# RAIRO Modélisation mathématique et analyse numérique

# D. N. ARNOLD

## F. Brezzi

### Mixed and nonconforming finite element methods : implementation, postprocessing and error estimates

*RAIRO – Modélisation mathématique et analyse numérique*, tome 19, nº 1 (1985), p. 7-32.

<http://www.numdam.org/item?id=M2AN\_1985\_\_19\_1\_7\_0>

© AFCET, 1985, tous droits réservés.

L'accès aux archives de la revue « RAIRO – Modélisation mathématique et analyse numérique » implique l'accord avec les conditions générales d'utilisation (http://www.numdam.org/legal.php). Toute utilisation commerciale ou impression systématique est constitutive d'une infraction pénale. Toute copie ou impression de ce fichier doit contenir la présente mention de copyright.

# $\mathcal{N}$ umdam

Article numérisé dans le cadre du programme Numérisation de documents anciens mathématiques http://www.numdam.org/ MATHEMATICAL MODELLING AND NUMERICAL ANALYSIS MODÉLISATION MATHÉMATIQUE ET ANALYSE NUMÉRIQUE

(vol 19, nº 1, 1985, p 7 à 32)

### MIXED AND NONCONFORMING FINITE ELEMENT METHODS : IMPLEMENTATION, POSTPROCESSING AND ERROR ESTIMATES (\*)

by D. N. ARNOLD (<sup>1</sup>) and F. BREZZI (<sup>2</sup>)

Communicated by E MAGENES

Abstract — We discuss a technique of implementing certain mixed finite elements based on the use of Lagrange multipliers to impose interelement continuity. The matrices arising from this implementation are positive definite. Considering some well-known mixed methods, namely the Raviart-Thomas methods for second order elliptic problems and the Hellan-Hermann-Johnson method for biharmonic problems, we show that the computed Lagrange multipliers may be exploited in a simple prostprocess to produce better approximation of the original variables. We further establish an equivalence between the mixed methods and certain modified versions of well-known nonconforming methods, notably the Morley method in the case of the biharmonic problem. The equivalence is exploited to provide error estimates for both the mixed and nonconforming methods

Résumé — Nous étudions ici une technique d'implémentation de certains éléments finis mixtes bases sur l'utilisation des multiplicateurs de Lagrange utilisés pour imposer la continuite à la traversée des éléments Les matrices qui apparaissent sont definies positives Considérant quelques méthodes d'éléments finis mixtes bien connues telles que les méthodes de Raviart-Thomas pour des problèmes elliptiques du second ordre et la méthode de Hellan-Hermann-Johnson pour les problèmes biharmoniques, nous voyons que les multiplicateurs de Lagrange calcules peuvent être exploites dans un post-traitement simple pour produire une meilleure approximation des variables originales. Nous etablissons en outre une équivalence entre les methodes mixtes et certaines versions modifiées de méthodes non conformes bien connues, en particulier la méthode de Morley pour le probleme biharmonique Cette equivalence est exploitee pour obtenir des estimations d'erreur pour les méthodes mixtes et non conformes a la fois.

#### 0. INTRODUCTION

The discretization of linear elliptic boundary value problems by mixed finite element methods typically leads to linear systems of the form

$$\begin{array}{c} A\sigma + Bu = f_1, \\ B^T \sigma = f_2. \end{array}$$
 (0.0)

<sup>(\*)</sup> Received in October 1983

<sup>(1)</sup> Dept of Mathematics, University of Maryland, College Park, MD 20742, USA

 $<sup>(^2)</sup>$  Dip Meccanica Strutturale Univ di Pavia and Istituto di Analisi Numerica del CNR di Pavia, Italy

The first author was supported as a North Atlantic Treaty Organisation Postdoctoral Fellow and by NSF Grant MCS-83-13247

M<sup>2</sup> AN Modélisation mathematique et Analyse numérique, 0399-0516/85/01/7/26/\$ 4,60 Mathematical Modelling and Numerical Analysis © AFCET-Gauthier-Villars

In order to fix ideas we shall think of  $\sigma$  and u as approximations to the stress field and displacement field respectively in a problem of linear elasticity. The choice of a numerical method to solve the system (0 0) is restricted by the fact that it is generally indefinite. However for many of the most widely used mixed methods this drawback is frequently circumvented by an implementational technique which leads to a positive definite system [8]. The technique applies, essentially, whenever the finite elements approximating the stress field are subject to continuity constraints only at points interior to the interelement boundaries, but *not* at element vertices. Then one may eliminate the continuity instead through Lagrange multipliers defined on the interelement boundaries. Denoting by  $\lambda$  the multipliers, which enter the discrete system as additional unknowns, the resulting system has the form

$$\begin{split} \tilde{A}\overline{\sigma} + \tilde{B}\overline{u} + C\lambda &= f_1 , \\ \tilde{B}^T \overline{\sigma} &= f_2 , \\ C^T \overline{\sigma} &= 0 \end{split}$$
 (0 1)

The third equation in (0.1) imposes the required continuity on the stress field, and — if the space of multipliers is chosen appropriately — then  $\overline{\sigma}$ , although a priori sought in a space of finite elements without interelement continuity constraints, will belong to the original finite element space for  $\sigma$ , and consequently  $\overline{\sigma}$  and  $\sigma$  will coincide. Moreover the displacement field  $\overline{u}$ defined by (0.1) coincides with u defined by (0 0).

The advantage of the system (0.1) is that the matrix corresponding to the operator  $\tilde{A}$  is *block-diagonal*, with each block corresponding to the stress unknowns in a single element. Hence  $\tilde{A}$  may be inverted easily and inexpensively at the element level, leading to the expression

$$\overline{\sigma} = \tilde{A}^{-1}(f_1 - \tilde{B}\overline{u} - C\lambda) \tag{0.2}$$

for the approximate stress field in terms of the other unknowns Substituting (0.2) into the second and third equations of (0.1) yields the linear system

$$\tilde{B}^{I} \tilde{A}^{-1} \tilde{B}\bar{u} + \tilde{B}^{I} \tilde{A}^{-1} C\lambda = \tilde{B}^{I} \tilde{A}^{-1} f_{1} - f_{2}, 
C^{T} \tilde{A}^{-1} \tilde{B}\bar{u} + C^{T} \tilde{A}^{-1} C\lambda = C^{T} \tilde{A}^{-1} f_{1},$$
(0.3)

which is symmetric positive definite. One may thus solve this system and then recover the stress field from (0.2) by a simple element-by-element post-process.

M<sup>2</sup> AN Modelisation mathematique et Analyse numerique Mathematical Modelling and Numerical Analysis

This technique may be (and sometimes is) regarded as a purely implementational trick, i.e., as a computationally convenient way to determine the solution of the original system (0.0). Still, one computes in this way, in addition to  $\sigma$  and u, also the Lagrange multiplier field,  $\lambda$ , which frequently admits a physical interpretation such as a displacement and is often so used [8]. The mathematical convergence theory for mixed methods, however, generally considers only the approximation of the original fields, neglecting the further information offered by the multipliers.

In this paper we consider two typical examples of mixed methods, the method of Raviart and Thomas [16] for membrane (and other second order elliptic) problems, and the Hellan-Hermann-Johnson method [9], [10], [12], for plate bending (and other fourth order elliptic) problems. The technique described above applies to both methods and in each case we show that the multipliers so obtained can be used in the reconstruction of an approximate displacement field which is asymptotically *more accurate* than the approximation furnished by the original field u. Our analysis further reveals that each of these mixed methods is equivalent (in the sense of leading to the same approximate solution) to a certain nonconforming displacement method, which in each case is an apparently slight modification of a well-known method and a modified  $\mathfrak{P}^k$ -nonconforming method, and between the Hellan-Hermann-Johnson method and a modification of the Morley method [13], [14], [15].)

As we show by example, this dual mode of regarding the methods (either as mixed or nonconforming displacement) is useful for deriving error estimates as well as for implementating the methods. However it raises the question of whether the modifications which render a displacement method equivalent to a corresponding mixed method actually improve the method in some sense. We cannot answer this question simply by comparing the asymptotic rates of convergence (which are generally not affected by the modification), although in one case we do show superior robustness of the modified method with respect to the regularity of the data. A general answer must await further analysis or numerical experimentation.

#### **1. ERROR ESTIMATES FOR THE LAGRANGE MULTIPLIER**

### A. The Raviart-Thomas elements

Let us recall the main features of the Raviart-Thomas method [16]. Let  $\Omega$  denote a bounded domain in  $\mathbb{R}^2$ , which, for the sake of simplicity, we sup-

pose to be a polygon. Let f be a given function in  $L^2(\Omega)$  and  $a_{z}$  a sufficiently smooth two by two matrix-valued function on  $\Omega$ . We assume that there exists  $\alpha > 0$  such that

$$\sum_{i,j} a_{ij}(\underline{x}) \,\xi_i \,\xi_j \ge \alpha \,\|\,\underline{\xi}\,\|^2 \quad \forall \underline{x} \in \Omega \quad \forall \underline{\xi} \in \mathbb{R}^2 \,. \tag{1.0}$$

Consider the boundary value problem :

$$\left. \begin{array}{c} \text{find } u \in H^{1}(\Omega) \text{ such that} \\ - \operatorname{div}\left(\underset{u = 0}{\operatorname{g} \operatorname{grad}} u\right) = f \quad \text{in } \Omega, \\ u = 0 \quad \text{on } \partial\Omega. \end{array} \right\}$$

$$(1.1)$$

It is well known that problem (1.1) has a unique solution. In the following we implicitly assume that u(x) has, at each step, the regularity required by the context. The exact requirements are easily obtained from inspection of the arguments. Note that if  $\Omega$  is convex and  $f \in H^s(\Omega)$  for some s > 0, then  $u \in H^r(\Omega)$  for some number r > 2 which depends on s and  $\Omega$ . In order to state a mixed formulation of (1.1) we define the space

$$H(\operatorname{div}; \Omega) = \left\{ \mathfrak{t} \mid \mathfrak{t} \in (L^2(\Omega))^2, \operatorname{div} \mathfrak{t} \in L^2(\Omega) \right\}$$
(1.2)

with the usual graph norm

$$\| \mathfrak{t} \|_{H(\operatorname{div};\Omega)}^{2} = \sum_{i=1}^{2} \| \tau_{i} \|_{L^{2}(\Omega)}^{2} + \| \operatorname{div} \mathfrak{t} \|_{L^{2}(\Omega)}^{2}, \qquad (1.3)$$

and set

$$\sigma = - \underset{z}{u} \operatorname{grad} u \in H(\operatorname{div}; \Omega) . \tag{1.4}$$

A mixed formulation of (1.1) is then

$$\begin{cases} find \ (\mathfrak{g}, u) \in H(\operatorname{div}; \Omega) \times L^{2}(\Omega) \text{ such that} \\ \int_{\Omega} \mathfrak{c}_{\mathfrak{g}}^{\mathfrak{g}} \mathfrak{c}_{\mathfrak{T}} d\underline{x} - \int_{\Omega} u \operatorname{div} \mathfrak{t} d\underline{x} = 0 \quad \forall \mathfrak{t} \in H(\operatorname{div}; \Omega) , \\ \int_{\Omega} v \operatorname{div} \mathfrak{g} d\underline{x} = \int_{\Omega} fv \, d\underline{x} \quad \forall v \in L^{2}(\Omega) , \end{cases}$$

$$\end{cases}$$

$$(1.5)$$

where  $c_{\tilde{z}} := a_{\tilde{z}}^{-1}$  is the compliance tensor. Problem (1.5) is obviously equivalent to (1.1) and (1.4).

