

AD-A148 345

A FEEDBACK FINITE ELEMENT METHOD WITH A POSTERIORI  
ERROR ESTIMATION PART 1. (U) MARYLAND UNIV COLLEGE PARK  
LAB FOR NUMERICAL ANALYSIS I BABUSKA ET AL. OCT 84

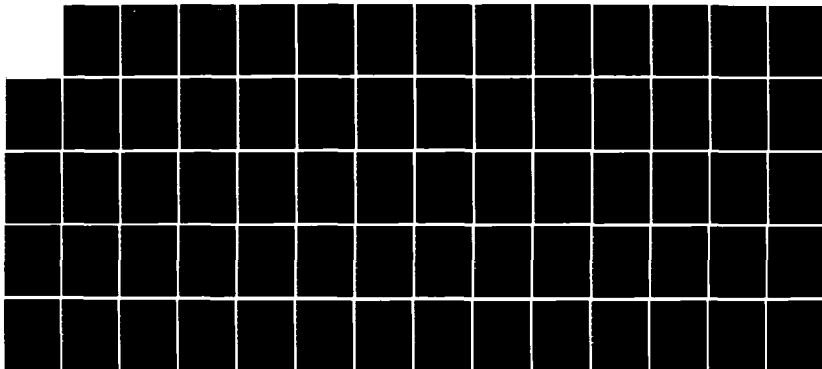
1/1

UNCLASSIFIED

BN-1031 N00014-77-C-0623

F/G 12/1

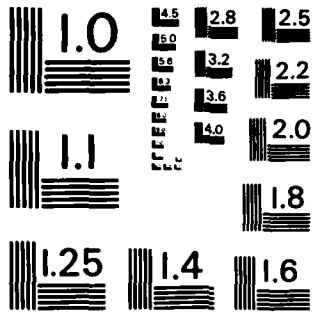
NL



END

FILED

DTIC



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



INSTITUTE FOR PHYSICAL SCIENCE  
AND TECHNOLOGY

15

Laboratory for Numerical Analysis

Technical Note EN-1031

AD-A148 345

A FEEDBACK FINITE ELEMENT METHOD WITH A POSTERIORI ERROR ESTIMATION

PART I. THE FINITE ELEMENT METHOD AND SOME BASIC PROPERTIES  
OF THE A POSTERIORI ERROR ESTIMATOR

by

I. Babuška

Institute for Physical Science and Technology  
University of Maryland

A. Miller

Centre for Mathematical Analysis  
The Australian National University

DTIC FILE COPY

This document has been approved  
for public release and its  
distribution is unlimited.

RECEIVED  
OCT 24 1984  
A

October 1984



UNIVERSITY OF MARYLAND

84 11 20 029

SECURITY CLASSIFICATION OF THIS PAGE (When Data Entered)

REPORT DOCUMENTATION PAGE		READ INSTRUCTIONS BEFORE COMPLETING FORM
1. REPORT NUMBER	2. GOVT ACCESSION NO.	3. RECIPIENT'S CATALOG NUMBER
4. TITLE (and Subtitle) A Feedback Finite Element Method with A Posteriori Error Estimation - Part I. The Finite Element Method and Some Basic Properties of the A Posteriori Error Estimator		5. TYPE OF REPORT & PERIOD COVERED Final life of the contract
7. AUTHOR(s) I. Babuška A. Miller		6. PERFORMING ORG. REPORT NUMBER
9. PERFORMING ORGANIZATION NAME AND ADDRESS Institute for Physical Science and Technology University of Maryland College Park, MD 20742		8. CONTRACT OR GRANT NUMBER(s) ONR N00014-77-C-0623
11. CONTROLLING OFFICE NAME AND ADDRESS Department of the Navy Office of Naval Research Arlington, VA 22217		10. PROGRAM ELEMENT, PROJECT, TASK AREA & WORK UNIT NUMBERS
14. MONITORING AGENCY NAME & ADDRESS (if different from Controlling Office)		12. REPORT DATE October 1984
		13. NUMBER OF PAGES 68
		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 is the first in a series of three in which we discuss some theoretical and practical aspect of a feedback finite element method for solving systems of linear second order elliptic partial differential equations (with particular interest in classical linear elasticity). In this first part we introduce some nonstandard finite element spaces, though based on the usual square bilinear elements, permit local mesh refinement. The algebraic structure of these spaces and their approximation properties are analyzed. An "equivalent estimator" for the $H^1$ finite element error is developed.		

DD FORM 1473  
1 JAN 73

EDITION OF 1 NOV 65 IS OBSOLETE

S N 0102-LF-014-6601

SECURITY CLASSIFICATION OF THIS PAGE (When Data Entered)

**A FEEDBACK FINITE ELEMENT METHOD WITH A POSTERIORI ERROR ESTIMATION**

**PART I. THE FINITE ELEMENT METHOD AND SOME BASIC PROPERTIES  
OF THE A POSTERIORI ERROR ESTIMATOR**

I. Babuška

A. Miller

Technical Note BN-1031

October 1984

This research was partially supported by ONR Contract N00014-77-C-0623.

A handwritten form with a signature and various fields. The form is tilted and has a grid-like structure. The signature "A1" is written in the bottom left corner. The form contains several lines of text, including "BIRG" and "BIRG".

Abstract

This paper is the first in a series of three in which we discuss some theoretical and practical aspect of a feedback finite element method for solving systems of linear second order elliptic partial differential equations (with particular interest in classical linear elasticity). In this first part we introduce some nonstandard finite element spaces, though based on the usual square bilinear elements, permit local mesh refinement. The algebraic structure of these spaces and their approximation properties are analysed. An equivalent estimator for the  $H^1_\Gamma$  finite element error is developed. In the second paper we shall discuss the asymptotic properties of this estimator. In the third paper we shall also report on some computational experience with the FEARS program which uses this estimator as part of a feedback loop to control mesh refinement and some of its programming features.

## §0. Introduction

The practical success or failure of many finite element computations often depends critically on the user's choices of finite element mesh and element type. As a simple illustration of this, consider the boundary value problems that arise in classical plane linear elasticity. For such problems it is well known that in the neighbourhood of certain critical boundary points (e.g. angular boundary points, or points where the boundary conditions change between specified tractions and specified displacements), the stresses exhibit some form of singular behaviour. Unless such critical points are handled carefully, the resulting finite element solution may have disappointing accuracy.

One way in which such critical points may be treated is to employ an appropriate mesh refinement strategy in the neighbourhood of the point. Broadly speaking, two kinds of refinement strategy can be identified. On the one hand, there are a priori refinement techniques which grade the mesh in a manner governed either by earlier experience with similar kinds of problems or by the results of some a priori analysis of the nature of the singularity. For many problems in linear elasticity an asymptotic representation of the solution in the vicinity of the critical point is available. Using such representations it is often possible to derive a sequence of graded meshes which can be shown to converge in the energy norm at an optimal rate with respect to the number of degrees-of-freedom of the resulting discrete system. From a practical point of view however, a weakness of such a priori methods is that the analysis or experience they are based upon is asymptotic in nature and is seldom discriminating enough to tell whether a singularity, though present in theory, is going to cause significant problems at the level of accuracy that one is working. As an example consider the stress singularity

$Kr^{-1/2}g(\theta)$  typically associated with cracks in plane elasticity. The potential of this singularity to affect the accuracy of a finite element approximation will depend on  $K$ . Usually an a priori analysis can give little insight into the value of  $K$ . If  $K$  is small enough (compared to the overall level of accuracy desired), then no harm will be done by employing, say, a uniform mesh near the crack tip. Indeed, were a refined mesh used, the extra degrees of freedom would not lead to any significant improvement in accuracy over the uniform case. On the other hand, any (fixed) mesh refined near the crack tip will give a far from optimal mesh as  $K \rightarrow \infty$ . (Optimal here indicates a sense of minimum error for a given number of degrees of freedom). This kind of phenomenon becomes more pronounced when, as occurs in most practical problems, there are a number of critical boundary points in the region of interest, and a decision must be made on how the corresponding refinements should be "weighted."

The other kind of refinement strategy referred to earlier is a posteriori in character. In this approach an initial finite element solution is calculated using some mesh. This solution is then examined in some fashion, and based upon this examination a refinement of the initial mesh is decided upon. Using this new mesh, a new finite element solution is computed. This process can obviously be iterated until some stopping criterion is satisfied. This kind of feedback technique could conceivably avoid the problems mentioned above, since the computational significance of the singularities present may be able to be ascertained during the feedback process. Everything, of course, depends on the a posteriori examination carried out after each step and the refinement decision arising from it.

So far we have only mentioned mesh refinement necessitated by some form of singular behaviour of the solution. However in many other situations



proper mesh refinement can also be crucial. As a typical example consider the use of curved elements for solving problems in elasticity. In general such elements cannot exactly represent rigid body motion. Because of this, proper mesh refinement may again be essential for satisfactory results, even though the solution is very smooth. It would seem very difficult to predict the pattern of such refinement a priori.

In this series of papers we shall describe and provide an analysis of one such feedback approach. The approach we shall deal with has been implemented in a practical form in the FEARS program [see [1]]. This algorithm is based upon an a posteriori examination which involves the calculation of an a posteriori error estimator for the energy (or similar) norm of the error. This estimator is composed of elementwise error "indicators", and the refinement decision at each stage is made on the basis of the distribution of these indicators. The theoretical analysis of this method is far from complete. At the moment there are many conjectures, etc., which though convincingly demonstrated by many numerical experiments can either not be proved rigorously, or not be proved in the generality that practical experience would indicate they hold.

Our theoretical analysis will concentrate upon the error indicators and estimators, in particular, upon clarifying in what sense they "estimate" the energy norm of the error of the finite element solution. We shall show under quite weak assumptions that the estimator,  $E$  say, is an equivalent estimator for the energy norm  $|e|$  of the error  $e$ , that is

$$C^{-1}E < |e| < CE \quad (1)$$

for some constant  $C > 0$  uniformly over large class of meshes; and moreover under some further restrictions that  $E$  is also an asymptotically exact estimator for  $|e|$ , that is

$$|e| = \mathcal{O}(1+o(|e|)) \quad (2)$$

as  $|e| \rightarrow 0$ .

The properties (1) and especially (2) suggest that the estimator can be reliably used for stopping the refinement process once some desired accuracy is achieved. Experience with FEARS confirms this.

To be considered worthwhile the process of successive construction of meshes should lead to a sequence of meshes whose rate of convergence is comparable to that of the theoretically optimal mesh grading (at least when dealing with practical problems). They should also have other "optimal" properties. When a feedback process has such optimal properties, it is called an adaptive process. (For more about this see [2] [3] [4].) Experience indicates that the FEARS program implements an adaptive process.

In §1.1 of this paper we shall describe the kinds of boundary value problems that we wish to consider (essentially those related to linear elasticity). §1.2 contains a description of the finite element discretization to be employed. In §1.3-1.5 we give some properties of the corresponding finite element spaces for use later in the paper. In particular, we introduce the important concept of a  $K$ -mesh in §1.4. Although we describe a rather general set up in §1.1 and §1.2, the subsequent analysis is carried out in a more restricted setting. This has been done for the sake of clarity and simplicity of notation. The analysis of the general case can be done analogously. In §2 we derive an equivalent estimator for the energy norm of the finite element error. An important step in this is the basic error estimate of §2.1. Using this result the error is able to be localized to a small number of elements. Some technical lemmas are proved in §2.2 and §2.3, while §2.4 contains the main results of this section.

In the second paper of this series we shall deal with the asymptotic

exactness of the estimator, as well as discussing the overall performance of the algorithm.

The third part will deal more specifically with design of the FEARS program and analyse its performance in the light of the developed theory.

## 1. Formulation of the problem and its finite element solution

1.1. We shall consider the boundary value problem

$$L_i(u) \equiv - \sum_{j,k,\ell=1}^2 D_k(a_{ijkl} D_\ell u_j) + \sum_{j=1}^2 c_{ij} u_j = f_i \quad \text{in } \Omega$$

$$u_i = 0 \quad \text{on } \partial^0 \Omega, \quad (1.1)$$

$$\sum_{k=1}^2 \sigma_{ik}(u) n_k \equiv \sum_{k=1}^2 \left( \sum_{j,\ell=1}^2 a_{ijkl} D_\ell u_j \right) n_k = t_i \quad \text{on } \partial^1 \Omega \quad (i = 1, 2)$$

on a bounded domain  $\Omega \subset \mathbb{R}^2$  whose boundary  $\partial\Omega$  is made up of two disjoint parts  $\partial^0\Omega \neq \emptyset$  and  $\partial^1\Omega$ ;  $\hat{n} = (n_1, n_2)$  is the outward pointing unit normal on  $\partial\Omega$ . This problem can be cast in a variational form. If  $\partial^0\Omega$  is sufficiently smooth (piecewise smooth, say) we may define the trace of  $H^1(\Omega)$  functions on  $\partial^0\Omega$ . Let us write

$$H = \{(v_1, v_2): v_i \in H^1(\Omega)\}$$

$$H_0 = \{(v_1, v_2): v_i \in H^1(\Omega), v_i = 0 \text{ on } \partial^0\Omega\}.$$

and  $H_0$  are Hilbert spaces with respect to the norm

$$\|v\|_{1,\Omega} = \left( \sum_{i=1}^2 \|v_i\|_{1,\Omega}^2 \right)^{1/2}, \quad \text{with } \|v_i\|_{1,\Omega} \text{ being the usual (scalar) } H^1(\Omega)$$

Sobolev norm. The boundary value problem (1.1) can now be posed as: Find  $u \in H_0$  such that

$$b(u, v) = \int_{\Omega} \sum_{i=1}^2 f_i v_i \, dx + \int_{\partial^1 \Omega} \sum_{i=1}^2 t_i v_i \, ds, \quad \forall v \in H_0 \quad (1.2)$$

where the bilinear form  $b: H_0 \times H_0 \rightarrow \mathbb{R}$  is defined by

$$b(w, z) = \int_{\Omega} \left( \sum_{i,j,k,\ell=1}^2 a_{ijkl} D_\ell w_j D_k z_i + \sum_{i,j=1}^2 c_{ij} w_j z_i \right) dx.$$

We shall assume that  $\Omega$  can be represented as a collection of transformed unit squares. To this end we make the following assumption on  $\Omega$ : There is a finite number of subdomains  $\Omega_d \subseteq \Omega$  ( $d = 1, \dots, N$ ) such that

- (i)  $\Omega_d \cap \Omega_g = \emptyset$  ( $d, g = 1, \dots, N; d \neq g$ )
- (ii)  $\bar{\Omega} = \bigcup_{d=1}^N \bar{\Omega}_d$
- (iii) There is an invertible transformation  $\phi_d: \bar{\Omega}_d \rightarrow [0, 1]^2$  which, together with its inverse  $\phi_d^{-1}$ , is sufficiently smooth (bounded derivatives of all orders, say)
- (iv)  $\Omega_d = \phi_d^{-1}((0, 1)^2)$ .

The image under  $\phi_d^{-1}$  of the closed edges and corners of  $[0, 1]^2$  will be referred to as the edges and corners respectively of  $\Omega_d$ .

- (v)  $\Omega_d$  and  $\Omega_e$  ( $d \neq e$ ) can only have a single edge or a single corner in common (or else  $\bar{\Omega}_d$  and  $\bar{\Omega}_e$  are disjoint)
- (vi)  $\partial^0 \Omega$  is the union of a number of (complete) edges of subdomains.  
We shall write  $\partial_d^0(0, 1)^2 = \phi_d(\partial^0 \Omega \cap \partial \Omega_d)$  and  $\partial_d^1(0, 1)^2 = \phi_d(\partial^1 \Omega \cap \partial \Omega_d)$ .

An edge common to  $\Omega_d$  and  $\Omega_e$  will be called an interface between  $\Omega_d$  and  $\Omega_e$ .

- (vii) If  $\Gamma$  is an interface between  $\Omega_d$  and  $\Omega_e$ , then  $\phi_d$  and  $\phi_e$  must "agree" on  $\Gamma$  in the following sense: if  $s$  is an arclength measured along  $\Gamma$ , and  $s_d$  and  $s_e$  are arclengths measured along  $\phi_d(\Gamma)$  and  $\phi_e(\Gamma)$  respectively, then regarding  $s_d$  and  $s_e$  as functions of  $s$ , either  $s_d(s) = s_e(s)$  or  $s_d(s) = 1 - s_e(s)$ .

Figure 1 shows a possible partitioning of a circle into subdomains. For a suitable choice for the mappings  $\phi_d$  in the case, see [5]. If  $\Gamma$  is an

interface between  $\Omega_d$  and  $\Omega_e$ , then we can naturally identify  $\Omega_d(\Gamma)$  and  $\Omega_e(\Gamma)$  as illustrated in Fig. 1. By virtue of (vii) above this identification takes a particularly simple form when expressed in terms of local arclength.

Figure 1. Partitioning of the circle into subdomains.

This representation of  $\Omega$  induces a natural correspondence between  $H$  and the set  $M$  of  $N$  tuples of the form

$$v = (v^{(1)}, \dots, v^{(d)}, \dots, v^{(N)})$$

where

$$(i) \quad v^{(d)} = (v_1^{(d)}, v_2^{(d)}), \quad v_1^{(d)} \in H^1((0,1)^2) \quad (i = 1, 2, d = 1, \dots, N)$$

(ii) if  $\Gamma$  is an interface between  $\Omega_d$  and  $\Omega_e$  ( $d, e = 1, \dots, N; d \neq e$ )

$$v_j^{(d)} \circ \phi_d = v_j^{(e)} \circ \phi_e \quad \text{on } \Gamma. \quad (j = 1, 2)$$

This correspondence is defined by:

$$v = (v_1, v_2) \in H \rightarrow V = (\dots, (v_1 \circ \phi_d^{-1}, v_2 \circ \phi_d^{-1}), \dots) \in M,$$

and

$$W = (W^{(1)}, \dots, W^{(N)}) \in M \rightarrow w = (w_1, w_2) \in H$$

where

$$w_i = w_i^{(d)} \circ \phi_d \quad \text{on } \Omega_d, \quad (i = 1, 2, d = 1, \dots, N)$$

( $w_i \in H^1(\Omega)$  by virtue of the interface continuity requirement (ii) above).

