

## Analysis of Some Mixed Finite Element Methods for Plane Elasticity Equations

By J. Pitkäranta and R. Stenberg

**Abstract.** We analyze some mixed finite element methods, based on rectangular elements, for solving the two-dimensional elasticity equations. We prove error estimates for a method proposed by Taylor and Zienkiewicz and for some new variants of the known equilibrium methods. A numerical example is given demonstrating the performance of the various algorithms considered.

**1. Introduction.** In the numerical solution of problems of continuum mechanics, the stresses are normally of primary interest in the elastic region. It is therefore natural to design the numerical algorithms so that the stresses can be obtained directly without first computing the displacements. Such methods can be derived from the dual variational formulation of the elasticity problem. The corresponding finite element algorithms are usually formulated as mixed methods where both the displacements and the stresses are first approximated, and the displacements are then eliminated from the discrete equations. In many cases the elimination can be rather effectively done using penalty/perturbation techniques or their iterative variants; cf. [3], [11], [12].

The best known finite element methods of the above type are the so-called equilibrium methods, first proposed by Fraeijis de Veubeke [17] (cf. also [14], [16], [18]) and analyzed theoretically by Johnson and Mercier [9] (cf. also [8]). In these methods, one uses specific composite elements which allow the equilibrium condition between the stresses and the volume load to be satisfied exactly in the case where the volume load is zero.

The main drawback of the equilibrium methods proposed so far is the relatively high number of free parameters as compared with displacement methods of the same order of accuracy. For example, if the composite quadrilateral element of [17], [9] is used on a regular rectangular grid, one has eight degrees of freedom per each interior node of the grid (after the local condensation of three extra degrees of freedom per node, cf. [9]) and the convergence rate  $O(h^2)$  for the stresses in  $L_2$  [9]. On the other hand, using the displacement method with reduced biquadratic elements (cf. [6]), one has the same convergence rate with six parameters per node, so the displacement method seems superior.

It is clear from the above example that the mixed or equilibrium methods should be further developed if they are desired to be competitive with displacement

---

Received August 18, 1982.  
1980 *Mathematics Subject Classification.* Primary 65N30.

methods. This is the motivation of the present paper. In particular, we try to find mixed or equilibrium methods which are simpler than those considered in [9] and still preserve the quadratic convergence rate of the stresses in  $L_2$ . We analyze in detail two candidates for such methods. In the first method, called Method I below, the stresses are approximated by continuous piecewise bilinear functions on a rectangular grid, and the displacements are taken to be piecewise constant on the same grid. This method (which probably is the simplest possible mixed method one can think of) was proposed recently by Taylor and Zienkiewicz [15]. The convergence rate of this method, however, does not seem to be quadratic. We are able to prove, under various restrictive assumptions, that the stresses converge with the rate  $O(h^{3/2})$  if the exact solution is sufficiently smooth. That this result is actually optimal is confirmed numerically.

As another alternative, called Method II below, we consider a class of algorithms based on the composite quadrilateral element of [17], [9]. We show that many of the degrees of freedom can be eliminated without affecting the convergence rate. In particular, we derive a method which contains only the average of four free parameters per node and still gives quadratic convergence rate for stresses in  $L_2$ .

We consider only the case of a uniform rectangular mesh on a rectangular domain in this paper. The assumption on mesh uniformity seems essential for Method I, but for Method II the results can very likely be extended to more general quadrilateral meshes.

The plan of the paper is as follows. In Section 2 we state the problem and its finite element discretization in a general form. In Sections 3 and 4 we analyze the two methods, and in Section 5 we present some results of numerical computations with both methods.

**2. Notation and Preliminaries.** Let us recall the basic problem of linear elasticity in two dimensions (plane stress or plane strain): given  $f = (f_1, f_2)$  find a symmetric stress tensor  $\sigma = \{\sigma_{ij}\}$ ,  $i, j = 1, 2$  and a displacement  $(u_1, u_2)$  satisfying

$$(2.1) \quad \begin{aligned} \varepsilon(u) &= \lambda \operatorname{tr}(\sigma)\delta + \mu\sigma && \text{in } \Omega, \\ \operatorname{div} \sigma + f &= 0 && \text{in } \Omega, \end{aligned}$$

subject to the boundary conditions

$$\begin{aligned} u &= 0 && \text{on } \Gamma_1, \\ \sigma \cdot n &= 0 && \text{on } \Gamma_2. \end{aligned}$$

Here

$$\varepsilon(u) = \{\varepsilon_{ij}(u)\}, \quad \varepsilon_{ij} = \frac{1}{2} \left( \frac{\partial u_i}{\partial x_j} + \frac{\partial u_j}{\partial x_i} \right), \quad i, j = 1, 2,$$

is the deformation tensor,

$$\begin{aligned} \operatorname{tr}(\sigma) &= \sigma_{11} + \sigma_{22}, \\ \operatorname{div} \sigma &= \left\{ \frac{\partial \sigma_{11}}{\partial x_1} + \frac{\partial \sigma_{12}}{\partial x_2}, \frac{\partial \sigma_{21}}{\partial x_1} + \frac{\partial \sigma_{22}}{\partial x_2} \right\}, \\ \sigma \cdot n &= (\sigma_{n,1}, \sigma_{n,2}) = (n_1\sigma_{11} + n_2\sigma_{12}, n_1\sigma_{21} + n_2\sigma_{22}), \\ \delta &= \{\delta_{ij}\}, \quad \delta_{ij} = \begin{cases} 1 & \text{if } i = j, \\ 0 & \text{if } i \neq j, \end{cases} \quad i, j = 1, 2, \end{aligned}$$

$\lambda$  and  $\mu$  are constants satisfying  $\mu > 0$  and  $2\lambda + \mu > 0$ ,  $\Omega$  is a bounded region in the plane with boundary  $\Gamma = \Gamma_1 \cup \Gamma_2$ , and  $n = (n_1, n_2)$  denotes the unit outward normal vector on  $\Gamma$ . We have assumed homogeneous boundary conditions for simplicity.

We assume below that  $\Omega$  is a rectangle:  $\Omega = \{x = (x_1, x_2) \in \mathbf{R}^2, 0 < x_1 < a_1, 0 < x_2 < a_2\}$ . Each side of  $\Omega$  is assumed to be fully contained in either  $\Gamma_1$  or  $\Gamma_2$ , i.e., the boundary conditions can change only at the vertices of  $\Omega$ . We also assume that  $\Gamma_1$  is nonempty.

Introducing the spaces

$$\begin{aligned} V &= V(\Omega) = [L_2(\Omega)]^2, \\ Y &= Y(\Omega) = \{\tau = (\tau_{ij}) : \tau_{ij} \in L_2(\Omega), \tau_{ij} = \tau_{ji}, i, j = 1, 2\}, \\ H &= H(\Omega) = \{\tau \in Y(\Omega) : \operatorname{div} \tau \in V(\Omega), \tau \cdot n = 0 \text{ on } \Gamma_2\}, \end{aligned}$$

the elasticity problem can be given the following variational formulation: Find  $(\sigma, u) \in H \times V$  such that

$$(2.2) \quad \begin{cases} a(\sigma, \tau) + (u, \operatorname{div} \tau) = 0, & \tau \in H, \\ (\operatorname{div} \sigma, v) + (f, v) = 0, & v \in V, \end{cases}$$

where  $(\cdot, \cdot)$  denotes the scalar product in  $V$  and

$$a(\sigma, \tau) = \int_{\Omega} (\lambda \operatorname{tr}(\sigma) \operatorname{tr}(\tau) + \mu \sigma \cdot \tau) dx,$$

where

$$\sigma \cdot \tau = \sum_{i,j=1}^2 \sigma_{ij} \tau_{ij}.$$

Since  $\Gamma_1$  is nonempty, (2.2) has a unique solution; cf. [7].