We now introduce the Raviart-Thomas discretization of (1.5). We shall use the notations

$$\mathfrak{P}^{k}(S) := Polynomials of degree \leq k \text{ on } S, \qquad (1.6)$$

$$\mathfrak{P}^{k}(S) := \mathfrak{P}^{k}(S) \times \mathfrak{P}^{k}(S), \qquad (1.7)$$

$$\Re \mathfrak{T}^{k}(S) := \left\{ f \mid f(\underline{x}) = p(\underline{x}) + \underline{x}q(\underline{x}), p \in \mathfrak{P}^{k}(S), q \in \mathfrak{P}^{k}(S) \right\}, \quad (1.8)$$

for any integer  $k \ge 0$  and any domain  $S \subseteq \mathbb{R}^n$ ,  $n \ge 1$ . We consider now a regular sequence of decompositions  $\{\mathfrak{T}_h\}_h$  of  $\Omega$  into triangles (see [5]) and define

$$RT_{-1}^{k}(\mathfrak{T}_{h}) = \left\{ \mathfrak{t} \mid \mathfrak{t} \in (L^{2}(\Omega))^{2}, \mathfrak{t} \mid_{T} \in \mathfrak{RT}^{k}(T) \quad \forall T \in \mathfrak{T}_{h} \right\},$$
(1.9)

$$M_{-1}^{k}(\mathfrak{T}_{h}) = \left\{ v \mid v \in L^{2}(\Omega), v \mid_{T} \in \mathfrak{P}^{k}(T) \quad \forall T \in \mathfrak{T}_{h} \right\},$$

$$(1.10)$$

$$RT_{0}^{k}(\mathfrak{T}_{h}) = \{ \mathfrak{T} \mid \mathfrak{T} \in RT_{-1}^{k}(\mathfrak{T}_{h}), \text{ the normal component of } \tau \text{ is } continuous across the interelement boundaries }$$
$$= RT_{-1}^{k}(\mathfrak{T}_{h}) \cap H(\operatorname{div}; \Omega) .$$
(1.11)

For k a fixed nonnegative integer, the kth order Raviart-Thomas mixed method now reads as follows :

$$\begin{cases} \text{find} \quad (\mathfrak{T}_{h}, u_{h}) \in RT_{0}^{k}(\mathfrak{T}_{h}) \times M_{-1}^{k}(\mathfrak{T}_{h}) \quad \text{such that} \\ \int_{\Omega} \mathfrak{L}_{0}^{\infty} \mathfrak{T}_{0} d\mathfrak{X}_{0} - \int_{\Omega} u_{h} \operatorname{div} \mathfrak{T}_{0} d\mathfrak{X}_{0} = 0 \quad \forall \mathfrak{T} \in RT_{0}^{k}(\mathfrak{T}_{h}) , \\ \int_{\Omega} v \operatorname{div} \mathfrak{T}_{h} d\mathfrak{X}_{0} = \int_{\Omega} fv \, d\mathfrak{X}_{0} \quad \forall v \in M_{-1}^{k}(\mathfrak{T}_{h}) . \end{cases}$$

$$(1.12)$$

The following results are known ([16], [7], [6]).

**THEOREM** 1.1: For any  $k \ge 0$ , problem (1.12) has a unique solution. Moreover there exists  $\gamma > 0$ , independent of h, such that

$$\| \operatorname{\mathfrak{g}} - \operatorname{\mathfrak{g}}_{h} \|_{0} \leq \gamma | h |^{k+1} \| \operatorname{\mathfrak{g}} \|_{k+1}, \qquad (1.13)$$

$$\| u - u_h \|_0 \leq \gamma \| h \|^{k+1} \| u \|_{r-1}, \quad (1.14)$$

$$|| P_{h} u - u_{h} ||_{0} \leq \gamma |h|^{k+2} || u ||_{r}, \qquad \int \gamma - \max(n+2,3), \qquad (1.15)$$

where we denote by |h| the maximum diameter of the triangles of  $\mathfrak{T}_h$  and by  $P_h$  the orthogonal projection of  $L^2$  onto  $M_{-1}^k(\mathfrak{T}_h)$ .

Note that the linear system associated with (1.12) has the structure (0.0). We now introduce Lagrange multipliers on the interelement boundaries and

so obtain a system with the structure (0.1). To this end we require some further notation. Let  $\mathfrak{E}_h$  denote the set of *edges* of triangles in  $\mathfrak{T}_h$  and set

$$\mathfrak{E}_{h}^{\partial} = \{ e \in \mathfrak{E}_{h} \mid e \subset \partial \Omega \}, \quad \mathfrak{E}_{h}^{0} = \mathfrak{E}_{h} \setminus \mathfrak{E}_{h}^{\partial}. \tag{1.16}$$

For  $T \in \mathfrak{T}_h$ ,  $e \in \mathfrak{G}_h$ , denote by  $h_T$  and  $h_e$  their respective diameters. Let  $\mathfrak{n}_T$  denote the exterior unit normal to T and  $\mathfrak{n}_e$  one of the unit vectors normal to e. The space of multipliers we shall use is the space  $M_{-1}^k(\mathfrak{G}_h^0)$  of all functions on  $\cup \mathfrak{G}_h$  which restrict to polynomial functions of degree at most k on each  $e \in \mathfrak{G}_h^0$  and vanish on  $\cup \mathfrak{G}_h^0$ . Now if  $\tau \in \mathfrak{RT}^k(T)$  then  $\tau.\mathfrak{n}_e \in \mathfrak{P}^k(e)$  for each edge e of T. The following lemma is an immediate consequence.

LEMMA 1.2 : If 
$$\underline{\tau} \in RT_{-1}^{k}(\underline{\mathfrak{I}}_{h})$$
, then  $\underline{\tau} \in RT_{0}^{k}(\underline{\mathfrak{I}}_{h})$  iff  

$$\sum_{T \in \underline{\mathfrak{I}}_{h}} \int_{\partial T} \underline{\tau} \cdot \underline{n}_{T} \ \mu \ de = 0 \quad \forall \mu \in M_{-1}^{k}(\underline{\mathfrak{C}}_{h}^{0}) \ . \tag{1.17}$$

Now consider the extended problem,

$$\begin{aligned} & find \ (\overline{\mathfrak{G}}_{h}, \overline{u}_{h}, \lambda_{h}) \in RT_{-1}^{k}(\mathfrak{T}_{h}) \times M_{-1}^{k}(\mathfrak{T}_{h}) \times M_{-1}^{k}(\mathfrak{G}_{h}^{0}) \ such \ that \\ & (\mathbf{i}) \ \int_{\Omega} \underline{c}_{\underline{\mathfrak{S}}} \overline{\mathfrak{G}}_{h} \cdot \underline{\tau} \ d\underline{\chi} - \sum_{T} \left\{ \int_{T} \overline{u}_{h} \operatorname{div} \underline{\tau} \ d\underline{\chi} \\ & - \int_{\partial T} \lambda_{h} \ \underline{\tau} \cdot \underline{n}_{T} \ de \right\} = 0 \quad \forall \underline{\tau} \in RT_{-1}^{k}(\mathfrak{T}_{h}) , \\ & (\mathbf{i}) \ \sum_{T} \int_{T} v \ \operatorname{div} \ \overline{\mathfrak{Q}}_{h} \ d\underline{\chi} = \int_{\Omega} fv \ d\underline{\chi} \quad \forall v \in M_{-1}^{k}(\mathfrak{T}_{h}) . \end{aligned}$$

$$(1.18)$$

$$(\mathbf{i}) \ \sum_{T} \int_{\partial T} \mu \overline{\mathfrak{Q}} \cdot \underline{n}_{T} \ de = 0 \quad \forall \mu \in M_{-1}^{k}(\mathfrak{G}_{h}^{0}) . \end{aligned}$$

The proof of the following lemma is immediate.

**LEMMA** 1.3 : Problem (1.18) has a unique solution  $(\overline{\mathfrak{g}}_h, \overline{u}_h, \lambda_h)$ . Moreover  $\overline{\mathfrak{g}}_h = \mathfrak{g}_h$  and  $\overline{u}_h = u_h$ , where  $(\mathfrak{g}_h, u_h)$  is the unique solution of (1.12).

This allows us to identify  $(\overline{\mathfrak{G}}_h, \overline{u}_h)$  and drop the upper bars in (1.18). Note that the equation  $\overline{\mathfrak{G}}_h = \mathfrak{G}_h$  is an identity among vector-valued functions. Clearly the corresponding coefficient arrays on the computer will not be equal to each other (they have different dimensions !). Note also that problem (1.18) has the form (0.1) and that, if a basis for  $RT_{-1}^k(\mathfrak{T}_h)$  is assembled from bases for the  $\mathfrak{RT}^k(T)$  in the obvious way, the compliance matrix  $\overline{A}$  corres-

ponding to  $\int \underset{k}{\varepsilon} \mathfrak{S}_{h} \mathfrak{T} d\mathfrak{X}$  is block diagonal as required. It is also easily proved in the present case that the final matrix corresponding to the form (0.3) is positive definite.

