

AD-A136 018

FEEDBACK AND ADAPTIVE FINITE ELEMENT SOLUTION OF  
ONE-DIMENSIONAL BOUNDARY..(U) MARYLAND UNIV COLLEGE  
PARK INST FOR PHYSICAL SCIENCE AND TECH..

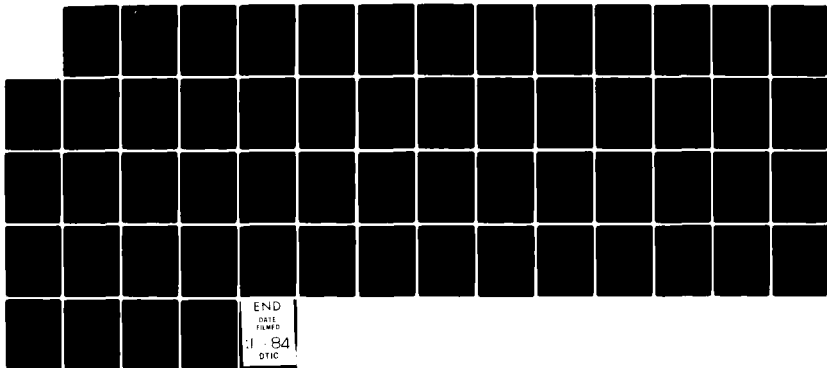
1/0

UNCLASSIFIED

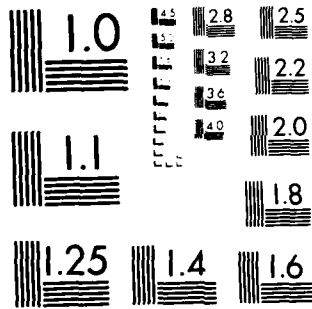
I BABUSKA ET AL. OCT 83 BN-1006

F/G 12/1

NL



END  
DATE  
FILMED:  
11-84  
DTIC



MICROCOPY RESOLUTION TEST CHART  
NATIONAL BUREAU OF STANDARDS-1963-A



INSTITUTE FOR PHYSICAL SCIENCE  
AND TECHNOLOGY

Laboratory for Numerical Analysis

Technical Note BN-1006

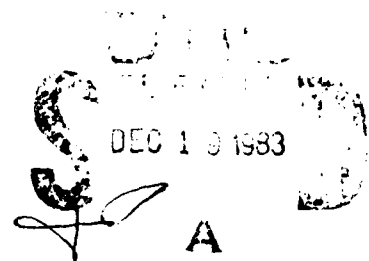
13

AD - A136018

FEEDBACK AND ADAPTIVE FINITE ELEMENT SOLUTION  
OF ONE-DIMENSIONAL BOUNDARY VALUE PROBLEMS

by

I. Babuška  
M. Vogelius



DTIC FILE COPY

October 1983

This document is available for unlimited distribution and use.



UNIVERSITY OF MARYLAND

83 12 16 091

SECURITY CLASSIFICATION OF THIS PAGE (When Data Entered)

REPORT DOCUMENTATION PAGE		READ INSTRUCTIONS BEFORE COMPLETING FORM	
1. REPORT NUMBER Technical Note BN-1006	2. GOVT ACCESSION NO. AD P116018	3. RECIPIENT'S CATALOG NUMBER	
4. TITLE (and Subtitle) Feedback and Adaptive Finite Element Solution of One-Dimensional Boundary Value Problems		5. TYPE OF REPORT & PERIOD COVERED Final life of the contract	
		6. PERFORMING ORG. REPORT NUMBER	
7. AUTHOR(s) I. Babuška and M. Vogelius		8. CONTRACT OR GRANT NUMBER(s) ONR N00014-77-C-0623	
9. PERFORMING ORGANIZATION NAME AND ADDRESS Institute for Physical Science & Technology University of Maryland College Park, MD 20742		10. PROGRAM ELEMENT, PROJECT, TASK AREA & WORK UNIT NUMBERS	
11. CONTROLLING OFFICE NAME AND ADDRESS Department of the Navy Office of Naval Research Arlington, VA 22217		12. REPORT DATE October 1983	
		13. NUMBER OF PAGES 54	
14. MONITORING AGENCY NAME & ADDRESS (if different from Controlling Office)		15. SECURITY CLASS. (of this report)	
		15a. DECLASSIFICATION/DOWNGRADING SCHEDULE	
16. DISTRIBUTION STATEMENT (of this Report)  Approved for public release: distribution unlimited			
17. DISTRIBUTION STATEMENT (of the abstract entered in Block 20, if different from Report)			
18. SUPPLEMENTARY NOTES			
19. KEY WORDS (Continue on reverse side if necessary and identify by block number)			
20. ABSTRACT (Continue on reverse side if necessary and identify by block number) This paper examines the concepts of feedback and adaptivity for the finite element method. The model problem concerns $C^0$ elements of arbitrary, fixed degree for a one-dimensional two-point boundary value problem. Three different feedback methods are introduced and a detailed analysis of their adaptivity is given.			

FEEDBACK AND ADAPTIVE FINITE ELEMENT SOLUTION  
OF ONE-DIMENSIONAL BOUNDARY VALUE PROBLEMS

I. Babuška  
M. Vogelius

Department of Mathematics  
and  
Institute for Physical Science and Technology  
University of Maryland  
College Park, MD 20742

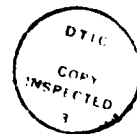
Technical Note BN-1006

Dedicated to F. L. Bauer

This research was partially supported by the Office of Naval Research  
under grant number N00014-77-C-0623.

ABSTRACT

This paper examines the concepts of feedback and Adaptivity for the Finite Element Method. The model problem concerns  $e^0$  elements of arbitrary, fixed degree for a one-dimensional two-point boundary value problem. Three different feedback methods are introduced and a detailed analysis of their adaptivity is given.



## 1. INTRODUCTION

In recent years there has been a growing interest in designing "adaptive" methods for the numerical solution of differential equations. "Adaptive" methods for numerical quadrature and integration of initial value problems for ODE's are common in available codes. The term "adaptive" however is generally used in a very loose sense.

In this paper we shall attempt to give a more precise analysis of the concepts of feedback and adaptivity as they relate to the numerical solution of differential equations. The definitions of feedback and adaptivity that we use are similar to those found in [11] (see also [7]): a feedback method is characterized by the fact that it produces a sequence (a trajectory) of approximate solutions, with each new solution depending on information available from the preceding solutions; a feedback method is called adaptive if it is optimal with respect to some performance measure, more precisely if its performance with respect to this measure is not worse than that of any other trajectory.

We restrict the discussion to finite element methods for a one-dimensional two-point boundary value problem; the elements are  $C^0$  piecewise polynomials of arbitrary, fixed degree. The performance measures we consider are related to the convergence, respectively the quasi-optimality, of the approximate solutions. We believe the approaches and the theorems discussed in this paper to be valid in dimensions higher than one; a feedback method of the same type as those introduced here has been implemented for a second order elliptic system of two equations on a two dimensional domain in the program FEARS (cf. [10]). Both the available theoretical results (cf. [2] [3] [6] and the extensive numerical experimentation strongly support our belief.

The feedback methods constructed here show some similarity to the algorithm for approximation introduced in [12], as well as the ODE-solver studied in [9] in the sense that they try to achieve an equilibration of certain (locally computable) quantities associated with the error. See also [8] [13] for a treatment of elliptic P.D.E.

The organization of this paper is as follows: In section 2 we introduce the model two point boundary value problem, its finite element solution and define the error-indicators and the error estimator on which all the feedback methods are based; the definitions are similar to those in [5]. Sections 3 through 5 make precise our notions of feedback and adaptivity and provide three examples of practically interesting feedback finite element methods. Theorem 6.1 is our first result concerning adaptivity of feedback methods, it shows that for a very large class of feedback methods the corresponding approximate solutions will always converge in energy to the exact solution. In section 7 we examine the question of adaptivity with respect to the stronger performance measure of quasi-optimality. The Theorems 7.1 and 7.2 verify that the three examples from sections 4 and 5 all produce approximate solutions that are quasi-optimal, provided the exact solution satisfies some additional assumptions; we go to some length to show that these assumptions are very mild and "almost" necessary. In order to compare two different feedback methods that are both adaptive, a relevant quantity is the cost of the computation. In section 8 we give a very simple-minded comparison of the feedback methods introduced in Sections 4 and 5. Finally the appendix contains the proof of a lemma that establishes the equivalence of our error indicators and a local interpolation error. This lemma was used in verifying quasi-optimality

and the proof of it is based on a stability estimate, Lemma 9.1, which is a slight extension of a result contained in [4].

## 2. PRELIMINARIES

Consider the two-point boundary value problem

$$(2.1a) \quad Lu = -\frac{d}{dx} a(x) \frac{du}{dx} + b(x)u = f \text{ in } I = (0,1)$$

$$(2.1b) \quad u(0) = u(1) = 0,$$

where  $a, b \in L_{\infty}(I)$  with

$$(2.2a) \quad 0 < a_1 \leq a(x) \leq a_2$$

$$(2.2b) \quad 0 \leq b(x) \leq a_2, \text{ and}$$

$$(2.3) \quad f \in H^{-1}(I).$$

The solution of (2.1a-b) is to be understood in the usual weak form, i.e.,  $u \in \overset{\circ}{H}^1(I)$  and

$$B(u, v) = \int_0^1 f v dx$$

for any  $v \in \overset{\circ}{H}^1(I)$ , where  $B$  denotes the bilinear form

$$(2.4) \quad B(u, v) = \int_0^1 \left( a \frac{du}{dx} \frac{dv}{dx} + buv \right) dx.$$

The space  $\overset{\circ}{H}^1(I)$  consists of all functions that have first derivatives in  $L_2(I)$  and that vanish at 0 and 1.  $H^{-1}(I)$  denotes the dual of  $\overset{\circ}{H}^1(I)$ .

We define the energy norm

$$\|u\|_E = [B(u,u)]^{1/2}.$$

For any mesh

$$\Delta : 0 = x_0^\Delta < x_1^\Delta \cdots < x_{N(\Delta)}^\Delta = 1$$

let

$$I_j^\Delta = (x_{j-1}^\Delta, x_j^\Delta), \quad h_j^\Delta = h(I_j^\Delta) = x_j^\Delta - x_{j-1}^\Delta, \quad j = 1, \dots, N(\Delta)$$

and  $h^\Delta = \max_j h_j^\Delta$ ;  $x_j^\Delta$  are called the nodal points,  $I_j^\Delta$ , the elements and  $N(\Delta)$  the cardinality of the mesh  $\Delta$ . A mesh may be identified by either the set of nodal points or the set of elements, we shall thus write  $\Delta = \{x_j^\Delta\}$ ,  $x \in \Delta$  or  $\Delta = \{I_j^\Delta\}$ ,  $I \in \Delta$  depending on the context.

Let  $p \geq 1$  be an integer. By  $S^p(\Delta)$  we denote the linear subspace of  $H^1(I)$  which consists of functions that are piecewise polynomials of degree  $\leq p$  on each  $I \in \Delta$ .

The finite element solution  $u(\Delta) \in S^p(\Delta)$  is defined in the usual way,

$$(2.5) \quad u(\Delta) = P(\Delta)u,$$

where  $u$  is the solution to (2.1a-b) and  $P(\Delta)$  is the elliptic projection of  $H^1(I)$  onto  $S^p(\Delta)$ . By this we mean  $P(\Delta)u \in S^p(\Delta)$  and

$$(2.6) \quad B(P(\Delta)u, v) = B(u, v)$$

for any  $v \in S^P(\Delta)$ . The expression

$$(2.7) \quad e(\Delta) = u - u(\Delta)$$

is the error of the finite element solution.