Below we consider finite element methods of the form: Find  $(\sigma_h, u_h) \in H_h \times V_h$  such that

$$(2.3a) \quad a(\sigma_h, \tau) + (u_h, \operatorname{div} \tau) = 0, \quad \tau \in H_h,$$

$$(2.3b) \quad (\operatorname{div} \sigma_h, v) + (f, v) = 0, \quad v \in V_h,$$

where  $H_h$  and  $V_h$  are finite element subspaces of  $H$  and  $V$ , respectively. In practice, Eqs. (2.3) are often solved by introducing a small perturbation parameter  $\varepsilon > 0$  and replacing (2.3b) by

$$(2.3b') \quad -\varepsilon(u_h, v) + (\operatorname{div} \sigma_h, v) + (f, v) = 0, \quad v \in V_h.$$

If  $\pi_h$  denotes the orthogonal projection of  $[L_2(\Omega)]^2$  onto  $V_h$ , (2.3b') may be written as

$$u_h = \frac{1}{\varepsilon} \pi_h(\operatorname{div} \sigma_h + f).$$

Upon substituting this into (2.3a) one obtains

$$(2.4) \quad a(\sigma_h, \tau) + \frac{1}{\varepsilon} (\pi_h(\operatorname{div} \sigma_h + f), \pi_h \operatorname{div} \tau) = 0, \quad \tau \in H_h.$$

This is a modified penalty method for approximately solving (2.1). The operator  $\pi_h$  corresponds frequently to the use of "selective reduced integration", i.e., some low-order quadrature rule for computing the integral in the penalty term; cf. [3], [12].

Below we associate the spaces  $V_h$  and  $H_h$  to a uniform rectangular partitioning of  $\Omega$  defined by

$$(2.5) \quad \mathcal{C}^h = \{K_{ij}, i = 1, \dots, m_1, j = 1, \dots, m_2\},$$

where

$$K_{ij} = \{x \in \mathbf{R}^2: (i-1)h_1 < x_1 < ih_1, (j-1)h_2 < x_2 < jh_2\}.$$

Here  $m_i h_i = a_i$ ,  $i = 1, 2$ , and  $h_1$  and  $h_2$  are associated to the mesh parameter  $h \equiv h_1$  in such a way that  $h_1/h_2$  is bounded from both above and below by a positive constant independent of  $h$ .

The set of nodal points in the grid induced by the partitioning  $\mathcal{C}^h$  is denoted by  $\mathfrak{N}$ :

$$\mathfrak{N} = \{P_{ij} = (ih_1, jh_2), i = 0, \dots, m_1, j = 0, \dots, m_2\}.$$

In all of the methods considered in this paper, the functions in  $V_h$  are fully discontinuous along the mesh lines, so that each  $v \in V_h$  can be expressed in the form

$$v(x) = \sum_{i=1}^{m_1} \sum_{j=1}^{m_2} \sum_{k=1}^{\nu} \alpha_{ijk} \phi_{ijk}(x), \quad \alpha_{ijk} \in \mathbf{R},$$

where  $\phi_{ijk}$  vanishes outside  $K_{ij} \in \mathcal{C}^h$ . Similarly, the space  $H_h$  is chosen to consist of functions of the form

$$\tau(x) = \sum_{P \in \mathfrak{N}} \sum_{k=1}^{\nu_P} \alpha_P^k \xi_P^k(x), \quad \alpha_P^k \in \mathbf{R},$$

where the basis function  $\xi_P^k \in H_h$  vanishes in any  $K \in \mathcal{C}^h$  that does not have a vertex at  $P$ . For this type of choice of the subspaces, the projection  $\pi_h$  in (2.4) can be computed locally in each rectangle  $K \in \mathcal{C}^h$ . It is also obvious that in (2.4) there is coupling only between the nodes that are vertices of the same rectangle  $K \in \mathcal{C}^h$ .

We shall denote by  $|\cdot|_{s,p,T}$  and  $\|\cdot\|_{s,p,T}$ , respectively, the seminorm and the norm of the Sobolev space  $[W^{s,p}(T)]^\alpha$ ,  $s$  and  $\alpha$  integers,  $1 \leq p \leq \infty$ . For nonintegral  $s$ ,  $s \geq 0$ ,  $\|\cdot\|_{s,p,T}$  is defined as usual by interpolation. If  $p = 2$ , we set  $\|\cdot\|_{s,2,T} = \|\cdot\|_{s,T}$ . The subscript  $T$  is dropped if  $T = \Omega$ .

In the analysis below we need frequently the partial integration formula

$$(2.6) \quad \int_T \operatorname{div} \tau \cdot v \, dx = \int_{\partial T} (\tau \cdot n) \cdot v \, ds - \int_T \tau \cdot \varepsilon(v) \, dx,$$

where  $n$  is the unit outward normal vector to  $\partial T$ . (2.6) is valid if  $\tau$  is a symmetric tensor such that  $\tau_{ij} \in L_2(T)$  and  $\operatorname{div} \tau \in [L_2(T)]^2$ , and  $v \in [H^1(T)]^2$ . We also recall that since  $\mu > 0$  and  $2\lambda + \mu > 0$ , we have (cf. [7])

$$(2.7) \quad a(\tau, \tau) \geq C \|\tau\|_0^2, \quad \tau \in Y.$$

Here and below,  $C$  denotes a positive constant which may depend on  $\Omega$  and on the parameters  $\lambda$  and  $\mu$  in (2.1) but is independent of other parameters unless indicated explicitly. We shall finally denote by  $P_k(T)$ ,  $T \subset \mathbf{R}^2$ , the set of polynomials in two variables of degree  $\leq k$  defined on  $T$ .

**3. Analysis of Method I.** We consider here the method proposed in [15]. Let  $S_h$  be the set of continuous piecewise bilinear functions on  $\mathcal{C}_h$  and let  $Q_h$  be the set of functions that are piecewise constant on  $\mathcal{C}^h$ . Then the method of [15] is equivalent to choosing the subspace  $H_h$  and  $V_h$  in (2.2) as

$$H_h = H(\Omega) \cap [S_h]^4, \quad V_h = [Q_h]^2.$$

Obviously any  $\tau \in H_h$  is determined uniquely by the values of  $\tau_{i,j}(P)$ ,  $P \in \mathfrak{N}$ ,  $P \notin \bar{\Gamma}_2$ . If  $P \in \bar{\Gamma}_2$ , then  $(\tau \cdot n)(P) = 0$ ,  $\tau \in H_h$ .

Let  $(\sigma, u)$  be the solution of (2.2), let  $\tilde{\sigma} \in H_h$  be the interpolant of  $\sigma$  defined by

$$\tilde{\sigma}(P) = \sigma(P), \quad P \in \mathfrak{N},$$

(assuming that  $\sigma$  is continuous) and let  $\tilde{u}$  be the  $L_2$ -projection of  $u$  into  $V_h$ . We need the following three lemmas.

**LEMMA 3.1.** *Assume that  $\Gamma_1$  contains two adjacent sides of  $\Omega$ . Then there is a constant  $C$  such that for all  $v \in V_h$*

$$\sup_{\tau \in H_h} \frac{(\operatorname{div} \tau, v)}{\|\tau\|_0} \geq C \|v\|_0. \quad \square$$

**LEMMA 3.2.** *For all  $v \in V_h$ ,*

$$|(\operatorname{div}(\sigma - \tilde{\sigma}), v)| \leq Ch^k |\sigma|_{k+1} \|v\|_0, \quad k = 2, 3. \quad \square$$

**LEMMA 3.3.** *For all  $\tau \in H_h$ ,*

$$|(u - \tilde{u}, \operatorname{div} \tau)| \leq C(h^2 |u|_3 + h^{3/2} |u|_{2,\infty}) \|\tau\|_0. \quad \square$$

*Remark.* Lemma 3.1 states a weak Babuška-Brezzi-type stability estimate (cf. [1],[5]) for the method (2.3). Without the additional assumption on the boundary condition, we have only been able to prove that

$$\sup_{\tau \in H_h} \frac{(\operatorname{div} \tau, v)}{\|\tau\|_0} \geq Ch \|v\|_0, \quad v \in V_h. \quad \square$$

Before proving the lemmas, let us show that they imply the following error estimate.