Our aim is now to derive error bounds for  $\lambda_h - u$ , which is defined on  $\bigcup \mathfrak{E}_h^0$ . The use of  $\lambda_h$  to approximate u within an element will be discussed in the next section (in the case of k even). Defining the norms on  $M_{-1}^k(\mathfrak{E}_h^0)$ 

$$\|\mu_{h}\|_{0,h}^{2} = \sum_{e \in \mathfrak{E}_{h}^{0}} \|\mu_{h}\|_{0,e}^{2}, \qquad (1.19)$$

$$\|\mu_{h}\|_{-1,2,h}^{2} = \sum_{e \in \mathfrak{E}_{h}^{0}} h_{e} \|\mu_{h}\|_{0,e}^{2}, \qquad (1.20)$$

we now compare  $\lambda_h$  with  $\Pi_h u$ , defined to be the orthogonal projection of  $u \mid_{\bigcup \mathfrak{E}_h^0}$  onto  $M_{-1}^k(\mathfrak{E}_h^0)$  in the norm (1.19).

**THEOREM** 1.4 : There exist constants  $\gamma_1, \gamma_2$  independent of u and h such that, for every  $T \in \mathfrak{T}_h$  and every edge e of T,

$$\|\lambda_{h} - \Pi_{h} u\|_{0,e} \leq \gamma_{1}(h_{T}^{1/2} \| \mathfrak{g} - \mathfrak{g}_{h} \|_{0,T} + h_{T}^{-1/2} \| P_{h} u - u_{h} \|_{0,T}), \quad (1.21)$$

$$|\lambda_{h} - \prod_{h} u|_{-1/2,h} \leq \gamma_{2}(|h| \| \mathfrak{g} - \mathfrak{g}_{h} \|_{0,\Omega} + \| P_{h} u - u_{h} \|_{0,\Omega}).$$
(1.22)

*Proof* : Clearly (1.22) is an immediate consequence of (1.21). In order to prove (1.21) let us consider  $T \in \mathfrak{T}_h$  and  $e \subset \partial T$ . It is proved in [16] that there exists a unique  $\overline{\mathfrak{t}} \in \mathfrak{RL}^k(T)$  such that

$$\overline{\overline{\tau}} \cdot \underline{n}_{e} = \lambda_{h} - \Pi_{h} u \quad \text{on} \quad e, \overline{\overline{\tau}} \cdot \underline{n}_{T} = 0 \quad \text{on} \quad \partial T \setminus e, \overline{\overline{\tau}} \perp \mathfrak{P}^{k-1}(T) \quad \text{in} \quad L^{2}(T).$$

$$(1.23)$$

Then a simple scaling argument shows that

$$h_{T} \| \overline{z} \|_{1,T} + \| \overline{z} \|_{0,T} \leq \gamma h_{T}^{1/2} \| \lambda_{h} - \Pi_{h} u \|_{0,e}.$$
 (1.24)

We may now choose  $\tau$  in (1.18i) such that

$$\mathfrak{t} = \overline{\mathfrak{t}} \quad \text{in } T, \qquad \mathfrak{t} = 0 \quad \text{in } \Omega \setminus T,$$
(1.25)

which gives, using (1.23),

$$\int_{T} \underbrace{c}_{\widetilde{z}} \underbrace{\sigma}_{h} \cdot \overline{t} d\underline{x} - \int_{T} u_{h} \operatorname{div} \overline{t} d\underline{x} + \int_{e} \lambda_{h} (\lambda_{h} - \Pi_{h} u) de = 0. \quad (1.26)$$

On the other hand from (1.4) we have

$$\underbrace{c\sigma}_{\widetilde{z}} = -\operatorname{grad} u, \qquad (1.27)$$

so that Green's formula implies

$$\int_{T} \underbrace{c}_{\widetilde{\Sigma}} \overline{\sigma} \cdot \overline{\overline{\tau}} \, d\underline{x} - \int_{T} u \operatorname{div} \overline{\overline{\tau}} \, d\underline{x} + \int_{e} u(\lambda_{h} - \Pi_{h} u) \, de = 0 \, . \qquad (1.28)$$

Subtracting (1.28) from (1.26) and using the fact that  $\operatorname{div} \overline{\overline{\tau}} \in \mathfrak{P}^k(T)$  we have

$$\|\lambda_{h} - \Pi_{h} u\|_{0,e}^{2} = \int_{e} (\lambda_{h} - u) (\lambda_{h} - \Pi_{h} u) de$$
$$= -\int_{T} \underbrace{c}_{\tilde{u}}(\underline{\sigma}_{h} - \underline{\sigma}) \cdot \overline{\underline{\tau}} d\underline{x} + \int_{T} (u_{h} - P_{h} u) \operatorname{div} \overline{\underline{\tau}} d\underline{x} . \quad (1.29)$$

Finally (1.29) and (1.24) give (1.21).

COROLLARY 1.5 : We have

$$|\lambda_{h} - \Pi_{h} u|_{-1,2,h} \leq \gamma |h|^{k+2} ||u||_{r}, \quad r = \max(k+2,3), \quad (1.30)$$

with  $\gamma$  independent of u and h.

The proof is immediate from (1.13), (1.15), and (1.22).

*Remark* : The norm (1.20) may be interpreted as an  $L^2(\Omega)$  norm of a suitable extension of  $\mu_h$  to the whole  $\Omega$ . In this sense, the estimate (1.30) may not seem better than (1.15) itself. However, this is not the case. Consider for instance the simplest case k = 0 : the estimate (1.15) gives superconvergence of  $0(h^2)$  at the center of gravity of each element (but nothing better), but (1.30) implies, as we shall see later on, that the  $\mathfrak{P}_1$  nonconforming extension of  $\lambda_h$  has a distance  $0(h^2)$  from u in  $L^2(\Omega)$ . This kind of argument will be developed in detail in the next section.

### B. The Hellan-Herrmann-Johnson element.

We consider, for the sake of simplicity, a very special model problem. However it is quite easy to check that all the results hold unchanged, for instance, for a general plate bending problem with constant (or piecewise cons-

M<sup>2</sup> AN Modélisation mathématique et Analyse numérique Mathematical Modelling and Numerical Analysis

tant) coefficients. Minor changes allow the treatment of the case with variable coefficients. Our model problem will be the following :

$$\begin{cases} \text{find } \psi \in H^2(\Omega) \text{ such that} \\ \Delta^2 \psi = f \text{ in } \Omega, \\ \psi = \frac{\partial \psi}{\partial n} = 0 \text{ on } \partial\Omega, \end{cases}$$

$$(1.31)$$

which we may write in variational form as

find 
$$\psi \in H_0^2(\Omega)$$
 such that  

$$\int_{\Omega} D_{\varepsilon}^2 \psi : D_{\varepsilon}^2 \phi \, dx = \int_{\Omega} f\phi \, dx \quad \forall \phi \in H_0^2(\Omega) ,$$
(1.32)

where  $D_{\tilde{z}}^2 \varphi = (\partial^2 \varphi / \partial x_i \partial x_j)_{ij}$  is the tensor of second partials and the colon denotes the scalar product of tensors. We shall analyze here the lowest order case of the family of H-H-J elements. For more information see [9], [10], [12], [4], [7], [2]. We maintain the assumptions and notations of the last section concerning the domain  $\Omega$  and the triangulations { $\mathfrak{T}_h$ }. The mixed discretization is based on a factorization of (1.31) into the equations

$$\sigma_{\widetilde{z}} = D_{\widetilde{z}}^2 \psi, \qquad (1.33)$$

$$\sum_{i,j} \frac{\partial^2 \sigma_{ij}}{\partial x_i \, \partial x_j} = f \,, \tag{1.34}$$

and seeks to approximate  $\underset{\approx}{\sigma}$  and  $\psi$  simultaneously. To define the finite element space we define first

$$\mathfrak{H}_{-1}^{0} = \{ \mathfrak{t}_{z} \mid \tau_{12} = \tau_{21} \text{ and } \tau_{ij} \in M_{-1}^{0}(\mathfrak{T}_{h}) \text{ for } i, j = 1, 2 \}$$
 (1.35)

and set

$$M_n(\underline{z}) = \underbrace{\tau}_{\underline{z}} \underbrace{n}_{e} \cdot \underbrace{n}_{e} \quad \text{on} \quad e, e \in \mathfrak{E}_h.$$
(1.36)

Now we can define the finite element spaces of Hellan-Herrmann-Johnson method as

$$\mathfrak{H}_{0}^{0}(\mathfrak{T}_{h}) = \left\{ \mathfrak{T} \in \mathfrak{H}_{-1}^{0}(\mathfrak{T}_{h}) \mid M_{n}(\mathfrak{T}) \quad is \ continuous \ at \ the \\ interelement \ boundaries \right\}, \qquad (1.37)$$

$$\tilde{M}_{0}^{1}(\mathfrak{I}_{h}) = M_{-1}^{1}(\mathfrak{I}_{h}) \cap H_{0}^{1}(\Omega) . \qquad (1.38)$$

The discretized problem may be written as follows :

(i) 
$$\begin{aligned} & find \ (\mathfrak{g}_{h}, \psi_{h}) \in \mathfrak{H}_{0}^{0}(\mathfrak{T}_{h}) \times M_{0}^{1}(\mathfrak{T}_{h}) \ such \ that \\ & \int_{\Omega} \mathfrak{g}_{h}^{0} : \mathfrak{T}_{\mathfrak{g}} d\mathfrak{X} + \sum_{T} \int_{\partial T} M_{n}(\mathfrak{T}) \frac{\partial \psi_{h}}{\partial n} de = 0 \quad \forall \mathfrak{T}_{\mathfrak{g}} \in \mathfrak{H}_{0}^{0}(\mathfrak{T}_{h}) , \end{aligned}$$

(ii) 
$$\sum_{T} \int_{\partial T} M_{n}(\underline{\mathfrak{g}}_{h}) \frac{\partial \varphi}{\partial n} de = - \int_{\Omega} f \varphi d\underline{\mathfrak{x}} \qquad \forall \varphi \in \mathring{M}_{0}^{1}(\mathfrak{T}_{h}) .$$

The following results are known (see [12], [4], [7], [2]).

**THEOREM** 1.6 : Problem (1.39) has a unique solution. Moreover if the solution  $\psi$  of (1.31) belongs to  $H^3(\Omega)$  we have

$$\| \underbrace{\sigma}_{\underline{v}} - \underbrace{\sigma}_{\underline{v}}_{h} \|_{0} + \| \psi - \psi_{h} \|_{1} \leq \gamma | h | \| \psi \|_{3}$$
(1.40)

with  $\gamma$  independent of  $\psi$  and h.

Problem (1.39) has again the structure (0.0) and hence we may again introduce Lagrange multipliers  $\lambda$  at the interelement boundaries in order to eliminate the condition of continuity of  $M_n(\underline{\sigma})$ . The proof of the following lemma is immediate.

