corrections to produce  $\tilde{a}$ , provided the function  $\hat{a}$  is in the space  $N$  associated with the local construction.

**4. Asymptotic Error Estimates.** In this section, we shall combine the results of Sections 2 and 3 with some approximation theory and some elliptic error estimates to derive asymptotic estimates of the errors that result when (2.3) and (2.4) are used with  $M$  and  $N$  chosen from particular families of spaces. First, elliptic error estimates will be made using general approximation assumptions to provide the needed bounds on terms involving  $\eta$  and  $\partial\eta/\partial t$  in Theorems 2.2 and 2.4. Second, the special case of Example 4 of Section 3 will be discussed. Next a special case of Example 2 will be presented.

In looking at asymptotic estimates, the following definition will be useful [5]. A family  $\{M_h\}_{0 < h \leq 1}$  of finite-dimensional subspaces of  $H^1(\Omega)$  is an  $S_{h,m}$  family if there is a constant  $C$  such that for all  $V \in H^s(\Omega)$  with  $1 \leq s \leq m$

$$(4.1) \quad \inf_{\chi \in M_h} (\|V - \chi\| + h\|V - \chi\|_1) \leq Ch^s \|V\|_s.$$

The elliptic error estimates we shall develop here are very similar to others which can be found in the literature [9], [10], [7], [4]; however, the previous results are not in quite the form needed here. Elliptic regularity is crucial in deriving these results. For the necessary regularity to hold, it is sufficient that all the second derivatives of  $a(x, u(x, t))$  be bounded in  $\bar{\Omega} \times [0, T]$  and that  $\partial\Omega$  be a  $C^3, (p-1)$ -dimensional manifold regularly imbedded in  $\mathbb{R}^p$ . We shall assume throughout this section that these conditions hold. However, it should be noted that certain corners can be tolerated; in particular, if  $\Omega$  is a rectangular parallelepiped, the elliptic regularity we use is still valid. In the special case of  $p = 1$  and  $\Omega$  a bounded interval, the regularity is trivial.

**LEMMA 4.1.** *Suppose that  $\{M_h\}_{0 < h \leq 1}$  is a  $S_{h,m}$  family with  $m \geq 3$ . There is a constant  $C$  such that if  $\eta = u - W$ , where  $u$  is the solution of (2.1) and  $W$  is defined by (2.17) with  $M = M_h$ , then*



$$(4.2) \quad \|\eta\|_{L^\infty(L^2)} + \left\| \frac{\partial \eta}{\partial t} \right\|_{L^2(H^{-1})} \leq Ch^m \left[ \|u\|_{L^\infty(H^m)} + \left\| \frac{\partial u}{\partial t} \right\|_{L^2(H^{m-1})} \right].$$

*Proof.* Since for each  $t \in [0, T]$

$$(4.3) \quad B(u, \eta, V) = (a(u)\nabla\eta, \nabla V) + (\eta, V) = 0, \quad V \in M,$$

we see that there is a constant  $C$  such that

$$(4.4) \quad \|\eta(t)\|_1 \leq C \inf_{\chi \in M} \|u(t) - \chi\|_1 \leq Ch^{m-1} \|u(t)\|_m.$$

If  $\phi$  is the unique element of  $H^1(\Omega)$  such that

$$(4.5) \quad B(u, \phi, V) = (\eta, V), \quad V \in H^1(\Omega),$$

then  $\phi \in H^2(\Omega)$  and there is a  $C$ , depending only on  $\Omega, \alpha_0$ , and a bound for  $\nabla a(u)$ , such that

$$(4.6) \quad \|\phi\|_2 \leq C\|\eta\|.$$

From (4.5) and (4.3) it follows that for  $\chi \in M$

$$\|\eta\|^2 = B(u, \eta, \phi - \chi) \leq C\|\eta\|_1 \|\phi - \chi\|_1.$$

Thus, taking the infimum over  $\chi \in M$ , we see that

$$\|\eta\|^2 \leq C\|\eta\|_1 h \|\phi\|_2 \leq C\|\eta\|_1 h \|\eta\|.$$

Hence, for each  $t \in [0, T]$

$$(4.7) \quad \|\eta(t)\| \leq Ch^m \|u(t)\|_m.$$

If we differentiate (4.3) with respect to  $t$ , we see that

$$(4.8) \quad B(u, \eta_t, V) = (-a_t \nabla \eta, \nabla V), \quad V \in M,$$

where  $\eta_t = \partial \eta / \partial t$  and  $a_t = (\partial / \partial t) a(x, u(x, t))$ . First note that (4.8) gives