**THEOREM 3.1.** *Assume that  $\Gamma_1$  contains two adjacent sides of  $\Omega$ . Let  $(\sigma, u)$  be the solution of (2.2) and  $(\sigma_h, u_h)$  the solution of (2.3), where the subspaces  $H_h$  and  $V_h$  are defined as above. Further let  $\tilde{u} \in V_h$  be the  $L_2$ -projection of  $u$ . Then we have the error estimate*

$$\|\sigma - \sigma_h\|_0 + \|u_h - \tilde{u}\|_0 \leq Ch^{3/2} \|u\|_{7/2}. \quad \square$$

**PROOF.** By (2.2) and (2.3) we have the identity

$$(3.1) \quad \mathfrak{B}(\sigma_h - \tilde{\sigma}, u_h - \tilde{u}; \tau, v) = \mathfrak{B}(\sigma - \tilde{\sigma}, u - \tilde{u}; \tau, v), \quad (\tau, v) \in H_h \times V_h,$$

where

$$(3.2) \quad \mathfrak{B}(\sigma, u; \tau, v) = a(\sigma, \tau) + (u, \operatorname{div} \tau) - (v, \operatorname{div} \sigma).$$

By a standard argument (cf. [1],[5]) (2.7) and Lemma 3.1 imply the existence of  $(\tau, v) \in H_h \times V_h$  satisfying

$$(3.3a) \quad \|\tau\|_0 + \|v\|_0 \leq C,$$

$$(3.3b) \quad \mathfrak{B}(\sigma_h - \tilde{\sigma}, u_h - \tilde{u}; \tau, v) \geq \|\sigma_h - \tilde{\sigma}\|_0 + \|u_h - \tilde{u}\|_0.$$

Combining (3.3b) with (3.1), we see that

$$(3.4) \quad \|\sigma_h - \tilde{\sigma}\|_0 + \|u_h - \tilde{u}\|_0 \leq |a(\sigma - \tilde{\sigma}, \tau)| + |(u - \tilde{u}, \operatorname{div} \tau)| + |(v, \operatorname{div}(\sigma - \tilde{\sigma}))|.$$

Using (3.3a) and standard approximation theory, the first term on the right side of (3.4) is estimated as

$$|a(\sigma - \tilde{\sigma}, \tau)| \leq C\|\sigma - \tilde{\sigma}\|_0\|\tau\|_0 \leq C_1h^2\|\sigma\|_2 \leq C_2h^2\|u\|_3.$$

In the second term we apply Lemma 3.3, (3.3a) and Sobolev imbedding to obtain

$$|(u - \tilde{u}, \operatorname{div} \tau)| \leq Ch^{3/2}(\|u\|_3 + \|u\|_{2,\infty})\|\tau\|_0 \leq C_1h^{3/2}\|u\|_{7/2}.$$

Finally, interpolating in Lemma 3.2 and using (3.3a), we have

$$|(\operatorname{div} \sigma - \tilde{\sigma}, v)| \leq Ch^{3/2}\|\sigma\|_{5/2} \leq C_1h^{3/2}\|u\|_{7/2}.$$

Upon combining these inequalities with (3.4), using the triangle inequality and recalling the estimate for  $\|\sigma - \tilde{\sigma}\|_0$ , the asserted estimate follows.  $\square$

*Remark.* Using the triangle inequality and the standard estimate  $\|u - \tilde{u}\|_0 \leq Ch\|u\|_1$ , it follows from Theorem 3.1 that  $\|u - u_h\|_0 \leq Ch\|u\|_{7/2}$ .  $\square$

*Remark.* Without the extra assumption on the boundary condition we can only prove the estimate

$$\|\sigma - \sigma_h\|_0 + \|u - u_h\|_0 \leq Ch^{1/2}\|u\|_{7/2}$$

(see the remark following Lemma 3.3).  $\square$

*Proof of Lemma 3.1.* Assume first that  $m_1 \leq 2$  in (2.5). Then, since  $m_2 = a_2(h_2/h_1)m_1$ ,  $m_2$  is bounded by a constant independent of  $h$ . Thus, also  $\dim(H_h)$  and  $\dim(V_h)$  are bounded by constants independent of  $h$ . By the equivalence of norms in a finite-dimensional space, the assertion then follows if one can show that

$$(3.5) \quad (\operatorname{div} \tau, v) = 0 \quad \forall v \in H_h \Rightarrow v = 0.$$

To see that this is valid, let  $K \in \mathcal{C}^h$  be such that  $K$  has two sides on  $\Gamma_1$ , and choose  $\tau \in H_h$  so that  $\tau$  is nonzero only in  $K$ . It is easy to see that one can have  $(\operatorname{div} \tau, v) = 0$  for all such  $\tau$  only if  $v = 0$  on  $K$ . Repeating this argument, (3.5) follows easily. We omit the details.

The case  $m_2 \leq 2$  can be handled as above, so let us assume that  $m_1, m_2 \geq 3$ . Let  $r_i \geq 0$  be the largest integer such that  $3 \leq m_i - 3r_i \leq 5$ ,  $i = 1, 2$ , and let  $\tilde{\mathcal{C}}^h$  be a coarser subdivision of  $\Omega$  into rectangles  $\tilde{K}_{\nu\mu}$ ,  $\nu = 1, \dots, r_1 + 1$ ,  $\mu = 1, \dots, r_2 + 1$ , where

$$\tilde{K}_{\nu\mu} = \{x \in \mathbf{R}^2: d_{1,\nu-1} < x_1 < d_{1\nu}, d_{2,\mu-1} < x_2 < d_{2\mu}\},$$

with

$$d_{i\nu} = \begin{cases} 3\nu h_i, & 0 \leq \nu \leq r_i, \\ a_i, & \nu = r_i + 1, i = 1, 2. \end{cases}$$

Thus, each  $\tilde{K} \in \tilde{\mathcal{C}}^h$  consists of  $l_1 \times l_2$  rectangles of  $\mathcal{C}^h$  with  $3 \leq l_1, l_2 \leq 5$ .

For each  $\tilde{K} \in \tilde{\mathcal{C}}^h$ , we define the following finite-dimensional spaces

$$\begin{aligned} U &= U(\tilde{K}) = \left\{ \tau_{|\tilde{K}} : \tau \in H_h, \tau = 0 \text{ on } \Omega \setminus \tilde{K} \right\}, \\ N &= \left\{ v \in V_{h|\tilde{K}} : \int_{\tilde{K}} \operatorname{div} \tau \cdot v \, dx = 0, \forall \tau \in U \right\}, \\ W &= \left\{ w \in V_{h|\tilde{K}} : \int_{\tilde{K}} w \cdot v \, dx = 0, \forall v \in N \right\}. \end{aligned}$$

Since  $\tilde{K}$  consists of at most  $5 \times 5$  rectangles of  $\mathcal{C}^h$ , we have  $\dim(U) \leq 27$  and  $\dim(W) \leq 50$ . Therefore, using a scaling argument and the equivalence of norms in a finite-dimensional space, we conclude the existence of a positive constant  $C$  independent of  $\tilde{K}$  such that

$$(3.6) \quad \sup_{\tau \in U(\tilde{K})} \frac{\int_{\tilde{K}} \operatorname{div} \tau \cdot v \, dx}{\|\tau\|_{0,\tilde{K}}} \geq Ch^{-1} \|v\|_{0,\tilde{K}} \quad \forall v \in W(\tilde{K}).$$

Now let  $v \in V_h$  be given, and write

$$(3.7) \quad v = v_0 + v_1, \quad (v_0)|_{\tilde{K}} \in W(\tilde{K}), \quad (v_1)|_{\tilde{K}} \in N(\tilde{K}), \quad \tilde{K} \in \tilde{\mathcal{C}}^h.$$

By (3.6), there exists for any  $\tilde{K} \in \tilde{\mathcal{C}}^h$  a function  $\tau_{\tilde{K}} \in U(\tilde{K})$  such that, for some constant  $C$ ,

$$\|\tau_{\tilde{K}}\|_{0,\tilde{K}} \leq Ch^{-1} \|v_0\|_{0,\tilde{K}}, \quad \int_{\tilde{K}} \operatorname{div} \tau_{\tilde{K}} \cdot v_0 \, dx \geq h^{-2} \|v_0\|_{0,\tilde{K}}^2.$$

Let  $\tau_0$  be defined on  $\Omega$  so that

$$\tau_0(x) = \tau_{\tilde{K}}(x), \quad x \in \tilde{K} \in \tilde{\mathcal{C}}^h.$$

Then  $\tau_0 \in H_h$  (since  $\tau_{\tilde{K}} = 0$  on  $\partial\tilde{K}$ ). Moreover,  $\|\tau_0\|_0 \leq Ch^{-1} \|v_0\|_0$ , and

$$(\operatorname{div} \tau_0, v) = \sum_{\tilde{K} \in \tilde{\mathcal{C}}^h} \int_{\tilde{K}} \operatorname{div} \tau_0 \cdot v \, dx = \sum_{\tilde{K} \in \tilde{\mathcal{C}}^h} \int_{\tilde{K}} \operatorname{div} \tau_0 \cdot v_0 \, dx \geq h^{-2} \|v_0\|_0^2.$$