LEMMA 1.7 : If 
$$\underline{\tau} \in \mathfrak{H}_{-1}^{0}(\mathfrak{T}_{h})$$
, then  $\underline{\tau} \in \mathfrak{H}_{0}^{0}(\mathfrak{T}_{h})$  iff  

$$\sum_{T \in \mathfrak{T}_{h}} \int_{\partial T} M_{n}(\underline{\tau}) \mu \underline{n}_{T} \cdot \underline{n}_{e} de = 0 \quad \forall \mu \in M_{-1}^{0}(\mathfrak{E}_{h}^{0}). \quad (1.41)$$

Hence we may consider the extended problem,

$$\begin{aligned} & find \ (\overline{\mathfrak{g}}_{h}, \overline{\psi}_{h}, \lambda_{h}) \in \mathfrak{H}_{-1}^{0}(\mathfrak{T}_{h}) \times \mathring{M}_{0}^{1}(\mathfrak{T}_{h}) \times M_{-1}^{0}(\mathfrak{E}_{h}^{0}) \ such \ that \\ & (\mathbf{i}) \quad \int_{\Omega} \overline{\mathfrak{o}}_{h} : \mathfrak{t} \ d\mathfrak{x} + \sum_{T} \left\{ \int_{\partial T} M_{n}(\mathfrak{t}) \ \frac{\partial \overline{\psi}_{h}}{\partial n} \ de \ - \int_{\partial T} M_{n}(\mathfrak{t}) \ \lambda_{h} \ \mathfrak{g}_{T} \cdot \mathfrak{g}_{e} \ de \right\} = 0 \\ & \quad \forall \mathfrak{r} \in \mathfrak{H}_{-1}^{0}(\mathfrak{T}_{h}) , \end{aligned}$$

$$\end{aligned}$$

$$\end{aligned}$$

$$\begin{aligned} & (\mathbf{ii}) \sum_{T} \int_{\partial T} M_{n}(\overline{\mathfrak{g}}_{h}) \ \frac{\partial \varphi}{\partial n} \ de \ = \ - \int f \varphi \ d\mathfrak{x} \quad \forall \varphi \in \mathring{M}_{0}^{1}(\mathfrak{T}_{h}) , \end{aligned}$$

$$\end{aligned}$$

$$\end{aligned}$$

$$\end{aligned}$$

$$\end{aligned}$$

$$\end{aligned}$$

$$\end{aligned}$$

$$\begin{aligned} & (\mathbf{iii}) \sum_{T} \int_{\partial T} M_{n}(\overline{\mathfrak{g}}_{h}) \ \mu \mathfrak{g}_{T} \cdot \mathfrak{g}_{e} \ de \ = \ 0 \qquad \forall \mu \in M_{-1}^{0}(\mathfrak{E}_{h}^{0}) . \end{aligned}$$

$$\end{aligned}$$

$$\end{aligned}$$

$$\end{aligned}$$

$$\end{aligned}$$

It is easy to prove the following lemma.

**LEMMA** 1.8 : Problem (1.42) has a unique solution  $(\overline{\mathfrak{g}}_h, \overline{\psi}_h, \lambda_h)$ . Moreover  $\overline{\mathfrak{g}}_h = \mathfrak{g}_h$  and  $\overline{\psi}_h = \psi_h$  where  $(\mathfrak{g}_h, \psi_h)$  is the unique solution of (1.39).

This allows us to identify  $\overline{\mathfrak{G}}_h$  with  $\mathfrak{G}_h$  and  $\overline{\psi}_h$  with  $\psi_h$  and drop the bars in (1.42). Following the pattern of part A of this section, we could now prove the convergence of  $\lambda_h$  (to  $\partial \psi / \partial n_e$ , in this case) and give a priori error bounds. We shall do this, but with a different technique. More precisely we shall introduce a nonconforming displacement method (a slight modification of the well-known Morley [13], [14], [15] method) and show that it is equivalent to (1.42). To this end we define

$$M_*^2(\mathfrak{T}_h) = \{ \varphi \in M_{-1}^2(\mathfrak{T}_h) \mid \varphi \text{ is continuous at the vertices and vanishes} \\ at the vertices of \partial\Omega, \partial\varphi/\partial n \text{ is continuous at the midpoint of} \\ each edge and vanishes at the midpoints of the edges in \partial\Omega \}. (1.43)$$

For a given  $\varphi \in H_0^2(\Omega) \cup M_*^2(\mathfrak{T}_h)$  we let  $\varphi^I$  be the *interpolant of*  $\varphi$  in  $\mathring{M}_0^1(\mathfrak{T}_h)$ , that is,  $\varphi^I$  is piecewise linear and continuous and coincides with  $\varphi$  at the vertices. We can now define our nonconforming displacement method as follows :

find 
$$\overset{*}{\Psi}^{h} \in M^{2}_{*}(\mathfrak{T}_{h})$$
 such that  

$$\sum_{T} \int_{T} \mathcal{D}^{2}_{\mathfrak{T}} \overset{*}{\Psi}^{h} : \mathcal{D}^{2}_{\mathfrak{T}} \varphi \, d\mathfrak{X} = \int_{\Omega} f\varphi^{l} \, d\mathfrak{X} \quad \forall \varphi \in M^{2}_{*}(\mathfrak{T}_{h}) .$$
(1.44)

Note that this method differs from the usual Morley method only by the presence of the interpolation operator in the right hand side. We shall now prove that the modified Morley method is actually *equivalent* with the method (1.42) (which in turn is *equivalent* with the original H-H-J method (1.39)). For this we need the following simple consequence of Green's formula :

$$if \tau_{ij} \in \mathfrak{P}^{0}(T), \quad \tau_{12} = \tau_{21} \text{ and } \varphi \in H^{2}(T) \text{ then}$$

$$\int_{T} \underbrace{\tau}_{z} : \underbrace{D}_{z}^{2} \varphi \, dx = \int_{\partial T} M_{n}(\underbrace{\tau}_{z}) \frac{\partial \varphi}{\partial n} \, de + \int_{\partial T} M_{nt}(\underbrace{\tau}_{z}) \frac{\partial \varphi}{\partial t} \, de$$

$$(1.45)$$

where  $M_n(\underline{z})$  is defined in (1.36) and  $M_m(\underline{z}) = \underline{z} \underline{n}_e \cdot t_e$  with  $t_e = (-n_e^2, n_e^1)$  denoting the unit vector tangent to e. We deduce as a consequence of (1.45) that

$$if_{\widetilde{z}} \in \mathfrak{H}_{0}^{0} \text{ and } \varphi \in M_{*}^{2}(\mathfrak{I}_{h}) + H_{0}^{2}(\Omega) \text{ then} \\ \sum_{T} \int_{T} \mathfrak{t}_{\widetilde{z}} : D_{\widetilde{z}}^{2} \varphi \, dx = \sum_{T} \int_{\partial T} M_{n}(\mathfrak{t}) \frac{\partial \varphi}{\partial t} \, de \, .$$

$$(1.46)$$

Finally we also note that

if 
$$\tau_{ij} \in \mathfrak{P}^{0}(T)$$
 and  $\varphi \in H^{2}(T)$  with  $\varphi = 0$  at the corners then,  

$$\int_{\partial T} M_{ni}(\underline{\tau}) \frac{\partial \varphi}{\partial t} de = 0.$$
(1.47)

Let now  $\hat{\Psi}^h$  be the solution of (1.44). We associate with it the functions

$$\underset{\widetilde{\mathfrak{S}}}{\overset{*}{\mathfrak{S}}} \in \mathfrak{H}_{-1}^{0}(\mathfrak{T}_{h}) \quad \text{defined by} \quad \underset{\widetilde{\mathfrak{S}}}{\overset{*}{\mathfrak{S}}} = \underbrace{D}_{\widetilde{\mathfrak{S}}}^{2} \overset{*}{\mathfrak{V}}^{h} \text{ in each } T , \qquad (1.48)$$

$$\stackrel{*}{\chi} \in \overset{\circ}{M_0^1}(\mathfrak{T}_h)$$
 defined by  $\stackrel{*}{\chi} = \underbrace{(\stackrel{*}{\Psi}{}^h)^I}_{*},$  (1.49)

$$\overset{*}{\lambda} \in M^{0}_{-1}(\mathfrak{E}^{0}_{h}) \text{ defined by } \overset{*}{\lambda} = \frac{\partial \psi^{h}}{\partial n_{e}} \text{ on } e \quad \forall e \in \mathfrak{E}^{0}_{h}, \qquad (1.50)$$

where, for  $\varphi \in M^2_*(\mathfrak{T}_h)$ , we have set

$$\frac{\overline{\partial \varphi}}{\partial n_e} := \text{ the value of } \frac{\partial \varphi}{\partial n_e} \text{ at the midpoint of } e, \quad e \in \mathfrak{E}_h^0.$$
(1.51)

**THEOREM** 1.9 : The triple  $(\xi, \chi, \lambda)$  defined in (1.48)-(1.50) is the solution of (1.42).

*Proof*: Using (1.46) and (1.45) we have

$$\int_{\Omega} \overset{\bullet}{\underset{\varepsilon}{\mathfrak{T}}} : \underset{\varepsilon}{\mathfrak{T}} d \underset{\varepsilon}{\mathfrak{T}} = \sum_{T} \left\{ \int_{\partial \mathfrak{T}} M_{n}(\underset{\varepsilon}{\mathfrak{T}}) \frac{\partial \overset{\bullet}{\psi}^{h}}{\partial n} d e + \int_{\partial \mathfrak{T}} M_{m}(\underset{\varepsilon}{\mathfrak{T}}) \frac{\partial \overset{\bullet}{\psi}^{h}}{\partial t} d e \right\}. \quad (1.52)$$

From (1.50) we get

$$\int_{\partial T} M_n(\underline{\tau}) \frac{\partial \underline{\psi}^h}{\partial n} de = \int_{\partial T} M_n(\underline{\tau}) \overset{*}{\lambda} n_T \cdot n_e de , \qquad (1.53)$$

and from (1.49), (1.47), and (1.45),

$$\int_{\partial T} M_{nt}(\underline{z}) \frac{\partial \underline{\psi}^{h}}{\partial t} de = \int_{\partial T} M_{nt}(\underline{z}) \frac{\partial \underline{\chi}}{\partial t} de = - \int_{\partial T} M_{n}(\underline{z}) \frac{\partial \underline{\chi}}{\partial n} de. \quad (1.54)$$