Let  $R(\Delta)$  and  $R(I \setminus I)$  be the elliptic projections of  $\mathring{H}^1(I)$  onto the spaces

$$\{v \in \mathring{H}^1(I) : v(x_j^\Delta) = 0, \quad j = 1, \dots, N(\Delta)\}$$

and

$$\{v \in \mathring{H}^1(I) : v(x) = 0, \quad x \in I \setminus I\}$$

respectively, with this notation

$$(2.8) \quad R(\Delta)v = \sum_{j=1}^{N(\Delta)} R(I \setminus I_j^\Delta)v.$$

We define the error indicators  $\eta_j^\Delta(I_j^\Delta) = \eta_j^\Delta$  and the error estimator  $\varepsilon(\Delta)$  as follows

$$(2.9) \quad \eta_j^\Delta = \|R(I \setminus I_j^\Delta)e(\Delta)\|_E$$

$$(2.10) \quad \varepsilon(\Delta) = \|R(\Delta)e(\Delta)\|_E.$$

Due to (2.8) we conclude that

$$(2.11) \quad \varepsilon(\Delta) = \left[ \sum_{j=1}^{N(\Delta)} (\eta_j^\Delta)^2 \right]^{1/2}.$$

The functions  $w_j(\Delta) = R(I \setminus I_j^\Delta)e(\Delta)$ , which enter into the definition

(2.9) of the error indicators, are characterized by

$$(2.12a) \quad Lw_j(\Delta) = r_j^\Delta = f - Lu(\Delta) \quad \text{in } I_j^\Delta,$$

$$(2.12b) \quad w_j(\Delta) \equiv 0 \quad \text{in } I \setminus I_j^\Delta.$$

To simplify the analysis in this paper we shall assume that the  $w_j(\Delta)$ 's are directly computable. In practice one derives explicit, approximate formulae based on the residuals  $r_j^\Delta$  (cf. [5]).

An important result is (see [1] [5]):

Lemma 2.1. There exists  $C \geq 1$  depending on  $\alpha_1$  and  $\alpha_2$  of (2.2a-b) but independent of  $a, b, f$  and  $\Delta$  such that

$$(2.13) \quad \varepsilon(\Delta) \leq \|e(\Delta)\|_E \leq C\varepsilon(\Delta).$$

For completeness we provide a short proof. Let

$$a_j^\Delta = \frac{1}{h_j^\Delta} \int_{x_{j-1}^\Delta}^{x_j^\Delta} a(x) dx, \quad j = 1, \dots, N(\Delta)$$

and define

$$(2.14) \quad \tilde{B}_\Delta(u, v) = \sum_{j=1}^{N(\Delta)} \int_{x_{j-1}^\Delta}^{x_j^\Delta} a_j^\Delta \frac{du}{dx} \frac{dv}{dx} dx$$

on  $\hat{H}^1(I) \times \hat{H}^1(I)$ .  $\hat{P}(\Delta)$  denotes the orthogonal projection of  $\hat{H}^1(I)$  onto  $S^P(\Delta)$  using the inner-product (2.14). It is well known that

$$\hat{P}(\Delta)v(x_j^{\Delta}) = v(x_j^{\Delta}), \quad j = 0, \dots, N(\Delta)$$

for any  $v \in H^1(T)$ . Thus

$$\begin{aligned} (2.15) \quad \|e(\Delta)\|_E^2 &= B(e(\Delta), e(\Delta)) \\ &= B(e(\Delta), e(\Delta) - \hat{P}(\Delta)e(\Delta)) \\ &= B(R(\Delta)e(\Delta), e(\Delta) - \hat{P}(\Delta)e(\Delta)) \\ &\leq C \|R(\Delta)e(\Delta)\|_E \|e(\Delta)\|_E. \end{aligned}$$

In (2.15) we used the fact that

$$e(\Delta) - \hat{P}(\Delta)e(\Delta) = 0 \quad \text{at the nodal points,}$$

we furthermore used the equivalence of the energy norms associated to  $B$  and  $\tilde{B}_\Delta$ :

$$\frac{1}{C} \tilde{B}_\Delta(v, v) \leq B(v, v) \leq C \tilde{B}_\Delta(v, v)$$

with  $C$  only depending on  $\alpha_1$  and  $\alpha_2$ . The inequality (2.15) immediately leads to the last inequality of (2.13). The first inequality in (2.13) follows from the fact that  $R(\Delta)$  is an orthogonal projection in the inner-product  $B$ .  $\square$

With additional assumptions on the coefficients  $a$  and  $b$  we may prove that not only is the estimator equivalent to the finite element error, it is indeed asymptotically exactly equal to the error.

Lemma 2.2. If in addition to satisfying (2.2a-b)  $a$  and  $b$  have the property that

$$\left(\frac{d}{dx}\right)^k a, \left(\frac{d}{dx}\right)^{k-1} b \in L_\infty(I)$$

for some integer  $k \geq 1$ , then

$$e(\Delta)_E = \varepsilon(\Delta)(1 + O(h_j^\Delta)^{2\mu})$$

with  $\mu = \min(p, k)$ . The term  $O(h_j^\Delta)$  is bounded by  $Ch_j^\Delta$  independently of  $f$  and  $\Delta$ .

Proof. Let  $R^1(\Delta)$  be the elliptic projection onto the orthogonal complement of  $R(\Delta)H^1(I)$ . Then

$$e(\Delta) = R(\Delta)e(\Delta) + R^1(\Delta)e(\Delta) = e_1(\Delta) + e_2(\Delta),$$

and

$$(2.16) \quad Le_2(\Delta) = 0 \quad \text{on each element } I_j^\Delta.$$

From (2.16) and the regularity assumptions about  $a$  and  $b$  it follows that

$$\|e_2(\Delta)\|_{H^{k+1}(I_j^\Delta)} \leq C \|e_2(\Delta)\|_{H^1(I_j^\Delta)}$$

where  $H^k$  denotes the standard Sobolev space of functions with derivatives up to and including order  $k$  in  $L^2(I_j^\Delta)$  (the constant  $C$  is independent of  $I_j^\Delta$ ). Thus there exists  $v \in S^p(\Delta)$  such that

$$e_2(\Delta) = v + z \quad \text{and}$$

$$\|z\|_{H^1(I_j^\Delta)} \leq C(h_j^\Delta)^\mu \|e_2(\Delta)\|_{H^1(I_j^\Delta)}, \quad \text{with } \mu = \min(p, k),$$

in particular we see that

$$(2.17) \quad \|z\|_E \leq C(h^\Delta)^{\mu} \|e_2(\Delta)\|_E.$$

Since  $0 = B(e(\Delta), v) = B(e_1(\Delta), v) + B(e_2(\Delta), v)$ , we compute that

$$\begin{aligned} B(e_2(\Delta), e_2(\Delta)) &= B(e_2(\Delta), v) + B(e_2(\Delta), z) \\ &= -B(e_1(\Delta), v) + B(e_2(\Delta), z) \\ &= -B(e_1(\Delta), e_2(\Delta)) + B(e_1(\Delta), z) \\ &\quad + B(e_2(\Delta), z) \\ &= B(e_1(\Delta), z) + B(e_2(\Delta), z) \end{aligned}$$

(in the last equality we used that  $e_1(\Delta)$  and  $e_2(\Delta)$  are orthogonal).

As a consequence of this and (2.17) we obtain

$$\|e_2(\Delta)\|_E^2 \leq C \|e_1(\Delta)\|_E \|e_2(\Delta)\|_E (h^\Delta)^{\mu} + C \|e_2(\Delta)\|_E^2 (h^\Delta)^{\mu},$$

and so for  $h^\Delta$  sufficiently small

$$\|e_2(\Delta)\|_E \leq C \|e_1(\Delta)\|_E (h^\Delta)^{\mu}.$$

This leads to

$$\begin{aligned} \|e(\Delta)\|_E^2 &= \|e_1(\Delta)\|_E^2 + \|e_2(\Delta)\|_E^2 \\ &\leq \|e_1(\Delta)\|_E^2 (1 + o(h^\Delta)^{2\mu}) \\ &= \varepsilon(\Delta) (1 + o(h^\Delta)^{2\mu}), \end{aligned}$$

thus proving the desired result for  $h^\Delta$  sufficiently small. For large  $h^\Delta$  this lemma is a simple consequence of Lemma 2.1.  $\square$

### 3. FEEDBACK AND ADAPTIVITY FOR FINITE ELEMENT METHODS

A standard finite element method is based on a single given mesh  $\Delta$ , whereas a feedback finite element method (feedback f.e.m.) produces a sequence of meshes  $\Delta_i$  and corresponding finite element solutions  $u_i$ ,  $i = 1, 2, \dots$ . The construction of the meshes  $\Delta_i$  is governed by a transition operator  $A$ .

$$\Delta_{i+1} = A(\Delta_1, \dots, \Delta_i; u_1, \dots, u_i).$$

A feedback f.e.m. may be terminated through some stopping criterion, e.g.,  $\Delta_i$ ,  $u_i$  are computed until

$$(3.1) \quad \varepsilon(\Delta_i) \leq \tau \|u_i\|_E$$

where  $\tau$  is a given tolerance and  $\varepsilon(\Delta)$  is the error estimator from before.

If a measure of performance is provided then a feedback f.e.m. shall be called adaptive if it is optimal with respect to this measure. We do not elaborate on the most general definitions of feedback, performance measure and adaptivity, instead we refer the interested reader to [11].

In this paper we analyze some concrete methods related to the problem (2.1a-b) and our performance measures evaluate the asymptotic convergence of the  $u_i$ . In practice, when analyzing the performance of a feedback f.e.m., a most relevant measure is the computational cost in satisfying a certain stopping criterion—this item is briefly discussed in section 8. One important question that enters into the evaluation of the performance of a terminated feedback f.e.m. is the effectivity of the

stopping criterion. We shall not address this issue in the present paper, but only point out that the result of Lemma 2.2 is very crucial to the effectivity of the stopping criterion (3.1), based on  $\epsilon(\Delta)$ .

For simplicity we consider only binary meshes, i.e., the nodal points all have the form  $k/2^j$  for some  $0 \leq j$ ,  $0 \leq k \leq 2^j$ . The transition operators never delete nodal points, so that the meshes  $\Delta_i$  form an increasing sequence. We shall call a transition operator simple if  $\Delta_{i+1}$  only depends on  $\Delta_i$  and  $u_i$ , all other transition operators are referred to as composite.

## 4. TWO SIMPLE FEEDBACK FINITE ELEMENT METHODS

Let  $\Delta_i = \{I_j^i\}$  and let  $\eta_j^i$ ,  $j = 1, \dots, N(\Delta_i)$  be the error indicators (2.9). Below are the definitions of two very useful transition operators  $A^{(1)}$  and  $A^{(2)}$ ; in each case we list the sets  $\Sigma_i$  of elements of  $\Delta_i$  that are bisected in passing to  $\Delta_{i+1}$

$$\Sigma_i^{(1)} : \{I_j^i : \eta_j^i \geq \beta \max_{I \in \Delta_i} \eta^i(I)\}, \quad 0 \leq \beta \leq 1.$$

$$\Sigma_i^{(2)} : \Sigma_i^{(1)} \cup \{I_{j(\ell)}^i\}_{\ell=1}^k, \quad \text{where}$$

$$\eta_{j(1)}^i \geq \eta_{j(2)}^i \geq \dots \geq \eta_{j(N(\Delta_i))}^i$$

is an ordering of the indicators by size, and  $k = [\gamma \cdot N(\Delta_i)]$   
 $0 \leq \gamma \leq 1.$

The transition operator  $A^{(1)}$  bisects all elements whose indicators are larger than or equal to  $\beta$  times the largest indicator.  $A^{(2)}$  bisects at least the same set but in addition always bisects the  $[\gamma \cdot N(\Delta_i)]$  elements with the largest indicators (approximately the fraction  $\gamma$ ). If  $\beta = 0$  both of these transition operators bisect all elements. Note also that  $A^{(1)}$  is just a special case of  $A^{(2)}$ , namely corresponding to  $\gamma = 0$ .