Let us assume for a while that there also exists  $\tau_1 \in H_h$  satisfying

$$(3.8a) \quad \|\tau_1\|_0 \leq C \|v_1\|_0,$$

$$(3.8b) \quad (\operatorname{div} \tau_1, v_1) \geq \|v_1\|_0^2.$$

Then, setting  $\tau = \tau_0 + \gamma\tau_1$ , where  $\gamma \in (0, 1]$  will be chosen below, we have  $\tau \in H_h$  and

$$(3.9) \quad \|\tau\|_0 \leq C(h^{-2} \|v_0\|_0^2 + \|v_1\|_0^2)^{1/2}.$$

Moreover, using the inverse inequality

$$\|\operatorname{div} \tau_1\|_0 \leq Ch^{-1} \|\tau_1\|_0$$

together with the above inequalities, we have

$$(\operatorname{div} \tau, v) \geq h^{-2} \|v_0\|_0^2 + \gamma \|v_1\|_0^2 + \gamma (\operatorname{div} \tau_1, v_0) \geq h^{-2} (1 - C\gamma) \|v_0\|_0^2 + \frac{\gamma}{2} \|v_1\|_0^2.$$

Thus, choosing  $\gamma = \min\{1, 1/2C\}$ , we have

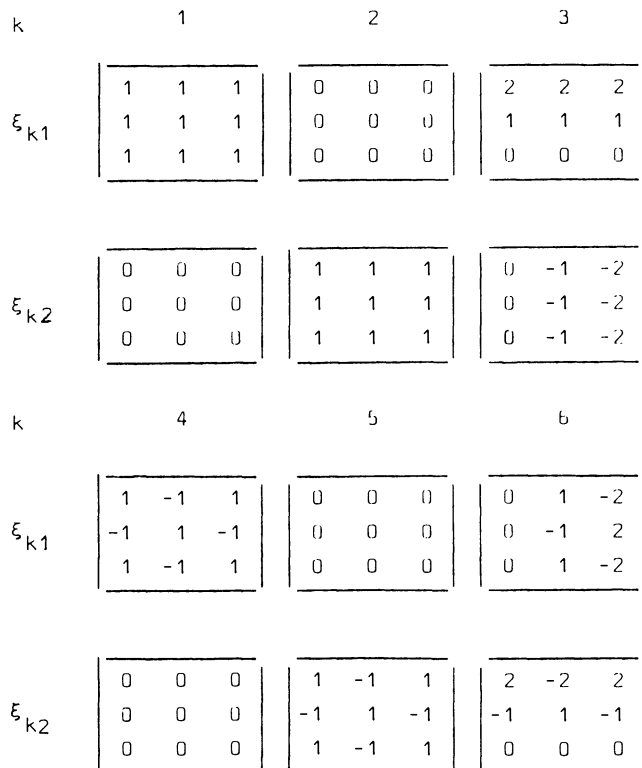
$$(\operatorname{div} \tau, v) \geq C(h^{-2} \|v_0\|_0^2 + \|v_1\|_0^2).$$

Together with (3.9) this proves the assertion of the lemma, so the remaining task is to construct  $\tau_1 \in H_h$  which satisfies (3.8).

We begin the construction of  $\tau_1$  by characterizing the spaces  $N(\tilde{K})$ . First, if  $\tilde{K}$  is a macroelement consisting of  $2 \times 3 = 6$  rectangles of  $\mathcal{C}^h$ , one has  $\dim(U(\tilde{K})) = 6$  and  $\dim(V_{H_{\tilde{K}}}) = 12$ . Hence, one expects  $N(\tilde{K})$  to be six-dimensional in this case, and this can indeed be verified by a straightforward computation. Let  $\nu$  and  $\mu$  be such that  $\tilde{K}$  contains the rectangles  $K_{\nu+i, \mu+j} \in \mathcal{C}^h$  for  $i = 1, 2$  and  $j = 1, 2, 3$ . Then we may choose for the basis of  $N(\tilde{K})$  the set  $\{\xi_1, \dots, \xi_6\}$ , where the functions  $\xi_k$  are defined on the subrectangles  $K_{\nu+i, \mu+j} \subset \tilde{K}$  as

$$\begin{aligned}
 (3.10) \quad & \xi_1(x) = (1, 0), \quad \xi_2(x) = (0, 1), \quad \xi_3(x) = (i - 1, 1 - j), \\
 & \xi_4(x) = ((-1)^{i+j}, 0), \quad \xi_5(x) = (0, (-1)^{i+j}), \\
 & \xi_6(x) = ((1 - j)(-1)^{i+j}, (i - 1)(-1)^{i+j}), \\
 & x \in K_{\nu+i, \mu+j}, \quad i, j \geq 1, \quad K_{\nu+i, \mu+j} \subset \tilde{K}.
 \end{aligned}$$

We omit the details of showing that if  $\tilde{K}$  is any macroelement consisting of  $l_1 \times l_2$  rectangles of  $\mathcal{C}^h$ , the space  $N(\tilde{K})$  is always six-dimensional and is spanned by the functions  $\xi_k$  defined by (3.10), provided that  $\min\{l_1, l_2\} \geq 2$  and  $\max\{l_1, l_2\} \geq 3$ . (This can be shown by splitting  $\tilde{K}$  into smaller subrectangles and using a construction similar to that given below.) In particular, this is the case for any  $\tilde{K} \in \tilde{\mathcal{C}}^h$ . The values of the basis functions  $\xi_k = (\xi_{k1}, \xi_{k2})$  on the subrectangles of  $\tilde{K}$  are shown in Figure 1 in the case where  $\tilde{K}$  consists of  $3 \times 3$  subrectangles.



**FIGURE 1**  
The basis of  $N(\tilde{K})$



Note that the first three functions represent the physical degrees of freedom in  $N(\tilde{K})$ , i.e., the rigid translations along the coordinate axes ( $k = 1, 2$ ) and rotation ( $k = 3$ ). The remaining degrees of freedom represent purely numerical “zero energy modes” (cf. [10] for similar modes in other mixed methods).

It follows from the above considerations that the function  $v_1$  in (3.7) can be written as

$$v_1 = \sum_{\nu=1}^{r_1+1} \sum_{\mu=1}^{r_2+1} \sum_{k=1}^6 \alpha_{\nu\mu k} \xi_{\nu\mu k}, \quad \alpha_{\nu\mu k} \in \mathbf{R},$$

where  $\xi_{\nu\mu k} = \xi_k(\tilde{K}_{\nu\mu})$ ,  $\tilde{K}_{\nu\mu} \in \tilde{\mathcal{C}}^h$ . Consider now a given  $\tilde{K}_{\nu\mu} \in \tilde{\mathcal{C}}^h$ , and let  $A, B, C, D \in \mathfrak{N}$  be nodes of the finite element grid located on the sides of  $\tilde{K}_{\nu\mu}$  as in Figure 2.