From (1.52)-(1.54) we see that (1.42i) is satisfied with  $(\overline{\mathfrak{g}}_h, \overline{\psi}_h, \lambda_h)$  replaced by  $(\overset{*}{\mathfrak{g}}, \overset{*}{\chi}, \overset{*}{\lambda})$ . Next, for each  $\mu \in M^0_{-1}(\mathfrak{E}^0_h)$  we define  $\varphi = \varphi(\mu) \in M^2_*(\mathfrak{E}_h)$  by

$$\frac{\partial \varphi}{\partial n_e} = \mu \quad \forall e \in \mathfrak{E}_h^0 \quad (\text{see } (1.51)), \qquad (1.55)$$

$$\varphi^{l} = 0$$
. (1.56)

Hence we have for all  $\mu \in M^0_{-1}(\mathfrak{E}^0_h)$  that

$$\sum_{T} \int_{\partial T} M_{n}(\overset{*}{\mathfrak{G}}) \mu_{\widetilde{\mathfrak{O}}_{T}} \cdot \underline{n}_{e} de = \sum_{T} \int_{\partial T} M_{n}(\overset{*}{\mathfrak{G}}) \frac{\partial \varphi}{\partial n} de . \qquad (1.57)$$

Now note that (1.45), (1.56), and (1.47) give

$$\int \overset{*}{\underset{\simeq}{\mathfrak{S}}} : \underbrace{D}_{\underset{\sim}{\mathfrak{S}}}^{2} \varphi \, d_{\underset{\sim}{\mathfrak{S}}} = \sum_{T} \int_{\partial T} M_{n}(\overset{*}{\underset{\simeq}{\mathfrak{S}}}) \frac{\widehat{c}\varphi}{\widehat{c}n} \, de \, . \tag{1.58}$$

From (1.57), (1.58), (1.44), and (1.56) we have now, for all  $\mu \in M_{-1}^{0}(\mathfrak{E}_{h}^{0})$ , that

$$\sum_{T} \int_{\partial T} M_n(\overset{*}{\underset{\sim}{\odot}}) \, \mu_{n_T} \cdot \underline{n}_e \, de = 0 \,, \qquad (1.59)$$

which is (1.42 iii); hence  $\mathfrak{F} \in \mathfrak{H}_0^0(\mathfrak{T}_h)$ . It remains to prove (1.42 ii) with  $\overline{\mathfrak{g}}_h$  replaced by  $\mathfrak{F}_0$ .

Associate, to each  $\varphi \in \overset{\circ}{M}_{0}^{1}(\mathfrak{T}_{h}), \zeta = \zeta(\varphi) \in M^{2}_{*}(\mathfrak{T}_{h})$  such that

$$\sigma_{\rm p}^{\rm d} = \varphi \,, \qquad (1.60)$$

$$\frac{\partial \zeta}{\partial n_e} = 0 \quad \text{on each} \quad e \in \mathfrak{E}_h^0 . \tag{1.61}$$

The using (1.45), (1.60), (1.47), and then again (1.45), we get for each  $T \in \mathfrak{T}_h$  that

$$\int_{\partial T} M_n(\overset{*}{\underline{\delta}}) \frac{\partial \varphi}{\partial n} de = - \int_{\partial T} M_{nt}(\overset{*}{\underline{\delta}}) \frac{\partial \varphi}{\partial t} de = - \int_{\partial T} M_{nt}(\overset{*}{\underline{\delta}}) \frac{\partial \zeta}{\partial t} de =$$
$$= \int_{\partial T} M_n(\overset{*}{\underline{\delta}}) \frac{\partial \zeta}{\partial n} de - \int_T \overset{*}{\underline{\delta}} : D^2 \zeta d\underline{x} . \quad (1.62)$$

Summing (1.62) over T, using the fact that  $\check{\mathfrak{F}} \in \mathfrak{H}_0^0(\mathfrak{T}_h)$  and applying (1.61), and finally using (1.48), (1.44), and (1.60), we obtain (1.42 ii). This completes the proof of theorem 1.9.

Note that the equivalence proved in the theorem can be used in both directions : from the solution  $\psi^h$  of (1.44) one can deduce the solution of (1.42) through (1.48)-(1.50), and, on the other hand, knowing the solution ( $\overline{\mathfrak{g}}_h$ ,  $\psi_h$ ,  $\lambda_h$ ) of (1.42) one can reconstruct  $\psi_h$ , the solution of (1.44), by

As a matter of fact (1.63) shows that such a reconstruction is much easier if (1.39) is solved in the equivalent form (1.42). The equivalence proved in theorem 1.9 can be a useful tool at the implementation level : according to the circumstances any of the formulations (1.39), (1.42), or (1.44) might be easiest to implement, although in our opinion (1.44) will usually be superior. We want to show now that the equivalence is indeed a very useful tool also in the asymptotic error analysis. First, we have as an immediate consequence of theorem 1.6, lemma 1.8, and theorem 1.9 the following error estimate for the modified Morley method.

**THEOREM** 1.10 : If  $\psi$ , the solution of (1.31), is in  $H^3(\Omega)$  and if  $\overset{*}{\psi}{}^h$  is the solution of the modified Morley method (1.44), then

$$\| \psi - \mathring{\psi}^{h} \|_{2,h} := \left( \sum_{T} \| D_{\approx}^{2} (\psi - \mathring{\psi}^{h}) \|_{0,T}^{2} \right)^{1/2} \leq \gamma \| h \| \| \psi \|_{3} \quad (1.64)$$

with  $\gamma$  independent of  $\psi$  and h.

Note that if  $\Omega$  is convex then (1.64) may be replaced by

 $\| \Psi - \mathring{\Psi}^{h} \|_{2,h} \leq \gamma | h | \| f \|_{-1} .$  (1.65)

As a consequence of theorem 1.10 we may deduce in a new way the known error estimate for the *usual* Morley method :

find 
$$w_h \in M^2_*(\mathfrak{T}_h)$$
 such that  

$$\sum_T \int_T \mathcal{D}^2_{\infty} w_h : \mathcal{D}^2_{\infty} \varphi \, d\mathfrak{X} = \int f\varphi \, d\mathfrak{X} \quad \forall \varphi \in M^2_*(\mathfrak{T}_h) .$$
(1.66)

COROLLARY 1.11 : If  $f \in L^2(\Omega)$  and if the solution  $\psi$  of (1.31) belongs to  $H^3(\Omega)$ , then

 $\| w_{h} - \psi \|_{2,h} \leq \gamma | h | (\| \psi \|_{3} + | h | \| f \|_{0})$  (1.67)

with  $\gamma$  independent of  $\psi$  and h.

*Proof* : Subtracting (1.44) from (1.66) with  $\varphi = \overset{*}{\psi}{}^{h} - w_{h}$  we have

$$\| \overset{*}{\Psi}{}^{h} - w_{h} \|_{2,h}^{2} = \int_{\Omega} f[(\overset{*}{\Psi}{}^{h} - w_{h})^{I} - (\overset{*}{\Psi}{}^{h} - w_{h})] dx . \qquad (1.68)$$

Since for all  $\varphi \in M^2_*(\mathfrak{T}_h)$  we have by scaling that

$$\| \phi' - \phi \|_{0} \leq \gamma | h |^{2} \| \phi \|_{2,h}, \qquad (1.69)$$

we easily obtain (1.67) from (1.64), (1.68). (1.69). and the triangle inequality.

For the original proof of (1.67) see [14]. For more information on the usual Morley method see [11], [5, p. 374], [15]. A comparison between (1.65) and (1.67) shows that the modified Morley method is *superior* to the usual one, at least with respect to the required regularity on f. We explicitly point out that a result of the form (1.64) or (1.65) cannot be true for the usual Morley method. This is obvious for (1.65), since  $M_*^2(\mathfrak{T}_h) \neq H_0^1(\Omega)$  so that the method (1.66) cannot be applied for a general  $f \in H^{-1}(\Omega)$ . Assume now that, for a fixed h, f is defined as a bounded linear functional on  $H_0^1(\Omega) + M_*^2(\mathfrak{T}_h)$  and set

$$\gamma_*(f,h) := \sup_{\varphi \in M^2_*(\mathfrak{T}_h)} \langle f, \varphi \rangle / \| \varphi \|_{1,h}, \qquad (1.70)$$

where

$$\| \phi \|_{1,h}^2 := \sum_T \| \underset{\mathcal{T}}{\text{grad}} \phi \|_{0,T}^2.$$
 (1.71)

From (1.68) we may easily deduce as in corollary 1.11 that

$$\|\psi - w_h\|_{2,h} \leq \gamma |h| (\|\psi\|_3 + \gamma_*(f,h)).$$
 (1.72)

However this is not of the form (1.64) unless we allow  $\gamma$ , in (1.64), to depend on *f*. More precisely, assume that we had

$$\|\psi - w_{h}\|_{2,h} \leq \gamma \|h\| \|\psi\|_{3}$$
(1.73)

with  $\gamma$  independent of h and  $\psi$  (hence independent of f) for all f defined on  $H_0^1(\Omega) + M_h^2(\mathfrak{T}_h)$ . For fixed h take  $\overline{f} \in H^{-1}(\Omega)$  such that

$$\sup_{\varphi \in \mathcal{M}^{2}_{*}(\mathfrak{T}_{h})} \int_{\Omega} \overline{f} \varphi \, d\underline{x} / \| \varphi \| = + \infty \qquad (1.74)$$

and take now a sequence  $f^{(n)} \in L^2(\Omega)$  such that  $f^{(n)} \to \overline{f}$  in  $H^{-1}(\Omega)$ . Assuming that  $\Omega$  is convex, the corresponding solutions,  $\psi^{(n)}$  of (1.31) are uniformly bounded in  $H^3(\Omega)$ . However (1.74) implies that  $|| w_h^{(n)} || \to \infty$  as  $n \to \infty$ and hence (1.73) must be false. Note that this example does not contradict (1.72) since  $\gamma_*(f_n, h)$  will also tend to infinity with n.

Our next goal is to use the equivalence of (1.42) and (1.44) to prove a new duality estimate.

**THEOREM** 1.12 : If  $\Omega$  is convex and  $f \in L^2(\Omega)$  then

$$\| \Psi - \overset{*}{\Psi}{}^{h} \|_{1,h} \leq \gamma |h|^{2} (\| \Psi \|_{3} + \| f \|_{0})$$
(1.75)

with  $\gamma$  independent of  $\psi$  and h.