There is of course no reason why we could not have considered dividing each interval into  $n = 2^m$  subintervals, we have chosen bisection for simplicity.

## 5. A COMPOSITE METHOD

Let  $I_j^{\Delta_i}$  be an element of the mesh  $\Delta_i$ . A direct predecessor of  $I_j^{\Delta_i}$  is an element  $I \in \Delta_k$ ,  $1 \leq k \leq i-1$  with the property that  $I_j^{\Delta_i}$  arises as one of the (two) subintervals in a bisection of  $I$ . We define

$$\tilde{I}_j = \begin{cases} \tilde{I}_j^{\Delta_i} & \text{if } I_j^{\Delta_i} \in \Delta_1 \\ \text{the direct predecessor of } I_j^{\Delta_i}, & \text{otherwise.} \end{cases}$$

For any  $I_j^{\Delta_i} \in \Delta_i$  the predictor  $\tilde{\eta}(I_j^{\Delta_i})$  is given by

$$\tilde{\eta}(I_j^{\Delta_i}) = \min \left( \frac{[\eta^{\Delta_i}(I_j^{\Delta_i})]^2}{\eta^{\Delta_k}(\tilde{I}_j)}, \kappa \eta^{\Delta_i}(I_j^{\Delta_i}) \right),$$

where  $k$  is the largest integer,  $1 \leq k \leq i$ , such that  $\tilde{I}_j \in \Delta_k$  and  $0 < \kappa < 1$  is a constant. The  $\tilde{\eta}$  is called a predictor since it uses past experience to predict the effect of a bisection of  $I_j^{\Delta_i}$  on the value of the indicator.

The transition operator  $A^{(3)}$  bisects all elements of

$$\Sigma_i^{(3)} : \{I_j^{\Delta_i} : \eta^{\Delta_i}(I_j^{\Delta_i}) \geq \beta \max_{I \in \Delta_i} \tilde{\eta}(I)\}$$

for some fixed  $0 \leq \beta \leq 1$ .  $A^{(3)}$  is clearly composite. The transition operator used in the FEARS program (cf. [10]) is a direct extension of this operator, including the so-called short passes: whenever the increase in the number of elements is less than a preset fraction of the total number of elements of  $\Delta_i$  no finite element solution is computed on

$\Delta_{i+1}$ . Instead indicators for the elements of  $\Delta_{i+1}$  are constructed directly from the  $\tilde{n}(I_j^{\Delta_i})$  and the operator  $A^{(3)}$  is applied to  $\Delta_1, \dots, \Delta_i, \Delta_{i+1}$  with this set of indicators for  $\tilde{n}^{\Delta_{i+1}}(I_j^{\Delta_{i+1}})$ . This process continues until the increase in the number of elements is larger than the preset fraction of the total number of elements in  $\Delta_i$ , at which point the short pass is concluded and the finite element solution is computed on the resulting mesh  $\Delta_{i+k}$ . This approach avoids the computation of  $u(\Delta_i)$  after only a small increase in the number of elements and makes the refinement sufficiently strong.

## 6. CONVERGENCE OF FEEDBACK FINITE ELEMENT METHODS

We shall refer to a transition operator (and the corresponding feedback f.e.m.) as  $\delta$ -regular if it bisects at least one element  $I^*$  with  $\eta^{\Delta_i}(I^*) \geq \delta \max_{I \in \Delta_i} \eta^{\Delta_i}(I)$ . Since the transition operators  $A^{(1)}$  through  $A^{(3)}$  all bisect an interval  $I^*$  with  $\eta^{\Delta_i}(I^*) = \max_{I \in \Delta_i} \eta^{\Delta_i}(I)$  these are indeed  $\delta$ -regular.

Theorem 6.1. Let  $u_i$ ,  $i = 1, 2, \dots$  be the finite element solutions computed with a  $\delta$ -regular feedback f.e.m.,  $0 < \delta \leq 1$ . Then

$$\|e_i\|_E = \|u - u_i\|_E \rightarrow 0 \text{ as } i \rightarrow \infty$$

provided only (2.2a-b) and (2.3) hold.

The following two lemmas are essential to the proof of Theorem 6.1.

Lemma 6.1. Let  $H$  be a Hilbert space and let  $S_1 \subseteq S_2 \subseteq \dots \subseteq S_i \subseteq S_{i+1} \subseteq \dots \subseteq H$  be a nested sequence of closed subspaces. If  $P_i$  denotes the orthogonal projection onto  $S_i$  and  $P$  denotes the orthogonal projection onto

$S = \overline{\bigcup_{i=1}^{\infty} S_i}$ , then for any  $u \in H$

$$P_i u \rightarrow Pu \text{ as } i \rightarrow \infty.$$

Proof: Since  $Pu \in \overline{\bigcup_{i=1}^{\infty} S_i}$  there exist  $v_i \in \bigcup_{j=1}^i S_j = S_i$  such that

$$v_i \rightarrow Pu \text{ as } i \rightarrow \infty.$$

$P_i Pu$  is by definition the closest element to  $Pu$  in  $S_i$ , and thus

$$(6.1) \quad P_i Pu \rightarrow Pu \quad \text{as } i \rightarrow \infty.$$

It is not difficult to see that  $P_i Pu = P_i u$ , and this in combination with (6.1) verifies the lemma.  $\square$

Lemma 6.2. Let  $H$  be a Hilbert space and let  $H \supseteq S_1 \supseteq S_2 \supseteq \dots \supseteq S_i \supseteq S_{i+1} \supseteq \dots$  be a nested sequence of closed subspaces. If  $P_i$  denotes the orthogonal projection onto  $S_i$  and  $P$  denotes the orthogonal projection onto  $S = \bigcap_{i=1}^{\infty} S_i$ , then for any  $u \in H$

$$P_i u \rightarrow Pu \quad \text{as } i \rightarrow \infty.$$

Proof: Consider the sequence  $S_1^\perp \subseteq S_2^\perp \subseteq \dots \subseteq S_i^\perp \subseteq S_{i+1}^\perp \subseteq \dots \subseteq H$ . The orthogonal projection onto  $S_i^\perp$  is  $\text{Id} - P_i$  and the orthogonal projection onto  $\bigcup_{i=1}^{\infty} S_i^\perp = S^\perp$  is  $\text{Id} - P$  ( $\text{Id}$  denotes the identity operator). The previous lemma proves that

$$(\text{Id} - P_i)u \rightarrow (\text{Id} - P)u \quad \text{as } i \rightarrow \infty,$$

and thus

$$P_i u \rightarrow Pu \quad \text{as } i \rightarrow \infty. \quad \square$$

Proof of Theorem 6.1: The spaces  $S^P(\Delta_1) \subseteq S^P(\Delta_2) \subseteq \dots$  form a nested sequence of closed subspaces. The elliptic projection  $P(\Delta_i)$  is the orthogonal projection onto  $S^P(\Delta_i)$  in the space  $\dot{H}(I)$ , equipped with the innerproduct  $B(\cdot, \cdot)$ . From Lemma 6.1 we conclude that

$$u_i = P(\Delta_i) u \rightarrow Pu,$$

where  $P$  is the elliptic projection onto  $\overline{\bigcup_{i=1}^{\infty} S^D(\Delta_i)}$ , i.e.,

$$(6.2) \quad e_i = u - u_i = u - Pu + \xi_i, \quad \xi_i \rightarrow 0 \quad \text{as } i \rightarrow \infty.$$

As before let  $R(\Delta)$  denote the elliptic projection onto the space

$$\{v \in \mathring{H}^1(I) : v(x_j^{\Delta}) = 0, \quad j = 1, \dots, N(\Delta)\}.$$

According to (6.2) and Lemma 6.2

$$(6.3) \quad R(\Delta_i)e_i \rightarrow R \circ (\text{Id} - P)u,$$

where  $R$  is the elliptic projection onto

$$\{v \in \mathring{H}^1(I) : v(x) = 0 \quad \text{at all nodal points } x \text{ of } \bigcup_{i=1}^{\infty} \Delta_i\}.$$

For each  $i \geq 1$ , let  $I_{j(i)} \in \Delta_i$  be an element which is bisected in the transition from  $\Delta_i$  to  $\Delta_{i+1}$ . Since  $\bigcup_{i=1}^{\infty} \Delta_i$  only contains finitely many different elements above a given positive size and since repetitions in the sequence  $\{I_{j(i)}\}_{i=1}^{\infty}$  are excluded, it follows that  $h(I_{j(i)}) \rightarrow 0$  as  $i \rightarrow \infty$ . The corresponding indicators

$$\begin{aligned} n^{\Delta_i}(I_{j(i)}) &= \|R(I \setminus I_{j(i)})e_i\|_E \\ &= \left[ \int_{I_{j(i)}} \left( a \left( \frac{d}{dx} + 1 \right)^2 + b(\rho_i)^2 \right) dx \right]^{1/2}, \end{aligned}$$

with

$$\varepsilon_i = R(\Delta_i)e_i$$

must converge to 0, due to the fact that  $h(I_{j(i)}) \rightarrow 0$  and  $R(\Delta_i)e_i \rightarrow R \circ (\text{Id} - P)u$  in  $H^1(I)$ . The feedback f.e.m. is assumed to be  $\delta$ -regular,  $0 < \delta \leq 1$ , and the above argument may thus be used to conclude that

$$(6.4) \quad \max_{I \in \Delta_i} \eta^{\Delta_i}(I) \rightarrow 0 \quad \text{as } i \rightarrow \infty.$$

From (6.3) and (6.4) we get that

$$(6.5) \quad R \circ (\text{Id} - P)u = 0.$$

If, to the contrary, there existed a point  $x^*$  such that  $R \circ (\text{Id} - P)u(x^*) \neq 0$ , then  $|R(\Delta_i)e_i(x^*)| \geq \gamma_0 > 0$  for  $i$  sufficiently large, because of the convergence (6.3); let  $I_j^{\Delta_i} \in \Delta_i$  denote the interval that contains  $x^*$ , then

$$\gamma_0 \leq |R(\Delta_i)e_i(x^*)| \leq Ch_j^{1/2} \eta^{\Delta_i}(I_j^{\Delta_i}),$$

and since  $h_j = h_j^{\Delta_i} \leq 1$  this would imply that

$$0 < C \leq \max_{I \in \Delta_i} \eta^{\Delta_i}(I),$$

a clear contradiction to (6.4). The identity (6.5) in combination with (6.3) and (2.10) leads to

$$\varepsilon(\Delta_i) \rightarrow 0 \quad \text{as } i \rightarrow \infty,$$

and an application of Lemma 2.1 now gives the desired result.  $\square$

Remark 6.1. Note that Theorem 6.1 is slightly different from most finite element convergence results in that it guarantees convergence without any assumption that the maximal element size,  $h^{\Delta_i}$ , approaches 0. Indeed if  $u$  is a piecewise polynomial of degree  $\leq p$  in part of the interval  $I$  then  $h^{\Delta_i}$  will not tend to 0.

Remark 6.2. A convergence result similar to Theorem 6.1 may also be proven in two dimensions, see [ 2 ].

## 7a. QUASI-OPTIMALITY, PRELIMINARY RESULTS

As already pointed out in section 3 we distinguish between a feedback f.e.m. and an adaptive feedback f.e.m.; an adaptive method is a feedback method which is optimal with respect to some given performance measure.