$$(4.9) \quad \begin{aligned} \|\eta_t\|_1 &\leq C \left[ \|\eta\|_1 + \inf_{\chi \in M} \|u_t - \chi\|_1 \right] \\ &\leq Ch^{m-2} [\|u\|_{m-1} + \|u_t\|_{m-1}]. \end{aligned}$$

Next, let  $\Psi \in H^1(\Omega)$  and take  $\phi$  such that

$$(4.10) \quad B(u, \phi, V) = (\Psi, V), \quad V \in H^1(\Omega).$$

Then  $\phi \in H^3(\Omega)$  and

$$(4.11) \quad \|\phi\|_3 \leq C\|\Psi\|_1,$$

where  $C$  depends on bounds for the  $x$ -derivatives of  $a(x, u(x, t))$  through order 2. We

can use (4.10) to estimate the  $H^{-1}$  norm of  $\eta_t$ . Note that, for appropriate  $\chi \in M$ ,

$$\begin{aligned} (\eta_t, \Psi) &= B(u, \eta_t, \phi) = B(u, \eta_t, \phi - \chi) + (a_t \nabla \eta, \nabla(\phi - \chi - \phi)) \\ &\leq C[\|\eta_t\|_1 + \|\eta\|_1] h^2 \|\phi\|_3 - (a_t \nabla \eta, \nabla \phi) \\ &\leq C[\|u\|_{m-1} + \|u_t\|_{m-1}] h^m \|\Psi\|_1 + (\eta, \nabla \cdot a_t \nabla \phi). \end{aligned}$$

There are no boundary terms that result from the integration by parts since the normal derivative of  $\phi$  is zero on the boundary. From the above, we see that, for  $t \in [0, T]$ ,

$$(4.12) \quad \|\eta_t\|_{-1} \leq Ch^m [\|u\|_m + \|u_t\|_{m-1}].$$

The result (4.2) follows easily from (4.7) and (4.12).

We are now ready to produce asymptotic error estimates for Example 4 of Section 3. Let the space  $M$  of (3.33) be  $M_h$  for  $h = 1/J$ ,  $J = 2, 3, \dots$ . It is well known that this gives a  $S_{h,4}$  family; this can be seen easily using  $J$  of (3.39) and the Peano kernel theorem. It is also easily checked that there is a constant  $C$  such that, if  $V \in H^s(\Omega)$  with  $2 \leq s \leq 4$ , then, for  $2 \leq p \leq \infty$  and  $l = 0, 1$ ,

$$(4.13) \quad \|(d/dx)^l(V - JV)\|_{L^p(\Omega)} \leq Ch^{s-l-1/2+1/p} \|V\|_s.$$

Let  $u$  be the solution of (2.1') and take  $W$  to be defined by (2.17). In order to apply Theorems 2.2 and 2.4, we need to know that  $\nabla W = W_x$  is bounded uniformly for  $J = 2, 3, \dots$  and  $(x, t) \in \Omega \times [0, T]$ . Note that for each  $t \in [0, T]$

$$\begin{aligned} \|W_x\|_{L^\infty(\Omega)} &\leq \|(W - Ju)_x\|_{L^\infty(\Omega)} + \|(u - Ju)_x\|_{L^\infty(\Omega)} + \|u_x\|_{L^\infty(\Omega)} \\ &\leq C(h^{-1/2} \|(W - Ju)_x\| + h^{1/2} \|u\|_2) + \|u_x\|_{L^\infty(\Omega)} \\ (4.14) \quad &\leq Ch^{-1/2} [\|(W - u)_x\| + \|(u - Ju)_x\|] + Ch^{1/2} \|u\|_2 + \|u_x\|_{L^\infty(\Omega)} \\ &\leq (Ch^{1/2} + \sqrt{2}) \|u\|_2. \end{aligned}$$

Thus, if  $u \in L^\infty(H^2)$ , we see that  $W_x$  is uniformly bounded. Hence, we obtain the following theorem from Theorem 2.2 and Lemma 3.7.