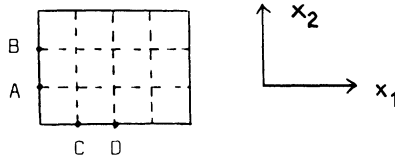


FIGURE 2

Let  $\tau \in H_h$  be such that  $\tau_{12} = 0$  and  $\tau_{11}$  and  $\tau_{22}$  vanish at all nodes except at  $A, B, C$  and  $D$ . Then if none of these four nodes is on the boundary  $\Gamma$ , we find that

$$(3.11) \quad (\operatorname{div} \tau, v_1) = h_2(\tau_{11}(A) + \tau_{11}(B))(\alpha_{\nu-1,\mu 1} - \alpha_{\nu\mu 1}) \\ + \frac{1}{2}h_2(\tau_{11}(A) + 3\tau_{11}(B))(\alpha_{\nu-1,\mu 3} - \alpha_{\nu\mu 3}) \\ + h_2(\tau_{22}(A) - \tau_{22}(B))(\alpha_{\nu-1,\mu 5} + \alpha_{\nu\mu 5}) \\ + \frac{1}{2}h_2(\tau_{22}(A) - 3\tau_{22}(B))(\alpha_{\nu-1,\mu 6} + \alpha_{\nu\mu 6}) \\ + h_1(\tau_{22}(C) + \tau_{22}(D))(\alpha_{\nu,\mu-1,2} - \alpha_{\nu\mu 2}) \\ + \frac{1}{2}h_1(\tau_{22}(C) + 3\tau_{22}(D))(\alpha_{\nu\mu 3} - \alpha_{\nu,\mu-1,3}) \\ + h_1(-\tau_{11}(C) + \tau_{11}(D))(\alpha_{\nu,\mu-1,4} + \alpha_{\nu\mu 4}) \\ + \frac{1}{2}h_1(\tau_{11}(C) - 3\tau_{11}(D))(\alpha_{\nu,\mu-1,6} + \alpha_{\nu\mu 6}).$$

By our assumption on the boundary conditions, we may assume that  $\Gamma_1$  contains the sides of  $\Omega$  at  $x_1 = 0$  and  $x_2 = 0$ . Then the values of  $\tau_{ij}$  at the nodes  $A, B, C, D$  are free parameters in  $H_h$ , and hence (3.11) holds even if some of the nodes are on the boundary if we set  $\alpha_{0\mu k} = \alpha_{\nu 0k} = 0$ .

Let us now choose  $\tilde{\tau} \in H_h$  in such a way that if  $A, B, C, D$  are nodes located on the boundary of some  $\tilde{K}_{\nu\mu} \in \tilde{\mathcal{C}}^h$  as in Figure 2, then

$$\tilde{\tau}_{11}(A) + \tilde{\tau}_{11}(B) = h^{-1}(\alpha_{\nu-1,\mu 1} - \alpha_{\nu\mu 1}), \\ \tilde{\tau}_{11}(A) + 3\tilde{\tau}_{11}(B) = h^{-1}(\alpha_{\nu-1,\mu 3} - \alpha_{\nu\mu 3}), \\ \tilde{\tau}_{22}(A) - \tilde{\tau}_{22}(B) = h^{-1}(\alpha_{\nu-1,\mu 5} + \alpha_{\nu\mu 5}), \\ \tilde{\tau}_{22}(A) - 3\tilde{\tau}_{22}(B) = h^{-1}(\alpha_{\nu-1,\mu 6} + \alpha_{\nu\mu 6}), \\ \tilde{\tau}_{22}(C) + \tilde{\tau}_{22}(D) = h^{-1}(\alpha_{\nu,\mu-1,2} - \alpha_{\nu\mu 2}),$$

$$\begin{aligned} \tilde{\tau}_{22}(C) + 3\tilde{\tau}_{22}(D) &= h^{-1}(\alpha_{\nu\mu 3} - \alpha_{\nu,\mu-1,3}), \\ -\tilde{\tau}_{11}(C) + \tilde{\tau}_{11}(D) &= h^{-1}(\alpha_{\nu,\mu-1,4} + \alpha_{\nu\mu 4}), \\ \tilde{\tau}_{11}(C) - 3\tilde{\tau}_{11}(D) &= h^{-1}(\alpha_{\nu,\mu-1,6} + \alpha_{\nu\mu 6}). \end{aligned}$$

At the remaining nodes we set  $\tilde{\tau} = 0$ . Then  $\tilde{\tau}$  is uniquely determined and  $\|\tilde{\tau}\|_0 \leq C\|v_1\|_h$ , where

$$\begin{aligned} \|v_1\|_h^2 &= \sum_{\nu=1}^{r_1+1} \sum_{\mu=1}^{r_2+1} \left[ \sum_{k=2,3} (\alpha_{\nu\mu k} - \alpha_{\nu,\mu-1,k})^2 \right. \\ &\quad + \sum_{k=4,6} (\alpha_{\nu\mu k} + \alpha_{\nu,\mu-1,k})^2 + \sum_{k=1,3} (\alpha_{\nu\mu k} - \alpha_{\nu-1,\mu k})^2 \\ &\quad \left. + \sum_{k=5,6} (\alpha_{\nu\mu k} + \alpha_{\nu-1,\mu k})^2 \right], \\ &\quad (\alpha_{0\mu k} = \alpha_{\nu 0k} = 0). \end{aligned}$$

Moreover, it follows from (3.11) that  $(\operatorname{div} \tilde{\tau}, v_1) \geq C\|v_1\|_h^2$ . Using the inequality

$$\alpha_1^2 + \sum_{i=2}^n (\alpha_i \pm \alpha_{i-1})^2 \geq n^{-2} \sum_{i=1}^n \alpha_i^2$$

and noting that  $r_i \leq Ch^{-1}$ ,  $i = 1, 2$ , we see that  $\|v_1\|_h \geq C\|v\|_0$ . Combining this with the above inequalities we see that  $\tau_1 = (\|v_1\|_0/C\|v_1\|_h)\tilde{\tau}$  satisfies (3.8) for  $C$  sufficiently large, and so the proof is complete.  $\square$

*Proof of Lemma 3.2.* The proof is based on the following result which can be verified by direct computation; cf. also [10].

**LEMMA 3.4.** *If  $\sigma \in [P_2(K)]^4$ ,  $K \in \mathcal{C}^h$  and if  $\tilde{\sigma}$  is the bilinear interpolant of  $\sigma$  on  $K$ , then*

$$\int_K \operatorname{div}(\sigma - \tilde{\sigma}) \, dx = 0. \quad \square$$

Using Lemma 3.4, it follows from the Bramble-Hilbert lemma [4] that, with  $\tilde{\sigma}$  as in Lemma 3.2,

$$\left| \int_K \operatorname{div}(\sigma - \tilde{\sigma}) v \, dx \right| \leq Ch^k |\sigma|_{k+1,K} \|v\|_{0,K}, \quad k = 2, 3, v \in V_h, K \in \mathcal{C}^h.$$

By summing over  $K \in \mathcal{C}^h$ , the asserted estimate follows.  $\square$

*Proof of Lemma 3.3.* We use the following easy-to-prove result. Note that this is only valid for a uniform mesh.

**LEMMA 3.5.** *Let  $\tau \in H_h$  be such that  $\tau$  vanishes at all nodes except at a given node  $P \in \mathfrak{N}$ . Denote by  $T$  the support of  $\tau$ , let  $u \in [L_2(T)]^2$  and let  $\tilde{u}$  be the  $L_2$ -projection of  $u$  into the subspace  $V_{hT}$ . Then*

$$\int_T (u - \tilde{u}) \cdot \operatorname{div} \tau \, dx = 0,$$

if either (a)  $P$  is an interior point of  $\Omega$  and  $u \in [P_2(T)]^2$  or (b)  $P$  is a boundary point but not a vertex of  $\Omega$  and  $u \in [P_1(T)]^2$ .  $\square$

Now let  $\tilde{u}$  be as in Lemma 3.3, let  $P \in \mathfrak{N}$  be a node interior to  $\Omega$  and let  $\tau$  satisfy the assumptions of Lemma 3.5. Then Lemma 3.5 and the Bramble-Hilbert lemma [4] imply that

$$(3.12) \quad \left| \int_T (u - \tilde{u}) \operatorname{div} \tau \, dx \right| \leq Ch^2 |u|_{3,T} \|\tau\|_{0,T}.$$

If  $P$  is a boundary node but not a vertex of  $\Omega$ , we have by the same argument

$$(3.13) \quad \left| \int_T (u - \tilde{u}) \cdot \operatorname{div} \tau \, dt \right| \leq Ch^2 |u|_{2,\infty,T} \|\tau\|_{0,T}.$$

Finally if  $P$  is a vertex of  $\Omega$ , we have three possibilities:

(i) The neighborhood of  $P$  is contained in  $\Gamma_2$ . In this case the boundary conditions imply that  $\tau(P) = 0$  for all  $\tau \in H_h$ .

(ii) The neighborhood of  $P$  is contained in  $\Gamma_1$ . In this case we have  $u = 0$  on the two sides meeting at  $P$ , so by Taylor expansion

$$|u(x)| \leq Ch^2 |u|_{2,\infty,T}, \quad x \in T.$$

Using this we have

$$\left| \int_T (u - \tilde{u}) \cdot \operatorname{div} \tau \, dx \right| \leq C \|u\|_{0,\infty,T} \|\tau\|_{1,1,T} \leq C_1 h^2 |u|_{2,\infty,T} \|\tau\|_{0,T},$$

so (3.13) is valid also in this case.