If the performance measure is 1 whenever the finite element solutions converge, i.e., whenever  $u_i \rightarrow u$  in  $\dot{H}^1(I)$ , and 0 otherwise, then Theorem 6.1 verifies that any  $\delta$ -regular method,  $0 < \delta \leq 1$ , is adaptive with respect to this measure for any  $u \in \dot{H}^1(I)$ .

In this section and the next we shall study adaptivity with respect to a stronger measure, which is 1 whenever the finite element solutions are quasi-optimal and 0 otherwise. The optimal approximation is defined through the expression

$$(7.1) \quad \phi(u, N) = \inf_{\substack{\Delta: \text{binary mesh} \\ N(\Delta) = N}} \min_{v \in S^P(\Delta)} \|u - v\|_E,$$

and the aforementioned measure will be 1 iff

$$(7.2) \quad \|u - u_i\|_E \leq C\phi(u, N(\Delta_i)),$$

$C$  may depend on  $u$ , but is independent of  $i$ .

Consider the following simple example:

Example 7.1

Let  $u \in \dot{H}^1(I) \setminus \{0\}$  be a piecewise polynomial of degree  $\leq P$  on the mesh  $0 = x_0 < x_1 < x_2 = 1$ , where  $x_1$  is not of the form  $k/2^j$ . Then  $\|u - u_i\|_E > 0$  for any binary mesh, but at the same time  $\phi(u, N) = 0$ , for any  $N \geq 2$ , and hence (7.2) can never be satisfied.

As is evident from this example we cannot expect to verify (7.2) for any of our feedback f.e.m. unless we require that  $\delta(u, N)$  not become too small. We shall make the following assumption about  $u$ ,

$$(A.1) \quad \begin{cases} \exists \epsilon > 0 \text{ and } C_1 > 0 \text{ such that} \\ \delta(u, N) \geq C_1 N^{-\epsilon} \quad \forall N \geq 1. \end{cases}$$

The indicator  $\eta_j^\Delta$  corresponds to the finite element solution of (2.1a-b) and it is therefore not entirely determined only from the behaviour of  $u$  in  $I_j^\Delta$ . Assumptions about  $\eta_j^\Delta$  would in practice be very difficult to verify. Instead we introduce an indicator  $\xi_j^\Delta$  which corresponds to an interpolant of  $u$  and which therefore is entirely determined based on knowledge of  $u$  in  $I_j^\Delta$ ; in the appendix it is shown that the  $\xi$ 's and the  $\eta$ 's are intimately connected. Let  $u_*(\Delta) = P_*(\Delta)u \in S^p(\Delta)$  denote the solution to (2.6) in case  $a = 1$ ,  $b = 0$ . Note that  $\frac{d}{dx} u_*$  on the interval  $I_j^\Delta$  is simply the  $L_2$  projection of  $\frac{d}{dx} u$  onto polynomials of degree  $\leq p-1$  and that  $u_*$  interpolates  $u$  at the nodal points of  $\Delta$ ;  $u_*$  is thus locally determined from  $u$ .

Define

$$(7.3) \quad \begin{aligned} \xi_j^\Delta &= \xi(I_j^\Delta) = \left[ \int_{I_j^\Delta} \left( \frac{d}{dx} (u - u_*(x)) \right)^2 dx \right]^{1/2} \\ &= \min_q \left[ \int_{I_j^\Delta} \left( \frac{d}{dx} u - q \right)^2 dx \right]^{1/2}, \end{aligned}$$

where the minimum is taken over all polynomials of degree  $\leq p-1$ .

A common feature of the transition operators  $A^{(1)} - A^{(3)}$  for

$0 < \beta \leq 1$  is that they bisect all intervals with indicators that are comparable (within a constant) to the maximal indicator. One might expect that this will lead to an equilibration of the set of indicators, and the intuition is now that this equilibration of error-indicators ensures quasi-optimality. In order to make this plausibility argument rigorous we naturally have to make the assumption that there exists a sequence of meshes for which a certain approximation converges at optimal rate and for which at the same time the errors element by element are equilibrated. Our exact assumption about  $u$  is

(A.2) For any  $N \geq 1$  there exists  $\Delta^{(N)}$  such that

- (i)  $N(\Delta^{(N)}) = N$
- (ii)  $\left[ \sum_{I \in \Delta^{(N)}} (\xi(I))^2 \right]^{1/2} \leq C_2 \xi(u, N)$  and
- (iii)  $\max_{I \in \Delta^{(N)}} \xi(I) \leq C_2 \min_{I \in \tilde{\Delta}^{(N)}} \xi(I)$   
for some subset of elements  $\tilde{\Delta}^{(N)} \subseteq \Delta^{(N)}$ , with the property that  $\Delta^{(N)} \setminus \tilde{\Delta}^{(N)}$  consists of at most  $K$  elements.

The constants  $C_2$  and  $K$  are independent of  $N$ .

Consider the following example:

Example 7.2

Let  $v_k : I \rightarrow \mathbb{R}$  denote the function given by

$$v_k(x) = \begin{cases} 1 & \text{for } 0 \leq x < 2^{-k-1} \\ (-1)^j & \text{for } (j - \frac{1}{2})2^{-k} \leq x < (j + \frac{1}{2})2^{-k}, \quad 1 \leq j \leq 2^k - 1 \\ 1 & \text{for } 1 - 2^{-k-1} \leq x \leq 1 \end{cases}$$

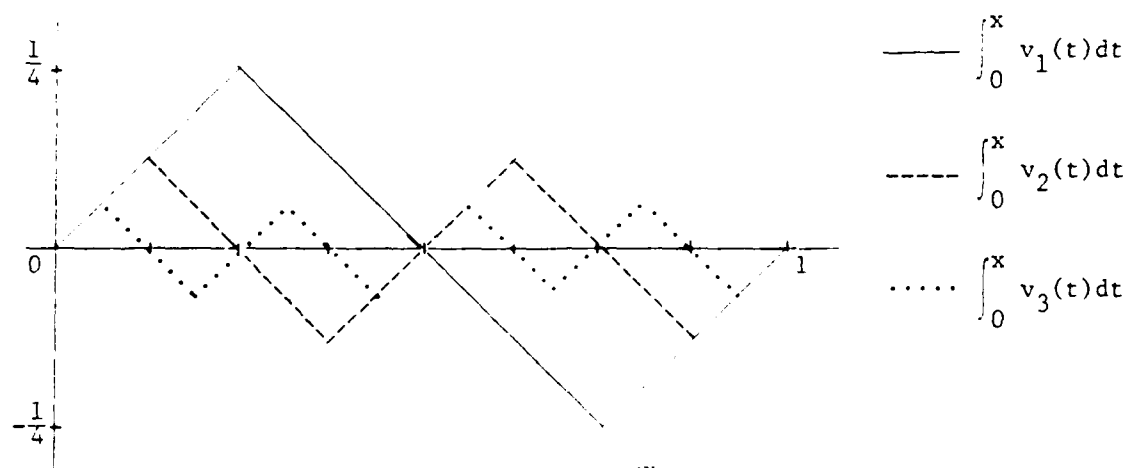


Fig. 1. The graphs of  $\int_0^x v_k(t) dt$ ,  $k = 1, 2, 3$

and define

$$\begin{aligned} u(x) &= \sum_{k=1}^{\infty} a_k \int_0^x v_k(t) dt \\ &= \int_0^x \sum_{k=1}^{\infty} a_k v_k(t) dt \in \dot{H}^1(I), \end{aligned}$$

$\{a_k\}$  square summable. Let  $\Delta_\ell$  denote the mesh  $\{j \cdot 2^{-\ell}\}_{j=0}^{2^\ell}$ , and let  $u_*^\ell$  denote  $u_*(\Delta_\ell)$  in the case  $p = 1$ . It is not difficult to see that the  $v_k$ 's,  $1 \leq k \leq \ell - 1$ , are constant on the elements of  $\Delta_\ell$  and that the  $v_k$ 's,  $\ell \leq k$ , on each element of  $\Delta_\ell$  are orthogonal (in  $L_2$ ) to the constants; we thus get

$$(7.4) \quad \frac{d}{dx} (u - u_*^\ell) = \sum_{k=\ell}^{\infty} a_k v_k.$$

The right hand side of (7.4) is symmetric in the nodal points of  $\Delta_\ell$  as a consequence

$$(7.5) \quad \xi_j^{\Delta_\ell} = \xi_{\ell-j}^{\Delta_\ell}, \text{ independently of } j.$$

From (7.4) and the fact that the  $v_k$ 's are indeed orthonormal in  $L_2(I)$ , it follows that

$$(7.6) \quad \int_0^1 \left( \frac{d}{dx} (u - u_\star^\ell) \right)^2 dx = \sum_{k=\ell}^{\infty} \alpha_k^2.$$

Pick

$$\alpha_k = 2^{-k} \text{ for } k = 2^j, \quad j = 0, 1, 2, \dots$$

$$\alpha_k = 0 \text{ otherwise,}$$

(7.6) then gives that

$$(7.7) \quad \int_0^1 \left( \frac{d}{dx} (u - u_\star^\ell) \right)^2 dx \geq CN(\Delta_\ell)^{-2}, \quad \ell = 2^m.$$

For  $\ell = 2^m$ , we define  $\ell' = \frac{\ell}{2} + 1$  and let  $\Delta_\ell'$  denote the mesh with nodal points

$$\left\{ \left( j - \frac{1}{2} \right) 2^{-\ell} \right\}_{j=1}^{2^\ell} \cup \left\{ j 2^{-\ell'} \right\}_{j=0}^{2^{\ell'}},$$

for which clearly  $N(\Delta_\ell')/N(\Delta_\ell) \rightarrow 1$  as  $m \rightarrow \infty$ . The mesh  $\Delta_\ell'$  is constructed exactly so that the  $v_k$ 's,  $k = 2^j$ ,  $0 \leq j \leq m$ , are constant in each element of  $\Delta_\ell'$ , and from the selection of the  $\alpha_k$ 's we thus get

$$(7.8) \quad \begin{aligned} \phi(u, N(\Delta'_k))^2 &\leq \sum_{k=1}^r \alpha_k^2 \\ &\leq C2^{-4k} \leq CN(\Delta'_k)^{-4}, \quad k = 2^m. \end{aligned}$$

We may add  $N(\Delta'_k) - N(\Delta_k)$  nodal points to the first element of  $\Delta_k$  without upsetting (7.7); instead of (7.5) we now get

$$(7.9) \quad \xi_j^{\Delta_k} = \xi_j^{\Delta'_k} \quad \text{for at least } \nu N(\Delta_k) \text{ values of } j,$$

with  $\nu$  arbitrarily close to 1 (provided  $m$  is sufficiently large). It follows directly from (7.7) and (7.8) that the  $u_{*k}^p$  are not quasioptimal even though the indicators are equilibrated in the sense of (7.9).  $\square$

The problem with the preceding example is not that  $\phi$  dies off too fast, rather it has to do with the fact that there is a total lack of correlation between the length of an element and the associated interpolation error. To avoid this problem we shall assume that

$$(A.3) \quad \left[ \begin{array}{l} \text{There exists a finite set of point } \{z_k\}_{k=1}^K \text{ and a constant } \\ C_3 \text{ such that } \forall I = (x, y) \subseteq I: \\ \\ I \cap \{z_k\}_{k=1}^K = \emptyset \\ \parallel \\ \downarrow \\ \xi(I) \leq C_3 \min\{\xi(I_{1/2}^-), \xi(I_{1/2}^+)\} \\ \\ \text{where } I_{1/2}^- = (x, \frac{x+y}{2}) \text{ and } I_{1/2}^+ = (\frac{x+y}{2}, y). \end{array} \right.$$

The assumptions (A.1)-(A.3) are easily seen to be satisfied for sufficiently smooth  $u$  and they also hold for a wide class of singular  $u$ ,

such as  $u(x) = x^\alpha - x$ . In the following example we demonstrate this for the case  $u(x) = x^\alpha - x$  and  $p = 1$ .