Let  $M_0 \subseteq M$  be the set of tuples satisfying the additional condition

(iii) If  $\Gamma$  is an edge of  $\Omega_d$  and  $\Gamma \subseteq \partial^0 \Omega$  then

$$v_j^{(d)} \circ \phi_d = 0 \quad \text{on } \Gamma. \quad (j = 1, 2).$$

Clearly  $M_0 \leftrightarrow H_0$  under the above natural correspondence.  $M$  and  $M_0$  are Hilbert spaces with respect the norm

$$\|V\|_{1, (\Omega_1, \dots, \Omega_d)} = \left( \sum_{d=1}^N \sum_{i=1}^2 \|V_i^{(d)}\|_{1, (0,1)}^2 \right)^{1/2}.$$

The problem (1.2) can now be reformulated as: Find  $U \in M_0$  such that

$$B(U, V) = \sum_{d=1}^N \left( \int_{(0,1)^2} \sum_{i=1}^2 F_i^{(d)} v_i^{(d)} dx + \int_{\partial_d^1(0,1)^2} \sum_{i=1}^2 T_i^{(d)} v_i^{(d)} ds \right), \quad \forall V \in M_0 \quad (1.3a)$$

where the bilinear form  $B: M_0 \times M_0 \rightarrow \mathbb{R}$  is defined by

$$B(W, Z) = \sum_{d=1}^N \int_{(0,1)^2} \left( \sum_{i,j,k,\ell=1}^2 A_{ijkl}^{(d)} D_\ell w_j^{(d)} D_k z_i^{(d)} + \sum_{i,j=1}^2 C_{ij}^{(d)} w_j^{(d)} z_i^{(d)} \right) dx. \quad (1.3b)$$

In (1.3) we have used the notation

$$A_{ijkl}^{(d)} = \frac{1}{|E^{(d)}|} \sum_{s,t=1}^2 (a_{ijst} \circ \phi_d^{-1}) E_{tk}^{(d)} E_{sl}^{(d)},$$

$$C_{ij}^{(d)} = \frac{1}{|E^{(d)}|} (c_{ij} \circ \phi_d^{-1}),$$

$$F_i^{(d)} = \frac{1}{|E^{(d)}|} (f_i \circ \phi_d^{-1}),$$

each defined on  $(0,1)^2$ ; while on  $\partial_d^1(0,1)^2$

$$T_i^{(d)} = \frac{1}{E^{(d)}} (t_i \circ \phi_d^{-1}).$$

We have set  $E_{ij}^{(d)} = D_i(\phi_d)_j$ ,  $|E^{(d)}| = \text{determinant } (E_{ij}^{(d)}: i, j = 1, 2)$ , and

$E^{(d)} = \left( \sum_{i=1}^2 \left( \frac{d}{ds} (\phi_d)_i \right)^2 \right)^{1/2}$  with  $\frac{d}{ds}$  denoting differentiation with respect

to arclength along  $\partial_d^1(0,1)^2$ .

With regard to the coefficients and input data of (1.2) we shall suppose that the  $a_{ijkl}$ ,  $c_{ij}$  and  $f_i$  are sufficiently smooth on each  $\Omega_d$  separately (say, bounded derivatives of all orders), and that the  $t_i$  are sufficiently smooth on each subdomain edge contained in  $\partial^1\Omega$ .

We assume the symmetries

$$a_{ijkl} = a_{jilk} \quad (i, j, k, l = 1, 2),$$

so ensuring that  $b(\cdot, \cdot)$ , is a symmetric bilinear form. Additionally, we shall suppose that the bilinear form  $b(\cdot, \cdot)$  is coercive over  $H_0$ , that is, there exists  $\alpha > 0$  such that



$$|b(w,w)| > \alpha \|w\|_{1,\Omega}^2 \quad \forall w \in H_0. \quad (1.4a)$$

Further we will assume that for some  $C > 0$

$$\sum_{i,j=1}^2 a_{ijkl} \xi_i \xi_j > C(\xi_1^2 + \xi_2^2), \quad k = 1,2; \quad \xi \in \mathbb{R}^2. \quad (1.4b)$$

By virtue of our assumptions on the mappings  $\phi_d$ , the above properties transfer naturally to the transformed system. Let us explicitly note the symmetry condition  $B(U,V) = B(V,U) \quad \forall V,W \in M_0$ , the coercivity condition

$$|B(W,W)| > \alpha \|W\|_{1,(\Omega_1, \dots, \Omega_d)}^2 \quad \forall W \in M_0,$$

and

$$\sum_{i,j=1}^2 A_{ijkl} \xi_i \xi_j > C(\xi_1^2 + \xi_2^2), \quad k = 1,2; \quad \xi \in \mathbb{R}^2.$$

Note that in the case when the  $a_{ijkl}$  are discontinuous across an interface  $\Gamma$  between two subdomains, then the classical formulation (1.1) needs to be supplemented by an interface condition expressing the continuity of "tractions" across  $\Gamma$ . This condition is of course implicit in the variational formulations (1.2) and (1.3) of the problem.

The finite element approximation that we shall discuss is based upon the formulation (1.3) of the problem. If  $\tilde{M}_0 \subset M_0$ , the corresponding Galerkin approximation  $\tilde{U} \in \tilde{M}_0$  to  $U$  is defined by

$$B(\tilde{U},V) = \sum_{d=1}^N \left\{ \int_{(0,1)^2} \sum_{i=1}^2 F_i^{(d)} v_i^{(d)} dx + \int_{\partial_d^1(0,1)} \sum_{i=1}^2 T_i^{(d)} v_i^{(d)} ds \right\} \quad \forall V \in \tilde{M}_0.$$

The coercivity of  $B(\cdot, \cdot)$  ensures the existence and uniqueness of  $\tilde{U}$ . We have of course the projection property

$$\|\tilde{U}\|_{1,(\Omega_1, \dots, \Omega_d)} < C \|U\|_{1,(\Omega_1, \dots, \Omega_d)} \quad (1.5)$$

The finite element error  $U - \tilde{U}$  satisfies the usual orthogonality relation

$$B(U - \tilde{U}, V) = 0 \quad \forall V \in \tilde{M}_0 \quad (1.6a)$$

which leads to the standard kind of best approximation estimate

$$\|U - \tilde{U}\|_{1,(\Omega_1, \dots, \Omega_d)} < C \inf_{V \in \tilde{M}_0} \|U - V\|_{1,(\Omega_1, \dots, \Omega_d)} \quad (1.6b)$$

The finite element approximation  $\tilde{U}$  corresponds (under the natural correspondence between  $H$  and  $M$ ) to an approximation  $\tilde{u}$  to the solution  $u$  of (1.1), (1.2). In terms of the energy norms of the respective errors we have

$$(b(u - \tilde{u}, u - \tilde{u}))^{1/2} = (B(U - \tilde{U}, U - \tilde{U}))^{1/2}. \quad (1.7)$$

1.2. Suppose that a subdivision of  $\Omega$  into subdomains  $\Omega_1, \dots, \Omega_N$  with corresponding mappings  $\phi_1, \dots, \phi_N$  satisfying the conditions of §1.1 has been decided upon. We shall now define what we mean by a mesh  $\mathcal{D}$  on  $\Omega$ .

A mesh  $\mathcal{D}$  is an  $N$ -tuple  $(\mathcal{D}_1, \dots, \mathcal{D}_N)$  where each  $\mathcal{D}_d$  ( $d = 1, \dots, N$ ) is a partition of  $[0, 1]^2$  into closed squares  $\Delta$  (with edges parallel to the coordinate axes). Each  $\mathcal{D}_d$  is called a submesh on  $\Omega_d$  and either

$$(i) \mathcal{D}_d = \mathcal{D}_d^0 = \{[0, 1]^2\}$$

or

$$(ii) \mathcal{D}_d = \mathcal{D}_d^{(1)} \text{ is constructed from an existing submesh } \mathcal{D}_d^{(i-1)}, \text{ by replacing any } \Delta \in \mathcal{D}_d^{(i-1)} \text{ by the four congruent}$$

squares resulting from the simultaneous bisections of  $\Delta$  in the two coordinate directions (see Fig. 2).

$$\mathcal{D}^{(0)} = \{(0,1)^2\} \quad \mathcal{D}_1^{(1)} = \{\Delta_1, \dots, \Delta_4\} \quad \mathcal{D}_1^{(2)} = \{\Delta_2, \dots, \Delta_8\}$$

Figure 2. The construction of a submesh on the subdomain  $\Omega_1$ .

Each closed square  $\Delta \in \mathcal{D}_d$  is called an element of  $\mathcal{D}_d$ . We shall use  $|\Delta|$  to denote the length of a side of  $\Delta$ . Clearly  $|\Delta| = 2^{-s}$  with  $s > 0$  an integer. Further we denote  $h(\mathcal{D}_d) = \max_{\Delta \in \mathcal{D}_d} |\Delta|$ .

A point  $P$  is called a node of  $\mathcal{D}_d$  if either

- (i)  $P$  is a vertex of an element of  $\mathcal{D}_d$

or

- (ii) there is an interface between  $\Omega_d$  and  $\Omega_e$ , and  $\phi_e \circ \phi_d^{-1}(P)$  is a vertex of an element of  $\mathcal{D}_e$ .

Nodes  $P$  of  $\mathcal{D}_d$  are classified as  $\mathcal{D}$ -proper if either

- (i)  $P \in (0,1)^2$  and whenever  $P \in \Delta'$  for  $\Delta' \in \mathcal{D}_d$ , then  $P$  is a vertex of  $\Delta'$

or

- (ii)  $\phi_d^{-1}(P) \in \partial\Omega$

or

- (iii)  $\phi_d^{-1}(P)$  lies on an interface between  $\Omega_d$  and  $\Omega_e$ ,  $P$  is the vertex of an element of  $\mathcal{D}_d$ , and  $(\phi_e \circ \phi_d^{-1})(P)$  is the vertex of an element of  $\mathcal{D}_e$ .

element of  $\mathcal{D}_e$ .

If a node is not  $\mathcal{D}$ -proper, it is said to be  $\mathcal{D}$ -improper. The cases (i), (ii), (iii) above are clearly mutually exclusive, and we shall further classify the  $\mathcal{D}$ -proper nodes as interior, boundary or interface nodes depending upon whether (i), (ii) or (iii) applies. Fig. 3 shows an example of mesh with proper and improper nodes indicated.

Each of the straight line segments  $\gamma$

Figure 3. A mesh on  $\Omega$  ( $\leftrightarrow$  denotes the natural identification of points across interfaces).

forming the boundary of an element of  $\mathcal{D}_d$  will be called an edge. For definiteness an edge will be assumed to be closed, that is, it includes its endpoints.  $|\gamma|$  will denote the length of an edge. An edge is called a primitive edge if it contains no nodes other than its endpoints.

We associate with a mesh  $\mathcal{D}$  the finite element subspace  $\mathcal{N}(\mathcal{D})$  consisting of all tuples of the form

$$W = (W^{(1)}, \dots, W^{(N)})$$

where

$$(i) \quad W^{(d)} = (W_1^{(d)}, W_2^{(d)}), \quad W_j^{(d)} \in H^1((0,1)^2) \quad (j=1,2; d=1, \dots, N)$$

$$(ii) \quad W_j^{(d)} \text{ is bilinear on each element } \Delta \in \mathcal{D}_d \quad (j=1,2; d=1, \dots, N)$$

(iii) if  $\Gamma$  is an interface between  $\Omega_d$  and  $\Omega_e$

$$W_j^{(d)} \circ \phi_d = W_j^{(e)} \circ \phi_e \quad \text{on } \Gamma, \quad (j=1,2; d,e=1, \dots, N).$$

Clearly  $\tilde{M}(\mathcal{D}) \subseteq M$ . Define  $\tilde{M}_0(\mathcal{D}) = \tilde{M}(\mathcal{D}) \cap M_0$ .

The finite element subspaces we have introduced are non-standard since improper nodes are permitted. (Note however that the spaces remain conforming). These spaces permit local mesh refinement, yet maintain many of the desirable programming characteristics of the more usual square or rectangular elements.

Having defined a general framework, we shall in the analysis for the remainder of this paper only consider the particular case of one subdomain ( $N=1$ ) and  $\phi$  the identity mapping (so  $\Omega = (0,1)^2$ ). This restriction is for notational simplicity only. Our results extend quite naturally to the case  $N > 1$  and  $\phi$  general, though with a considerable growth in notation. In the light of this simplification we shall from now on suppress the subdomain index and just write  $\Omega$ ,  $M$ ,  $\partial^0(0,1)^2$ , etc. There is also no need now to distinguish between the original problem (1.2) and the transformed problem (1.3). For definiteness we shall from now on use the notation of the original problem.

To further contain notation, where no confusion is possible we shall not notationally distinguish between vector valued functions and their components. In such instances all operations, relations, etc., are to be understood in a componentwise sense.

1.3. On a matter of terminology we shall say that a line  $[P_1, P_2]$  is a binary segment of another line  $[Q_1, Q_2]$  if

$$P_1 = Q_1 + \frac{k}{2^m} (Q_2 - Q_1)$$

$$m = 0, 1, \dots; k = 0, 1, \dots, 2^m - 1$$

$$P_2 = Q_1 + \frac{k+1}{2^m} (Q_2 - Q_1).$$

Lemma 3.1.

(a) If  $\Delta$ ,  $\Delta^*$  are distinct elements of  $\mathcal{D}$  with  $|\Delta| < |\Delta^*|$ , then one and only one of the following holds:

- (i)  $\Delta \cap \Delta^* = \emptyset$
- (ii)  $\Delta$  and  $\Delta^*$  share only one point. This point being a proper node (and hence also a vertex of both  $\Delta$  and  $\Delta^*$ ).
- (iii)  $\Delta \cap \Delta^*$  is an edge of  $\Delta$ . This edge is a binary segment of an edge of  $\Delta^*$ .

(b) If  $Q$  is an improper node, then  $Q$  lies in the interior (i.e., is not an endpoint) of an edge of a unique element  $\Delta^*$ . Furthermore, if  $Q$  is the vertex of  $\Delta \in \mathcal{D}$  then  $|\Delta| < |\Delta^*|$ .

Proof. (a) These results follow readily by an induction based on the refinement process used to construct  $\mathcal{D}$ .

(b) Since  $Q$  is improper, there is an element  $\Delta^*$  with  $Q \in \Delta^*$  but  $Q$  is not a vertex of  $\Delta^*$ . It is readily seen that  $\Delta^*$  is unique. Case (iii) of (a) applies. Assuming that  $|\Delta| \geq |\Delta^*|$  leads to a contradiction since  $Q$  is not a vertex of  $\Delta^*$ .  $\square$

Lemma 3.2. Suppose  $\mathcal{D}$  is a mesh and that  $z \in L_2((0,1)^2)$  is bilinear on each  $\Delta \in \mathcal{D}$ . If  $z$  has a well defined limiting value at each interior node (i.e.,  $z$  is continuous in these points) then  $z \in \tilde{M}(\mathcal{D})$ .

Proof. It suffices to show that  $z \in C^0((0,1)^2)$ . Let  $x \in (0,1)^2$ . If  $x$  lies in the interior of an element then  $z$  is obviously continuous at  $x$ . The other possibility is that  $x$  is a common boundary point of (at least) two elements,  $\Delta$  and  $\Delta^*$  say. There are two cases to consider here: (i)  $x$  is a vertex of  $\Delta$  or  $\Delta^*$ . In this case the hypothesis guarantees continuity. (ii)  $x$  is not a vertex of  $\Delta$  or  $\Delta^*$ . Case (iii) of Lemma 3.1(a) must then apply. Quite generally we may suppose that  $|\Delta| < |\Delta^*|$ , so  $\Delta^* \supset \Delta$  is an edge,  $[P_1, P_2]$  say, of  $\Delta$  and is contained in an edge of  $\Delta^*$ . The limit of  $z$  on  $[P_1, P_2]$  from within  $\Delta$  is a linear function, as is the limit from within  $\Delta^*$ . But, by our hypothesis, at  $P_1$  and  $P_2$  these limits must agree. Thus they must agree throughout  $[P_1, P_2]$ , and so in particular at  $x$ . The continuity of  $z$  at  $x$  follows.  $\square$

Theorem 3.3.

(a) For any proper node  $P$  there is a  $\phi_P \in \tilde{M}(\mathcal{D})$  satisfying

$$\phi_P(Q) = \begin{cases} 1 & \text{if } Q = P \\ 0 & \text{if } Q \text{ is a proper node, } Q \neq P. \end{cases} \quad (3.1)$$

In addition  $\phi_P > 0$  everywhere.

(b) If  $\varphi \in \tilde{M}(\mathcal{D})$  and  $\varphi(P) = 0$  for all proper nodes  $P$ , then  $\varphi = 0$

(In particular, the  $\phi_P$  are unique).

(c)  $\{\phi_P: P \text{ a proper node}\}$  is a basis for  $\tilde{M}(\mathcal{D})$ . In fact, for any  $\varphi \in \tilde{M}(\mathcal{D})$ ,  $\varphi = \sum_P \varphi(P)\phi_P$  where  $\sum_P$  denotes summation over all proper nodes  $P$  of  $\mathcal{D}$ .