*Proof*: Let us set  $\vartheta := \psi - \psi^{*}, \ \vartheta^{I} :=$  interpolant of  $\vartheta$  in  $\mathring{M}_{0}^{1}(\mathfrak{T}_{h})$ , and  $q = -\Delta \vartheta^{I} \in H^{-1}(\Omega)$ .

Consider the auxiliary Dirichlet problem for the biharmonic :

find 
$$\zeta \in H_0^2(\Omega)$$
 such that  

$$\int_{\Omega} \mathcal{D}_{\varepsilon}^2 \zeta : \mathcal{D}_{\varepsilon}^2 \varphi \, dx = \langle q, \varphi \rangle \quad \forall \varphi \in H_0^2(\Omega) .$$
(1.76)

Clearly we have

$$\| \zeta \|_{3} \leq c \| q \|_{-1} \leq c \| \vartheta' \|_{1}.$$
 (1.77)

On the other hand we have

$$\| \mathfrak{G}^{l} \|_{1}^{2} = \langle q, \mathfrak{G}^{l} \rangle$$

$$= \left[ \sum_{T} \int_{T} \mathcal{D}_{z}^{2} \zeta : \mathcal{D}_{z}^{2} \mathfrak{G} d_{\tilde{x}} \right] - \left[ \sum_{T} \int_{T} \mathcal{D}_{z}^{2} \zeta : \mathcal{D}_{z}^{2} \mathfrak{G} d_{\tilde{x}} - \langle q, \mathfrak{G}^{l} \rangle \right]$$

$$:= E_{1} - E_{2} . \qquad (1.78)$$

We bound  $E_1$  and  $E_2$  separately. Let  $\zeta_h$  be the usual interpolant of  $\zeta$  in  $M^2_*(\mathfrak{T}_h)$ . We have that

$$E_{1} = \sum_{T} \int_{T} D_{z}^{2} \zeta : D_{z}^{2} (\psi - \psi^{h}) dx = \int_{\Omega} f\zeta dx - \int_{\Omega} f\zeta_{h}^{I} dx - \sum_{T} \int_{T} D_{z}^{2} (\zeta - \zeta_{h}) : D_{z}^{2} \psi^{h} dx . \quad (1.79)$$

However theorem 1.9 together with (1.46) and (1.47) implies that

$$\sum_{T} \int_{T} \sum_{z}^{2} (\zeta - \zeta_{h}) : \sum_{z}^{2} \psi^{h} dx = 0, \qquad (1.80)$$

so from (1.79) we get

$$|E_{1}| \leq \gamma |h|^{2} ||f||_{0} ||\zeta||_{2}.$$
(1.81)

#### M<sup>2</sup> AN Modélisation mathématique et Analyse numérique Mathematical Modelling and Numerical Analysis

22

Now set  $\mathfrak{t} := D^2_{\mathfrak{s}} \zeta$ . We have that

$$E_{2} = \sum_{T} \int_{T} \underbrace{\mathfrak{r}}_{\Xi} : \underbrace{D}_{\Xi}^{2} \vartheta \, d\underline{\mathfrak{x}} - \langle q, \vartheta' \rangle$$
  
$$= \sum_{T} \left\{ -\int_{T} \operatorname{div} \underbrace{\mathfrak{r}}_{\Xi} : \underbrace{\operatorname{grad}}_{\Xi} \vartheta \, d\underline{\mathfrak{x}} + \int_{\partial T} \left[ M_{n}(\underbrace{\mathfrak{r}}_{\Xi}) \frac{\partial \vartheta}{\vartheta n} + M_{n'}(\underbrace{\mathfrak{r}}_{\Xi}) \frac{\partial \vartheta}{\vartheta t} \right] de \right\} - \langle q, \vartheta' \rangle$$
  
$$= -\sum_{T} \int_{T} \operatorname{div} \underbrace{\mathfrak{r}}_{\Xi} : \underbrace{\operatorname{grad}}_{\Xi} \vartheta \, d\underline{\mathfrak{x}} - \langle q, \vartheta' \rangle + E_{3}, \qquad (1.82)$$

where  $E_3$  is defined as the sum of the integrals over the element boundaries appearing in this equation.

However, since it is easily seen that

$$-\sum_{T} \int_{T} \operatorname{div}_{\widetilde{z}} \operatorname{grad}_{\widetilde{z}} \vartheta' \, d\widetilde{z} = \langle q, \vartheta' \rangle, \qquad (1.83)$$

(1.82) becomes

$$E_2 = -\sum_T \int_T \operatorname{div}_{\widetilde{z}} \operatorname{grad}_{\widetilde{z}} \left(\vartheta - \vartheta'\right) dx + E_3, \qquad (1.84)$$

and therefore

$$|E_{2}| \leq |E_{3}| + ||\zeta||_{3} |h| ||\vartheta||_{2,h}.$$
(1.85)

Next, to estimate  $E_3$  we note that the jumps of  $\frac{\partial \vartheta}{\partial n}$  and  $\frac{\partial \vartheta}{\partial t}$  have zero mean value on each interelement boundary. Setting  $\overline{M} :=$  projection of M onto  $M_{-1}^0(\mathfrak{E}_h^0)$  in the  $L^2$  norm we have

$$E_{3} = \sum_{T} \int_{\partial T} \left[ (M_{n}(\underline{\tau}) - \overline{M}_{n}) \frac{\partial \vartheta}{\partial n} + (M_{nt}(\underline{\tau}) - \overline{M}_{nt}) \frac{\partial \vartheta}{\partial t} \right] de \qquad (1.86)$$

so that

$$|E_{3}| \leq c |h| || \zeta ||_{3} || \vartheta ||_{2,h}.$$
(1.87)

Combining (1.78), (1.81), (1.85), (1.87), and (1.77) we get

$$\| \mathfrak{G}' \|_{1}^{2} \leq c \| \mathfrak{G}' \|_{1} \left( |h|^{2} \| f \|_{0} + |h| \| \mathfrak{G} \|_{2,h} \right).$$
(1.88)

which implies

$$\| \vartheta' \|_{1} \leq c(\|h\|^{2} \| f \|_{0} + \|h\| \| \vartheta \|_{2,h}).$$
 (1.89)

Finally we have from the triangle inequality and (1.89) that

$$\| \mathfrak{g} \|_{1,h} \leq \| \mathfrak{g} - \mathfrak{g}' \|_{1,h} + \| \mathfrak{g}' \|_{1} \leq c \| h \| \| \mathfrak{g} \|_{2,h} + \| \mathfrak{g}' \|_{1}$$
  
 
$$\leq c(\| h \|^{2} \| f \|_{0} + h \| \mathfrak{g} \|_{2,h}).$$
(1.90)

Using (1.90) and (1.64) we get (1.75).

*Remark* : In the proof of theorem 1.12 we see that the term  $|| f ||_0$  in (1.75) appears only in the estimate (1.79)-(1.81) of  $E_1$  through the bound :

$$\int_{\Omega} (\zeta - \zeta'_h) f \, dx \leq \gamma \parallel f \parallel_0 |h|^2 \parallel \zeta \parallel_2.$$
 (1.91)

Hence a slight modification of the proof gives for instance

$$\| \Psi - \tilde{\Psi}_h \|_{1,h} \leq \gamma_{\varepsilon} h^{2-\varepsilon} \| \Psi \|_{3-\varepsilon} \quad (\varepsilon > 0)$$
 (1.92)

when f is the Dirac measure at a vertex of the decomposition.

*Remark*: Setting  $\vartheta := w_h - \overset{*}{\psi}{}^h$  (where  $w_h$  is the solution of (1.66)) and repeating the arguments of theorem 1.12 one sees that nothing changes except for the estimate of  $E_1$ . This now reads :

$$E_{1} = \sum_{T} \int_{T} \sum_{\tilde{z}}^{2} \zeta : \sum_{\tilde{z}}^{2} (w_{h} - \tilde{\psi}^{h}) d\tilde{z} = \sum_{T} \int_{T} \sum_{\tilde{z}}^{2} (\zeta - \zeta_{h}) : \sum_{\tilde{z}}^{2} (w_{h} - \tilde{\psi}^{h}) d\tilde{z} + \int_{\Omega} f(\zeta_{h} - \zeta_{h}^{I}) d\tilde{z}$$
(1.93)

so that

$$|E_{1}| \leq \gamma(|h| || \zeta ||_{3} || \vartheta ||_{2,h} + || f ||_{0} |h|^{2} || \zeta ||_{2}).$$
(1.94)

Hence one may show that

$$\|w_{h} - \tilde{\Psi}^{h}\|_{1,h} \leq \gamma \|h\|^{2} (\|\psi\|_{3} + \|f\|_{0})$$
(1.95)

which joined to (1.75) gives

$$\| \Psi - w_h \|_{1,h} \leq \gamma \| h \|^2 \left( \| \Psi \|_3 + \| f \|_0 \right).$$
 (1.96)

This is a new error estimate for the usual Morley method.

#### 2. DISPLACEMENT FORMULATION AND POST-PROCESSING OF THE RAVIART-THOMAS MIXED METHODS

In the previous section we presented a different implementation technique for two mixed methods, making use of Lagrangian multipliers  $\lambda_h$  at the interelement boundaries. In theorem 1.4 we gave an estimate for the distance  $\lambda_{\rm b} - u$  on  $\mathfrak{E}_{\rm b}^0$  in the case of the Raviart-Thomas methods. Then we shifted to the Hellan-Herrmann-Johnson methods and proved in theorem 1.9 that this is equivalent to a slight modification of a classical nonconforming displacement method, the Morley method. The equivalence proved to be very fruitful not only from the point of view of implementation but also from the point of view of error analysis : in particular the known error estimate (1.40) for the mixed allowed the very simple proof of the estimate (1.67) for the Morley method. On the other hand, using the displacement formulation we proved the duality estimate (1.75) which one cannot naturally derive from the original mixed formulation (1.39), based as it is on piecewise linear displacements. Of course, a posteriori, this can be done; we would claim, however, that such an estimate does not come in mind looking at the formulation (1.39). In turn (1.75) was employed to prove the duality estimate (1.96) for the usual Morley method.

Our next goal is to do something in this direction for the Raviart-Thomas methods of section 1-A. However this time we shall first deduce an error estimate in  $L^2(\Omega)$  from the estimate " on the edges " (1.21); for the sake of simplicity we do this only in the case of k even (which includes the lowest order case k = 0); the case of odd k presents more technical difficulties as we discuss. On the other hand we shall stick to the case of variable coefficients which complicates the equivalence between the mixed method and a nonconforming displacement method (as in theorem 1.9 for the H-H-J method). Therefore we separately remark on the simplest case of k = 0, a constant.