(iii) The boundary condition changes at  $P$ . Assume, for example, that  $P = (0, 0)$  and that the sides at  $x_i = 0$  are contained in  $\Gamma_i$ ,  $i = 1, 2$ . If  $u \in [W^{2,\infty}(\Omega)]^2$ , it follows from (2.1) and from the boundary conditions that  $u(P) = 0$ ,  $\partial u_i(P)/\partial x_2 = 0$ ,  $i = 1, 2$  and  $\partial u_2(P)/\partial x_1 = 0$ . Therefore,  $u$  admits the Taylor expansion

$$u(x) = (0, \alpha x_2) + v(x), \quad x \in T,$$

where  $\alpha = \partial u_2(P)/\partial x_2$  and  $\|v\|_{0,\infty,T} \leq Ch^2 |u|_{2,\infty,T}$ . On the other hand,  $\tau_{12}(P) = \tau_{22}(P) = 0$  by the boundary conditions, so that  $\operatorname{div} \tau = (\partial \tau_{11}/\partial x_1, 0)$ . Combining these observations it is easy to see that (3.13) is again valid.

Now let  $\tau \in H_h$  be arbitrary and let  $\tau_p$ ,  $P \in \mathfrak{N}$ , be such that  $\tau_p(P) = \tau(P)$  and  $\tau_p(P') = 0$  for  $P' \in \mathfrak{N}$ ,  $P' \neq P$ . Denote the support of  $\tau_p$  by  $T_p$ . Then, by (3.12) and (3.13),

$$\begin{aligned} |(u - \tilde{u}, \operatorname{div} \tau)| &= \left| \sum_{P \in \mathfrak{N}} (u - \tilde{u}, \operatorname{div} \tau_p) \right| \\ &\leq Ch^2 \sum_{\substack{P \in \mathfrak{N} \\ P \notin \Gamma}} |u|_{3,T_p} \|\tau_p\|_0 + Ch^2 \sum_{\substack{P \in \mathfrak{N} \\ P \in \Gamma}} |u|_{2,\infty,T_p} \|\tau_p\|_0 \\ &\leq C_1 h^2 \left\{ \sum_{\substack{P \in \mathfrak{N} \\ P \notin \Gamma}} |u|_{3,T_p}^2 \right\}^{1/2} \|\tau\|_0 + C_1 h^2 \left\{ \sum_{\substack{P \in \mathfrak{N} \\ P \in \Gamma}} |u|_{2,\infty,T_p}^2 \right\}^{1/2} \|\tau\|_0 \\ &\leq C_2 (h^2 |u|_3 + h^{3/2} |u|_{2,\infty}) \|\tau\|_0, \end{aligned}$$

which proves the assertion.  $\square$

**4. Method II.** In this section we consider a class of methods based on (2.3), where the subspace  $H_h$  is defined in terms of a composite rectangular element described in [9]. Let  $K \in \mathcal{C}^h$ , and let  $K$  be subdivided into four triangles  $T_i$  as in Figure 3.

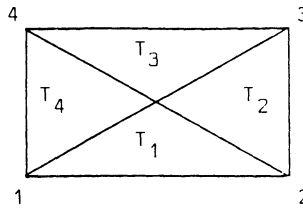


FIGURE 3  
Subdivision of  $K \in \mathcal{C}^h$

We introduce the spaces

$$R_K = \left\{ \tau \in Y(K) : \tau \in [P_1(T)]^4, \tau \cdot n = 0 \text{ on the diagonal } 1-3 \right. \\ \left. \text{and } \tau = 0 \text{ on } K \setminus T \right\},$$

$$R'_K = \left\{ \tau \in Y(K) : \tau \in [P_1(T')]^4, \tau \cdot n' = 0 \text{ on the diagonal } 2-4 \right. \\ \left. \text{and } \tau = 0 \text{ on } K \setminus T' \right\},$$

$$S_K = Y(K) \cap [P_1(K)]^4 \quad \text{and} \quad G_K = S_K \oplus R_K \oplus R'_K,$$

where  $n$  and  $n'$  are normals to the diagonals 1–3 and 2–4, respectively,  $T = T_1 \cup T_2$  and  $T' = T_1 \cup T_4$ .

Setting  $V_K = [P_1(K)]^2$  and  $H_K = G_K$ ,  $K \in \mathcal{C}^h$  and defining the spaces  $V_h$  and  $H_h$  in (2.3) as

$$(4.1a) \quad V_h = \left\{ v \in [L_2(\Omega)]^2 : v|_K \in V_K, K \in \mathcal{C}^h \right\},$$

$$(4.1b) \quad H_h = \left\{ \tau \in H(\Omega) : \tau|_K \in H_K, K \in \mathcal{C}^h \right\},$$

we obtain the method analyzed in [9]. Here we have (see below)

$$\dim(H_h) = 11m_1m_2 + O(h^{-1}),$$

i.e., there are about 11 degrees of freedom per node in  $H_h$ . In what follows we consider methods where  $V_h$  and  $H_h$  are chosen as proper subspaces of those defined by (4.1).

The first modification we consider is as follows. Let  $V_K$  be the space of rigid displacements on  $K$ , i.e.,

$$V_K = \left\{ v(x) = (\alpha_1 + \alpha_3x_2, \alpha_2 - \alpha_3x_1), x \in K, \alpha_i \in \mathbf{R} \right\}.$$

Further, let

$$N_K = \left\{ v \in [P_1(K)]^2 : \int_K v \cdot w \, dx = 0 \, \forall v \in V_K \right\}$$

and

$$(4.2) \quad H_K = \left\{ \tau \in G_K : \int_K \operatorname{div} \tau \cdot v \, dx = 0 \, \forall v \in N_K \right\},$$

where  $G_K$  is as above. With this new definition of  $V_K$  and  $H_K$ , define now  $V_h$  and  $H_h$  by (4.1).

To see how  $H_K$  in (4.2) is constructed, we recall from [9] that any  $\tau \in G_K$  is defined uniquely by the following 19 degrees of freedom:

- (i) the values of  $\tau \cdot n$  at two points on each side of  $K$ ,
- (ii)  $\int_K \tau_{ij} dx, \quad i, j = 1, 2$ .

Let  $Q = (q_1, q_2)$  be the midpoint of  $K$ . Then if  $\tau \in G_K$ , we have  $\tau \in H_K$  if and only if

$$(4.3) \quad \int_K \operatorname{div} \tau \cdot v_{ij} dx = 0, \quad i, j = 1, 2,$$

where  $v_{ij} \in N_K$  are defined by

$$\begin{aligned} v_{11}(x) &= (x_1 - q_1, 0), \\ v_{12}(x) &= v_{21}(x) = \frac{1}{2}(x_2 - q_2, x_1 - q_1), \\ v_{22}(x) &= (0, x_2 - q_2), \quad x \in K. \end{aligned}$$

Applying (2.6) and noting that  $\tau \cdot \epsilon(v_{ij}) = \tau_{ij}$ , we see that (4.3) is equivalent to

$$(4.4) \quad \int_K \tau_{ij} dx = \int_{\partial K} (\tau \cdot n) \cdot v_{ij} ds, \quad i, j = 1, 2.$$

Thus,  $H_K$  is constructed by eliminating from  $G_K$  the inner degrees of freedom (ii) using (4.4).

If  $H_K$  is defined by (4.2), the inclusion  $H_K \supset S_K$  obviously remains valid, since  $\int_K \operatorname{div} \tau \cdot v dx = 0$  if  $v \in N_K$  and  $\tau \in [P_1(K)]^4$ . Therefore the space  $H_K$  has the same approximation properties as  $G_K$  [9]: If  $\sigma \in Y(K) \cap [H^2(K)]^4$  and  $\tilde{\sigma} \in H_K$  is the interpolant of  $\sigma$ , i.e.,  $\tilde{\sigma} \cdot n = \sigma \cdot n$  at two points on each side of  $K$ , then

$$(4.5) \quad \|\sigma - \tilde{\sigma}\|_{0,K} \leq Ch^2 |\sigma|_{2,K}.$$

In the sequel we choose for the degrees of freedom of  $H_K$  the limiting values of  $(\tau \cdot n)(x)$  on each side of  $K$  as  $x$  approaches a vertex of  $K$ . Thus, there are eight degrees of freedom in  $H_h$  associated to each interior node. In order to further reduce the space  $H_h$ , we need some notation. For  $\tau \in H_h$  and  $P$  an interior node, let  $A_P(\tau)$  be a  $2 \times 4$  matrix with coefficients

$$a_{ij} = \lim_{\substack{P' \rightarrow P \\ P' \in l_j}} (\tau \cdot n_j)_i, \quad i = 1, 2, \quad j = 1, \dots, 4,$$

where  $l_j, j = 1, \dots, 4$ , denote the semi-infinite mesh lines starting at  $P$ , numbered as in Figure 4.