**THEOREM 4.1.** *Assume that*

$$\|a(u)\|_{L^\infty(H^4)} + \|u\|_{L^\infty(H^4)} + \|u_t\|_{L^2(H^3)}$$

*is finite. Let  $\tilde{a}$  be given by (3.40) with  $\hat{a}$  defined as in Example 4 of Section 3. Let  $U_0(X) = W(x, 0)$ , and let  $U$  be defined by (2.3). Then there is a constant  $C$ , independent of  $h$ , such that*

$$(4.15) \quad \|U - u\|_{L^\infty(L^2)} \leq Ch^4.$$

We shall use (4.15) to show that, for  $h$  sufficiently small,  $\hat{a}(U) = \tilde{a}(U)$ . All that

we need to show is that for  $t \in [0, T]$

$$(4.16) \quad \|\hat{a}(U) - a(u)\|_{L^\infty} \leq \frac{1}{2} \alpha_0.$$

Recall that we have assumed that  $a(x, u(x, t))$  is boundedly twice differentiable in  $\bar{\Omega} \times [0, T]$ . Note that

$$(4.17) \quad \begin{aligned} \|\hat{a}(U) - a(u)\|_{L^\infty(\Omega)} &\leq \|\hat{a}(U) - \hat{a}(u)\|_{L^\infty(\Omega)} + \|\hat{a}(U) - J(a(u))\|_{L^\infty(\Omega)} \\ &\quad + \|J(a(u)) - a(u)\|_{L^\infty(\Omega)} \\ &\leq Ch^{-1/2} [\|\hat{a}(U) - \hat{a}(u)\|_p + \|\hat{a}(u) - J(a(u))\|_p] + Ch^{3/2} \|a(u)\|_2, \end{aligned}$$

where we used the Peano kernel theorem to bound  $J(a(u)) - a(u)$ . The first term is bounded as follows:

$$(4.18) \quad \begin{aligned} \|\hat{a}(U) - \hat{a}(u)\|_p &\leq C \|U - u\|_p \leq C [\gamma_0 \|U - J(u)\| + \|J(u) - u\|_p] \\ &\leq C [\|U - u\| + \|u - J(u)\| + \|J(u) - u\|_p] \leq Ch^4. \end{aligned}$$

The second term is estimated as

$$(4.19) \quad \|\hat{a}(u) - J(a(u))\|_p \leq 2 \|a(u) - J(a(u))\|_p \leq Ch^2 \|a(u)\|_2.$$

From these estimates, it is clear that (4.16) holds, for  $h$  sufficiently small.

We also have the following result for the discrete-time Galerkin approximation.

**THEOREM 4.2.** *Assume that*

$$\|a(u)\|_{L^\infty(H^4)} + \|u\|_{L^\infty(H^4)} + \|u_t\|_{L^2(H^3)} + \|u_{tt}\|_{L^2(H^1)} + \|u_{ttt}\|_{L^2(H^{-1})}$$

*is finite. Take  $U_0$  and  $\tilde{a}$  as in Theorem 4.1 and let  $\{U_n\}_{n=1}^M$  be defined by (2.4) with  $\Delta t_n = T/M = \Delta t$ . Then there is a constant  $C$ , independent of  $h$  and  $\Delta t$ , such that*

$$(4.20) \quad \|U - u\|_{L_{\Delta t}^\infty(L^2)} \leq C(h^4 + (\Delta t)^2).$$

Note that if we take  $h$  and  $\Delta t$  to zero in such a fashion that  $h^{-1/2}(\Delta t)^2$  goes to zero, then  $\tilde{a}(U_n) = \hat{a}(U_n)$ , for  $h$  and  $\Delta t$  sufficiently small. In particular, this holds if the natural choice  $\Delta t \approx h^2$  is used.

In order to illustrate possible applications of Example 2 of Section 3, we shall derive asymptotic error estimates for a family of spaces which are built from piecewise polynomials on triangulations of a bounded domain  $\Omega \subset \mathbf{R}^2$ .