Proof. (a) We shall construct a function  $z \in L_2(0,1)^2$  which is bilinear on each  $\Delta \in \mathcal{D}$ . Arrange the elements of  $\mathcal{D}$  in order  $\Delta_1, \dots, \Delta_m$  where  $|\Delta_j| > |\Delta_{j+1}|$  ( $j = 1, \dots, m-1$ ). Suppose that  $z > 0$  has been defined on each  $\Delta_j$  for  $j < n$ . To define  $z$  on  $\Delta_n$  it suffices to specify the limiting values of  $z$  from within  $\Delta_n$  at each vertex  $Q$  of  $\Delta_n$ :

(i) If  $Q$  is a proper node, then set

$$z(Q) = \begin{cases} 1 & \text{if } Q = P \\ 0 & \text{if } Q \neq P \end{cases}$$

(ii) If  $Q$  is improper, then by Lemma 3.1(b),  $Q$  lies in the interior of an edge,  $[P_1, P_2]$  say, of some unique element  $\Delta^*$  with  $|\Delta^*| > |\Delta_n|$ . Thus  $\Delta^* = \Delta_j$  for some  $j < n$  and we may define  $z(Q) = z^*(Q) > 0$  when  $z^*(Q)$  is the limiting value of  $z$  at  $Q$  from within  $\Delta^*$ .

Thus we are able to define  $z > 0$  on  $\Delta_n$ . It is clear from this method of construction that  $z$  has a well defined limiting value at each node of  $\mathcal{D}$ , and so by Lemma 3.2,  $z \in \tilde{M}(\mathcal{D})$ . Obviously if we set  $\phi_P = z$ , then  $\phi_P > 0$  and (3.1) is satisfied.

(b) Again arrange the elements of  $\mathcal{D}$  in order of nonincreasing size. Suppose that  $\phi = 0$  on all  $\Delta_j$  for  $j < n$ . Consider any vertex  $Q$  of  $\Delta_n$ . If  $Q$  is proper then  $\phi(Q) = 0$ . On the other hand, if  $Q$  improper, then again by Lemma 3.1(b),  $Q$  lies in the interior of an edge,  $[P_1, P_2]$  say, of some unique element  $\Delta^*$  with  $|\Delta^*| > |\Delta_n|$ . Thus  $\Delta^* = \Delta_j$  with  $j < n$ . But  $\phi = 0$ , on  $\Delta_j$ , and so  $\phi(Q) = 0$ . Thus  $\phi = 0$  at all vertices of  $\Delta_n$ , and therefore  $\phi = 0$  on  $\Delta_n$ .

(c) This follows readily from (a) and (b).



Corollary 3.4.

- (a)  $\{\phi_P: P \notin \partial^0(0,1)^2\}$  is a basis for  $M_0(\mathcal{D})$ .
- (b)  $\sum_P \phi_P = 1, \quad |\phi_P| < 1.$
- (c) On any  $\Delta \in \mathcal{D}, \quad |D_j \phi_P| < \frac{1}{|\Delta|}. \quad (j = 1, 2).$

Proof. This follows readily from the Theorem 3.

Theorem 3.5.

(a)  $\text{Supp}(\phi_P) = \text{closure } \{x \in [0,1]^2: \phi_P(x) \neq 0\}$  is the union of a (whole) number of (closed) elements.

(b)  $\text{Supp}(\phi_P)$  is "connected" in the sense that if  $\Delta', \Delta'' \subseteq \text{supp } \phi_P$  then there exists a sequence of elements  $\Delta' = \Delta_0, \dots, \Delta_s = \Delta''$  such that:

- (i)  $\Delta_j \subseteq \text{supp } \phi_P \quad (j = 0, \dots, s)$
- (ii) The pair  $\Delta_j, \Delta_{j+1}$  share an edge of the smaller of the pair  $(j = 0, \dots, s-1).$

Proof. (a) This follows from the simple observation that if a bilinear function on an element  $\Delta$  vanishes on  $K \subseteq \Delta$ , then either  $K = \Delta$  or  $K$  is a (one dimensional) curve.

(b) Let  $\text{supp } \phi_P = \bigcup_{k=1}^r \Delta_k$ . Assume that  $|\Delta_k| > |\Delta_{k+1}|, k = 1, \dots, r-1$ . It will suffice to show the result in the case  $\Delta' = \Delta_1$ .

First we show that one of the vertices of  $\Delta_1$  must be  $P$ . Suppose that this is not the case, and let  $Q$  be any vertex of  $\Delta_1$ . If  $Q$  is proper,  $\phi_P(Q) = 0$ . So at least one vertex,  $Q$  say, of  $\Delta_1$  must be improper with  $\phi_P(Q) \neq 0$ . By Lemma 3.1(b), there therefore exists an element  $\Delta^*$  with  $Q$  lying inside an edge of  $\Delta^*$  and  $|\Delta^*| > |\Delta_1|$ . Since  $Q \in \Delta^*$  with  $\phi_P(Q) \neq 0$ , it follows that  $\Delta^* \subseteq \text{supp } \phi_P$ . But this contradicts the maximality of  $\Delta_1$ .

Next notice that any element with  $P$  as a vertex can be "connected" to  $\Delta_1$  either directly (if (iii) of Lemma 3.1(a) applies) or indirectly by way of one intermediate element (if (ii) of Lemma 3.1(a) applies).

We now prove the result by induction on  $k$ . Suppose that we can "connect"  $\Delta_k$  to  $\Delta_1$  for all  $k < j$ , and consider  $\Delta_{j+1}$ . If  $P$  is a vertex of  $\Delta_{j+1}$ , then by what was said above,  $\Delta_{j+1}$  can be "connected" to  $\Delta_1$ . If  $P$  is not a vertex then at least one vertex,  $Q$  say, of  $\Delta_{j+1}$  must be improper with  $\phi_P(Q) \neq 0$ . From Lemma 3.1(b) it follows that there exists  $\Delta^* \in \mathcal{D}$  with  $Q \in \Delta^*$ ,  $|\Delta^*| > |\Delta_{j+1}|$  and  $\Delta^* \cap \Delta_{j+1}$  an edge of  $\Delta_{j+1}$ . Clearly  $\Delta^* \subseteq \text{supp } \phi_P$ , so  $\Delta^* = \Delta_k$  for some  $k < j$ . Thus  $\Delta_{j+1}$  can be "connected" to  $\Delta_1$  by appending  $\Delta_{j+1}$  to the connecting chain for  $\Delta_k$ .  $\square$

For any proper node  $P$  we shall refer to  $\Omega_P = \text{Interior}(\text{supp}(\phi_P))$  as the star of  $P$ . We shall also make use of the following notation:

$$\mathcal{D}_P = \{\Delta \in \mathcal{D} : \Delta \subseteq \text{supp}(\phi_P)\},$$

$$\Gamma_P^* = \cup\{\gamma : \gamma \text{ and edge of some } \Delta \in \mathcal{D}_P, \gamma \not\subseteq \partial\Omega_P\},$$

$$\Gamma_P^1 = \partial\Omega_P \cap \partial^1(0,1)^2,$$

$$\Gamma_P^0 = \partial\Omega_P - \Gamma_P^1,$$

$$\Gamma_P = \Gamma_P^* \cup \Gamma_P^1.$$

The three tuple  $\langle P, \mathcal{D}_P, \Gamma_P^0 \rangle$  will be called the star tuple of  $P$ . Note that the star uniquely determines  $\Omega_P$ ,  $\Gamma_P^*$ ,  $\Gamma_P^1$  and  $\Gamma_P$ . For instance

$$\Omega_P = \text{Int}\left(\bigcup_{\Delta \in \mathcal{D}_P} \Delta\right),$$

$$\Gamma_P^* = \bigcup_{\Delta \in \mathcal{D}_P} (\partial\Delta - \partial\Omega_P).$$

Furthermore observe that whether or not an edge contained in  $\Gamma_P$  is primitive can be determined from  $\mathcal{D}_P$  alone.

Let us also introduce the notation

$$\Gamma^* = \cup \{ \gamma : \gamma \text{ and edge of some } \Delta \in \mathcal{D}, \gamma \not\subseteq \partial(0,1)^2 \},$$

$$\Gamma = \Gamma^* \cup \partial^1 \Omega.$$

Lemma 3.6.

(a) If  $\Delta \in \mathcal{D}$  (and  $\mathcal{D}$  contains more than one element), then  $\Delta \in \mathcal{D}_P$  for some interior proper node  $P$ .

$$(b) \quad \Gamma^* = \bigcup_P \Gamma_P^* \\ \text{(all proper nodes)}$$

$$(c) \quad \partial^1 \Omega = \bigcup_P \Gamma_P^1 \\ \text{(all proper nodes)}$$

Proof. (a) An easy induction on the refinement process used to construct the mesh shows that for any element at least one vertex is an interior proper node.

(b) Clearly  $\Gamma_P^* \subseteq \Gamma^*$ , so  $\bigcup_P \Gamma_P^* \subseteq \Gamma^*$ . For the reverse inclusion: Suppose  $\gamma$  is an edge of  $\Delta$  with  $\gamma \not\subseteq \partial \Omega$ . Let  $x$  be the midpoint of  $\gamma$ . By Corollary 4(b),  $\sum_P \phi_P(x) = 1$ . So for at least one proper node  $P$ ,  $\phi_P(x) \neq 0$ . It follows that  $\Delta \in \mathcal{D}_P$  and  $\gamma \not\subseteq \partial \Omega_P$ .

(c) Again it is obvious that  $\bigcup_P \Gamma_P^1 \subseteq \partial^1 \Omega$ . To show the reverse. Let  $\gamma$  be an edge of  $\Delta \in \mathcal{D}$ . Then by (a) of the lemma,  $\Delta \in \mathcal{D}_P$  for some proper node  $P$ .  $\square$

Suppose  $\gamma$  is a straight line segment which when extended to infinity in either direction divides  $\mathbb{R}^2$  into half planes  $\pi^{(+)}$  and  $\pi^{(-)}$ . Let  $x \in \gamma$

and suppose  $v = (v_1, v_2)$  is a sufficiently smooth function defined in  $S^+ = \pi^+ \cap \{y: |x-y| < \epsilon\}$  and  $S^- = \pi^- \cap \{y: |x-y| < \epsilon\}$ . Let  $\hat{n}^+$  and  $\hat{n}^{(-)}$  be the unit normals to  $\gamma$  pointing out of  $\pi^{(+)}$  and  $\pi^{(-)}$  respectively. We shall define  $[[v \cdot \hat{n}]]$  at  $x$  by

$$[[v \cdot \hat{n}]] = \lim_{\substack{y \rightarrow x \\ y \in S^{(+)}}} v(y) \cdot \hat{n}^{(+)} + \lim_{\substack{y \rightarrow x \\ y \in S^{(-)}}} v(y) \cdot \hat{n}^{(-)}.$$

Lemma 3.7. Suppose  $z \in C^0(\Omega_P)$  is bilinear on each  $\Delta \in \mathcal{D}_P$ . If

$$[[\nabla z \cdot \hat{n}]] = 0 \text{ on each } \gamma \in \Gamma_P^* \text{ then } z \text{ is bilinear in } \Omega_P.$$

Proof. Let  $\Delta'$  be any element of  $\mathcal{D}_P$ . The function  $z$  is bilinear on  $\Delta'$ . Extend this bilinear function to all of  $\Omega_P$ , and call this extension  $z^*$ . We claim that  $z = z^*$  on all elements of  $\mathcal{D}_P$ . To show this it suffices by Theorem 3.5(b) to show that if  $z = z^*$  on  $\Delta_1$  and the pair  $\Delta_1, \Delta_2$  share an edge,  $\gamma$  say, of the smaller of the pair, then  $z = z^*$  on  $\Delta_2$ . But this follows readily since (i)  $z$  is bilinear on  $\Delta_2$ , (ii)  $z$  is continuous across  $\gamma$ , and (iii)  $\gamma \in \Gamma_P^*$  and so the directional derivative of  $z$  normal to  $\gamma$  is continuous across  $\gamma$ .  $\square$

For any  $\Delta \in \mathcal{D}$  define  $\Omega_\Delta$ , the influence region for  $\Delta$  by

$$\Omega_\Delta = \text{Interior} \left( \bigcup_{\substack{P \\ (\Delta \in \mathcal{D}_P)}} \text{supp } \phi_P \right).$$

Introduce the notation  $\mathcal{D}_\Delta = \{\Delta^* \in \mathcal{D}: \Delta^* \in \bigcup_{\substack{P \\ \Delta \in \mathcal{D}_P}} \text{supp } \phi_P\}$  and  $\Gamma_\Delta = \{\gamma: \gamma \text{ an edge of some } \Delta^* \in \mathcal{D}_\Delta, \gamma \subseteq \partial^0 \Omega\}$ . The three tuple  $\langle \Delta, \mathcal{D}_\Delta, \Gamma_\Delta \rangle$  will be called the influence tuple of  $\Delta$ .

1.4. We shall now introduce a restriction on the "spread" of  $\text{supp}(\phi_p)$ . This restriction will be essential for our analysis later on.

For any set  $S \subseteq \mathbb{R}^2$  let  $|S| = \max_{i=1,2} \sup_{x,y \in S} |x_i - y_i|$ .

Let  $K > 1$ . A mesh  $\mathcal{D}$  is called a  $K$ -mesh if for each proper node  $P$

$$|\text{supp}(\phi_p)| < K \min_{\Delta \in \mathcal{D}_p} |\Delta|. \quad (4.1)$$

For any  $K > 1$  there are certainly meshes that are not  $K$ -meshes. For instance, for meshes of the type shown in Fig. 4(a),  $\text{supp}(\phi_p)$  is always the

Figure 4(a). Example of a mesh which is not a  $K$ -mesh.

shaded region, yet the minimum element size can be made arbitrarily small by continuing the refinement process sufficiently far. The  $K$ -mesh property is, in some sense, only a local property as it permits many natural forms of mesh grading. For instance, the grading pattern of Figs. 3(b) can be continued indefinitely without violating (4.1) with  $K = 2$ .

Figure 4(b). Example of a K-mesh with refinement.

[It is conjectured that the above definition of a K-mesh is equivalent to: There is  $L > 1$  such that if  $\Delta, \Delta' \in \mathcal{D}$  and  $\Delta \cap \Delta' \neq \emptyset$  then

$$L^{-1} < |\Delta'|/|\Delta| < L. \quad (4.2)$$

By equivalence here we mean that if (4.1) holds then (4.2) holds with  $L = L(K)$ ; and if (4.2) holds then (4.1) also holds with  $K = K(L)$ .]

In what follows we shall always assume that the mesh  $\mathcal{D}$  is a K-mesh. We now state some important properties of K-meshes.

Lemma 4.1. Suppose  $\mathcal{D}$  is a K-mesh, then there exist  $C = C(K) > 0$  such that

- (i) If  $\Delta \in \mathcal{D}_p$ ,  $C|\Delta| > \max_{\Delta \in \mathcal{D}_p} |\Delta|$ .
- (ii) If  $P$  is a proper node,  $\text{card } \mathcal{D}_p < C$ .
- (iii) If  $\Delta \in \mathcal{D}$ ,  $\text{card}\{P: \Delta \in \mathcal{D}_p\} < C$ .
- (iv) If  $\gamma$  is an edge of an element of  $\mathcal{D}$ ,  $\text{card}\{P: \gamma \subseteq \Gamma_p\} < C$ .

- (v) If  $P$  is a proper node,  $\text{card}\{Q: D_P \cap D_Q \neq \emptyset\} < C$ .
- (vi) If  $\Delta^* \in \mathcal{D}$ ,  $\text{card}\{\Delta: \Delta^* \in \mathcal{D}_\Delta\} < C$ .
- (vii) If  $\gamma$  is an edge of an element of  $\mathcal{D}$ , then there are at most  $C$  nodes on  $\gamma$ .

Proof.