Example 7.3

Let  $u(x) = x^\alpha - x \in \tilde{H}^1(I)$ ,  $1/2 < \alpha$ . It is not difficult to see that (A.1) holds with  $c = 1$ .

Let  $\phi(x)$  denote the function  $\phi(x) = x^{(2\alpha-1)/3}$  and let  $\Delta^{(N)}$  be the mesh with nodal points  $x_j = \phi^{-1}(j/N) = (j/N)^{3/(2\alpha-1)}$ . It is not difficult to see that  $(\xi(I_j))^2$ ,  $1 \leq j \leq N$ , is of the order

$$\left| \left( \frac{d}{dx} \right)^2 u(x_j) \right|^2 h_j^3$$

which is bounded from above and below by a constant times  $N^{-3}$ . Since  $\phi(u, N) \geq C_1 N^{-1}$  this leads us to conclude that assumption (A.2) is also satisfied.

Let  $I = (x, y) \subseteq I$  be any interval; it is not difficult to see that  $(\xi(I))^2$  is of the order

$$\left| \left( \frac{d}{dx} \right)^2 u(y) \right|^2 (y-x)^3$$

and hence

$$\begin{aligned} (\xi(I))^2 &\leq C \left| \left( \frac{d}{dx} \right)^2 u(y) \right|^2 (y-x)^3 \\ &= 8C \min \left\{ \left| \left( \frac{d}{dx} \right)^2 u\left(\frac{x+y}{2}\right) \right|^2, \left| \left( \frac{d}{dx} \right)^2 u(y) \right|^2 \right\} \left(\frac{y-x}{2}\right)^3 \\ &\leq C' \min \{ \xi(I_{1/2}^-)^2, \xi(I_{1/2}^+)^2 \} \end{aligned}$$

which verifies (A.3).  $\square$

Definition 7.1. Let  $\Delta_i$ ,  $i = 1, 2, \dots$  be a sequence of meshes. We shall call  $\{\Delta_i\}_{i=1}^{\infty}$  semi-equilibrated if there exist constants  $C$ ,  $\nu > 0$  and subsets  $\tilde{\Delta}_i \subset \Delta_i$  with at least  $\nu N(\Delta_i)$  elements, such that

$$\max_{I \in \Delta_i} \xi(I) \leq C \min_{I \in \tilde{\Delta}_i} \xi(I).$$

The following lemma justifies the intuitive feeling that equilibration implies quasi-optimality.

Lemma 7.1. Assume that  $u \in \mathring{H}^1(I)$  satisfies (A.1)-(A.3). If  $\Delta_i$ ,  $i = 1, 2, \dots$  is a semi-equilibrated sequence of meshes and  $u_i = u(\Delta_i)$  is the associated finite element projection, defined in (2.5)-(2.6), then

$$\|e_i\|_E = \|u - u_i\|_E \leq C\phi(u, N(\Delta_i)).$$

Proof: For any mesh  $\Delta$  we have

$$\|u - u(\Delta)\|_E^2 \leq C \sum_{I \in \Delta} (\xi(I))^2,$$

where  $C$  depends only on  $\alpha_2$  of (2.2a-b). It thus suffices to verify that

$$\sum_{I \in \Delta_i} (\xi(I))^2 \leq C \sum_{I \in \Delta'_i} (\xi(I))^2$$

where  $\Delta'_i = \Delta^{(N(\Delta_i))}$  is the mesh, of same cardinality as  $\Delta_i$ , that occurs in the assumption (A.2).

Let  $\nu$ ,  $\tilde{\Delta}_i$  be as in definition 7.1 and let  $\tilde{\tilde{\Delta}}_i \subseteq \tilde{\Delta}_i$  consist of those elements that do not contain any of the points  $\{z_j\}_{j=1}^K$  of the

assumption (A.3).  $\tilde{\Delta}_i$  must contain at least  $vN(\Delta_i) - K$  elements and since it is only necessary to prove this lemma for  $N(\Delta_i)$  sufficiently large we may assume that  $\tilde{\Delta}_i$  contains at least  $vN(\Delta_i)$  elements,  $0 < v' < v$ . Divide each interval of  $\tilde{\Delta}_i$  into  $2^k$  subintervals of equal size, where  $k$  is picked such that  $2^k v' \geq 1$ . Since there are at least  $2^k v' N(\Delta_i) \geq N(\Delta_i)$  of these subintervals and exactly  $N(\Delta_i)$  elements in  $\Delta_i'$ , it follows that at least one of the subintervals  $I_0$  is entirely contained in an element  $I_0' \in \Delta_i'$ . From the definition of  $\tilde{\Delta}_i$ , it follows that

$$\begin{aligned}
 (7.10) \quad \max_{I \in \Delta_i} \xi(I) &\leq C\xi(J) \\
 &\leq C C_3^k \xi(I_0) \\
 &\leq C C_3^k \xi(I_0'),
 \end{aligned}$$

where  $J$  is the element of  $\tilde{\Delta}_i$ , which was bisected  $k$  times to give  $I_0$ , and  $C_3$  is the same constant as in (A.3). Based on (7.10) and the assumption (A.2) we conclude that

$$(7.11) \quad \xi(I) \leq C\xi(I') \quad \forall I \in \Delta_i, \quad I' \in \tilde{\Delta}_i'$$

where  $\tilde{\Delta}_i' = \tilde{\Delta}^{(N(\Delta_i'))}$  is the same as in (A.2). Since  $\Delta_i$  and  $\tilde{\Delta}_i'$  have a comparable number of elements (7.11) leads to

$$\begin{aligned}
 \sum_{I \in \Delta_i} (\xi(I))^2 &\leq C \sum_{I \in \tilde{\Delta}_i'} (\xi(I))^2 \\
 &\leq C \sum_{I \in \Delta_i'} (\xi(I))^2,
 \end{aligned}$$

exactly as desired.  $\square$

## 7b. QUASI-OPTIMALITY

The indicators  $\xi$  used in the statements of the preliminary result in section 7a were all associated to the interpolant  $u_*$  (cf.(7.3)). The transition operators  $A^{(1)}-A^{(3)}$  defined in sections 4 and 5 were all based on the "actually computed" indicators (cf.(2.9)). In order to establish quasi-optimality of any of these feedback f.e.m. we thus need a lemma relating the  $\xi$ 's and the  $\eta$ 's.

Lemma 7.2. Assume that  $a$  in addition to satisfying (2.2a-b) is Hölder continuous with exponent  $1/2$ . Then there exists  $h_0$  (independent of  $u$ ) such that for any mesh  $\Delta$  with  $h^\Delta \leq h_0$  one has

$$\begin{aligned} (a_j^\Delta)^{1/2} \xi_j^\Delta &= C(h_j^\Delta)^{1/2} \max_{I \in \Delta} \xi(I) \\ &\leq \eta_j^\Delta \leq (a_j^\Delta)^{1/2} \xi_j^\Delta + C(h_j^\Delta)^{1/2} \max_{I \in \Delta} \xi(I), \end{aligned}$$

$1 \leq j \leq N(\Delta)$ , where  $a_j^\Delta = a\left(\frac{x_{j-1}^\Delta + x_j^\Delta}{2}\right)$  and  $C$  is independent of  $u$  and  $\Delta$ .

The proof of this lemma is given in the appendix; at this point we proceed to verify quasi-optimality of the feedback method associated with  $A^{(1)}$ .

Theorem 7.1. Assume that  $a$  is Hölder continuous with exponent  $1/2$  and satisfies (2.2a-b). Assume that  $u \in \mathring{H}^1(I)$  satisfies (A.1)-(A.3) and that  $\frac{d}{dx} u \in L^r(I)$  for some  $r > 2$ . Let  $\Delta_i, u_i, i = 1, 2, \dots$  be the meshes and the finite element solutions generated with the feedback

finite element method based on the transition operator  $A^{(1)}$ ,  $0 < \epsilon \leq 1$ . Assume that  $h_i \leq h_0$  where  $h_0$  is as given in Lemma 7.2. Then

$$\|e_i\|_E = \|u - u_i\|_E \leq C\phi(u, N(\Delta_i)).$$

Proof: Let  $\Sigma_i^{(1)}$  denote those elements of  $\Delta_i$  which are bisected in passing from  $\Delta_i$  to  $\Delta_{i+1}$ . Since there are at most finitely many different elements in  $\bigcup_{i=1}^{\infty} \Delta_i$  above a given positive size and since the sets  $\Sigma_i^{(1)}$  are mutually disjoint, it follows easily that

$$(7.12) \quad \max_{I \in \Sigma_i^{(1)}} h(I) \rightarrow 0 \quad \text{as } i \rightarrow \infty.$$

Let  $I_i^* \in \Delta_i$  denote an element with

$$\xi(I_i^*) = \max_{I \in \Delta_i} \xi(I).$$

Due to Theorem 6.1 we know that

$$(7.13) \quad \xi(I_i^*) \rightarrow 0 \quad \text{as } i \rightarrow \infty,$$

and the assumption (A.1) guarantees that

$$(7.14) \quad \xi(I_i^*) \neq 0 \quad \forall i.$$

A combination of (7.13) and (7.14) gives

$$(7.15) \quad h(I_i^*) \rightarrow 0 \quad \text{as } i \rightarrow \infty.$$

Based on Lemma 7.2 and (7.15) we conclude that

$$(7.16) \quad \max_{I \in \Delta_i} \xi(I) \leq C \max_{I \in \Delta_i} n^{\Delta_i}(I)$$

for  $i$  sufficiently large. Similarly Lemma 7.2, (7.12), (7.16) and the definition of  $\Sigma_i^{(1)}$  show that for any  $I_j^{\Delta_i} \in \Sigma_i^{(1)}$

$$(7.17) \quad \beta' \max_{I \in \Delta_i} \xi(I) \leq \xi(I_j^{\Delta_i}),$$

with  $0 < \beta'$ , provided  $i$  is sufficiently large.

The meshes  $\Delta_i$ ,  $i \geq i_0$  are thus according to (7.17) generated by always bisecting elements that satisfy  $\beta' \max_{I \in \Delta_i} \xi(I) \leq \xi(I_j^{\Delta_i})$ . Since the

indicator  $\xi(I)$  does not increase through bisection of an interval, we may without loss of generality assume that each transition from  $\Delta_i$  to  $\Delta_{i+1}$  consists in the bisection of only one interval, satisfying (7.17) (if this is not the case we add the intermediate meshes to the sequence  $\{\Delta_j\}_{j \geq i_0}$ ).

We divide the elements of  $\Delta_i$  into three distinct categories

$$\Delta_i^S : \text{those } \tilde{I} \in \Delta_i \text{ for which } \xi(\tilde{I}) < \beta' C_3^{-1} \max_{I \in \Delta_i} \xi(I),$$

$$\Delta_i^Z : \text{those } \tilde{I} \in \Delta_i \setminus \Delta_i^S \text{ that contain one (or more) of the point } \{z_k\}_{k=1}^K \text{ of the assumption A.3,}$$

$$\Delta_i^R : \Delta_i \setminus (\Delta_i^S \cup \Delta_i^Z),$$

the constant  $\varepsilon'$  is the same as in (7.17) and  $C_3$  is the constant of the assumption (A.3). We shall prove that  $\Delta_i^r$  contains at least  $\nu N(\Delta_i)$  elements, for some  $\nu > 0$ . According to the definition 7.1 this ensures that the sequence of meshes  $(\Delta_i)_{i \geq i_0}$  is semi-equilibrated and the conclusion of this theorem then follows directly from Lemma 7.1.