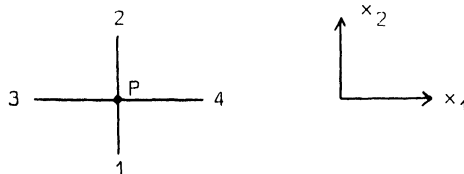


FIGURE 4

Local numbering of mesh lines

If  $P$  is located on  $\Gamma$ , we think the mesh to be extended outside  $\Omega$  and set  $a_{ij} = 0, i = 1, 2$ , if  $l_j \cap \Omega = \emptyset$ .

Let  $H^P$  denote the subspace of  $H_h$  consisting of those functions that vanish at all nodes except at  $P$ . If  $P$  is an interior node, then  $\dim(H^P) = 8$ , and we may choose

for the basis of  $H^P$  the set  $\{\tau_i^P, i = 1, \dots, 8\}$  such that

$$\begin{aligned} A(\tau_1^P) &= \begin{pmatrix} 1 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{pmatrix}, & A(\tau_2^P) &= \begin{pmatrix} 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 1 \end{pmatrix}, \\ A(\tau_3^P) &= \begin{pmatrix} 0 & 0 & 1 & 1 \\ 1 & 1 & 0 & 0 \end{pmatrix}, & A(\tau_4^P) &= \begin{pmatrix} 1 & -1 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{pmatrix}, \\ A(\tau_5^P) &= \begin{pmatrix} 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & -1 \end{pmatrix}, & A(\tau_6^P) &= \begin{pmatrix} 0 & 0 & -1 & -1 \\ 1 & 1 & 0 & 0 \end{pmatrix}, \\ A(\tau_7^P) &= \begin{pmatrix} 0 & 0 & 0 & 0 \\ 1 & -1 & 0 & 0 \end{pmatrix}, & A(\tau_8^P) &= \begin{pmatrix} 0 & 0 & 1 & -1 \\ 0 & 0 & 0 & 0 \end{pmatrix}. \end{aligned}$$

If  $P$  is a boundary node, we define  $\tau_i^P$  on  $\Omega$  as above.

With the above notation, the space  $H_h$  can be defined as

$$(4.6) \quad H_h = \left\{ \tau = \sum_{P \in \mathfrak{N}} \sum_{i=1}^8 \alpha_i^P \tau_i^P, \alpha_i^P \in \mathbf{R}, \tau \cdot n = 0 \text{ on } \Gamma_2 \right\}.$$

We note that if  $\sigma \in H(\Omega)$  is continuous and  $\tilde{\sigma} \in H_h$  is the interpolant of  $\sigma$  defined by

$$(4.7) \quad A_P(\tilde{\sigma}) = A_P(\sigma), \quad P \in \mathfrak{N},$$

then

$$(4.8) \quad \tilde{\sigma} = \sum_{P \in \mathfrak{N}} [\sigma_{11}(P)\tau_1^P + \sigma_{22}(P)\tau_2^P + \sigma_{12}^P(\tau_3^P)].$$

Thus, only the functions  $\tau_i^P, i = 1, 2, 3, P \in \mathfrak{N}$ , are required in the interpolation. In particular, if  $\sigma \in H(\Omega) \cap [H^2(\Omega)]^4$  and if  $\tilde{\sigma}$  is defined by (4.8), then (4.5) holds and so, by summing over  $K \in \mathcal{C}^h$  in (4.5),

$$(4.9) \quad \|\sigma - \tilde{\sigma}\|_0 \leq Ch^2 \|\sigma\|_2.$$

From now on we denote by  $H_{h,1}$  the subspace defined by (4.6) and by  $H_{h,0}$  the space

$$(4.10) \quad H_{h,0} = \left\{ \tau = \sum_{P \in \mathfrak{N}} \sum_{i=1}^3 \alpha_i^P \tau_i^P, \alpha_i^P \in \mathbf{R}, \tau \cdot n = 0 \text{ on } \Gamma_2 \right\}.$$

We shall see below that with  $H_h = H_{h,1}$  and  $V_h$  as above, we have the quasioptimal estimate

$$(4.11) \quad \|\sigma - \sigma_h\|_0 \leq Ch^2 \|\sigma\|_2,$$

where  $\sigma$  and  $\sigma_h$  are solutions to (2.2) and (2.3), respectively. Our aim is now to construct smaller subspaces  $H_h$ , with  $H_{h,0} \subset H_h \subset H_{h,1}$ , so that (4.11) still holds.

Let us first formulate a stability criterion for the method (2.3), assuming that  $V_h$  is defined by (4.1a) and  $H_{h,0} \subset H_h \subset H_{h,1}$ . We follow here the lines of Babuška [1] and Brezzi [5] (cf. also [2], [10]). Let  $\tau \in H(K), K \in \mathcal{C}^h$ , and let  $v \in V^h$ . Then, by (2.6) and since  $\varepsilon(v) = 0$  on  $K$ ,

$$\int_K \operatorname{div} \tau \cdot v \, dx = \int_{\partial K} (\tau \cdot n) \cdot v \, ds.$$

Summing this over all  $K \in \mathcal{C}^h$ , we obtain

$$(4.12) \quad (\operatorname{div} \tau, v) = \int_{\Gamma_h} (\tau \cdot n) \cdot (v^- - v^+) ds + \int_{\Gamma_1} (\tau \cdot n) \cdot v ds,$$

$$\tau \in H(\Omega), v \in V_h,$$

where  $\Gamma_h$  denotes the union of mesh lines in the interior of  $\Omega$ ,  $n$  is a normal to  $\Gamma_h$  or an exterior normal to  $\Gamma$  and

$$v^\pm(x) = \lim_{\epsilon \rightarrow 0^\pm} v(x + \epsilon n), \quad x \in \Gamma_h.$$

By (4.12),

$$(4.13) \quad |(\operatorname{div} \tau, v)| \leq |\tau|_{0,h} |v|_{1,h}, \quad \tau \in H(\Omega), v \in V_h,$$

where

$$|\tau|_{0,h}^2 = h \int_{\Gamma_h \cup \Gamma_1} |\tau \cdot n|^2 ds,$$

$$|v|_{1,h}^2 = h^{-1} \int_{\Gamma_h} |v^+ - v^-|^2 ds + h^{-1} \int_{\Gamma_1} |v|^2 ds.$$

By the following lemma, the seminorms  $|\cdot|_{0,h}$  and  $|\cdot|_{1,h}$  are norms on  $H_h$  and  $V_h$ , respectively.