- (i)  $\max_{\Delta \in \mathcal{D}_P} |\Delta| < |\text{supp } \phi_P| < K \min_{\Delta \in \mathcal{D}_P} |\Delta|$ .
- (ii)  $\text{card}(\mathcal{D}_P) < \frac{|\text{supp } \phi_P|^2}{\min_{\Delta \in \mathcal{D}_P} |\Delta|^2} < K^2$ .
- (iii) If  $\Delta \in \mathcal{D}_P$ , then  $|\text{supp } \phi_P| < K \min_{\Delta^* \in \mathcal{D}_P} |\Delta^*| < K|\Delta|$ . So  
 $\bigcup_P (\text{supp } \phi_P) \subseteq Q_\Delta$  where  $Q_\Delta$  is a square with centre at the  
 $(\Delta \in \mathcal{D}_P)$   
 center of  $\Delta$  and side  $2(K - 1/2) |\Delta|$ .
- Now if  $\Delta' \in \bigcup_P \mathcal{D}_P$ , then  $\Delta'$  and  $\Delta$  are elements of the  
 $(\Delta \in \mathcal{D}_P)$   
 same  $\mathcal{D}_P$  for at least one proper node  $P$ . So

$$|\Delta'| > \frac{1}{K} |\text{supp } \phi_P| > \frac{1}{K} |\Delta|.$$

$$\text{Thus } \text{card}\left(\bigcup_P \mathcal{D}_P\right) < \frac{|Q_\Delta|^2}{\left(\frac{1}{K} |\Delta|\right)^2} < 4K^2 \left(K - \frac{1}{2}\right)^2.$$

$$(\Delta \in \mathcal{D}_P)$$

But certainly  $\text{card}\left(\bigcup_P \mathcal{D}_P\right) > \frac{1}{4} \text{card}\{P: \Delta \in \mathcal{D}_P\}$  since  $P$   
 $(\Delta \in \mathcal{D}_P)$

must be the vertex of at least one element of  $\mathcal{D}_P$  and, of course,  
 no element has more than four vertices. Thus

$$\text{card}\{P: \Delta \in \mathcal{D}_P\} < 16 K^2 \left(K - \frac{1}{2}\right)^2.$$

- (iv)  $\gamma$  can be an edge of at most two elements. By (iii) each of these elements is contained in  $\mathcal{D}_P$  for at most  $C$  proper nodes  $P$ . Thus  $\gamma$  can be contained in  $\Gamma_P$  for at most  $2C$  proper nodes  $P$ .
- (v) By (ii) there are at most  $C$  elements in  $\mathcal{D}_P$ , while, by (iii), each such element in turn belongs to a  $\mathcal{D}_Q$  for no more than  $C$  proper nodes  $Q$ . Thus  $\mathcal{D}_P \cap \mathcal{D}_Q \neq \emptyset$  for at most  $C^2$  proper nodes  $Q$ .
- (vi) By (iii),  $\Delta^*$  is contained in  $\text{supp } \phi_P$  for at most  $C$  proper nodes  $P$ . For each such node in turn,  $\text{supp } \phi_P$  can contain at most  $C$  elements  $\Delta$ . Thus,  $\text{card}\{\Delta: \Delta^* \in \mathcal{D}_\Delta\} < C^2$ .
- (vii) If  $\gamma$  is on  $\partial(0,1)^2$  then the endpoints of  $\gamma$  are the only nodes on  $\gamma$ . Otherwise suppose  $\gamma$  is an edge of  $\Delta$ . Since  $\sum_P \phi_P = 1$ , at, say, the midpoint of  $\gamma$ , then for at least one proper node  $P$ ,  $\phi_P > 0$  at the midpoint of  $\gamma$ . In fact it follows that  $\phi_P > 0$  everywhere on  $\gamma$ , except perhaps at an endpoint. Any node on  $\gamma$  must then be a vertex of an element of  $\mathcal{D}_P$ . By (ii) the number of such elements is bounded.  $\square$

We would expect that many stars  $\Omega_P$  are essentially identical except for a translation and scaling. Indeed, as we shall see in Theorem 3, there are in fact only a finite number (depending on  $K$ ) of stars up to translation and rescaling. To set the scene for the result let us consider an arbitrary proper node  $P$ . Define an affine transformation  $\tau_P: \mathbb{R}^2 \rightarrow \mathbb{R}^2$ ,

$$\tau_P(x) = \frac{1}{\left(\max_{\Delta \in \mathcal{D}_P} |\Delta|\right)} (x - P).$$



It is readily verified that  $\tau_P(P) = 0$  and  $\max_{\Delta \in \mathcal{D}_P} |\tau_P(\Delta)| = 1$ . We shall consider  $\tau_P$  to map the star of  $P$  onto some "standard" star. The mapping  $\tau_P$  will be referred to as the star transformation of  $P$ . Let us set

$$\Omega_0 = \tau_P(\Omega_P),$$

$$\mathcal{D}_0 = \{\tau_P(\Delta) : \Delta \in \mathcal{D}_P\},$$

$$\Gamma_0^* = \tau_P(\Gamma_P^*),$$

$$\Gamma_0^0 = \tau_P(\Gamma_P^0),$$

$$\Gamma_0^1 = \tau_P(\Gamma_P^1),$$

$$\Gamma_0 = \Gamma_0^* \cup \Gamma_0^1,$$

$$\psi_0 = \psi_P \circ \tau_P^{-1}.$$

The two-tuple  $\langle \mathcal{D}_0, \Gamma_0^0 \rangle$  will be called the standard star tuple of  $P$ . With a harmless abuse of notation we shall refer to members of  $\mathcal{D}_0$  as elements and their (closed) sides as edges. Notice that  $\Omega_0$ ,  $\Gamma_0^*$ ,  $\Gamma_0^1$  and  $\Gamma_0$  can be reconstructed from  $\langle \mathcal{D}_0, \Gamma_0^0 \rangle$  in much the same way as  $\Omega_P$ ,  $\Gamma_P^*$ ,  $\Gamma_P^1$  and  $\Gamma_P$  can be expressed in terms of star tuple of  $P$ . Somewhat less obvious is the following result.

Lemma 4.2.  $\psi_0$  is completely determined by  $\langle \mathcal{D}_0, \Gamma_0^0 \rangle$ .

Proof. This is clearly the same as saying that if  $\mathcal{D}$  is a mesh and  $P$  is a proper node of  $\mathcal{D}$ , then  $\psi_P$  is completely determined by  $\langle P, \mathcal{D}_P, \Gamma_P^0 \rangle$ . To show this, we shall construct a function  $z \in \tilde{M}(\mathcal{D})$  using only our knowledge of  $\langle P, \mathcal{D}_P, \Gamma_P^0 \rangle$ , and then prove  $z = \psi_P$ :

Arrange the elements of  $\mathcal{D}_p$  in order of nonincreasing size,  $\Delta_1, \Delta_2, \dots, \Delta_n$ , say. Suppose that  $z$  has been defined on each  $\Delta$ ; for  $j < n$ . To define  $z$  on  $\Delta_n$  it suffices to specify the limiting values of  $z$  from within  $\Delta_n$  at each vertex  $Q$  of  $\Delta_n$ :

(I) If  $Q = P$ , then  $z(Q) = 1$ .

(II) If  $Q \neq P$ :

(i) If  $Q \in \partial\Omega_p$ , then  $z(Q) = 0$ .

(ii) If  $Q \in \Omega_p$  and  $Q$  is proper node (in the context of the  $\langle P, \mathcal{D}_p, \Gamma_p^0 \rangle$  this just means that whenever  $Q \in \Delta^*$  for  $\Delta^* \in \mathcal{D}_p$  then  $Q$  is a vertex of  $\Delta^*$ ) then  $z(Q) = 0$ .

(iii) If  $Q \in \Omega_p$  and  $Q$  is improper (just as in (ii), this can be determined from  $\mathcal{D}_p$  alone), then by Lemma 3.1(b)  $Q$  lies in the interior of an edge  $[P_1, P_2]$ , say, of some unique element  $\Delta^*$  with  $|\Delta^*| > |\Delta_n|$ . In our present case  $\Delta^* \in \mathcal{D}_p$  and so  $\Delta^* = \Delta_j$  for some  $j < n$  and we may define  $z(Q) = z^*(Q)$  where  $z^*(Q)$  is the limiting value of  $z$  at  $Q$  from within  $\Delta^*$ .

This enables us to define  $z$  on each  $\Delta \in \mathcal{D}_p$ . Define  $z$  on the remaining elements of  $\mathcal{D}$  to be identically zero. By the above construction  $z$  has a well defined limiting value at each node of  $\mathcal{D}$ , and so by Lemma 3.2,  $z \in \tilde{M}(\mathcal{D})$ . However  $z$  satisfies (3.1), and so by Theorem 3.3,  $z = \phi_p$ .  $\square$

**Theorem 4.3.** There are at most  $C = C(K)$  standard star tuples  $\langle \mathcal{D}_0, \Gamma_0^0 \rangle$ .

**Proof.** It will suffice to show that there are at most  $C$  possibilities for  $\mathcal{D}_0$ , since once  $\mathcal{D}_0$  is known there are only  $C$  possibilities for  $\Gamma_0^0$ . ( $\mathcal{D}_0$  has at most  $C$  elements and each such element has no more than four edges.  $\Gamma_0^0$  must be the union of some of these edges.) Since

$$\bigcup_{\Delta \in \mathcal{D}_0} |\Delta| < K \min_{\Delta \in \mathcal{D}_0} |\Delta| < K,$$

$\bigcup_{\Delta \in \mathcal{D}_0} \Delta$  must be contained in a square with centre 0 and side  $2K$ .

Let  $m$  be an integer such that  $2^m > K$ . Let  $Q$  be a square with centre 0 and side  $2^{m+1}$ . Define a uniform grid on  $Q$  of grid size  $2^{-m}$ . We claim that each  $\Delta \in \mathcal{D}_0$  consists of a (whole) number of grid squares. If we can show this, then the theorem will follow at once by simple combinatoric considerations.

Firstly note that for any  $\Delta \in \mathcal{D}_0$ , certainly  $|\Delta| = 2^{-\xi}$  for some integer  $\xi > 0$ . Moreover,

$$|\Delta| > \min_{\Delta' \in \mathcal{D}_0} |\Delta'| > \frac{1}{K} \bigcup_{\Delta' \in \mathcal{D}_0} |\Delta'| > 2^{-m}$$

and, of course,

$$|\Delta| < \max_{\Delta' \in \mathcal{D}_0} |\Delta'| = 1.$$

In particular, if  $(0,0)$  is a vertex of  $\Delta$ , then  $\Delta$  must consist of a whole number of grid squares. To prove this result for an arbitrary element of  $\mathcal{D}_0$ , it will suffice by Theorem 3.5(b) to show that whenever the pair  $\Delta^*, \Delta^{**} \in \mathcal{D}_0$  are connected in the sense of Theorem 3.5(b) and  $\Delta^*$  consists of a whole number of grid squares, then so also  $\Delta^{**}$ . But this follows readily on recalling (iii) of Lemma 3.1.  $\square$

We also shall need standard forms for the influence tuples. To this end for each  $\Delta \in \mathcal{D}$  define an affine transformation  $\tau_\Delta: \mathbb{R}^2 \rightarrow \mathbb{R}^2$

$$\tau_\Delta(x) = \frac{1}{|\Delta|} \left( x - x_\Delta + \left( \frac{1}{2}, \frac{1}{2} \right) \right)$$

where  $x_\Delta$  is the centre of  $\Delta$ . It is easily seen that  $\tau_\Delta(\Delta) = [0,1]^2$ . Just as for the case of star tuples we shall consider  $\tau_\Delta$  as mapping the influence tuple of  $\Delta$  onto some standard reference tuple. Let us set

$$\Omega_* = \tau_\Delta(\Omega_\Delta),$$

$$\mathcal{D}_* = \{\tau_\Delta(\Delta^*): \Delta^* \in \mathcal{D}_\Delta\},$$

$$\Gamma_* = \tau_\Delta(\Gamma_\Delta).$$

The two tuple  $\langle \mathcal{D}_*, \Gamma_* \rangle$  will be called the standard influence tuple of  $\Delta$ , and clearly  $\Omega_*$  can be reconstructed from the standard tuple:

$$\Omega_* = \text{Int} \left( \bigcup_{\Delta \in \mathcal{D}_*} \Delta \right).$$

Akin to Theorem 4.3 we have

Theorem 4.4. There are at most  $C = C(K)$  standard influence tuples  $\langle \mathcal{D}_*, \Gamma_* \rangle$ .

Proof. It will suffice to show that there are at most  $C$  possibilities for  $\mathcal{D}_*$ , since once  $\mathcal{D}_*$  is known there are only  $C$  possibilities for  $\Gamma_*$ . (By (ii) and (iii) of Lemma 4.1,  $\mathcal{D}_*$  has at most  $C$  elements;  $\Gamma_*$  must be some union of the edges of these elements.)

To prove the claim for  $\mathcal{D}_*$ , we apply an almost identical argument to that used in Theorem 4.3.  $\square$

1.5. For later use we shall need to know something about the approximation properties of the finite element subspaces  $\tilde{M}(\mathcal{D})$  and  $\tilde{M}_0(\mathcal{D})$ . In this section we shall prove some results in this direction. The methods of proof are more or less standard, the only real difference arising from the fact that the star

$\Omega_p$  may spread further than just those elements adjacent to  $P$ . The  $K$ -mesh assumption however restricts this spread and allows us to retrieve much of the standard theory.

Lemma 5.1. Let  $\langle \mathcal{D}_*, \Gamma_* \rangle$  be a standard influence tuple. There is a constant  $C = C(\mathcal{D}_*, \Gamma_*) > 0$  such that if  $z \in H^r(\Omega_*)$  ( $r = 1, 2$ ) and  $z = 0$  on  $\Gamma_*$  then

$$\inf_{q \in Q} \|z - q\|_{r, \Omega_*} < C \|z\|_{r, \Omega_*} \quad (5.1)$$

where  $Q$  is the set of all polynomials of degree  $r - 1$  on  $\Omega_*$  which vanish on  $\Gamma_*$ . By  $\|\cdot\|_{r, \Omega}(\|\cdot\|_{r, \Omega})$  we denote the usual Sobolev norm (seminorm).

Proof. Suppose not, then we can find  $z^{(n)}$  ( $n = 1, 2, \dots$ )

$$\inf_{q \in Q} \|z^{(n)} - q\|_{r, \Omega_*} > n \|z^{(n)}\|_{r, \Omega_*}. \quad (5.2)$$

Without loss of generality we may as well suppose that,

$$\inf_{q \in Q} \|z^{(n)} - q\|_{r, \Omega_*} = 1 \quad (5.3)$$

(by taking suitable multiples of the original  $z^{(n)}$  if necessary) and

$$\|z^{(n)}\|_{r, \Omega_*} < 2 \quad (5.4)$$

(by adding a suitable  $q^n \in Q$  to each of the original  $z^{(n)}$  if necessary).

From (5.4) we can conclude using Rellich's Lemma that a subsequence of the  $z^{(n)}$  converges in  $H^{r-1}(\Omega_*)$ . We can suppose that this subsequence is the entire sequence (by deleting members of the original sequence if necessary).

But from (5.2),  $\lim_{n \rightarrow \infty} \|z^{(n)}\|_{r, \Omega_*} = 0$ . Thus  $z^{(n)}$  converges in  $H^r(\Omega_*)$ .

Let  $z^{(\infty)}$  be the limit. Obviously,  $\|z^{(\infty)}\|_{r, \Omega_*} = 0$ , and as  $\Omega_*$  is connected

(this fact follows readily from Theorem 3.5(b)), the only possibility is for  $z^{(\infty)}$  to be a polynomial of degree  $r - 1$  on  $\Omega_*$ . Moreover,  $z^{(\infty)} = 0$  on  $\Gamma_*$  (since taking traces is a continuous operation in  $H^r(\Omega^*)$ ). Thus  $z^{(\infty)} \in Q$ . But this contradicts the limiting form of (5.3).  $\square$

**Theorem 5.2.** There is a constant  $C > 0$  such that for any  $z \in H^r(\Omega) \cap H_0$  ( $r = 1, 2$ ) and any  $K$ -mesh  $\mathcal{D}$ , there exists a function  $\pi z \in \tilde{M}_0(\mathcal{D})$  such that

$$|z - \pi z|_{s, \Delta} < C |\Delta|^{r-s} |z|_{r, \Omega_\Delta}, \quad (s = 0, 1; \Delta \in \mathcal{D}). \quad (5.1)$$

**Proof.** Let  $p(t)$  be a polynomial in one variable satisfying  $\int_0^1 p(t) dt = 1$  and  $\int_0^1 tp(t) dt = 0$ . For any  $\varepsilon > 0$  define

$$\phi_\varepsilon = \begin{cases} p(t/\varepsilon) & t \in [0, \varepsilon] \\ p(-t/\varepsilon) & t \in [-\varepsilon, 0] \\ 0 & |t| > \varepsilon. \end{cases}$$

For any proper node  $P$  of  $\mathcal{D}$  set

$$\varepsilon_P = \frac{1}{2} \min_{\tilde{\Delta} \in \mathcal{E}_P} |\tilde{\Delta}|$$

and  $S_P = \{x \in \Omega: |x_i - P_i| < \varepsilon_P, i = 1, 2\}$ . Our choice of  $\varepsilon_P$  guarantees  $S_P \subseteq \Omega_P$ . For any  $y \in H^1(\Omega)$  let

$$Y_P = \int_S y(x) \phi_{\varepsilon_P}(x_1 - P_1) \phi_{\varepsilon_P}(x_2 - P_2) dx / \int_S dx \quad (5.2)$$

and define

$$\pi y = \sum_{P \in \mathcal{D}} Y_P \phi_P.$$

Clearly (Corollary 3.4(a))  $\pi y \in \tilde{M}_0(\mathcal{D})$ .