Since  $1 \leq C_3$  and since only elements that satisfy (7.17) are bisected, it is clear that the element bisected in passing from  $\Delta_i$  to  $\Delta_{i+1}$  is not in  $\Delta_i^s$ . If the interval being bisected is in  $\Delta_i^r$  then it follows from (7.17), the assumption (A.3) and the definition of  $\Delta_i^r$  that the two resulting new elements are in  $\Delta_{i+1}^r$ , i.e., we obtain

$$(7.18) \quad \#\Delta_{i+1}^r \geq \#\Delta_i^r + 1$$

(# denotes the number of elements). For any  $i_0 + 1 \leq i$  define

$M_i$  = number of indices  $j$ ,  $i_0 \leq j \leq i-1$ , for which the transition from  $\Delta_j$  to  $\Delta_{j+1}$  consisted in bisecting an element of  $\Delta_j^r$ .

If we prove

$$(7.19) \quad \liminf_{i \rightarrow \infty} \frac{M_i}{i} > \nu > 0,$$

then we immediately get from (7.18) that

$$\#\Delta_i^r \geq \nu N(\Delta_i), \text{ for } i \text{ sufficiently large}$$

for some  $\nu > 0$ , which completes the proof of this theorem. The proof

of (7.19) proceeds by contradiction: assume that

$$\frac{M_i}{i} \rightarrow 0 \text{ for some (sub)sequence } i \rightarrow \infty,$$

this implies that there exist arbitrarily large  $i$  such that the transition from  $\Delta_j$  to  $\Delta_{j+1}$ ,  $i_0 \leq j \leq i-1$  consist at least  $i/2$  times in bisecting an element of  $\Delta_j^z$ . Since there are only  $K$  points  $z_k$ ,  $1 \leq k \leq K$ , in the assumption (A.3) we conclude that

$$\Delta_i^z \text{ contains an element } I'_i \text{ of size } 2^{-\tau}, \text{ with } \tau = i/2K.$$

On this interval  $I'_i$

$$\begin{aligned} \xi(I'_i)^2 &\leq \int_{I'_i} \left(\frac{du}{dx}\right)^2 dx \\ &\leq \left( \int_{I'_i} \left(\frac{du}{dx}\right)^r dx \right)^{\frac{2}{r}} h(I'_i)^{1-\frac{2}{r}} \\ &\leq C 2^{-\nu i}, \text{ with } \nu = (1-\frac{2}{r})/2K. \end{aligned}$$

From this and the fact that  $I'_i \in \Delta_i^z$  it follows that

$$\begin{aligned} \max_{I \in \Delta_i} \xi(I)^2 &\leq C 2^{-\nu i}, \text{ or} \\ \Phi(u, N(\Delta_i)) &\leq C [N(\Delta_i) 2^{-\nu N(\Delta_i)}]^{1/2}, \end{aligned}$$

for certain arbitrarily large  $i$ , and this clearly contradicts the assumption (A.1).  $\square$

Remark 7.1. The assumption  $h^i \leq h_0$  can obviously be replaced by the (weaker) assumption that  $h^i \leq h_0$  for  $i$  sufficiently large. This is satisfied, for example, when the solution  $u$  is not a polynomial of degree  $\leq p$  in any subinterval of  $I$ .

Remark 7.2. If we had assumed that the set of points  $\{z_k\}_{k=1}^K$  was empty, i.e.,

$$\varepsilon(I) \leq C_3 \min\{\varepsilon(I_{1/2}^-), \varepsilon(I_{1/2}^+)\}$$

for any  $I \subseteq I$ , then the set  $\Delta_i^z$  would also have been empty and we had not needed to assume that  $\frac{d}{dx} u \in L^r(I)$ ,  $r > 2$ , or  $\varphi(u, N) \geq CN^{-\sigma}$  (only that  $\varphi(u, N) > 0$ ).  $\square$

Remark 7.3. If we assume that the set of point  $\{z_k\}_{k=1}^K$  of the assumption (A.3) is empty, then it is not difficult to verify, based on Lemma 7.2, that

$$(7.20) \quad \max_{I \in \Delta_i} \eta^{\Delta_i}(I) \leq C \max_{I \in \Delta_i} \tilde{\eta}^{\Delta_i}(I), \quad i \geq i_0,$$

where  $\tilde{\eta}^{\Delta_i}(I)$  is the predictor defined in section 5, and  $\Delta_i$ ,  $i = 1, 2, \dots$  is a sequence of meshes generated by the transition operator  $A^{(3)}$ .

Consequently there exists  $0 < \beta'$  such that

$$\beta' \max_{I \in \Delta_i} \eta^{\Delta_i}(I) \leq \eta(I_j^{\Delta_i}), \quad i \geq i_0,$$

for any  $I_j^{\Delta_i} \in \mathcal{L}_i^{(3)}$  and we may now use the argument of the previous proof

to show that the finite element solutions computed, using  $A^{(3)}$ , are quasi-optimal.  $\square$

In order to verify quasi-optimality of the feedback f.e.m. associated with the transition operator  $A^{(2)}$  we shall need a slightly more restrictive assumption than (A.3).

$$(A.3)' \quad \left[ \begin{array}{l} \text{There exist } 1 < C_3 \text{ and } \beta_3 < 1 \text{ such that for any} \\ I = (x, y) \subseteq I \\ \\ \varepsilon(I) \leq C_3 \min\{\varepsilon(I_{1/2}^-), \varepsilon(I_{1/2}^+)\} \text{ and} \\ \\ \max\{\varepsilon(I_{1/2}^-), \varepsilon(I_{1/2}^+)\} \leq \beta_3 \varepsilon(I) \\ \\ \text{with } I_{1/2}^- = (x, \frac{x+y}{2}) \text{ and } I_{1/2}^+ = (\frac{x+y}{2}, y) \end{array} \right.$$

We could permit exceptional points  $\{z_k\}_{k=1}^K$  as in (A.3), but for clarity of exposition we have omitted these here.

Theorem 7.2. Assume that  $a$  is Hölder continuous with exponent  $1/2$  and satisfies (2.2a-b). Assume that  $u \in \hat{H}^1(I)$  satisfies (A.1)-(A.2) and (A.3)'. Let  $\Delta_i$ ,  $u_i$  be the meshes and the finite element solutions generated by the feedback finite element method based on the transition operator  $A^{(2)}$ ,  $0 \leq \gamma \leq 1$ ,  $0 < \beta < 1$ . There exists  $0 < \gamma_0$  such that

$$\|e_i\|_E = \|u - u_i\|_E \leq C(u, N(\Delta_i))$$

provided  $0 \leq \gamma \leq \gamma_0$  (both  $\gamma_0$  and  $C$  may depend on  $u$ ).

Assumption (A.3)' implies that  $u$  cannot be a polynomial of degree  $\leq p$  in any subinterval unless  $u$  is indeed a polynomial of degree  $\leq p$  in the entire interval  $I$ ; this latter possibility is excluded by assumption (A.1). From Theorem 6.1 we know that  $\sum_{I \in \Delta_i} \xi(I)^2 \rightarrow 0$  as  $i \rightarrow \infty$ , i.e., in particular  $\max_{I \in \Delta_i} \xi(I) \rightarrow 0$  as  $i \rightarrow \infty$ . Since there are only finitely many different elements in  $\bigcup_{i=1}^{\infty} \Delta_i$  above a given positive size, it follows that either  $h^{\wedge i} = \max_{I \in \Delta_i} h(I) \rightarrow 0$  as  $i \rightarrow \infty$ , or  $\xi(I) = 0$  for some subinterval  $I \in \Delta_i$ ,  $I \neq \emptyset$ ; the last is impossible due to the fact that  $u$  is nowhere a polynomial of degree  $\leq p$ , and we thus conclude

$$(7.21) \quad h^{\wedge i} = \max_{I \in \Delta_i} h(I) \rightarrow 0 \quad \text{as } i \rightarrow \infty.$$

Theorem 7.2 will be proven by help of the next two lemmas.

Lemma 7.3. Let  $C_4$  denote a constant

$$C_4 > \max(C_3, 1/\beta).$$

Let  $i_0$  be a given index. If  $0 < \nu_0$  is taken sufficiently small then there exists an index  $i_0 \leq i$  for which the cardinality of

$$\bar{\Delta}_i = \{I \in \Delta_i : a_I^{1/2} \xi(I) \geq C_4^{-3} \max_{I \in \Delta_i} a_I^{1/2} \xi(I)\}$$

is at least  $N(\Delta_i)/2$  for any  $0 \leq \nu \leq \nu_0$  ( $a_I$  denotes the value of  $a$  at the midpoint of the interval  $I$ ).

Proof: From (7.21) and Lemma 7.2 we see that

$$(7.22) \quad \left[ \begin{array}{l} \text{if } I \in \Delta_i, \text{ and} \\ r_i(I) \geq \mu \max_{I \in \Delta_i} r_i(I) \text{ or } a_I^{1/2 \xi(I)} \geq \mu \max_{I \in \Delta_i} a_I^{1/2 \xi(I)}, \mu > 0, \\ \text{then} \\ r_i(I) = a_I^{1/2 \xi(I)} (1 + \rho_i) \end{array} \right.$$

where  $\rho_i \rightarrow 0$  as  $i \rightarrow \infty$ , independently of  $I$ . Let  $0 < \delta < 1$  be the constant entering in the definition of  $A^{(2)}$  and select

$$\max(\delta, \delta_3) < \delta < 1,$$

then it is clear from the assumption (A.3)', (7.22) and the definition of the transition operator  $A^{(2)}$  that

$$(7.23) \quad \max_{I \in \Delta_{i+1}} a_I^{1/2 \xi(I)} < \delta \max_{I \in \Delta_i} a_I^{1/2 \xi(I)}$$

after each application of  $A^{(2)}$ , for  $i$  sufficiently large. Take a large fixed  $i_0 \leq i_1$  so that  $\rho_i, i_1 \leq i$ , are sufficiently small and (7.23) is satisfied for  $i_1 \leq i$  (to simplify notation we assume  $i_0 = i_1 = 1$ ) and define

$$C_0 = \max_{I \in \Delta_1} a_I^{1/2 \xi(I)} / \min_{I \in \Delta_1} a_I^{1/2 \xi(I)};$$

if  $C_0 \leq C_4^3$  then there is nothing to prove. If  $C_0 > C_4^3$  we consider the following sets  $S_i, i \geq 1$ :

$$S_i = \{I \in \Delta_i : a_I^{1/2 \xi(I)} \geq \min_{I \in \Delta_1} a_I^{1/2 \xi(I)}\};$$

of course  $S_1$  starts out to be  $\Delta_1$ . Define

$$J_i = \max_{I \in \Delta_i} a_I^{1/2} \xi(I) / \min_{I \in \Delta_1} a_I^{1/2} \xi(I)$$

then it follows directly from (7.23) that

$$(7.24) \quad J_{i+1} \leq \delta J_i.$$

Let  $k \geq 1$  be the first integer so that

$$J_{k+1} \leq C_4^3.$$

Notice that  $k \leq k_0$ , where  $k_0$  is the first integer to satisfy

$$\delta^{k_0} C_0 \leq C_4^3.$$

Since  $J_i > C_4^3$ ,  $1 \leq i \leq k$ , it follows that whenever the transition from  $\Delta_i$  to  $\Delta_{i+1}$ ,  $1 \leq i \leq k$ , consists in bisecting all intervals that satisfy

$$n^{\Delta_i}(I) \geq \varepsilon \max_{I \in \Delta_i} n^{\Delta_i}(I)$$

then the newly created elements will lie in  $S_{i+1}$ . Whenever  $A^{(2)}$  bisects a fixed fraction  $\gamma$  of the elements, the newly created intervals may all fall outside  $S_{i+1}$ . In summary we get that for  $\gamma$  sufficiently small

$$(7.25) \quad \left[ \begin{array}{l} S_{k+1} \text{ contains at least} \\ (1-\gamma - (1+\gamma)\gamma - (1+\gamma)^{k-1}\gamma) \frac{N(\Delta_{k+1})}{(1+\gamma)^k} \\ \text{elements, and} \\ J_{k+1} \leq C_4^3. \end{array} \right.$$