Let

$$(4.21) \quad \begin{aligned} B_1 &= \{(x, y): x \geq 0, y \geq 0, x + y \leq 3/2\}, \\ B_2 &= S_2 = \{(x, y): x + y \leq 1\} \cap B_1, \\ S_1 &= \{(x, y): x + y \leq 1/2\} \cap B_1. \end{aligned}$$

Fix an integer  $m \geq 3$  and let  $N_1 = N_2$  be the class of all polynomials in two variables of

degree less than  $m$ . For a sequence of positive  $h$ 's tending to zero, let

$$T_h = \{(\Gamma_j, T_j, k_j), j = 1, \dots, J_h\},$$

where  $k_j = 1$  or  $2$ ,  $\Gamma_j$  is a closed set such that  $S_{k_j} \subset \Gamma_j \subset B_{k_j}$ ,  $T_j$  is a one-to-one affine map on  $\mathbb{R}^2$ . Assume that  $\bar{\Omega}$  is the nonoverlapping union of the sets  $\sigma_j = T_j \Gamma_j$  as in Example 2. We shall further assume that there is a constant  $C$  independent of  $h$  such that

$$(4.22) \quad \|_j T'\| \leq C/h, \quad \|T'_j\| \leq Ch.$$

Also assume that  $\partial\Omega$  is contained in the union of the  $\sigma_j$ 's for which  $k_j = 1$  and that if  $k_j = 1$ ,  $\sigma_j \cap \partial\Omega$  is a smooth curve from  $T_j(1, 0)$  to  $T_j(0, 1)$ . The  $N_{T_h}$  defined by (3.22) is taken to be bounded independently of  $h$ ; with our assumption that  $\partial\Omega$  is  $C^3$ , it is clear that, for  $h$  sufficiently small, we can choose  $T_h$  such that  $N_{T_h} \leq 2$ . So that the piecewise polynomial functions on each  $\sigma_j$  fit together nicely, we assume that, if  $j_1 \neq j_2$ ,  $\sigma_{j_1} \cap \sigma_{j_2}$  is either void, a point, or  $T_{j_1}s_1 = T_{j_2}s_2$ , where  $s_1$  and  $s_2$  are sides of  $S_2$ .

Note that  $N_1$  and  $N_2$  have dimension  $m_1 = m(m+1)/2$  and that we can find  $m_1$  points  $Z_1, \dots, Z_{m_1}$  in  $S_2$  such that  $V \in N_1$  or  $N_2$  is determined by its values at these points. The points  $Z_1, \dots, Z_{m_1}$  can be chosen so as to include the vertices of  $S_2$ ,  $m-2$  evenly spaced points in the interior of each side, and  $m_1 - 3m + 3$  points in the interior of  $S_2$ . The space  $N$  is, as in Example 2, the space of all functions  $\phi$  such that  $\phi$  is a polynomial in two variables of degree less than  $m$  on each  $\Gamma_j$ .

Let  $M = N \cap H^1(\Omega)$ ; it is easily seen that  $M$  consists of those functions in  $N$  which are continuous on  $\bar{\Omega}$ . The functions in  $M$  can be represented by their values at the points  $T_j Z_k$ ,  $j = 1, \dots, J$ ,  $k = 1, \dots, m_1$ ; these points will not all be in  $\bar{\Omega}$  unless  $\Omega$  is convex, but we can use the values at these points of the natural extension of the polynomial on  $\sigma_j$ . Let  $J$  be the map of  $C(\mathbb{R}^2)$  into  $M$  such that  $V - JV = 0$  at each point  $T_j Z_k$ ,  $j = 1, \dots, J$ ,  $k = 1, \dots, m_1$ . Since  $\partial\Omega$  is smooth,  $\Omega$  has the restricted cone property and we can apply an argument very similar to the one used to prove Lemma 3.4 to show that this family of spaces  $M = M_h$  is a  $S_{h,m}$  family. In particular, if  $m \geq s \geq 2$ ,  $\phi \in H^s(\Omega)$  and  $\tilde{\phi} = E_s \phi \in H^s(\mathbb{R}^2)$  is the extension discussed in Example 2, we see that

$$(4.23) \quad \|\phi - J\tilde{\phi}\| + h\|\phi - J\tilde{\phi}\|_1 \leq Ch^s \|\phi\|_s.$$