We shall now prove (5.1). Suppose  $\Delta \in \mathcal{D}$  and let  $\langle \mathcal{D}_*, \Gamma_* \rangle$  be the standard influence tuple of  $\Delta$ , and, as usual, write  $\Omega_* = \text{Interior} \left( \bigcup_{\Delta \in \mathcal{D}_*} \Delta \right)$ . Then

$$\begin{aligned} |\pi y|_{s, \Delta} &< \sum_{\substack{P \\ P \not\subset \partial^0 \Omega}} |Y_P| |\phi_P|_{s, \Delta} \\ &< C \sum_{\substack{P \\ \Delta \in \mathcal{D}_P}} |Y_P| |\Delta|^{1-s} \end{aligned}$$

by Corollary 3.4 (b) and (c). However, from (5.2)

$$|Y_P| < \varepsilon_P^{-1} \|y\|_{0, S_P} < C |\Delta|^{-1} \|y\|_{0, S_P}$$

whenever  $\Delta \in \mathcal{D}_P$ . Thus

$$|\pi y|_{s, \Delta} < C \sum_{\substack{P \\ \Delta \in \mathcal{D}_P}} \|y\|_{0, S_P} |\Delta|^{-s} < C |\Delta|^{-s} \|y\|_{0, \Omega_\Delta} \quad (5.3)$$

by (iii) of lemma 4.1 and since  $S_P \subseteq \Omega_P \subseteq \Omega_\Delta$  when  $\Delta \in \mathcal{D}_P$ . Upon rescaling  $\Delta$  to the unit square, (5.3) becomes

$$|\pi y \circ \tau_\Delta^{-1}|_{s, (0,1)^2} < C \|y \circ \tau_\Delta^{-1}\|_{0, \Omega_*}$$

and so

$$\|y \circ \tau_\Delta^{-1} - \pi y \circ \tau_\Delta^{-1}\|_{s, (0,1)^2} < C \|y \circ \tau_\Delta^{-1}\|_{r, \Omega_*}. \quad (5.4)$$

If  $\zeta$  is a polynomial of degree  $r - 1$  on  $\Omega_\Delta$  which vanishes on  $\Gamma_\Delta$ , then  $\zeta$  can trivially be extended to all of  $\Omega$ . Direct calculation shows that  $\pi\zeta = \zeta$ . Thus

$$\begin{aligned} \|z \circ \tau_\Delta^{-1} - \pi z \circ \tau_\Delta^{-1}\|_{s, (0,1)^2} &= \inf_{\zeta} \|(z+\zeta) \circ \tau_\Delta^{-1} - \pi(z+\zeta) \circ \tau_\Delta^{-1}\|_{s, (0,1)^2} \\ &< \inf_{\zeta} \|(z+\zeta) \circ \tau_\Delta^{-1}\|_{r, \Omega_*} \end{aligned} \quad (5.5)$$

by (5.4). But clearly as  $\zeta$  ranges over all possibilities  $\zeta \circ \tau_\Delta^{-1}$  ranges over all polynomials of degree  $r - 1$  on  $\Omega_*$  which vanish on  $\Gamma_*$ . Thus by the result of Lemma 5.1,

$$\inf_{\zeta} \|z \circ \tau_\Delta^{-1} - \zeta \circ \tau_\Delta^{-1}\|_{r, \Omega_*} < C \|z \circ \tau_\Delta^{-1}\|_{r, \Omega_*}. \quad (5.6)$$

Combining (5.5) and (5.6), and rescaling back to  $\Delta$  gives the desired result (5.1) upon noting that the constant in (5.6) can be taken to be mesh independent by Theorem 4.4.  $\square$

Corollary 5.3. There is a constant  $C > 0$  such that for any  $K$ -mesh  $\mathcal{D}$ ,

(i) If  $z \in H^r(\Omega) \cap M_0$  ( $r = 1, 2$ ).

$$\|z - \pi z\|_{s, \Omega} < Ch(\mathcal{D})^{r-s} \|z\|_{r, \Omega}, \quad (s = 1, 2)$$

(ii) If  $z \in H^1(\Omega) \cap M_0$

$$\sum_P \|\phi_P(z - \pi z)\|_{1, \Omega_P}^2 < C \|z\|_{1, \Omega}^2.$$

Proof.

(i) It follows directly from the theorem that



$$\begin{aligned} \|z - \pi z\|_{s, \Omega} &< \text{Ch}(\mathcal{D})^{r-s} \left( \sum_{\Delta} \|z\|_{r, \Omega_{\Delta}}^2 \right)^{1/2} \\ &< \text{Ch}(\mathcal{D})^{r-s} \|z\|_{r, \Omega} \end{aligned}$$

by (vi) of Lemma 4.1.

(ii) If  $\Delta \in \mathcal{D}_P$ , then the theorem and Corollary 3.4(b) and (c) shows

$$\|\phi_P(z - \pi z)\|_{1, \Delta} < C \left( \max_{i=1,2} \sup_{x \in \Delta} |D_i \phi_P| \|z - \pi z\|_{0, \Delta} + \sup_{x \in \Delta} |\phi_P| \|z - \pi z\|_{1, \Delta} \right).$$

On the other hand, if  $\Delta \notin \mathcal{D}_P$  then obviously

$$\|\phi_P(z - \pi z)\|_{1, \Delta} = 0.$$

Fixing  $\Delta$ , squaring and summing over all proper nodes  $P$  gives

$$\sum_P \|\phi_P(z - \pi z)\|_{1, \Delta}^2 < C \|z\|_{1, \Omega_{\Delta}}^2$$

by (iii) of Lemma 4.1. The result now follows, just as in (i), from (vi) of Lemma 4.1.  $\square$

We shall also need a lower bound on the approximation power of  $M(\mathcal{D})$ . To this end we prove the following lemma.

Lemma 5.4. Suppose that  $z \in H^1(\Omega)$  and that there is an open disc  $D \subseteq \Omega$  where

$$\begin{aligned} (i) \quad & z \in C^{\infty}(\overline{D}), \\ (ii) \quad & \max_{i=1,2} \inf_{x \in \overline{D}} |D_i z(x)| > 0. \end{aligned} \tag{5.7}$$

Suppose that the mesh  $\mathcal{D}$  satisfies for some  $\kappa > 0$

$$\min_{\substack{\Delta \in \mathcal{D} \\ \Delta \subseteq \bar{D}}} |\Delta| > \kappa h(\mathcal{D}). \quad (5.8)$$

Then there are constants  $C, h_0$  such that if  $h(\mathcal{D}) < h_0$ ,

$$\inf_{v \in \tilde{M}(\mathcal{D})} \|z-v\|_{1,\Omega} > Ch(\mathcal{D}).$$

Proof. Certainly

$$\begin{aligned} \inf_{v \in \tilde{M}(\mathcal{D})} \|z-v\|_{1,\Omega}^2 &> \sum_{\Delta \in \mathcal{D}} \inf_w \|z-w\|_{1,\Delta}^2 \\ &\quad (\text{w bilinear on } \Delta) \\ &> \sum_{\substack{\Delta \in \mathcal{D} \\ \Delta \subseteq \bar{D}}} \inf_w \|z-w\|_{1,\Delta}^2. \end{aligned} \quad (5.9)$$

On each element  $\Delta \subseteq \bar{D}$ , we can expand  $z$  as a Taylor series about the centre,  $\bar{x}$  say, of  $\Delta$ ,

$$z(x) = \bar{w}(x) + \frac{1}{2} \sum_{i=1}^2 D_{ii} z(\bar{x}) (x_i - \bar{x}_i)^2 + r(x) \quad (x \in \Delta)$$

where  $\bar{w}$  is bilinear on  $\Delta$  and

$$|r(x)| < C|x-\bar{x}|^3, \quad |D_i r(x)| < C|x-\bar{x}|^2 \quad (i = 1,2)$$

for some constant  $C$ , independent of the particular  $\Delta$  under consideration. Writing

$$\bar{z}(x) = \bar{w}(x) + \frac{1}{2} \sum_{i=1}^2 D_{ii} z(\bar{x}) (x_i - \bar{x}_i)^2$$

we obtain using the finite dimension of the classes of functions involved

$$\inf_w \|\bar{z}-w\|_{1,\Delta} > C \left( \sum_{i=1}^2 |D_{ii} z(\bar{x})| \right) |\Delta|^2$$

w bilinear on  $\Delta$

$$> C^* |\Delta|^2$$

where  $C^* > 0$  is independent of  $\Delta$  by (5.7). Hence

$$\inf_w \|\bar{z}-w\|_{1,\Delta} > \inf_w \|\bar{z}-w\|_{1,\Delta} - \|r(x)\|_{1,\Delta}$$

(w bilinear on  $\Delta$ )

$$> C^* |\Delta|^2 - C |\Delta|^3$$

$$> C |\Delta|^2$$

for  $h(D) < h_0$ , and hence  $|\Delta|$ , sufficiently small. Therefore (5.9) enables us to say

$$\inf_{v \in \tilde{M}(D)} \|z-v\|_{1,\Omega}^2 > C \sum_{\substack{\Delta \in D \\ \Delta \in \bar{D}}} |\Delta|^4$$

$$> Ch(D)^2 \sum_{\substack{\Delta \in D \\ \Delta \in \bar{D}}} |\Delta|^2 \tag{5.10}$$

using (5.8). Note however, that for  $h(D)$  small enough,  $D^* \subseteq \bigcup_{\substack{\Delta \in D \\ \Delta \subseteq \bar{D}}} \Delta$  where

$D^*$  is a disc concentric with  $D$  but of radius half that of  $D$ . Thus

$$\sum_{\substack{\Delta \in D \\ \Delta \subseteq \bar{D}}} |\Delta|^2 > \text{area } D^*, \text{ and the lemma follows at once from (5.10).}$$

Let us remark that the hypotheses of the lemma are not very demanding at all.

## 2. The A-posteriori Error Estimate

2.1. For each proper node  $P$  of  $\mathcal{D}$ , define

$$M_P = \{v \in H^1(\Omega_P): v = 0 \text{ on } \Gamma_P^0\}.$$

(Clearly by extending functions in  $M_P$  by zero to the remainder of  $\Omega$  we can consider  $M_P \subseteq H_0$ ). We shall associate the following local subproblem with the proper node  $P$ : Find  $\eta_P \in M_P$  such that

$$b(\eta_P, v) = b(u - \tilde{u}, v), \quad v \in M_P. \quad (1.1)$$

After an integration by parts the right hand side of (1.1) can be written out explicitly in terms of the input data of (1.1.1):

$$\begin{aligned} b(u - \tilde{u}, v) &= \sum_{\Delta \in \mathcal{D}_P} \int_{\Delta} (f - L(\tilde{u}))v \, dx - \int_{\Gamma} \sigma(\tilde{u}) \cdot \hat{n} \, v \, ds \\ &+ \int_{\partial' \Omega} (\tau - \sigma(\tilde{u}) \cdot \hat{n})v \, ds. \end{aligned} \quad (1.2)$$

Note also that as long as  $P \in \partial^0(0,1)^2$ , then  $\phi_P \in \tilde{M}_0(\mathcal{D})$ , and that therefore

$$b(u - \tilde{u}, \phi_P) = 0 \quad (1.3)$$

by the usual finite element error orthogonality relation (1.1.5b).

We are now able to state and prove the fundamental error estimate of this paper. This result will allow us to estimate the global finite element error  $\|u - \tilde{u}\|_{1,\Omega}$  in terms of the solution  $\eta_P$  of the local subproblems (1.1). The basic ideas behind this estimate were first presented in [6].

**Theorem 1.1.** There is a constant  $C > 0$  such that for any  $K$ -mesh

$$C^{-1} \left( \sum_P \|\eta_P\|_{1, \Omega_P}^2 \right)^{1/2} < \|u - \tilde{u}\|_{1, \Omega} < C \left( \sum_P \|\eta_P\|_{1, \Omega_P}^2 \right)^{1/2} \quad (1.4)$$

where  $\sum_P$  denotes a summation over all proper nodes  $P$  of  $\mathcal{D}$ .

**Proof.** (1) The right hand inequality: By the coercivity of  $B$  (see 1.1.4c), and the finite element error orthogonality relation (1.1.6a), we have for any  $v \in \tilde{M}_0(\mathcal{D})$

$$\begin{aligned} \|u - \tilde{u}\|_{1, \Omega}^2 &< C b(u - \tilde{u}, u - \tilde{u}) \\ &< C b(u - \tilde{u}, u - \tilde{u} - v) \\ &< C b(u - \tilde{u}, \left( \sum_P \phi_P \right) (u - \tilde{u} - v)) \end{aligned}$$

since  $\sum_P \phi_P = 1$  by Corollary 1.3.4(b). Thus,

$$\|u - \tilde{u}\|_{1, \Omega}^2 < C \sum_P b(u - u, \phi_P (u - \tilde{u} - v)).$$

But  $\phi_P(u - \tilde{u} - v) \in M_P$  (certainly  $\phi_P(u - \tilde{u} - v) \in H^1(\Omega_P)$ ); for the trace behaviour note that if  $P \notin \partial(0, 1)^2$  then  $\phi_P = 0$  on  $\partial\Omega_P$ , whereas if  $P \in \partial\Omega$  then  $u - \tilde{u} - v = 0$  on  $\Gamma_P^0 \cap \partial^0(0, 1)^2$  while  $\phi_P = 0$  on  $\Gamma_P^0 \cap \Omega$ . Thus

$$b(u - \tilde{u}, \phi_P (u - \tilde{u} - v)) = b(\eta_P, \phi_P (u - \tilde{u} - v))$$

and we have

$$\begin{aligned} \|u - u\|_{1, \Omega}^2 &< C \sum_P |b(\eta_P, \phi_P (u - \tilde{u} - v))| \\ &< C \left( \sum_P \|\eta_P\|_{1, \Omega_P}^2 \right)^{1/2} \left( \sum_P \|\phi_P (u - \tilde{u} - v)\|_{1, \Omega_P}^2 \right)^{1/2} \end{aligned}$$

$$< C \left( \sum_P \|\eta_P\|_{1,\Omega_P}^2 \right)^{1/2} \|u-\tilde{u}\|_{1,\Omega}$$

on making the particular choice of  $V$  described in Corollary 1.5.3(ii).

The left hand inequality: Let us partition the set of proper nodes of  $\mathcal{D}$  into disjoint classes  $\chi_j$  ( $j = 1, \dots, J(\mathcal{D})$ ) say, with the property that whenever  $P, P' \in \chi_j$ , then  $\Omega_P \cap \Omega_{P'} = \emptyset$ . It follows from Lemma 1.4.1(v) that we may always ensure that  $J(\mathcal{D}) < C$ , for some mesh independent constant  $C > 0$ . Now

$$\begin{aligned} b(u-\tilde{u}, \sum_P \eta_P) &= \sum_P b(u-\tilde{u}, \eta_P) \\ &= \sum_P b(\eta_P, \eta_P) \\ &> C \sum_P \|\eta_P\|_{1,\Omega_P}^2 \end{aligned}$$

by the coercivity of  $B$  over  $M_P \subset M_0$  (1.1.4c). But

$$\begin{aligned} |b(u-\tilde{u}, \sum_P \eta_P)| &< C \|u-\tilde{u}\|_{1,\Omega} \left\| \sum_{j=1}^J \sum_{P \in \chi_j} \eta_P \right\|_{1,\Omega} \\ &< C \|u-\tilde{u}\|_{1,\Omega} \sum_{j=1}^J \left( \sum_{P \in \chi_j} \|\eta_P\|_{1,\Omega_P}^2 \right)^{1/2} \\ &< C \|u-\tilde{u}\|_{1,\Omega} \sum_{j=1}^J \left( \sum_P \|\eta_P\|_{1,\Omega_P}^2 \right)^{1/2} \\ &< C \|u-\tilde{u}\|_{1,\Omega} \left( \sum_P \|\eta_P\|_{1,\Omega_P}^2 \right)^{1/2}. \end{aligned}$$

The left hand inequality now follows.  $\square$

The significance of this result is that within the class of  $K$ -meshes the ratio  $\theta = (\sum_P |\eta_P|^2)^{1/2} / \|u - \tilde{u}\|_{1,\Omega}$  is bounded above and below (away from zero), independently of the global mesh refinement pattern (as long as the  $K$ -mesh property is maintained). Note also that Theorem 1 demands no assumptions on the regularity of the solution  $u$  (other than it lie in  $M_0$ , of course).

However the practicality (1.4) will depend on whether the solutions of the local subproblems (1.1) can in some sense be effectively estimated. This is the matter that shall concern us for the remainder of this section. Our main result will be Theorem 4.5 which will turn out to estimate the  $\eta_P$ 's in terms of the "jump residuals"  $[\sigma(\tilde{u}) \cdot \hat{n}]$  across the interelement boundaries.

2.2. In this and the following section we establish some preliminary results which will be the basis of our effective estimation of  $\|\eta_P\|_{1,\Omega}$ . We shall be working initially with a version of the subproblem (1.1) "standardized" to the standard star tuple  $\langle \mathcal{D}_0, \Gamma_0^0 \rangle$ . We now describe this standardized problem: Let

$$M_0 = \{u \in H^1(\Omega_0) : u = 0 \text{ on } \Gamma_0^0\}$$

and consider a continuous bilinear form  $\Lambda_0: M_0 \times M_0 \rightarrow \mathbb{R}$  which satisfies the conditions

$$\left. \begin{aligned} |\Lambda_0(u,v)| &< \beta \|u\|_{1,\Omega_0} \|v\|_{1,\Omega_0} \\ |\Lambda_0(u,u)| &> \beta^{-1} \|u\|_{1,\Omega_0}^2 \end{aligned} \right\} u, v \in M_0$$

for some  $\beta > 0$ . Corresponding to any  $R_0 \in L_2(\Omega_0)$  and  $r_0 \in L_2(\Gamma_0)$  let  $\eta_0 \in M_0$  be the solution of

$$\Lambda_0(\eta_0, v) = \int_{\Omega_0} R_0 v \, dx + \int_{\Gamma_0} r_0 v \, ds, \quad v \in M_0. \quad (2.1)$$

Let  $\varphi_0 \in M_0$  be a function with the properties

$$(i) \quad \varphi_0 > 0 \quad \text{on} \quad \Omega_0, \quad (2.2)$$

(ii) for some  $\delta > 0$ ,  $\varphi_0 > \delta$  on some disc contained in  $\Omega_0$ .

Finally let  $G_0 \subseteq L_2(\Omega_0)$  and  $g_0 \subseteq L_2(\Gamma_0)$  be finite dimensional subspaces.