We return now to the notations of section 1-A and to the estimate (1.21), which we now write as

$$\|\lambda_{h} - \Pi_{h}^{k} u\|_{0,e} \leq \gamma(h_{T}^{1/2} \| \mathfrak{g} - \mathfrak{g}_{h} \|_{0,T} + h_{T}^{-1/2} \| P_{h}^{k} u - u_{h} \|_{0,T}), \quad (2.0)$$

where

$$\Pi_h^k := L^2 \text{ projection onto } M_{-1}^k(\mathfrak{E}_h), \qquad (2.1)$$

$$P_h^k := L^2 \text{ projection onto } M_{-1}^k(\mathfrak{T}_h).$$
(2.2)

We have now two pieces of information at our disposal,  $\lambda_h$ , which is a polynomial of degree  $\leq k$  on each  $e \in \mathfrak{E}_h^0$ , and  $u_h$ , which is a polynomial of degree

 $\leq k$  in each  $T \in \mathfrak{T}_h$ . We shall use them in order to construct a new approximation,  $\tilde{u}_h$ , which is of degree k + 1 in each T and which converges to u faster than  $u_h$ . In order to define  $\tilde{u}_h$  we need the following lemma.

**LEMMA** 2.1. Let k be a nonnegative even integer and let  $T \in \mathfrak{T}_h$  be a triangle with edges  $e_1, e_2, e_3$ . Then for all  $p_i \in L^2(e_i)$  (i = 1, 2, 3) and  $q \in L^2(T)$  there exists a unique  $\chi = \chi(p_i, q) \in \mathfrak{P}^{k+1}(T)$  such that

$$\int_{e_i} (\chi - p_i) z \, de = 0 \quad \forall z \in \mathfrak{P}^k(e_i) \quad i = 1, 2, 3 , \qquad (2.3)$$

$$\int_{T} (\chi - q) z \, d\underline{x} = 0 \quad \forall z \in \mathfrak{P}^{k-2}(T) \,. \tag{2.4}$$

Moreover,

$$\|\chi\|_{0,T} \leq \gamma \left( \|q\|_{0,T} + h_T^{1/2} \sum_{i=1}^3 \|p_i\|_{0,e_i} \right),$$
 (2.5)

with  $\gamma$  depending only on k and on the minimum angle of T.

Proof : Clearly (2.3), (2.4) is a square linear system with

$$3(k + 1) + k(k - 1)/2 = (k + 2)(k + 3)/2$$

equations and unknowns. Hence for proving existence and uniqueness of  $\chi$ it is enough to consider the case q = 0,  $p_i = 0$  (i = 1, 2, 3) and show that  $\chi = 0$ is the unique solution. Conditions (2.3) with  $p_i = 0$  imply that  $\chi | e_i$  is a multiple of the Legendre polynomial of degree k + 1 on each  $e_i$ . Since k + 1is odd, this implies that  $\chi$  takes opposite values at the endpoints of each  $e_i$ . Hence the continuity of  $\chi$  on  $\partial T$  and the nonvanishing of the Legendre polynomial at the endpoints of the interval imply that  $\chi \equiv 0$  on  $\partial T$ , and therefore that  $\chi$  has the form  $\chi = \lambda_1 \lambda_2 \lambda_3 \tilde{z}$  where  $\lambda_i$  (i = 1, 2, 3) are the barycentric coordinates on T and  $\tilde{z} \in \mathfrak{P}^{k-2}(T)$  (for  $k \ge 2$ ; for k = 0 the condition  $\chi \equiv 0$ on  $\partial T$  clearly implies that  $\chi \equiv 0$  in T). Taking  $z = \tilde{z}$  and q = 0 in (2.4) yields  $\chi \equiv 0$  since  $\lambda_1 \lambda_2 \lambda_3 > 0$  in the interior of T. This proves the existence and uniqueness of  $\chi \in \mathfrak{P}^{k+1}(T)$  satisfying (2.3) and (2.4). The inequality (2.5) follows by simple scaling arguments.

In light of lemma 2.1 we now use  $\lambda_h$  and  $u_h$  to define our "better approximation",  $\tilde{u}_h \in M_{-1}^{k+1}(\mathfrak{T}_h)$  by

$$\Pi_h^k \, \tilde{u}_h = \lambda_h \,, \tag{2.6}$$

$$P_h^{k-2}(\tilde{u}_h - u_h) = 0 \quad \text{(for } k \ge 2).$$
(2.7)

Note that by lemma 2.1  $\tilde{u}_h$  is uniquely determined and that (2.6) implies some continuity of  $\hat{u}_h$  at the interelement boundaries, together with some vanishing on  $\partial\Omega$ . More precisely,  $\tilde{u}_h$  is continuous at the k + 1 Gauss points of each edge  $e \in \mathfrak{E}_h^0$  and vanishes at the Gauss points of each  $e \in \mathfrak{E}_h^{\hat{\theta}}$  (where  $\lambda_h \equiv 0$ ). However, in general  $\hat{u}_h \notin \mathring{H}^1(\Omega)$ . Hence  $\tilde{u}_h$  is a *nonconforming* approximation of u. We now prove that it indeed approximates u with a higher order of accuracy than  $u_h$ .

**THEOREM** 2.2 : Let u be the solution of (1.1) and  $(\mathfrak{g}_h, u_h, \lambda_h)$  the solution of (1.18) (for k even). Define  $\tilde{u}_h \in M_{-1}^{k+1}(\mathfrak{T}_h)$  by (2.6), (2.7). Then

$$\| u - \tilde{u}_h \|_0 \leq \gamma | h |^{k+2} (\| u \|_r + \| \mathfrak{g} \|_{k+1}) \quad r = \max (k+2,3) \quad (2.8)$$

with  $\gamma$  independent of u and h.

*Proof*: We first define  $\tilde{u}_h^* \in M_{-1}^{k+1}(\mathfrak{T}_h)$ , the nonconforming projection of u, by

$$\Pi_{h}^{k}(u - \tilde{u}_{h}^{*}) = 0, \qquad (2.9)$$

$$P_h^{k-2}(u - \tilde{u}_h^*) = 0, \quad (k \ge 2).$$
 (2.10)

Lemma 2.1 implies existence and uniqueness of  $\tilde{u}_h^*$ ; by standard arguments it is easily proved that

$$\| u - \tilde{u}_{h}^{*} \|_{0} \leq \gamma | h |^{k+2} \| u \|_{k+2}.$$
(2.11)

Note that from (2.6), (2.7) and (2.9), (2.10) we get

$$\Pi_h^k(\tilde{u}_h - \tilde{u}_h^*) = \lambda_h - \Pi_h^k u, \qquad (2.12)$$

$$P_h^{k-2}(\tilde{u}_h - \tilde{u}_h^*) = P_h^{k-2}(u_h - u) = P_h^{k-2}(u_h - P_h^k u). \qquad (2.13)$$

Using now lemma 2.1 with  $q := u_h - P_h^k u$ ,  $p_i := (\lambda_h - \Pi_h^k u)|_{e_i}$  we obtain for each T in  $\mathfrak{T}_h$  that

$$\| \tilde{u}_{h} - \tilde{u}_{h}^{*} \|_{0,T} \leq \gamma \left( \| u_{h} - P_{h}^{k} u \|_{0,T} + h_{T}^{1/2} \sum_{i=1} \| \lambda_{h} - \Pi_{h}^{k} u \|_{0,c_{i}} \right).$$
(2.14)

Combining (2.14) with (2.0) and then using (1.13), (1.15) we get

$$\| \tilde{u}_{h} - \tilde{u}_{h}^{*} \|_{0} \leq \gamma \| h \|^{k+2} (\| u \|_{r} + \| \mathfrak{g} \|_{k+1}), \quad r = \max(k+2,3), \quad (2.15)$$

and (2.8) follows from (2.11), (2.15).

The nature of the technical difficulties connected with the case of odd k should be clear now. What is needed is a nonconforming element which is vol. 19, n<sup>o</sup> 1, 1985

exact up to the degree k + 1 (i.e., reproduces exactly the polynomials of degree  $\leq k + 1$ ) and uses as degrees of freedom, some or all of the quantities

$$\int_{e} \chi z \, de \,, \quad z \in \mathfrak{P}^{k}(e) \,, \qquad \int_{T} \chi z \, d\mathfrak{x} \,, \quad z \in \mathfrak{P}^{k}(T) \,. \tag{2.16}$$

For instance for k = 1 we may define  $\tilde{u}_h \in M^2_{-1}(\mathfrak{T}_h)$  by the orthogonalities

$$\Pi_{h}^{0}(\tilde{u}_{h}-\lambda_{h})=0, \qquad P_{h}^{1}(\tilde{u}_{h}-u_{h})=0, \qquad (2.17)$$

which would give (2.8) with k = 1. For k = 3 the choice

$$\Pi_{h}^{2}(\tilde{u}_{h}-\lambda_{h})=0, \qquad P_{h}^{2}(\tilde{u}_{h}-u_{h})=0, \qquad (2.18)$$

works and again gives (2.8). Other *ad hoc* choices may be made for each particular odd k. However we didn't find an elegant general structure. For this reason the following *interpretation of* (1.18) *as a generalized displacement method* will be carried out in the case of k even.

We introduce the spaces

$$M_{NC}^{k+1}(\mathfrak{T}_{h}) = \left\{ v \in M_{-1}^{k+1}(\mathfrak{T}_{n}) \mid v \text{ is continuous at the } k+1 \text{ Gauss} \\ \text{points of each } e \in \mathfrak{E}_{h}^{0} \text{ and vanishes at the } k+1 \text{ Gauss} \\ \text{points of each } e \in \mathfrak{E}_{h}^{\partial} \right\}.$$

$$(2.19)$$

$$B^{k+3}(\mathfrak{I}_h) = \{ v \in M_0^{k+3}(\mathfrak{I}_h) \mid v \text{ vanishes on each } e \in \mathfrak{E}_h \}.$$

$$(2.20)$$

$$N^{k+1}(\mathfrak{I}_{h}) = M_{NC}^{k+1}(\mathfrak{I}_{h}) + B^{k+3}(\mathfrak{I}_{h}). \qquad (2.21)$$

Note that  $B^{k+3}(\mathfrak{T}_h)$  consists of bubble functions; hence the space defined in (2.21) is a classical nonconforming space augmented with bubbles.

With the same arguments as in lemma 2.1 one easily proves the following lemma.