This result is a straightforward application of the Bramble-Hilbert lemma and change of variables; the fact that  $\|T'_j\|, \|_j T'\|$  is bounded independently of  $h$  and  $j$  is used in estimating the error in the derivatives. It then follows from Lemma 9 of [2] that (4.1) holds. We shall also use the easily checked result that, for  $3 \leq s \leq m$  and  $\tilde{\phi} = E_s \phi$ ,

$$(4.24) \quad \|\phi - J\tilde{\phi}\|_{L^\infty(\Omega)} + h\|\nabla(\phi - J\tilde{\phi})\|_{L^\infty(\Omega)} \leq Ch^{s-1} \|\phi\|_{s^+}.$$

In order to show that, for  $W$  defined by (2.17),  $\nabla W$  is bounded on  $\bar{\Omega} \times [0, T]$ , assume that  $u \in L^\infty(H^3)$ . Then, with  $\tilde{u} = E_3 u$ , for each  $t \in [0, T]$ ,

$$\begin{aligned}
\|\nabla W\|_{L^\infty(\Omega)} &\leq \|\nabla(W - J\tilde{u})\|_{L^\infty(\Omega)} + \|\nabla(u - J\tilde{u})\|_{L^\infty(\Omega)} + \|\nabla u\|_{L^\infty(\Omega)} \\
&\leq C\{h^{-1}\|\nabla(W - J\tilde{u})\| + h\|u\|_3 + \|u\|_3\} \\
(4.25) \quad &\leq C\{h^{-1}\|\nabla(W - u)\| + h^{-1}\|\nabla(u - J\tilde{u})\| + \|u\|_3\} \\
&\leq C\|u\|_3.
\end{aligned}$$

A natural choice for  $\hat{a}$  is interpolation at  $m_1$  points in each set  $\sigma_j$ . In particular, if we choose

$$(4.26) \quad \|\phi\|_{\mathbb{P}_1}^2 = \sum_{l=1}^{m_1} \phi^2(\tfrac{1}{2}Z_l), \quad \|\phi\|_{\mathbb{P}_2}^2 = \sum_{l=1}^{m_1} \phi^2(Z_l),$$

then  $\hat{a}(V)$  is obtained by interpolating  $a(V)$  at the points  $\{T_j(\tfrac{1}{2}Z_l): l = 1, \dots, m_1\}$  or  $\{T_j(Z_l): l = 1, \dots, m_1\}$  for  $k_j = 1$  or 2, respectively. Assume that  $\hat{a}$  is given by (3.1) with  $(\cdot, \cdot)_p$  defined as in Example 2 using (4.26). Use the second construction of  $\tilde{a}$  from  $\hat{a}$  in Section 3. I.e.,  $\tilde{a}(V)$  is either  $\hat{a}(V)$  on  $\sigma_j$  or is  $a^*(V)$  on  $\sigma_j$ , where  $a^*$  is defined by (3.41), the choice being based on the truth or falsity of inequality (3.42).

Note that, in order to use the schemes (2.3) and (2.4), it is necessary to be able to compute integrals of the form  $(\tilde{a}(U)V_j, V_i)$ , where  $V_i$  and  $V_j$  are basis functions for  $M$ . In the interior of the region, this can be done exactly (up to rounding error). However, at the boundary, it will be necessary to build accurate approximations of integrals of the form  $(\tilde{V}_i V_j, \nabla V_i)$  where the  $\tilde{V}_i$ 's are basis functions for  $N$ . These are computed once a problem rather than once a time step. Thus, construction of these approximations is not extremely time consuming, even for very accurate approximations. We shall not consider here the effect of the errors made in constructing these integrals.

From Theorem 2.2 and Lemma 3.4, we obtain the following theorem.