Lemma 2.1. There is a constant  $C = C(\mathcal{D}_0, \Gamma_0^0) > 0$  such that for any  $v \in M_0$

$$\|v\|_{0, \Gamma_0} < C \|v\|_{1, \Omega_0},$$

and provided  $\Gamma_0^0 \neq \emptyset$

$$\|v\|_{0, \Omega_0} < C \|v\|_{1, \Omega_0}.$$

Proof. The first estimate is a standard trace result, while the second is a Poincaré-type bound on the  $L_2$  norm in terms of the  $H_1$  semi-norm.  $\square$

Lemma 2.2. There is a constant  $C = C(\beta, \mathcal{D}_0, \Gamma_0^0) > 0$  such that

$$\|\eta_0\|_{1, \Omega_0} < C (\|R_0\|_{0, \Omega_0} + \|r_0\|_{0, \Gamma_0}).$$

Proof.

$$\begin{aligned} \|\eta_0\|_{1, \Omega_0}^2 &< C \Lambda_0(\eta_0, \eta_0) \\ &< C \left| \int_{\Omega_0} R_0 \eta_0 \, dx + \int_{\Gamma_0} r_0 \eta_0 \, ds \right| \\ &< C (\|R_0\|_{0, \Omega_0} \|\eta_0\|_{0, \Omega_0} + \|\eta_0\|_{0, \Gamma_0} \|r_0\|_{0, \Gamma_0}). \end{aligned}$$



By Lemma 2.1,  $\|\eta_0\|_{0,\Gamma_0} < C\|\eta_0\|_{1,\Omega_0}$ , and of course  $\|\eta_0\|_{0,\Omega_0} < \|\eta_0\|_{1,\Omega_0}$ .  
The lemma follows immediately.  $\square$

**Lemma 2.3.** Suppose  $G_0$  and  $g_0$  are finite-dimensional subspaces of  $L_2(\Omega_0)$  and  $L_2(\Gamma_0)$  respectively, then there is a constant  $C = C(\beta, \mathcal{D}_0, \Gamma_0, G_0, g_0) > 0$  such that if  $R_0 \in G_0$  and  $r_0 \in g_0$  then

$$\|\eta_0\|_{1,\Omega_0} > C(\|R_0\|_{0,\Omega_0} + \|r_0\|_{0,\Gamma_0}).$$

**Proof.** Suppose the result does not hold, then there are sequences  $R_0^{(n)}$  and  $r_0^{(n)}$ , ( $n = 1, \dots$ ) from  $G_0$  and  $g_0$  respectively, and corresponding solutions  $\eta_0^{(n)}$  of (2.1) such that

$$\|\eta_0^{(n)}\|_{1,\Omega_0} < \frac{1}{n} (\|R_0^{(n)}\|_{0,\Omega_0} + \|r_0^{(n)}\|_{0,\Gamma_0}). \quad (2.3)$$

We may without loss of generality assume that

$$\|R_0^{(n)}\|_{0,\Omega_0} + \|r_0^{(n)}\|_{0,\Gamma_0} = 1. \quad (2.4)$$

In particular  $\|R_0^{(n)}\|_{0,\Omega_0} < 1$  and  $\|r_0^{(n)}\|_{0,\Gamma_0} < 1$ . By the finite dimensions of  $G_0$  and  $g_0$  we may suppose that

$$R_0^{(n)} \rightarrow R \in G_0 \text{ in } L_2(\Omega_0)$$

and

$$r_0^{(n)} \rightarrow r \in g_0 \text{ in } L_2(\Gamma_0).$$

Let  $\eta$  be the solution of (2.1) corresponding to  $R$  and  $r$ . By Lemma 2.2 applied to  $(\eta - \eta_0^{(n)})$ ,

$$\|\eta - \eta_0^{(n)}\|_{1, \Omega_0} < C(\|R - R_0^{(n)}\|_{0, \Omega_0} + \|r - r_0^{(n)}\|_{0, \Gamma_0}) \rightarrow 0;$$

but by (2.3),  $\eta_0^{(n)} \rightarrow 0$  in  $L^1(\Omega_0)$ , and so we conclude that  $\eta = 0$ .

Thus for any  $v \in M_0$

$$0 = \Delta_0(\eta, v) = \int_{\Omega_0} Rv \, dx + \int_{\Gamma_0} rv \, ds.$$

But this can only mean that  $R = 0$  and  $r = 0$ . However this clearly contradicts the limiting case of (2.4).  $\square$

Lemma 2.4. Let  $G_0$  and  $g_0$  be as in Lemma 2.3, then there is a constant  $C = C(\Omega_0, \Gamma_0, G_0, g_0, \phi_0)$  such that if  $R_0 \in G_0$ ,  $r_0 \in g_0$  and  $p \in P_0$  (the set of constant functions on  $\Omega_0$ )

$$\|R_0\|_{0, \Omega_0} < C(\|r_0\|_{0, \Gamma_0} + \|R_0 - p\|_{0, \Omega_0} + |\int_{\Omega_0} R_0 \phi_0 \, dx + \int_{\Gamma_0} r_0 \phi_0 \, ds|).$$

Proof. We shall again argue by contradiction. So suppose the result is not true, then there are sequences  $R_0^{(n)}$ ,  $r_0^{(n)}$  and  $p^{(n)}$  from  $G_0$ ,  $g_0$  and  $P_0$  respectively, such that

$$\|R_0^{(n)}\|_{0, \Omega_0} > n(\|r_0^{(n)}\|_{0, \Gamma_0} + \|R_0^{(n)} - p^{(n)}\|_{0, \Omega_0} + |\int_{\Omega_0} R_0^{(n)} \phi_0 \, dx + \int_{\Gamma_0} r_0^{(n)} \phi_0 \, ds|) \quad (2.5)$$

where we can suppose that

$$\|R_0^{(n)}\|_{0, \Omega_0} = 1. \quad (2.6)$$

By the finite dimension of  $G_0$  we may as well assume that  $R_0^{(n)} \rightarrow R \in G_0$  in  $L_2(\Omega_0)$ . But by (2.5),  $\|R_0^{(n)} - p^{(n)}\|_{0,\Omega_0} \rightarrow 0$ . Thus  $p^{(n)} \rightarrow R$  in  $L_2(\Omega_0)$ .

But  $P_0$  is certainly a closed subspace of  $L_2(\Omega_0)$  (since it is finite dimensional), and so  $R \in P_0$ , i.e.,  $R$  is a constant on  $\Omega_0$ . Now (2.5) also shows that  $r_0^{(n)} \rightarrow 0$  in  $L_2(\Gamma_0)$  and that

$$|\int_{\Omega_0} R_0^{(n)} \varphi_0 dx + \int_{\Gamma_0} r_0^{(n)} \varphi_0 ds| \rightarrow 0 \text{ as } n \rightarrow \infty. \text{ Thus we must have}$$

$$0 = \lim_{n \rightarrow \infty} \int_{\Omega_0} R_0^{(n)} \varphi_0 dx = \int_{\Omega_0} R \varphi_0 dx.$$

The only way this can happen when  $\varphi_0$  has the properties (2.2) and  $R$  is a constant, is for  $R = 0$ . But this contradicts the limiting case of (2.6).  $\square$

Lemma 2.5. There is a constant  $C = C(\Omega_0, \mathcal{D}_0) > 0$  such that for any  $z \in \mathcal{B}_0 = \{z \in C^0(\Omega_0) : z \text{ bilinear on each } \Delta \in \mathcal{D}_0\}$ ,

$$\inf_{p \in \mathcal{D}_0} \|Dz - p\|_{0,\Omega_0} < C \left\{ \left( \sum_{\Delta \in \mathcal{D}_0} \|Dz\|_{1,\Delta}^2 \right)^{1/2} + \left( \int_{\Gamma_0} [\nabla z \cdot \hat{n}]^2 ds \right)^{1/2} \right\} \quad (2.7)$$

where  $Dz$  denotes any (element-by-element) zeroth, first or second order derivative of  $z$ .

Proof. Let  $M = \{z \in \mathcal{B} : Dz \in P_0\}$ . Clearly the quotient space  $\mathcal{B}_0/M$  is finite dimensional, and (2.7) will follow if we can show that both sides of (2.7) define norms on  $\mathcal{B}_0/M$ . In both cases it clearly suffices to show that if the respective quantity in (2.7) vanishes then  $z \in M$ :

(i) The left side: If  $\inf_{p \in P_0} \|Dz - p\|_{0,\Omega_0} = 0$ , then  $Dz$  is a constant on

$\Omega_0$  since  $P_0$  is finite dimensional and hence closed.

(ii) The right side: If  $(\sum_{\Delta \in \mathcal{D}_0} |Dz|_{1,\Delta}^2)^{1/2} = 0$ , then  $Dz$  is constant on each  $\Delta \in \mathcal{D}_0$ . On the other hand, if  $\int_{\Gamma_0^0} \|\nabla z \cdot \underline{n}\|^2 ds = 0$ , then  $\|\nabla z \cdot \underline{n}\| = 0$  on  $\Gamma_0^0$ , and so by Lemma 1.3.7  $z$  is bilinear on  $\Omega_0$ . In particular,  $Dz \in C^0(\Omega_0)$ . The net result being that  $Dz \in P_0$ .  $\square$

We shall now relax the requirements of Lemmas 2.3 and 2.4 that  $R_0$  and  $r_0$  lie in the finite dimensional spaces  $G_0$  and  $g_0$ . Instead we shall require that they can be approximated from within these spaces.

Lemma 2.6. Let  $G_0$  and  $g_0$  be as in Lemma 2.3, and suppose that

$$\inf_{R^* \in G_0} \|R_0 - R^*\|_{0,\Omega_0} + \inf_{r^* \in g_0} \|r_0 - r^*\|_{0,\Gamma_0} < \varepsilon.$$

Then

(i) there is a constant  $C = C(\beta, \Omega_0, \Gamma_0, G_0, g_0) > 0$  such that

$$C^{-1} (\|R_0\|_{0,\Omega_0} + \|r_0\|_{0,\Gamma_0} - \varepsilon) < \|\eta_0\|_{1,\Omega_0} < C (\|R_0\|_{0,\Omega_0} + \|r_0\|_{0,\Gamma_0}) \quad (2.8)$$

(ii) there is a constant  $C = C(\Omega_0, \mathcal{D}_0, \Gamma_0, G_0, g_0, \psi_0) > 0$  such that if

$$\int_{\Omega_0} R_0 \psi_0 dx + \int_{\Gamma_0} r_0 \psi_0 ds = 0 \quad (2.9)$$

then

$$\|R_0\|_{0,\Omega_0} < C (\|r_0\|_{0,\Gamma_0} + \inf_{p \in P_0} \|R_0 - p\|_{0,\Omega_0} + \varepsilon). \quad (2.10)$$

Proof. Choose  $R^* \in G_0$  and  $r^* \in g_0$  such that

$\|R_0 - R^*\|_{0, \Omega_0} + \|r_0 - r^*\|_{0, \Gamma_0} < 2\varepsilon$ . Let  $\eta^*$  be the solution of (2.1) corresponding to  $r^*$  and  $R^*$ . Applying Lemma 2.2 to the difference  $(\eta_0 - \eta^*)$  we get

$$\|\eta_0 - \eta^*\|_{1, \Omega_0} < C(\|R_0 - R^*\|_{0, \Omega_0} + \|r_0 - r^*\|_{0, \Gamma_0}) < C\varepsilon.$$

Using Lemma 2.3,

$$\begin{aligned} \|R_0\|_{0, \Omega_0} + \|r_0\|_{0, \Gamma_0} &< \|R^*\|_{0, \Omega_0} + \|r^*\|_{0, \Gamma_0} + 2\varepsilon \\ &< C\|\eta^*\|_{1, \Omega_0} + 2\varepsilon \\ &< C(\|\eta_0\|_{1, \Omega_0} + \varepsilon), \end{aligned}$$

proving the left side of (2.8). The right side follows from Lemma 2.2.

For (2.10) we have

$$\|R_0\|_{0, \Omega_0} < \|R^*\|_{0, \Omega_0} + 2\varepsilon$$

while Lemma 2.4 gives for any  $p \in P_0$

$$\begin{aligned} \|R^*\|_{0, \Omega_0} &< C(\|r^*\|_{0, \Gamma_0} + \|R^{*-p}\|_{0, \Omega_0} + |\int_{\Omega_0} R^* \phi_0 \, dx + \int_{\Gamma_0} r^* \phi_0 \, ds|) \\ &< C(\|r_0\|_{0, \Gamma_0} + \|R_0^{-p}\|_{0, \Omega_0} + |\int_{\Omega_0} R_0 \phi_0 \, dx + \int_{\Gamma_0} r_0 \phi_0 \, ds| + \varepsilon) \\ &< C(\|r_0\|_{0, \Gamma_0} + \|R_0^{-p}\|_{0, \Omega_0} + \varepsilon) \end{aligned}$$

by virtue of the orthogonality relation (2.9).  $\square$

2.3. In the last section we considered a "standardized" form of the local subproblems (1.1). In this section we extend those results to the actual star tuple  $\langle P, \rho, \Gamma_P^0 \rangle$ .

For any proper node  $P$  with standard star tuple  $\langle \mathcal{D}_0, \Gamma_0^0 \rangle$ , let  $\tau_P: \mathbb{R}^2 \rightarrow \mathbb{R}^2$  be the star transformation of  $P$  (see §1.4), that is

$$\tau_P(x) = \frac{1}{\rho} (x-P)$$

where  $\rho = \rho_P = \max_{\Delta \in \mathcal{D}_P} |\Delta|$ . (Notice that  $\rho < 1$ .)

Corresponding to the finite dimensional spaces  $G_0$  and  $g_0$  of §2.2 let us define

$$G_P = \{v \in L_2(\Omega_P): v \circ \tau_P^{-1} \in G_0\} \subset L_2(\Omega_P),$$

$$g_P = \{v \in L_2(\Gamma_P): v \circ \tau_P^{-1} \in \tau_P^{-1} g_0\} \subset L_2(\Gamma_P).$$

At this stage it is convenient to consider a more general version of (1.1):

Find  $\eta \in M_P$  such that

$$b(\eta, v) = \int_{\Omega_P} Rv \, dx + \int_{\Gamma_P} rv \, ds, \quad \forall v \in M_P \quad (3.1)$$

where  $R \in L_2(\Omega_P)$  and  $r \in L_2(\Gamma_P)$ .

Lemma 3.1. Set

$$\epsilon = \rho_P \inf_{R \in G_P} \|R - \bar{R}\|_{0, \Omega_P} + \rho_P^{1/2} \inf_{r \in g_P} \|r - \bar{r}\|_{0, \Gamma_P} \quad (3.2)$$

then

(i) there is a constant  $C = C(D_0, \Gamma_0, G_0, g_0)$  such that

$$C^{-1}(\rho_P \|R\|_{0, \Omega_P} + \rho_P^{1/2} \|r\|_{0, \Gamma_P} - \varepsilon) < \|\eta\|_{1, \Omega_P} < C(\rho_P \|R\|_{0, \Omega_P} + \rho_P^{1/2} \|r\|_{0, \Gamma_P})$$

(ii) there is a constant  $C = C(D_0, \Gamma_0^0, G_0, g_0)$  such that if

$$\int_{\Omega_P} R \phi_P \, dx + \int_{\Gamma_P} r \phi_P \, ds = 0$$

then

$$\rho_P \|R\|_{0, \Omega_P} < C(\rho_P^{1/2} \|r\|_{0, \Gamma_P} + \rho_P \inf_{p \in P_0(\Omega_P)} \|R - p\|_{0, \Omega_P} + \varepsilon)$$

where  $P_0(\Omega_P)$  is the space of constant functions on  $\Omega_P$ .

Proof. It is clear that  $M_0 = \{v \in H^1(\Omega_0) : v \circ \tau_P \in M_P\}$ . Define a bilinear form  $\Lambda_0 : M_0 \times M_0 \rightarrow \mathbb{R}$  by

$$\Lambda_0(u, v) = b(u \circ \tau_P, v \circ \tau_P).$$

It follows from the boundedness of  $b(\cdot, \cdot)$  on  $M_P \times M_P$

$$\begin{aligned} |\Lambda_0(u, v)| &= |b(u \circ \tau_P, v \circ \tau_P)| \\ &< C \|u \circ \tau_P\|_{1, \Omega_P} \|v \circ \tau_P\|_{1, \Omega_P} \\ &< C(\rho^2 \|u\|_{0, \Omega_0}^2 + \|u\|_{1, \Omega_0}^2)^{1/2} (\rho^2 \|v\|_{0, \Omega_0}^2 + \|v\|_{1, \Omega_0}^2)^{1/2} \\ &< C \|u\|_{1, \Omega_0} \|v\|_{1, \Omega_0}. \end{aligned}$$

In addition, by the coercivity of  $b(\cdot, \cdot)$  over  $M_P$ ,

$$\begin{aligned}
|\Lambda_0(u,u)| &= |b(u \circ \tau_p, u \circ \tau_p)| > \alpha \|u \circ \tau_p\|_{1,\Omega_p}^2 \\
&> C \|u\|_{1,\Omega_0}^2 \\
&> C \|u\|_{1,\Omega_0}^2
\end{aligned}$$

using lemma 2.1, noting that  $\Gamma_0^0 \neq \emptyset$ . It is readily verified by a change of variable of integration that for any  $v \in M_0$

$$\begin{aligned}
\Lambda_0(\eta \circ \tau_p, v) &= b(\eta, v \circ \tau_p^{-1}) = \int_{\Omega_p} R(v \circ \tau_p^{-1}) dx + \int_{\Gamma_p} r(v \circ \tau_p^{-1}) ds \\
&= \int_{\Omega_0} (\rho^2 R \circ \tau_p) v dx + \int_{\Gamma_0} (\rho r \circ \tau_p) v ds.
\end{aligned}$$

So we are in a situation covered by the analysis of §2.2 with  $\eta_0 = \eta \circ \tau_p$ ,  $R_0 = \rho^2 R \circ \tau_p$  and  $r_0 = \rho r \circ \tau_p$ . Furthermore, for any  $\bar{R} \in G_p$  and  $\bar{r} \in g_p$ , we have  $\rho^2 \bar{R} \circ \tau_p \in G_0$  and  $\rho \bar{r} \circ \tau_p \in g_0$  with

$$\begin{aligned}
&\|\rho^2 R \circ \tau_p - \rho^2 \bar{R} \circ \tau_p\|_{0,\Omega_0} + \|\rho r \circ \tau_p - \rho \bar{r} \circ \tau_p\|_{0,\Gamma_0} \\
&= \rho \|R - \bar{R}\|_{0,\Omega_p} + \rho^{1/2} \|r - \bar{r}\|_{0,\Gamma_p}.
\end{aligned}$$

By virtue of (3.2) we are now able to apply Lemma 2.6 to infer from (2.8)

$$\begin{aligned}
C^{-1} (\|\rho^2 R \circ \tau_p\|_{0,\Omega_0} + \|\rho r \circ \tau_p\|_{0,\Gamma_0} - \varepsilon) &< \|\eta \circ \tau_p\|_{1,\Omega_0} \\
&< C (\|\rho^2 R \circ \tau_p\|_{0,\Omega_0} + \|\rho r \circ \tau_p\|_{0,\Gamma_0}).
\end{aligned} \tag{3.3}$$

Upon rescaling back to the actual star  $\Omega_p$ ,



$$\begin{aligned} \|\eta \circ \tau_P\|_{1, \Omega_0}^2 &= \|\eta \circ \tau_P\|_{0, \Omega_0}^2 + \|\eta \circ \tau_P\|_{1, \Omega_0}^2 = \frac{1}{\rho^2} \|\eta\|_{0, \Omega_P}^2 + \|\eta\|_{1, \Omega_P}^2 \\ &> \|\eta\|_{1, \Omega_P}^2; \end{aligned}$$

while by Lemma 2.1 we also have the opposing bound

$$\|\eta \circ \tau_P\|_{1, \Omega_0} < C \|\eta \circ \tau_P\|_{1, \Omega_0} < C \|\eta\|_{1, \Omega_P} < C \|\eta\|_{1, \Omega_P}.$$

Rescaling the other terms in (3.3) leads at once to (1) of the lemma.