If  $0 < \gamma_0$  is selected appropriately and  $0 \leq \gamma \leq \gamma_0$  then we immediately see from (7.25) that

$$S_{k+1} \text{ contains at least } N(\Delta_{k+1})/2 \text{ elements, and } J_{k+1} \leq C_4^3;$$

this verifies the lemma.  $\square$

Based on Lemma 7.3, it is now easy to prove

Lemma 7.4. Let  $C_4$  denote a constant

$$C_4 > \max(C_3, 1/\beta).$$

There exist integers  $i$  and  $k_0$  such that if  $0 < \gamma_0$  is taken sufficiently small and  $0 \leq \gamma \leq \gamma_0$  then the cardinality of

$$\bar{\Delta}_j = \{I \in \Delta_j : a_I^{1/2} \xi(I) \geq C_4^{-3} \max_{I \in \Delta_j} a_I^{1/2} \xi(I)\}$$

is at least  $N(\Delta_j)/2$ , for arbitrary  $\ell \in \mathbb{N}$  and some  $i+k_0 < j \leq i+k_0(\ell+1)$ .

Proof: Since the proof in many ways is very similar to the previous, we shall omit some of the details.

Pick  $0 < \gamma_0$  sufficiently small and pick  $i$  such that the assertion

of Lemma 7.3 holds. Let  $k$  be the smallest integer so that

$$(7.26) \quad \max_{I \in \Delta_{i+k}} a_I^{1/2} \xi(I) \leq C_4^{-1} \max_{I \in \Delta_i} a_I^{1/2} \xi(I).$$

Due to (7.23) it is clear that  $k \leq k_0$  where  $k_0$  satisfies

$$\delta^{k_0} \leq C_4^{-1}$$

The size of  $0 < \gamma_0$  is furthermore taken so small that

$$(7.27) \quad \frac{1}{2} - \gamma - (1+\gamma)\gamma - \dots - (1+\gamma)^{k-1} \gamma > 0, \quad 0 \leq \gamma \leq \gamma_0.$$

For  $i \leq j \leq i+k$  define

$$S_j = \{I \in \Delta_j : a_I^{1/2} \xi(I) \geq C_4^{-4} \max_{I \in \Delta_i} a_I^{1/2} \xi(I)\};$$

using (7.22), (7.26), (7.27), the definition of  $i$  and the assumption

(A.3)', we get that for  $0 \leq \gamma \leq \gamma_0$

$$(7.28) \quad S_j \text{ has cardinality at least } N(\Delta_j)/2, \quad i \leq j \leq i+k.$$

The definitions of  $S_{i+k}$  and  $\tilde{\Delta}_{i+k}$  in combination with (7.26) gives

$$S_{i+k} \subseteq \tilde{\Delta}_{i+k},$$

and the desired result follows now directly from (7.28), in the case

$\ell = 1$ ; the case of  $\ell \geq 2$  follows by induction.  $\square$

Proof of Theorem 7.2. From Lemma 7.4 we immediately conclude that there exist a  $C$  such that for any  $j$

$$\tilde{\Delta}_j = \{I \in \Delta_j : \xi(I) \geq C^{-1} \max_{I \in \Delta_j} \xi(I)\}$$

has cardinality at least  $N(\Delta_j)/2$ . According to Definition 7.1,  $\{\Delta_j\}_{j=1}^{\infty}$  is therefore semi-equilibrated, and the theorem follows from Lemma 7.1.  $\square$

## 8. A SIMPLE COST ANALYSIS

Let us assume that the cost of computing the function  $u(\Delta_i)$ , given the mesh  $\Delta_i$ , is proportional to  $N(\Delta_i)^\lambda$ ,  $\lambda > 0$  (the cost of computing the indicators  $\eta_j^{\Delta_i}$  and generating the mesh  $\Delta_{i+1}$  is considered negligible compared to this). The total cost of the feedback f.e.m., at arriving at  $u(\Delta_i)$ , may therefore be estimated by

$$(8.1) \quad \sum_{j=1}^i N(\Delta_j)^\lambda.$$

In the case of the transition operator  $A^{(2)}$  we have for any  $\gamma' < \gamma$  (and  $N(\Delta_j)$  sufficiently large) that

$$N(\Delta_j) (1 + \gamma') \leq N(\Delta_{j+1}),$$

and using the expression (8.1) we get the total cost estimated by

$$\sum_{j=0}^{i-1} (1 + \gamma')^{-j\lambda} N(\Delta_1)^\lambda$$

$$\gamma' \leq \frac{(1 + \gamma')^\lambda}{(1 + \gamma')^\lambda - 1} N(\Delta_1)^\lambda, \quad 0 < \gamma' < \gamma$$

The transition operator  $A^{(2)}$  with  $0 < \gamma \leq 1$  thus has the property that the total cost is comparable to the cost of computing the last finite element solution. The same is in general not true for the transition operators  $A^{(1)}$  and  $A^{(3)}$ : if the function  $u$  has a strong singularity, then a few indicators near this singularity will remain larger than the rest through successive refinements, and the total cost will be of a quite different magnitude than the cost of a single finite element

solution. In the program FEARS this problem has been solved for the transition operator  $A^{(3)}$  through the introduction of the aforementioned short passes (cf. section 5).

## 9. APPENDIX

The bilinear form  $B$  introduced in (2.4) coerces the norm on  $\mathring{H}^1(I)$ , actually by definition  $\|u\|_E = [B(u,u)]^{1/2}$ . Due to the form of  $B$  it is also very natural to consider it relative to the normed spaces  $\mathring{W}_\infty^1(I)$  and  $\mathring{W}_1^1(I)$ . As the norm on  $\mathring{W}_\infty^1(I)$  we choose  $\|u\|_{1,\infty} = \sup_{x \in I} \left| \frac{d}{dx} u(x) \right| + \sup_{x \in I} |u(x)|$  and as the norm on  $\mathring{W}_1^1(I)$ ,  $\|u\|_{1,1} = \int_0^1 \left| \frac{d}{dx} u \right| dx + \int_0^1 |u| dx$ .

Theorem 9.1. Let  $a$  and  $b$  satisfy (2.2a-b). Assume that  $a$  is Hölder continuous with exponent  $\theta$ ,  $0 < \theta \leq 1$ . There exist  $\delta$  and  $h_0$  such that

$$\inf_{\substack{u \in S^P(\Delta) \\ \|u\|_{1,\infty} = 1}} \sup_{\substack{v \in S^P \\ \|v\|_{1,1} = 1}} |B(u,v)| \geq \delta > 0$$

for any mesh with  $h^\Delta \leq h_0$ .

Before we give the proof of this theorem, let us show how the result implies Lemma 7.2.

Let  $u_*(\Delta) = P_*(\Delta)u \in S^P(\Delta)$  denote the solution to (2.6) in the special case  $a = 1$ ,  $b = 0$  (as in section 7). If  $e(\Delta)$  denotes  $u - u(\Delta)$  and  $e_*(\Delta) = u - u_*(\Delta)$ , then according to (2.9) and (7.3)

$$\eta_j^\Delta = \|R(I \setminus I_j^\Delta) e(\Delta)\|_E$$

$$\xi_j^\Delta = \left( \int_{I_j^\Delta} \left( \frac{d}{dx} e_*(\Delta) \right)^2 dx \right)^{1/2}.$$

For any  $v \in S^P(\Delta)$  we have that

$$\begin{aligned}
(9.1) \quad B(u(\Delta) - u_*(\Delta), v) &= B(e_*(\Delta) - e(\Delta), v) = B(e_*(\Delta), v) \\
&= \sum_{j=1}^{N(\Delta)} \int_{I_j^\Delta} a \frac{d}{dx} e_*(\Delta) \frac{d}{dx} v \, dx + \sum_{j=1}^{N(\Delta)} \int_{I_j^\Delta} b e_*(\Delta) v \, dx \\
&= \sum_{j=1}^{N(\Delta)} \int_{I_j^\Delta} (a - a_j) \frac{d}{dx} e_*(\Delta) \frac{d}{dx} v \, dx + \sum_{j=1}^{N(\Delta)} \int_{I_j^\Delta} b e_*(\Delta) v \, dx
\end{aligned}$$

where  $a_j^\Delta = a\left(\frac{x_{j-1}^\Delta + x_j^\Delta}{2}\right)$ ; in (9.1) we have used the fact that  $\int_{I_j^\Delta} \frac{d}{dx} e_*(\Delta) q \, dx$

= 0 for any  $j$ , and any polynomial  $q$  of degree  $\leq p - 1$ .

Since  $a$  is Hölder continuous with exponent  $1/2$ ,  $|a - a_j| \leq C(h_j^\Delta)^{1/2}$  on  $I_j^\Delta$ ,

so that based on (9.1)

$$\begin{aligned}
(9.2) \quad &|B(u(\Delta) - u_*(\Delta), v)| \\
&\leq C \left( \sum_{j=1}^{N(\Delta)} (h_j^\Delta)^{1/2} \sup_{x \in I_j^\Delta} \left| \frac{d}{dx} v(x) \right| \left( \int_{I_j^\Delta} \left| \frac{d}{dx} e_*(\Delta) \right| dx + \sum_{j=1}^{N(\Delta)} \sup_{x \in I_j^\Delta} |v(x)| \int_{I_j^\Delta} |e_*(\Delta)| dx \right) \right) \\
&\leq C \left( \sum_{j=1}^{N(\Delta)} \int_{I_j^\Delta} \left| \frac{d}{dx} v \right| dx \left( \int_{I_j^\Delta} \left| \frac{d}{dx} e_*(\Delta) \right|^2 dx \right)^{1/2} + \sum_{j=1}^{N(\Delta)} (h_j^\Delta)^{-1} \int_{I_j^\Delta} |v| dx \int_{I_j^\Delta} |e_*(\Delta)| dx \right) \\
&\leq C \sum_{j=1}^{N(\Delta)} \left( \int_{I_j^\Delta} \left| \frac{d}{dx} v \right| dx + (h_j^\Delta)^{1/2} \int_{I_j^\Delta} |v| dx \right) \varepsilon_j^\Delta.
\end{aligned}$$

In the last inequality we have used the fact that  $e_*(\Delta) = 0$  at the nodal points of  $\Delta$ . It follows now directly from Theorem 9.1 and the estimate

(9.2) that

$$\|u(\Delta) - u_*(\Delta)\|_{1,\infty} \leq C \max_{I \in \Delta} \xi(I).$$

Hence

$$\begin{aligned} \eta_j^\Delta &= \|R(I \setminus I_j^\Delta) e(\Delta)\|_E \\ &\geq \left[ \int_{I_j^\Delta} a \left( \frac{d}{dx} (u - u_*(\Delta)) \right)^2 dx + \int_{I_j^\Delta} b (u - u_*(\Delta))^2 dx \right]^{1/2} \\ &\quad - \left[ \int_{I_j^\Delta} a \left( \frac{d}{dx} (u_*(\Delta) - u(\Delta)) \right)^2 dx + \int_{I_j^\Delta} b (u_*(\Delta) - u(\Delta))^2 dx \right]^{1/2} \\ &\geq (a_j^\Delta)^{1/2} \xi_j^\Delta - C (h_j^\Delta)^{1/2} \max_{I \in \Delta} \xi(I), \end{aligned}$$

and similarly

$$\eta_j^\Delta \leq (a_j^\Delta)^{1/2} \xi_j^\Delta + C (h_j^\Delta)^{1/2} \max_{I \in \Delta} \xi(I).$$

In our proof of Theorem 9.1 we use the following observations.