**THEOREM 4.3.** *Assume that  $m \geq 3$  and that*

$$\|u\|_{L^\infty(H^m)} + \|a(u)\|_{L^\infty(H^m)} + \|u_t\|_{L^2(H^{m-1})}$$

*is finite. Let  $U_0(x) = W(x, 0)$  and let  $U$  be defined by (2.3). Then there is a constant  $C$ , independent of  $h$ , such that*

$$(4.27) \quad \|u - U\|_{L^\infty(L^2)} \leq Ch^m.$$

A computation that parallels that of (4.17), (4.18) and (4.19) shows that, for  $h$  sufficiently small,

$$(4.28) \quad \hat{a}(U) = \tilde{a}(U).$$

The analogous discrete-time result follows from Theorem 2.4.

**THEOREM 4.4.** *Suppose that  $m \geq 3$  and that*

$$\begin{aligned}
&\|a(u)\|_{L^\infty(H^m)} + \|u\|_{L^\infty(H^m)} + \|u_t\|_{L^2(H^{m-1})} \\
&+ \|u_t\|_{L^\infty(L^\infty)} + \|u_{tt}\|_{L^2(H^1)} + \|u_{ttt}\|_{L^2(H^{-1})}
\end{aligned}$$

is finite. Take  $U_0(x) = W(x, 0)$  and let  $\{U_n\}_{n=1}^M$  be defined by (2.4) with  $\Delta t_n = T/M = \Delta t$ . Then there is a constant  $C$ , independent of  $h$  and  $\Delta t$ , such that

$$\|U - u\|_{L_{\Delta t}^{\infty}(L^2)} \leq C(h^m + (\Delta t)^2).$$

In this case, if we take  $h$  and  $\Delta t$  to zero such that  $h^{-1}(\Delta t)^2$  goes to zero, then (4.28) holds for  $h$  and  $\Delta t$  sufficiently small. In particular, this is true if we use the natural choice  $(\Delta t)^2 \approx h^m$ .

Department of Mathematics  
University of Chicago  
Chicago, Illinois 60637

1. S. AGMON, *Lectures on Elliptic Boundary Value Problems*, Van Nostrand, Princeton, N. J., 1965. MR 31 #2504.
2. J. H. BRAMBLE, T. DUPONT & V. THOMÉE, "Projection methods for Dirichlet's problem in approximating polygonal domains with boundary-value corrections," *Math. Comp.*, v. 26, 1972, pp. 869–879.
3. J. H. BRAMBLE & S. R. HILBERT, "Bounds for a class of linear functionals with applications to Hermite interpolation," *Numer. Math.*, v. 16, 1970/71, pp. 362–369. MR 44 #7704.
4. J. H. BRAMBLE & J. E. OSBORN, "Rate of convergence estimates for nonselfadjoint eigenvalue approximations," *Math. Comp.*, v. 27, 1973, pp. 525–549.
5. J. H. BRAMBLE & A. H. SCHATZ, "Rayleigh-Ritz-Galerkin methods for Dirichlet's problem using subspaces without boundary conditions," *Comm. Pure Appl. Math.*, v. 23, 1970, pp. 653–675. MR 42 #2690.
6. J. DOUGLAS, JR. & T. DUPONT, "Galerkin methods for parabolic equations," *SIAM J. Numer. Anal.*, v. 7, 1970, pp. 575–626. MR 43 #2863.
7. J. DOUGLAS, JR. & T. DUPONT, "Galerkin methods for parabolic equations with non-linear boundary conditions," *Numer. Math.*, v. 20, 1973, pp. 213–237.
8. J. L. LIONS, *Équations Différentielles Opérationnelles et Problèmes aux Limites*, Die Grundlehren der math. Wissenschaften, Band 111, Springer-Verlag, Berlin, 1961. MR 27 #3935.
9. JOACHIM A. NITSCHKE, "Ein Kriterium für die Quasi-Optimalität des Ritzschen Verfahrens," *Numer. Math.*, v. 11, 1968, pp. 346–348. MR 38 #1823.
10. M. H. SCHULTZ, " $L^2$  error bounds for the Rayleigh-Ritz-Galerkin method," *SIAM J. Numer. Anal.*, v. 8, 1971, pp. 737–748. MR 45 #7967.
11. M. F. WHEELER, "A priori  $L_2$  error estimates for Galerkin approximations to parabolic partial differential equations," *SIAM J. Numer. Anal.*, v. 10, 1973, pp. 723–759.