**LEMMA** 2.3 : Let k be an even integer  $\geq 0$ . For any  $v_h \in M_{-1}^k(\mathfrak{T}_h)$  and any  $\mu_h \in M_{-1}^k(\mathfrak{T}_h^0)$  there exists a unique  $\chi \in N^{k+1}(\mathfrak{T}_h)$  such that

$$P_h^k \chi = v_h, \qquad \Pi_h^k \chi = \mu_h. \qquad (2.22)$$

Let us go back now to problem (1.18). We have the following result.

**LEMMA** 2.4 : Let  $(\mathfrak{T}_h, u_h, \lambda_h)$  be the unique solution of (1.18) and let  $\psi_h \in N^{k+1}(\mathfrak{T}_h)$  be defined by

$$P_h^k \psi_h = u_h, \qquad \Pi_h^k \psi_h = \lambda_h. \qquad (2.23)$$

Then  $(\sigma_h, \psi_h)$  is the unique solution of the following problem :

$$\begin{array}{l} \text{find } (\mathfrak{T}_{h}, \psi_{h}) \in RT_{-1}^{k}(\mathfrak{T}_{h}) \times N^{k+1}(\mathfrak{T}_{h}) \text{ such that} \\ (i) \quad \int_{\Omega} \underbrace{c}_{\mathfrak{T}} \mathfrak{T}_{h}, \underbrace{c}_{\mathfrak{T}} dx + \sum_{T} \int_{T} \mathfrak{T}_{\mathfrak{T}} \mathfrak{g}_{\mathfrak{r}} \operatorname{grad}_{\mathfrak{T}} \psi_{h} d\mathfrak{X} = 0 \quad \forall \mathfrak{T} \in RT_{-1}^{k}(\mathfrak{T}_{h}) , \\ (ii) \quad \sum_{T} \int_{T} \mathfrak{T}_{h}, \operatorname{grad}_{\mathfrak{T}} \chi d\mathfrak{X} = - \int_{\Omega} (P_{h}^{k} f) \chi d\mathfrak{X} \quad \forall \chi \in N^{k+1}(\mathfrak{T}_{h}) . \end{array} \right)$$

$$(2.24)$$

The proof is immediate by Green's formula and obvious properties of projection operators.

Now assume temporarily that the compliance tensor c is of the form  $c_{ij} = c(\underline{x}) \, \delta_{ij}$  with  $c(\underline{x})$  constant on each  $T \in \mathfrak{T}_h$ . In this case (2.24 i) clearly means that  $\underset{\mathbb{Z}}{c} \mathfrak{G}_h$  is the L<sup>2</sup>-projection of  $-\operatorname{grad} \psi_h$  onto  $RT_{-1}^k(\mathfrak{T}_h)$ . Denoting by  $P_{RT}^{k}$  this projection operator, we may write  $c \sigma_{h} = -P_{RT}^{k} (\text{grad } \psi_{h})$  and, since  $c_{z} := a_{z}^{-1}$ , problem (2.24) is now equivalent to the following problem :

find 
$$\psi_h \in N^{k+1}(\mathfrak{T}_h)$$
 such that  

$$\sum_T \int_T \underbrace{a}_T P_{RT}^k(\operatorname{grad} \psi_h) \cdot \operatorname{grad} \chi \, d\chi = \int_\Omega (P_h^k f) \chi \, d\chi \quad \forall \chi \in N^{k+1}(\mathfrak{T}_h) \, . \quad \left\{ \begin{array}{c} (2.25) \\ \end{array} \right\}$$

Let us briefly discuss the structure of (2.25) in the special case k = 0. We then have  $N^{1}(\mathfrak{T}_{h}) = M^{1}_{NC}(\mathfrak{T}_{h}) + B^{3}(\mathfrak{T}_{h})$ . Note now that for  $v_{h} \in M^{1}_{NC}(\mathfrak{T}_{h})$ we have  $P_{RT}^{0}(\text{grad } v_h) = \text{grad } v_h$ , a piecewise constant. Moreover the gradient of a bubble function has zero mean value on each T. It follows that the solution  $\psi_h$  of (2.25) may be determined as  $z_h + \zeta_h$  where  $(z_h, \zeta_h)$  is the unique solution to the problem

We remark that (2.26) is block diagonal, in that (2.26 i) and (2.26 ii) may be solved independently. Moreover (2.26 ii) gives rise to a diagonal matrix equation, and so  $\zeta_h$  can be computed independently in each triangle. The vol. 19, nº 1, 1985

1

system (2.26 i) on the other hand coincides with the usual  $\mathfrak{P}^1$  nonconforming method, except for the appearance of  $P_h^0 f$  rather than f in the right hand side. It appears to us that (2.26) offers the simplest implementation of the lowest order Raviart-Thomas element in the case of a piecewise constant diagonal coefficient matrix.

Let us go back now to the case of a general coefficient matrix. It will prove convenient to introduce the operator  $P_{RT,\underline{c}}^k$ , defined as the projection operator onto  $RT_{-1}^k$  with respect to the scalar product

$$[\mathfrak{g},\mathfrak{r}] := \int_{\Omega} \mathfrak{g}\mathfrak{g}.\mathfrak{r} \, d\mathfrak{x} \, . \tag{2.27}$$

Then writing in (2.24 i) grad  $\psi_h = \underset{\approx}{c} (\underset{\approx}{a} \operatorname{grad} \psi_h)$  we have  $\mathfrak{g}_h = -P_{RT,\underline{s}}^k (\underset{\approx}{\underline{s}} \operatorname{grad} \psi_h)$  and substituting this expression into (2.24 ii) we obtain the problem

find 
$$\psi_h \in N^{k+1}(\mathfrak{T}_h)$$
 such that  

$$\sum_T \int_T P_{RT,\underline{\xi}}^k (\underbrace{a}_{\underline{\xi}} \operatorname{grad} \psi_h) \cdot \operatorname{grad} \chi \, d\underline{\chi} = \int_{\Omega} (P_h^k f) \chi \, d\underline{\chi} \quad \forall \chi \in N^{k+1}(\mathfrak{T}_h) \,. \quad (2.28)$$

This is the displacement version of (1.18) in the case of a general coefficient matrix (and even k). The usual nonconforming method for this problem, on the other hand, reads

$$\begin{cases} \text{find } \overline{\psi}_h \in M_{NC}^{k+1}(\mathfrak{T}_h) \text{ such that} \\ \sum_T \int_T \underbrace{a}_{\widetilde{z}} \operatorname{grad} \overline{\psi}_h \cdot \operatorname{grad} \chi \, d\underline{x} = \int_{\Omega} f\chi \, d\underline{x} \quad \forall \chi \in M_{NC}^{k+1}(\mathfrak{T}_h) \, . \end{cases}$$

$$(2.29)$$

Let us point out the differences between (2.28) and (2.29).

1) On the right hand side of (2.28)  $P_h^k f$  appears in place of f.

2) The space  $M_{NC}^{k+1}(\mathfrak{T}_h)$  of (2.29) is augmented with bubble functions to obtain  $N^{k+1}(\mathfrak{T}_h)$  in (2.28).

3) The gradients are projected onto the Raviart-Thomas space  $RT_{-1}^{k}(\mathfrak{T}_{h})$  in (2.28).

Note that the projection referred to in 3) is the orthogonal projection with respect to the scalar product (2.27).

We believe that this may account for a significant difference in the performance of the mixed and standard methods. Through this projection the weighted averages over the elements of  $c_{\tilde{c}}$ , the inverse of the coefficient matrix, enter the numerical scheme. This is in contrast to the standard method (2.29), which sees only weighted local averages of the coefficient matrix  $a_{\tilde{z}}$  itself Now in one dimension it is known that when a rough coefficient is to be replaced locally by a constant, the best value is the harmonic average, i.e., the inverse of the average of the inverse ([1], see also the literature on homogenization referenced in [3]) In higher dimensions the harmonic average is not the best strategy, but is nonetheless often still superior to the ordinary average This may be one of the main reasons for the good performance of mixed methods for rough coefficient problems

It would be very interesting to determine through numerical experiments the effects of each of the differences 1-3 on the numerical solution

#### REFERENCES

- I BABUSKA and J E OSBORN, Generalized finite element methods their performance and their relation to mixed methods, SIAM J Numer Anal 20 (1983), 510-536
- [2] I BABUSKA, J OSBORN and J PITKARANTA, Analysis of mixed methods using mesh dependent norms, Math Comput 35 (1980), 1039-1062
- [3] A BENSOUSSON, J L LIONS, G PAPANICOLAU, Asymptotic Analysis of Periodic Structures, North-Holland, Amsterdam, 1978
- [4] F BREZZI and P A RAVIART, Mixed finite element methods for 4th order elliptic equations, in Proc of the Royal Irish Academy Conference on Numerical Analysis, Academic Press, London, 1977
- [5] P G CIARLET, The Finite Element Method for Elliptic Equations, North-Holland, Amsterdam, 1978
- [6] J DOUGLAS and J E ROBERTS, Global estimates for mixed methods for second order elliptics, to appear in Math Comput
- [7] R S FALK and J E OSBORN, Error estimates for mixed methods, RAIRO Anal numer 14 (1980), 309-324
- [8] B FRAEJIS DE VEUBEKE, Displacement and equilibrium models in the finite element method, in Stress Analysis, O C Zienkiewicz and G Holister, eds, Wiley, New York, 1965
- K HELLAN, Analysis of elastic plates in flexure by a simplified finite element method, Acta Polytechnica Scandinavica, Ci 46, Trondheim, 1967
- [10] L HERRMANN, Finite element bending analysis for plates, J Eng Mech Div ASCE, a 3, EM5 (1967), 49-83
- [11] P LASCAUX and P LESAINT, Some nonconforming finite elements for the plate bending problem, R A I R O Anal numer 9 (1975), 9-53
- [12] C JOHNSON, On the convergence of a mixed finite element method for plate bending problems, Numer Math 21 (1973), 43-62
- [13] L. S. D. MORLEY, The triangular equilibrium element in the solution of plate bending problems, Aero. Quart. 19 (1968), 149-169
- vol 19, nº 1, 1985

- [14] R RANNACHER, Nonconforming finite element methods for eigenvalue problems in linear plate theory, Numer Math 33 (1979), 23-42
- [15] R RANNACHER, On nonconforming and mixed finite elements for plate bending problems The linear case RAIRO Anal numer 13 (1979), 369-387
- [16] P A RAVIART and J M THOMAS A mixed finite element method for second order elliptic problems in Mathematical Aspects of the Finite Element Method, Lecture Notes in Mathematics 606, Springer-Verlag, Berlin, 1977

#### 32