**Lemma 9.1.** Let  $v$  be a polynomial of degree  $\leq r$ . Let  $I = (y, z)$  and  $a(x) \in L_\infty$ ,  $0 < \alpha_1 < a(x) < \alpha_2 < \infty$ . Assume that  $v$  is given by the expression

$$v(x) = \sum_{i=0}^r c_i (x-y)^i.$$

Define the norms

$$\|v\|' = \left[ \sum_{i=0}^r c_i^2 h^{2i} \right]^{1/2}$$

$$\|v\|'' = \sum_{i=0}^r |c_i| h^i$$

$$\|v\|''' = \sum_{i=0}^r \frac{|\int_I a(x)v(x)(x-y)^i dx|}{h^i}$$

with  $h = z - y$ . There exists a constant  $C$  independent of  $I$ ,  $v$ ,  $a$ , (but depending on  $\alpha_1$  and  $r$ ) such that

$$1/C \|v\|'' \leq \|v\|_{L_\infty}(I) \leq C \|v\|'$$

and

$$\|v\|_{L_1}(I) \leq C \|v\|''.$$

The proof of this lemma consists of a simple scaling argument and the fact that all norms are equivalent on a finite dimensional space. That  $C$  depends only on  $\alpha_1$  and  $r$  follows directly from this argument.

Proof of Theorem 9.1. Let us first consider the case that  $b \equiv 0$  in the definition (2.4) of  $B$ .

Given  $u \in S^p(\Delta)$  let  $I = (x_{j_0-1}, x_{j_0})$  be an interval with the property that the maximum of  $|\frac{d}{dx} u|$  is attained in  $\bar{I}$ . Let

$$\frac{d}{dx} u(x) = \sum_{i=0}^{p-1} c_i (x - x_{j_0-1})^i, \quad x \in I,$$

and define

$$z(x) = \sum_{i=0}^{p-1} d_i (x-x_{j_0-1})^i, \quad x \in I$$

where the  $d_i$ 's are selected so that

$$\int_I a(x) z(x) (x-x_{j_0-1})^i = h^{2i} c_i, \quad 0 \leq i \leq p-1.$$

Using Lemma 9.1 we thus get

$$\begin{aligned} & \int_I a(x) \frac{d}{dx} u(x) z(x) dx \\ (9.3) \quad &= \sum_{i=1}^{p-1} c_i \int_I a(x) z(x) (x-x_{j_0-1})^i dx \\ &= \sum_{i=1}^{p-1} c_i^2 h^{2i} \geq \delta \left\| \frac{d}{dx} u \right\|_{L_\infty(I)}^2 = \delta \left\| \frac{d}{dx} u \right\|_{L_\infty(I)}^2, \end{aligned}$$

and furthermore

$$\begin{aligned} (9.4) \quad & \|z\|_{L_1(I)} \leq C \|z\|_{L_1(I)}''' \\ &= C \sum_{i=0}^{p-1} \frac{|\int_I a(x) z(x) (x-x_{j_0-1})^i dx|}{h^i} \\ &= C \sum_{i=0}^{p-1} |c_i| h^i = C \left\| \frac{d}{dx} u \right\|''' \\ &\leq C \left\| \frac{d}{dx} u \right\|_{L_\infty(I)} = C \left\| \frac{d}{dx} u \right\|_{L_\infty(I)}. \end{aligned}$$

Extend  $z$  to be zero outside of  $I$  and define

$$\phi(x) = A \int_0^x \frac{1}{a(t)} dt$$

where  $A$  is chosen such that

$$\phi(1) = \int_0^1 z(x) dx.$$

It is clear that

$$|A| \leq C \|z\|_{L_1(I)}.$$

Let  $\tilde{v}$  be the piecewise linear interpolant of  $\phi$  on  $\Delta$ , then

$$\begin{aligned} (9.5) \quad \left\| \frac{d}{dx} (\phi - \tilde{v}) \right\|_{L_\infty(I)} &\leq C |A| (h^\Delta)^{\frac{1}{2}} \leq C \|z\|_{L_1(I)} (h^\Delta)^{\frac{1}{2}} \\ &\leq C \left\| \frac{d}{dx} u \right\|_{L_\infty(I)} (h^\Delta)^{\frac{1}{2}}. \end{aligned}$$

Define

$$v(x) = \int_0^x z(t) dt - \tilde{v}(x),$$

it is clear that  $v \in S^p(\Delta)$  (the construction of  $\tilde{v}$  guarantees that  $v = 0$  at the endpoints of  $I$ ), and

$$\begin{aligned} \int_0^1 a \frac{d}{dx} u \frac{d}{dx} v dx &= \int_0^1 a \frac{d}{dx} u z dx - \int_0^1 a \frac{d}{dx} u \frac{d}{dx} \phi dx \\ &\quad - \int_0^1 a \frac{d}{dx} u \left( \frac{d}{dx} \tilde{v} - \frac{d}{dx} \phi \right) dx \end{aligned}$$

$$\begin{aligned}
&= \int_0^1 a \frac{d}{dx} u z \, dx - A \int_0^1 \frac{d}{dx} u \, dx \\
&\quad - \int_0^1 a \frac{d}{dx} u \frac{d}{dx} (\tilde{v} - z) \, dx.
\end{aligned}$$

Due to the estimates (9.3), (9.5) and the fact that  $u(0) = u(1) = 0$ , we now get

$$\begin{aligned}
(9.6) \quad \int_0^1 a \frac{d}{dx} u \frac{d}{dx} v &\geq \delta \left\| \frac{d}{dx} u \right\|_{L_\infty(I)}^2 - C \left\| \frac{d}{dx} u \right\|_{L_\infty(I)}^2 (h^\Delta)^{\frac{1}{2}} \\
&\geq \delta \left\| \frac{d}{dx} u \right\|_{L_\infty(I)}^2 \\
&\geq \delta \|u\|_{1,\infty}^2
\end{aligned}$$

for  $h^\Delta$  sufficiently small. From (9.4) and (9.5) we immediately see

$$\begin{aligned}
\|v\|_{1,1} &\leq C \left\| \frac{d}{dx} v \right\|_{L_1(I)} \\
&\leq C (\|z\|_{L_1(I)} + \left\| \frac{d}{dx} \tilde{v} \right\|_{L_1(I)}) \\
&\leq C \left\| \frac{d}{dx} u \right\|_{L_\infty(I)} \leq C \|u\|_{1,\infty},
\end{aligned}$$

and a combination of this and (9.6) leads to the desired result in case  $b \equiv 0$ .

We now consider the general case of a nonzero  $b$ . Let  $v$  be the same function as before and set  $v_1 = v + \tilde{w}$ , with  $\tilde{w}$  to be picked later. Then we have

$$\begin{aligned}
B(u, v_1) &= \int_0^1 a \frac{d}{dx} u \frac{d}{dx} v \, dx + \int_0^1 b u v \, dx \\
&\quad + \int_0^1 a \frac{d}{dx} u \frac{d}{dx} \tilde{w} \, dx + \int_0^1 b u \tilde{w} \, dx.
\end{aligned}$$

Let  $w$  denote the solution to the two point boundary value problem

$$\begin{aligned}
-\frac{d}{dx} \left( a \frac{d}{dx} w \right) + bw &= -bv \\
w(0) &= w(1) = 0;
\end{aligned}$$

it is not difficult to see that

$$\begin{aligned}
\left| \frac{d}{dx} w(x) - \frac{d}{dx} w(x+\varepsilon) \right| &\leq C \|v\|_{1,1} \varepsilon^\theta \\
&\leq C \|u\|_{1,\infty} \varepsilon^\theta.
\end{aligned}$$

We pick  $\tilde{w}$  to be the linear interpolant to  $w$  on the mesh  $\Delta$ , then  $v_1 = v + \tilde{w} \in S^P(\Delta)$  and

$$\begin{aligned}
(9.7) \quad B(u, v_1) &\geq \int_0^1 a \frac{d}{dx} u \frac{d}{dx} v - C \|u\|_{1,\infty} \|w - \tilde{w}\|_{1,\infty} \\
&\geq \delta \|u\|_{1,\infty}^2 - C \|u\|_{1,\infty}^2 (h^\Delta)^\theta \\
&\geq \delta \|u\|_{1,\infty}^2
\end{aligned}$$

for  $h^\Delta$  sufficiently small. Similarly as before we also have that

$$(9.8) \quad \|v_1\|_{1,1} \leq C \|u\|_{1,\infty}$$

and a combination of (9.7) and (9.8) leads to the desired result.  $\square$

## REFERENCES

- [1] Babuška, I., Rheinboldt, W. C., A posteriori error estimates for the finite element method., Int. J. Numer. Meth. Eng. 12 (1978), pp. 1597-1615.
- [2] Babuška, I., and Miller A., A-posteriori Error Estimates and Adaptive Techniques for the Finite Element Method, to appear.
- [3] Babuška, I., Miller, A. and Vogelius M., Adaptive Methods and Error Estimation for Elliptic Problems of Structural Mechanics, Proceedings ARO Workshop on Adaptive Computational Methods for Partial Differential Equations, to appear SIAM Publication.
- [4] Babuška, I., and Osborn, J., Analysis of Finite Element Method for Second Order Boundary Value Problems Using Mesh Dependent Norms, Num. Math. 34 (1980), pp. 41-62.
- [5] Babuška, I., and Rheinboldt, W. C., A-Posteriori Error Analysis of Finite Element Solutions for One-Dimensional Problems, SIAM J. Num. Anal. 18 (1981), pp. 565-589.
- [6] Babuška, I., and Rheinboldt, W. C., A Survey of A-Posteriori Error Estimators and Adaptive Approaches in the Finite Element Method, Proceedings of the China-France Symposium on Finite Element Method, Beijing, China, 1982, to appear.
- [7] Babuška, I., and Rheinboldt, W. C., Adaptive Finite Element Processes in Structural Mechanics, Proceedings of the Conference of Elliptic Solvers, Monterey, California, 1983, to appear.
- [8] Bank, R. E., and Weiser, A., Some A-Posteriori Error Estimators for Elliptic Partial Differential Equations, to appear.
- [9] Childs, B., Scott, M., Daniel, J. V., Denman, E., and Nelson, P., Codes for Boundary-Value Problems in Ordinary Differential Equations, Proceedings of Working Conference, May 14-17, 1978. Lecture Notes in Computer Science No. 76, Springer-Verlag, Berlin, 1979.
- [10] Mesztenyi, C., and Szymczak, W., FEARS User's Manual for UNIVAC 1100, University of Maryland, Inst. for Physical Sci. Tech., Note BN-991.
- [11] Rheinboldt, W. C., Feedback Systems and Adaptivity for Numerical Computations, Proceedings ARO Workshop on Adaptive Computational Methods for Partial Differential Equations, to appear, SIAM Publication.
- [12] Rice, J. R., Adaptive Approximation. J. of Approx. Theory 16 (1976), pp. 329-337.

- [13] Weiser, A., Local-mesh, Local-order Adaptive Finite Element Methods with A-Posteriori Error Estimators for Elliptic Partial Differential Equations, Tech. Rep. 213, Yale University, Dept. of Comp. Sci., 1981.

LMED  
8