To show (ii) note that  $\psi_0 = \psi_P \circ \tau_P$  satisfies (2.2) and that

$$\int_{\Omega_0} (\rho^2 R \circ \tau_P) \psi_0 + \int_{\Gamma_0} (\rho r \circ \tau_P) \psi_0 = \int_{\Omega_P} R \psi_P dx + \int_{\Gamma_P} r \psi_P ds = 0.$$

Part (ii) of Lemma 2.6 then gives

$$\|\rho^2 R \circ \tau_P\|_{0, \Omega_0} < C (\|\rho r \circ \tau_P\|_{0, \Gamma_0} + \inf_{P \in P_0} \|\rho^2 R \circ \tau_P - P\|_{0, \Omega_0} + \varepsilon),$$

where, since  $\psi_0$  is determined by the standard tuple  $\langle \Omega_0, \mathcal{D}_0, \Gamma_0 \rangle$ , the constant  $C$  of Lemma 2.6 (ii) is in fact  $C = C(\Omega_0, \mathcal{D}_0, \Gamma_0, G_0, g_0)$ . Rescaling to  $\Omega_P$  gives the result.  $\square$

2.4. We now return to the task of effectively estimating  $(\sum_P \|\eta_P\|_{1, \Omega_P}^2)^2$ . The local subproblems (1.1) fit into the pattern of §2.3 if we set (see (1.2))

$$\begin{aligned} R &= (f - L(\tilde{u})) && \text{on each } \Delta \\ \text{and} &&& \\ r &= \begin{cases} -[\sigma(\tilde{u}) \cdot \hat{n}] & \text{on } \Gamma^* \\ t - \sigma(\tilde{u}) \cdot \hat{n} & \text{on } \partial^1(0, 1)^2. \end{cases} \end{aligned} \quad (4.1)$$

To complete the framework required for applying the results of §2.3 we define

$$G_0 = \{v \in L_2(\Omega_0): v \text{ is constant on each } \Delta \in \mathcal{D}_0\},$$

$$g_0 = \{v \in L_2(\Gamma_0): v \text{ is linear on each primitive edge } \gamma \subseteq \Gamma_0\}.$$

It is clear that  $G_0$  and  $g_0$  are finite dimensional and that they are determined completely by the standard star tuple  $\langle \mathcal{D}_0, \Gamma_0^0 \rangle$ .

Lemma 4.1:

- (i)  $\sum_{\Delta \in \mathcal{D}} |\Delta|^2 \|\tilde{u}\|_{2,\Delta}^2 = o(1)$  as  $h(\mathcal{D}) \rightarrow 0$ .
- (ii)  $\sum_P \rho_P^2 \inf_{R \in G_P} \|R - \bar{R}\|_{0,\Omega_P}^2 = o(h(\mathcal{D})^2)$  as  $h(\mathcal{D}) \rightarrow 0$ .
- (iii)  $\sum_P \rho_P \inf_{r \in g_P} \|r - \bar{r}\|_{0,\Gamma_P}^2 = o(h(\mathcal{D})^2)$  as  $h(\mathcal{D}) \rightarrow 0$ .
- (iv) There exists  $r^* \in L_2(\Gamma)$  which is linear on each primitive edge  $\gamma \in \Gamma$ , which agrees with  $r$  at all nodes lying on  $\Gamma$ , and which satisfies

$$\sum_{\Delta} |\Delta| \|r - r^*\|_{\partial\Delta \cap \Gamma}^2 = o(h(\mathcal{D})^2) \quad \text{as } h(\mathcal{D}) \rightarrow 0.$$

Proof. To prove (i),

$$\sum_{\Delta \in \mathcal{D}} |\Delta|^2 \|\tilde{u}\|_{2,\Delta}^2 < h^2(\mathcal{D}) \|\tilde{u}\|_{1,\Omega}^2 + \sum_{\Delta \in \mathcal{D}} |\Delta|^2 \|\tilde{u}\|_{2,\Delta}^2.$$

By the projection property (1.1.5) of the finite element solution the first terms obviously tends to 0. For the second term

$$\begin{aligned}
\sum_{\Delta \in \mathcal{D}} |\Delta|^2 \|\tilde{u}\|_{2,\Delta}^2 &< C \sum_{\Delta \in \mathcal{D}} \inf_v \|\tilde{u}-v\|_{1,\Delta}^2 \\
&\quad (\text{v linear on } \Delta) \\
&< C \sum_{\Delta \in \mathcal{D}} (\|u-\tilde{u}\|_{1,\Delta}^2 + \inf_v \|u-v\|_{1,\Delta}^2) \\
&\quad (\text{v linear on } \Delta) \\
&< C (\|u-\tilde{u}\|_{1,\Omega}^2 + \sum_{\Delta} \inf_v \|u-v\|_{1,\Delta}^2) \\
&\quad (\text{v linear on } \Delta)
\end{aligned}$$

which certainly tends to zero as  $h(\mathcal{D}) \rightarrow 0$ .

Now turn to (ii) of the lemma. On each element  $\Delta \in \mathcal{D}_P$ ,  $R$  is a sum of terms each of the form  $\lambda(x)\xi(x)$  where  $\lambda$  has bounded derivatives of all orders on  $\Omega$  and  $\xi = 1, \tilde{u}, D_i \tilde{u}$  or  $D_{ij} \tilde{u}$  ( $i, j = 1, 2$ ). We can find a constant function  $\theta$  on  $\Delta$  such that

$$\begin{aligned}
\|\lambda\xi - \theta\|_{0,\Delta} &< C|\Delta| \|\lambda\xi\|_{1,\Delta} \\
&< C|\Delta| \|\xi\|_{1,\Delta} \\
&< C|\Delta| (\|\tilde{u}\|_{2,\Delta} + |\Delta|).
\end{aligned}$$

Thus

$$\inf_{\bar{R} \in G_P} \|\bar{R} - R\|_{0,\Omega_P} = C \left\{ \sum_{\Delta \in \mathcal{D}_P} (|\Delta|^2 \|\tilde{u}\|_{2,\Delta}^2 + |\Delta|^4) \right\}^{1/2},$$

and so

$$\begin{aligned}
\sum_P \rho_P^2 \inf_{\bar{R} \in G_P} \|\bar{R} - R\|_{0,\Omega_P}^2 &< C \sum_P \rho_P^2 \sum_{\Delta \in \mathcal{D}_P} (|\Delta|^2 \|\tilde{u}\|_{2,\Delta}^2 + |\Delta|^4) \\
&< C h(\mathcal{D})^2 \left( \sum_{\Delta \in \mathcal{D}} |\Delta|^2 \|\tilde{u}\|_{2,\Delta}^2 + h(\mathcal{D})^2 \right)
\end{aligned}$$

where we have made use of the fact that each element  $\Delta \in \mathcal{D}$  belongs to a  $\mathcal{D}_P$  for at least one, but no more than  $C = C(K)$ , proper nodes  $P$  (Lemmas 1.3.6 and 1.4.1). Part (ii) of the lemma now follows immediately on using (i).

For part (iii), observe that on each element edge  $\gamma$  making up  $\Gamma_P$ ,  $r$  is a sum of terms of the form  $\lambda(t)\xi(t)$  where  $\lambda$  has bounded derivatives of all orders (independent of  $\mathcal{D}$ ) and  $\xi = 1$  or a limiting value of  $D_1 \tilde{u}$  ( $i = 1, 2$ ) on the edge. There is no loss of generality in assuming that  $D_1 \tilde{u}$  is linear on this edge. We can find a linear function  $\beta$  on  $\gamma$  satisfying

$$\begin{aligned} \|\lambda\xi - \beta\|_{0,\gamma} &< C|\gamma|^2 \|\lambda\xi\|_{2,\gamma} \\ &< C|\gamma|^2 (\|\xi\|_{2,\gamma}) \\ &< C|\gamma|^2 \left( \sum_{i=1}^2 \|D_1 \tilde{u}\|_{1,\gamma} + |\gamma|^{1/2} \right) \end{aligned}$$

where  $D_1 \tilde{u}$  is the limiting value on  $\gamma$  from within an element  $\Delta \in \mathcal{D}_P$ . We have

$$\|D_1 \tilde{u}\|_{1,\gamma} < C|\Delta|^{-1/2} \|\tilde{u}\|_{2,\Delta}$$

and so

$$\|\lambda\xi - \beta\|_{0,\gamma} < C|\Delta|^{1/2} (|\Delta| \|\tilde{u}\|_{2,\Delta}^2 + |\Delta|^2). \quad (4.2)$$

Therefore

$$\inf_{r \in \mathcal{S}_P} \|r - \tilde{r}\|_{0,\Gamma_P} < C\rho_P^{1/2} \left\{ \sum_{\Delta \in \mathcal{D}_P} |\Delta|^2 \|\tilde{u}\|_{2,\Delta}^2 + |\Delta|^4 \right\}^{1/2}$$

and so, much as part (i), squaring and summing over all proper nodes

$$\begin{aligned} \sum_P \rho_P \inf_{r \in \mathcal{G}_P} \|r - \bar{r}\|_{0, \Gamma_P}^2 &< C \sum_P \rho_P^2 \sum_{\Delta \in \mathcal{D}_P} (|\Delta|^2 \|\tilde{u}\|_{2, \Delta}^2 + |\Delta|^4) \\ &< Ch(\mathcal{D})^2 \left( \sum_{\Delta} |\Delta|^2 \|\tilde{u}\|_{2, \Delta}^2 + h(\mathcal{D})^2 \right). \end{aligned}$$

Again, on using (i), part (iii) of the lemma follows.

Part (iv) of the lemma likewise follows from (4.2) and (4.1), since we may as well choose  $\beta$  to interpolate  $\lambda \xi$  at the endpoints of  $\lambda$ .  $\square$

We are now able to state and prove the main results of this paper concerning an equivalent estimator for the energy norm of the finite element error. Because of the coercivity of  $b(\cdot, \cdot)$  there is no need to distinguish between equivalent estimators for  $(b(u - \tilde{u}, u - \tilde{u}))^{1/2}$ , or  $\|u - \tilde{u}\|_{1, \Omega}$ . For definiteness we shall phrase our result in terms of the latter.

Theorem 4.2. There is a constant  $C > 0$  such that for any  $K$ -mesh

$$\|u - \tilde{u}\|_{1, \Omega} < CE \tag{4.4a}$$

where

$$E = \left\{ \sum_{\Delta \in \mathcal{D}} [|\Delta|^2 \|R\|_{0, \Delta}^2 + |\Delta| \|r\|_{0, \partial \Delta \cap \Gamma}^2]^{1/2} \right\}.$$

If in addition for some  $L > 0$

$$\|u - \tilde{u}\|_{1, \Omega} > Lh(\mathcal{D}) \tag{4.5}$$

then there are constants  $C > 0$ ,  $h_0 > 0$  such that for any  $K$ -mesh  $\mathcal{D}$  with  $h(\mathcal{D}) < h_0$

$$C^{-1}E < \|u - \tilde{u}\|_{1, \Omega}. \tag{4.4b}$$

Proof. Consider a particular proper node  $P$  and its star  $\Omega_P$ . By Lemma 3.1

$$C^{-1}(E_P - \rho_P \inf_{R \in G_P} \|R - \bar{R}\|_{0, \Omega_P} + \rho_P^{1/2} \inf_{r \in g_P} \|r - \bar{r}\|_{0, \Gamma_P}) < \|\eta_P\|_{1, \Omega_P} < C E_P \quad (4.6)$$

where

$$E_P^2 = \sum_{\Delta \in \mathcal{D}_P} [|\Delta|^2 \|R\|_{0, \Delta}^2 + |\Delta| \|r\|_{0, \partial\Delta}^2]_{\Gamma_P}.$$

Here we have used the fact that for any  $\Delta \in \mathcal{D}_P$ ,  $C^{-1}\rho_P < |\Delta| < \rho_P$  (see Lemma 1.4.1).

The constant in (4.6) depends only on the standard star tuple  $\langle \mathcal{D}_0, \Gamma_0^0 \rangle$  associated with the proper node  $P$ , but is otherwise independent of the mesh or the node  $P$ . However, by Theorem 1.4.3, for any  $K$  mesh there are at most  $C = C(K) < \infty$  possible standard star tuples. We may therefore regard the constant  $C$  in (4.6) as uniform for all proper nodes  $P$  and all meshes of the specified class.

We claim that  $C^{-1} E^2 < \sum_P E_P^2 < C E^2$ . To see this we need only recall

- (i) any element  $\Delta \in \mathcal{D}$  belongs to a  $\mathcal{D}_P$  for at least one but no more than  $C = C(K)$  proper nodes  $P$  (see Lemmas 1.3.6 and 1.4.1)
- (ii) any edge contained in  $\Gamma$  is contained in a  $\Gamma_P$  for at least one, but no more than  $C = C(K)$  proper nodes  $P$  (see Lemma 1.4.1).

Thus squaring (4.6), summing over all proper nodes  $P$ , and using (ii) and (iii) of lemma 4.1 gives

$$C^{-1}(E^2 - h(\mathcal{D})^2 o(1)) < \sum_P \|\eta_P\|_{1, \Omega_P}^2 < C E^2. \quad (4.7)$$

Theorem 1.1 allows us to conclude from (4.7)

$$C^{-1}(E^2 - h(\mathcal{D})^2 o(1)) < \|u - \tilde{u}\|_{1,\Omega}^2 < C E^2.$$

Making use of (4.5) gives the desired result.  $\square$

At first glance the estimate of (4.4a) may not seem too surprising since it bounds the error in terms of a "residual." Note however a simple integration by parts of (1.2) would only have yielded

$$\|u - \tilde{u}\|_{1,\Omega} < C E^*$$

$$\text{where } (E^*)^2 = \sum_{\Delta} [ \|R\|_{0,\Delta}^2 + |\Delta|^{-1} \|r\|_{0,\partial\Delta\cap\Gamma}^2 ].$$

To obtain the extra powers of  $|\Delta|$  in (4.4) it was essential that we were able to reduce the estimation of the global error  $\|u - \tilde{u}\|_{1,\Omega}$  to the consideration of a number of local subproblems.

The estimate of (4.4) is sharp in the sense that uniformly for all  $K$ -meshes

$$0 < C^{-1} < E / \|u - \tilde{u}\|_{1,\Omega} < C \quad (4.8)$$

(at least for  $h(\mathcal{D})$  sufficiently small). In contrast, for  $E^*$ , although we certainly have the estimate

$$C^{-1} < \frac{E^*}{\|u - \tilde{u}\|_{1,\Omega}},$$

no two-sided estimate similar to (4.8) can hold, since  $E^*/E \rightarrow \infty$  and so

$$\frac{E^*}{\|u - \tilde{u}\|_{1,\Omega}} \rightarrow \infty \quad \text{as } h(\mathcal{D}) \rightarrow 0.$$

Theorem 4.2 requires no extra assumptions on the regularity of the exact solution  $u$  other than (4.5). However this requirement is hardly a restriction at all being satisfied in all but trivial cases (see Lemma 1.5.4).

Although the estimate of Theorem 4.2 is sharp in the sense described above, it can be simplified even further. It turns out that the  $L_2$  residual contribution can be dropped from  $E$ , giving an expression solely in terms of "line residuals," while still maintaining the two-sided estimate (4.4). This will be shown in Theorem 4.5. It is natural to ask whether the "reverse" simplification is valid, that is, can the "line residuals" be dropped from  $E$  without affecting the estimate. It is not difficult to see that this cannot be true. Consider the simple case where  $L(\cdot)$  is the Laplacian. Then for bilinear elements,  $L(\tilde{u}) = 0$ , and the  $L_2$  residual becomes just

$$\left( \sum_{\Delta} |\Delta|^2 \int_{\Delta} |f|^2 dx \right)^{1/2} > Ch(\mathcal{D}) \|f\|_{0,\Omega}$$

if the mesh  $\mathcal{D}$  is quasiuniform. However, if the solution  $u$  has some singular behaviour so that  $u \notin H^2(\Omega)$  then  $\|u - \tilde{u}\|_{1,\Omega}$  must converge at a slower rate than  $O(h(\mathcal{D}))$ . So the  $L_2$  residual alone cannot suffice for  $E$  in (4.5) in this case.

For proving Theorem 5 we shall need the following lemmas:

**Lemma 4.3.** There is a mesh independent constant  $C > 0$  such that for any proper node  $P$  and any  $z \in \mathcal{B}(\Omega_P) = \{z \in C^0(\Omega_P) : z \text{ bilinear on each } \Delta \in \mathcal{D}_P\}$ ,

$$\inf_{P \in \mathcal{P}_0(\Omega_P)} \|Dz - p\|_{0,\Omega_P} < C \left\{ \rho_P \left( \sum_{\Delta \in \mathcal{D}_P} |Dz|_{1,\Delta}^2 \right)^{1/2} + \rho_P^s \left( \int_{\Gamma_P^*} \|\nabla z \cdot \hat{n}\|^2 ds \right)^{1/2} \right\}$$

where  $Dz$  denotes any (element-by-element) zeroth ( $s = 3/2$ ), first ( $s = 1/2$ ) or second ( $s = -1/2$ ) derivative of  $z$ .



Proof. This is just the rescaled version of Lemma 2.5, taking note again that there are only a finite number of standard star tuples.  $\square$

Lemma 4.4. Suppose  $\gamma \in \Omega$  is a straight horizontal or vertical line segment, then for any  $z \in \tilde{M}(\mathcal{D})$ ,

$$\int_{\gamma} \llbracket \nabla z \cdot \hat{n} \rrbracket^2 ds < C \int_{\gamma} \llbracket \sigma(z) \cdot \hat{n} \rrbracket^2 ds.$$

Proof. For definiteness suppose that  $\gamma$  is vertical and let superscript  $(-)$  and  $(+)$  indicate limits from the left and righthand sides of  $\gamma$  respectively. (See Fig. 5.) Since all the  $a_{ijk\ell}$  are continuous in  $\Omega$ ,

Figure 5. The scheme of notation.

$$\begin{aligned} \llbracket \sigma(z) \cdot \hat{n} \rrbracket &= \sum_{i,j,k,\ell} \{ (a_{ijk\ell} D_{\ell} z_j^{n_k})^{(+)} + (a_{ijk\ell} D_{\ell} z_j^{n_k})^{(-)} \} \\ &= \sum_{i,j,\ell} \{ a_{ij\ell\ell} [ -(D_{\ell} z_j)^{(+)} + (D_{\ell} z_j)^{(-)} ] \} \end{aligned}$$

since  $n_2 = 0$ ,  $n_1^{(+)} = -1$  and  $n_1^{(-)} = 1$ . But  $D_2 z_j$  is continuous across  $\gamma$ , so

$$\begin{aligned} \llbracket \sigma(z) \cdot \hat{n} \rrbracket &= \sum_{i,j} a_{ij11} [ -(D_1 z_j)^{(+)} + (D_1 z_j)^{(-)} ] \\ &= a^{11} \llbracket \nabla z \cdot \hat{n} \rrbracket \end{aligned}$$

where

$$a^{11} = \begin{pmatrix} a_{1111} & a_{1211} \\ a_{2111} & a_{2211} \end{pmatrix}.$$

But the matrix  $a^{11}$  is uniformly positive definite on  $\Omega$  (see 1.1.4b) and so  $\det(a^{11}) > C > 0$  on  $\Omega$ . It follows then that on  $\gamma$

$$|[\nabla z \cdot \hat{n}]| < C |[\sigma(z) \cdot \hat{n}]|,$$

and the lemma is proven in this case.

The case of  $\gamma$  horizontal is treated similarly.  $\square$

Theorem 4.5. Suppose that for some  $L > 0$

$$\|u - \tilde{u}\|_{1,\Omega} > Lh(\mathcal{D})$$

then there are constants  $C, h_0 > 0$  such that for any  $K$ -mesh  $\mathcal{D}$  with  $h(\mathcal{D}) < h_0$

$$C^{-1}E < \|u - \tilde{u}\|_{1,\Omega} < CE \tag{4.9a}$$

where

$$E = \left\{ \sum_{\Delta} |\Delta| \|r\|_{0,\partial\Delta \cap \Gamma}^2 \right\}^{1/2}. \tag{4.9b}$$

Proof. First note that the left hand side of (4.9) follows trivially from (4.4b) since  $E < E$ . To prove the right hand side it will suffice to show

$$\sum_{\Delta} |\Delta|^2 \|R\|_{0,\Delta}^2 < E^2 + o(h(\mathcal{D})^2) \tag{4.10}$$

with the  $o(1)$  term valid as  $h(\mathcal{D}) \rightarrow \infty$ .

For any interior proper node  $P$  we have (see (1.2), (1.3))

$$\int_{\Omega_P} R \phi_P \, dx + \int_{\Gamma_P} r \phi_P \, ds = 0$$

with  $R$  and  $r$  as in (4.1). Thus we may apply part (ii) of Lemma 3.1,

$$\rho_P \|R\|_{0, \Omega_P} < C(\rho_P^{1/2} \|r\|_{0, \Gamma_P} + \rho_P \inf_{p \in P_0(\Omega_P)} \|R-p\|_{0, \Omega_P} + \varepsilon_P), \quad (4.11)$$

where  $P_0(\Omega_P)$  is the space of constant functions on  $\Omega_P$  and

$$\varepsilon_P = \rho_P \inf_{\bar{R} \in G_P} \|R-\bar{R}\|_{0, \Omega_P} + \rho_P^{1/2} \inf_{r \in G_P} \|r-\bar{r}\|_{0, \Gamma_P}.$$

We shall concentrate on the term  $\inf_{p \in P_0(\Omega_P)} \|R-p\|_{0, \Omega_P}$  for a while. As

the coefficients of  $L(\cdot)$  are smooth on  $\Omega$ , let us write  $\bar{L}(\cdot)$  for  $L(\cdot)$

but with its coefficients replaced by their averages over  $\Omega_P$ . If we then set

$$\bar{R} = \sum_{\Delta \in \mathcal{D}_P} (F-\bar{L}(\tilde{u})),$$

we have

$$\begin{aligned} \inf_{p \in P_0(\Omega_P)} \|R-p\|_{0, \Omega_P} &< \|R-\bar{R}\|_{0, \Omega_P} + \inf_{p \in P_0(\Omega_P)} \|\bar{R}-p\|_{0, \Omega_P} \\ &< C\rho_P \|\tilde{u}\|_{2, \Omega_P} + \inf_{p \in P_0(\Omega_P)} \|\bar{R}-p\|_{0, \Omega_P} \\ &< C\rho_P (\|\tilde{u}\|_{2, \Omega_P} + |\Omega_P|) + \inf_{p \in P_0(\Omega_P)} \left\| \sum_{\Delta \in \mathcal{D}_P} L(\tilde{u}) - p \right\|_{0, \Delta} \end{aligned} \quad (4.12)$$

where  $\|\tilde{u}\|_{2, \Omega_P}$  is to be understood in an element-by-element sense.

Now  $L(\tilde{u})$  is just a sum of terms of the form  $\lambda(D\tilde{u})$  where  $\lambda$  is a constant

and  $D\tilde{u}$  denotes a zeroth, first or second derivative of the function

$\tilde{u} \in \mathcal{B}(\Omega_P)$ . (The constant  $\lambda$  is bounded independently of  $P$  and  $\mathcal{D}$ ). We may

clearly apply Lemma 4.3 to obtain

$$\begin{aligned}
& \inf_{P \in \mathcal{P}_0(\Omega_P)} \left\| \sum_{\Delta \in \mathcal{D}_P} \bar{L}(\tilde{u}) - p \right\|_{0,\Delta} \\
& < C(\rho_P \|\tilde{u}\|_{2,\Omega_P} + \rho_P^{-1/2} (\int_{\Gamma_P^*} [\nabla \tilde{u} \cdot \hat{n}]^2 ds)^{1/2}) \\
& < C(\rho_P \|\tilde{u}\|_{2,\Omega_P} + \rho_P^{-1/2} (\int_{\Gamma_P^*} [\sigma(\tilde{u}) \cdot \hat{n}]^2 ds)^{1/2})
\end{aligned} \tag{4.13}$$

by Lemma 4.4.

Substituting (4.13) in (4.12), and then this in turn into (4.11) gives

$$\rho_P \|R\|_{0,\Omega_P} < C(\rho_P^{1/2} \|r\|_{0,\Gamma_P} + \epsilon_P^*) \tag{4.14}$$

where

$$\epsilon_P^* < \epsilon_P + \rho_P^2 (\|\tilde{u}\|_{2,\Omega_P} + |\Omega_P|).$$

Squaring (4.14) and summing over all interior proper nodes  $P$  gives

$$\begin{aligned}
\sum_{\Delta \in \mathcal{D}} |\Delta|^2 \|R\|_{0,\Delta}^2 & < C \sum_P \rho_P^2 \|R\|_{0,\Omega_P}^2 \\
& \quad \text{(P an interior proper node)} \\
& < C(E^2 + \sum_P (\epsilon_P)^2 + h(\mathcal{D})^2 (\sum_{\Delta} |\Delta|^2 \|\tilde{u}\|_{2,\Delta}^2 + h(\mathcal{D})^2))
\end{aligned} \tag{4.15}$$

where, much as in the proof of Theorem 4.2, we have used the following properties of our meshes:

- (i) for  $\Delta \in \mathcal{D}_P$ ,  $C^{-1} \rho_P < |\Delta| < C \rho_P$  (see Lemma 1.4.1)
- (ii) any element  $\Delta \in \mathcal{D}$  belongs to  $\mathcal{D}_P$  for at least one interior proper node  $P$  (see Lemma 1.3.6) and at most  $C = C(K)$  proper nodes  $P$  (Lemma 1.4.1)

(iii) any edge contained in  $\gamma$  is contained in  $\Gamma_P$  for one more than  $C = C(K)$  proper nodes  $P$  (Lemma 1.4.1).  $\square$

The desired result (4.10) follows at once from (4.15) and the assumed lower bound on  $\|u-\tilde{u}\|_{1,\Omega}$  upon appealing to (i), (ii) and (iii) of Lemma 4.1.

A slightly more computationally convenient form of estimator than that given by Theorem 4.5 would be obtained if the integrals involved in  $\|r\|_{0,\partial\Delta\cap\Gamma}^2$  could be replaced by discrete sums.

If  $g$  is a (sufficiently smooth) function defined on an edge,  $\gamma$  say, of an element, and if  $Q_1, \dots, Q_k$  are the nodes which lie on  $\gamma$  then let us write

$$\|g\|_{0,\gamma}^* = (|\gamma| \sum_{j=1}^k |g(Q_j)|^2)^{1/2}.$$

If  $\Delta$  is an element and  $g$  is defined on  $\partial\Delta \cap \Gamma$  then we shall write

$$\|g\|_{0,\partial\Delta\cap\Gamma}^* = \left( \sum_{\substack{\gamma \\ (\gamma \text{ an edge of } \Delta) \\ \gamma \subset \Gamma}} \|g\|_{0,\gamma}^* \right)^{1/2}.$$

Corollary 4.6. Theorem 4.5 remains true if  $E$  is replaced by  $E^*$  where

$$E^* = \left\{ \sum_{\Delta} |\Delta| (\|r\|_{0,\partial\Delta\cap\Gamma}^*)^2 \right\}^{1/2}.$$

Proof. Let  $r^*$  be as in (iv) of Lemma 4.1. Since any edge  $\gamma$  can have no more than  $C = C(K)$  nodes on it (see Lemma 1.4.1)

$$C^{-1} \|r^*\|_{0,\gamma} < \|r\|_{0,\gamma}^* < C \|r^*\|_{0,\gamma}.$$

Thus

$$C^{-1}E^* < \left( \sum_{\Delta} |\Delta| \|r^*\|_{0, \partial\Delta \cap \Gamma}^2 \right)^{1/2} < CE^*.$$

But

$$\begin{aligned} \sum_{\Delta} |\Delta| \|r^*\|_{0, \partial\Delta \cap \Gamma}^2 &< C \sum_{\Delta} |\Delta| (\|r\|_{0, \partial\Delta \cap \Gamma}^2 + \|r-r^*\|_{0, \partial\Delta \cap \Gamma}^2) \\ &< C(E^2 + \sum_{\Delta} |\Delta| \|r-r^*\|_{0, \partial\Delta \cap \Gamma}^2) \\ &< C(E^2 + o(h(\mathcal{D})^2)) \end{aligned}$$

by (iv) of Lemma 4.1. Likewise

$$E^2 < C \left( \sum_{\Delta} |\Delta| \|r^*\|_{0, \partial\Delta \cap \Gamma}^2 + o(h(\mathcal{D})^2) \right).$$

The corollary now follows by virtue of the assumed lower bound on  $\|u-\tilde{u}\|_{1, \Omega}$  in the theorem.  $\square$

Since we are only concerned here with equivalent estimators for  $\|u-\tilde{u}\|_{1, \Omega}$  there are many other modifications that can be made to  $E$  or  $E^*$  without affecting their equivalent estimator property. Let us just mention one other such modification which we shall say a little more about in Part II. Using the same notation as above, define

$$\begin{aligned} \|g\|_{0, \gamma}^{**} &= (|\gamma| \max_{j=1, \dots, k} |g(Q_j)|^2)^{1/2} \\ \|g\|_{0, \partial\Delta \cap \Gamma}^{**} &= \left\{ \begin{array}{ll} \max_{\substack{\gamma \\ \gamma \subset \Gamma}} ( \|g\|_{0, \gamma}^{**} )^2 + & \max_{\substack{\gamma \\ \gamma \subset \Gamma}} ( \|g\|_{0, \gamma}^{**} )^2 \\ (\gamma \text{ a vertical edge of } \Delta) & (\gamma \text{ a horizontal edge of } \Delta) \end{array} \right\}^{1/2}. \end{aligned}$$

Now set

$$E^{**} = \left\{ \sum_{\Delta \in \mathcal{D}} |\Delta| (\|r\|_{0, \partial\Delta \cap \Gamma}^{**})^2 \right\}^{1/2}.$$

Since there are at most  $C = C(K)$  nodes on an edge of any element, it is immediate that  $C^{-1}E^* \leq E^{**} \leq CE^*$ . Thus  $E^{**}$  is an equivalent estimator under the same conditions that apply to  $E^*$ .

## REFERENCES

- [1] C. Mesztenyi, W. Szymczak, FEARS User's Manual for Univac 1100. Tech. Note BN-991, Institute for Physical Science and Technology, University of Maryland, October 1982.
- [2] W. C. Rheinboldt, Feedback Systems and Adaptivity for Numerical Computations, in "Adaptive Computational Methods for Partial Differential Equations," eds. I. Babuška, J. Chandra, J. E. Flaherty, SIAM Publications, Philadelphia, PA, 1983, pp. 3-19.
- [3] I. Babuška, M. Vogelius, Feedback and Adaptive Finite Element Solution of One-Dimensional Boundary Value Problems, Numerische Mathem. 44, 1984, pp. 75-103.
- [4] I. Babuška, A. Miller, M. Vogelius, Adaptive Method and Error Estimation for Elliptic Problems of Structural Mechanics, in "Adaptive Computational Method for Partial Differential Equations, eds. I. Babuška, J. Chandra, J. E. Flaherty, SIAM Publications, Philadelphia, PA, 1983, pp. 35-56.
- [5] C. Mesztenyi, A. Miller, W. Szymczak, FEARS Details of Mathematical Formulation, Tech. Note BN-994, Institute for Physical Science and Technology, University of Maryland, December 1982.
- [6] I. Babuška, W. C. Rheinboldt, Error estimates for adaptive finite element computations, SIAM, J. Numer. Anal. 15 (1978), pp. 736-754.



The Laboratory for Numerical Analysis is an integral part of the Institute for Physical Science and Technology of the University of Maryland, under the general administration of the Director, Institute for Physical Science and Technology. It has the following goals:

- To conduct research in the mathematical theory and computational implementation of numerical analysis and related topics, with emphasis on the numerical treatment of linear and nonlinear differential equations and problems in linear and nonlinear algebra.
- To help bridge gaps between computational directions in engineering, physics, etc. and those in the mathematical community.
- To provide a limited consulting service in all areas of numerical mathematics to the University as a whole, and also to government agencies and industries in the State of Maryland and the Washington Metropolitan area.
- To assist with the education of numerical analysts, especially at the postdoctoral level, in conjunction with the Interdisciplinary Applied Mathematics Program and the programs of the Mathematics and Computer Science Departments. This includes active collaboration with government agencies such as the National Bureau of Standards.
- To be an international center of study and research for foreign students in numerical mathematics who are supported by foreign governments or exchange agencies (Fulbright, etc.).

Further information may be obtained from Professor I. Babuška, Chairman, Laboratory for Numerical Analysis, Institute for Physical Science and Technology, University of Maryland, College Park, Maryland 20742.

**END**

**FILMED**

1-85

**DTIC**