LEMMA 4.1. *There is a constant  $C$  such that for all  $\tau \in H_h$  and  $v \in V_h$*

$$C^{-1} \|\tau\|_0 \leq |\tau|_{0,h} \leq C \|\tau\|_0,$$

and

$$|v|_{1,h} \geq \begin{cases} C^{-1} \|v\|_0, \\ C^{-1} (1 + |\log h|)^{-1/2} \|v\|_{0,\infty}. \end{cases} \quad \square$$

*Proof.* The equivalence of  $|\cdot|_{0,h}$  and  $\|\cdot\|_0$  on  $H_h$  follows easily from local scaling arguments (cf. also [2]), so let us only prove the lower bounds for  $|\cdot|_{1,h}$ . Let  $\tilde{\mathcal{C}}^h$  be another rectangular subdivision of  $\Omega$  such that the interior nodes in  $\tilde{\mathcal{C}}^h$  are at the midpoints of each  $K \in \mathcal{C}^h$ . For  $v \in V_h$  given, let  $\tilde{v}$  be a smoothing of  $v$  defined in terms of the bicubic Bogner-Fox-Schmidt element [6] as follows: For any  $\tilde{K} \in \tilde{\mathcal{C}}^h$ ,  $\tilde{v}|_{\tilde{K}}$  is a bicubic polynomial, and if  $\tilde{P}$  is a vertex of  $\tilde{K} \in \tilde{\mathcal{C}}^h$ , then

$$\tilde{v}(P) = v(P), \quad \frac{\partial}{\partial x_i} \tilde{v}(P) = \frac{\partial}{\partial x_i} v(P), \quad i = 1, 2,$$

$$\frac{\partial^2}{\partial x_1 \partial x_2} \tilde{v}(P) = \frac{\partial^2}{\partial x_1 \partial x_2} v(P).$$

As is well known (cf. [6]),  $\tilde{v}$  is uniquely determined and  $\tilde{v} \in H^1(\Omega)$ . Now if  $\tilde{K} \in \tilde{\mathcal{C}}^h$ , it is easy to see that we can have

$$\int_{\Gamma_h \cap \tilde{K}} |v^+ - v^-|^2 ds = 0$$

if and only if  $\|\epsilon(\tilde{v})\|_{0,\tilde{K}} = 0$ , i.e., if and only if  $v_{|\tilde{K}} = \tilde{v}_{|\tilde{K}}$  is a rigid displacement on  $\tilde{K}$ . Therefore, by a scaling argument and by the equivalence of norms in a finite-dimensional space, there is a positive constant  $C$  independent of  $\tilde{K}$  such that

$$h^{-1} \int_{\Gamma_h \cap \tilde{K}} |v^+ - v^-|^2 ds \geq C \|\epsilon(\tilde{v})\|_{0,\tilde{K}}^2.$$

Summing this over all  $\tilde{K} \in \tilde{\mathcal{C}}^h$ , we obtain

$$h^{-1} \int_{\Gamma_h} |v^+ - v^-|^2 ds \geq C \|\epsilon(\tilde{v})\|_0^2.$$

Noting also that, by the construction of  $\tilde{v}$ ,  $\|v\|_{0,\Gamma_1}^2 \geq C \|\tilde{v}\|_{0,\Gamma_1}^2$ , we obtain the inequality

$$\|v\|_{1,h}^2 \geq C(\|\epsilon(\tilde{v})\|_0^2 + \|\tilde{v}\|_{0,\Gamma_1}^2).$$

Further, since  $\Gamma_1$  contains at least one side of  $\Omega$ , we have by Korn's inequality (cf. [7], [13])

$$\|\epsilon(\tilde{v})\|_0^2 + \|\tilde{v}\|_{0,\Gamma_1}^2 \geq C \|\tilde{v}\|_1^2.$$

We finally note that, by the construction of  $\tilde{v}$ ,  $\|\tilde{v}\|_0 \geq C \|v\|_0$ . Upon combining the last three inequalities, the first part of the assertion follows.

For the second part of the assertion, we need the additional estimates

$$(4.14) \quad \|v\|_{0,\infty} \leq C \|\tilde{v}\|_{0,\infty} \leq C_1(1 + |\log h|)^{1/2} \|\tilde{v}\|_1.$$

Combining (4.14) with the above inequalities, the proof is completed, so it remains to prove (4.14).

The first inequality in (4.14) follows easily from the construction of  $\tilde{v}$ , so let us only prove the second part of (4.14). Let  $w$  be an extension of  $\tilde{v}$  to  $\mathbf{R}^2$  such that  $w$  vanishes outside some disc  $S$  of finite radius and satisfies

$$\|w\|_{H^1(\mathbf{R}^2)} \leq C \|\tilde{v}\|_1.$$

Let  $x \in \mathcal{N}$ ,  $x \notin \partial\Omega$ . Then  $\tilde{v}(x)$  may be written, by Green's formula, as

$$\tilde{v}(x) = \frac{1}{2\pi} \int_S \nabla(\log|x - y|) \cdot \nabla w(y) dy.$$

Let  $S_h$  be a disc centered at  $x$  and of radius  $\rho = \min\{h_1, h_2\}$ . Then  $S_h \subset \Omega$ , and we have by a scaling argument

$$\left| \int_{S_h} \nabla(\log|x - y|) \cdot \nabla w(y) dy \right| \leq \int_{S_h} \frac{1}{|x - y|} |\nabla \tilde{v}(y)| dy \leq C \|\tilde{v}\|_{1,S_h}.$$

Since, on the other hand,

$$\begin{aligned} & \left| \int_{S \setminus S_h} |\nabla \log(x - y) \cdot \nabla w(y) dy \right| \\ & \leq \left\{ \int_{S \setminus S_h} \frac{1}{|x - y|^2} dy \right\}^{1/2} \|w\|_{H^1(\mathbf{R}^2)} \leq C(1 + |\log h|)^{1/2} \|\tilde{v}\|_1, \end{aligned}$$

it follows that, for  $x \in \mathcal{N}$ ,  $x \notin \partial\Omega$ ,

$$(4.15) \quad |\tilde{v}(x)| \leq C(1 + |\log h|)^{1/2} \|\tilde{v}\|_1.$$



Finally, if  $x \in \bar{\Omega}$  is arbitrary, we find by a scaling argument that  $|\tilde{v}(x)| \leq |v(y)| + C|\tilde{v}|_{1,K}$ , where  $K \in \mathcal{C}^h$  is such that  $x \in \bar{K}$  and  $y$  is any vertex of  $K$ . Thus, (4.15) holds for any  $x \in \bar{\Omega}$ , and so the proof is complete.  $\square$

Using the above mesh-dependent norms, we can state a basic stability condition for method (2.3) as: There is a constant  $C$  such that

$$(4.16) \quad \sup_{\tau \in H_h} \frac{(\operatorname{div} \tau, v)}{|\tau|_{0,h}} \geq C|v|_{1,h}, \quad v \in V_h.$$

Note that, by Lemma 4.1, (4.16) is a stronger stability inequality than that proved for Method I (Lemma 3.1).

Let us first prove

LEMMA 4.2. *If  $H_h = H_{h,1}$ , then (4.16) holds with  $C = 1$ .*  $\square$

*Proof.* If  $H_h = H_{h,1}$ , the trace space  $\{\tau \cdot n : \tau \in H_h\}$  contains the trace space of  $V_h$  on  $\Gamma_h \cup \Gamma_1$ . Thus if  $v \in V_h$  is given, there exists  $\tau \in H_h$  such that  $\tau \cdot n = v^- - v^+$  on  $\Gamma_h$  and  $\tau \cdot n = v$  on  $\Gamma_1$ . Using (4.12) we see that  $(\operatorname{div} \tau, v) = |\tau|_{0,h}|v|_{1,h}$ , so the assertion follows.  $\square$

*Remark.* If one chooses  $H_h = H_{h,0}$ , it is easy to see that (4.16) can only be valid if  $C$  depends on  $h$ . In fact, it is not difficult to verify that, even if the space  $V_h$  is reduced to consist only of functions that are piecewise constant, (4.16) does not hold if  $H_h = H_{h,0}$ . Note that such a method has essentially the same properties as Method I above.  $\square$

Before showing examples where the stability assumption (4.16) holds with  $H_h \neq H_{h,1}$ , let us prove an error estimate for the method (2.3) assuming merely that  $H_{h,1} \supset H_h \supset H_{h,0}$  and that (4.16) holds. We need the following nonstandard interpolant in  $V_h$ : if  $u \in V$ , define  $\tilde{u} \in V_h$  by requiring

$$(4.17) \quad \int_K (u - \tilde{u}) \operatorname{div} \tau \, dx = 0, \quad \tau \in H_K, K \in \mathcal{C}^h,$$

where  $H_K$  is defined as above.

LEMMA 4.3. *The interpolant  $\tilde{u}$  is uniquely determined. Moreover, one has the estimates*

$$\|u - \tilde{u}\|_0 \leq Ch|u|_1 \quad \text{and} \quad |(u - \tilde{u})(Q_K)| \leq Ch^2|u|_{2,\infty,K}, \quad K \in \mathcal{C}^h,$$

where  $Q_K$  is the midpoint of  $K$ .  $\square$

*Proof.* Let  $w_K = \{\operatorname{div} \tau, \tau \in G_K\}$ , where  $G_K$  is as in (4.2). In [9] it is shown that  $w_K = \pi_K([P_1(K)]^2)$ , where  $\pi_K$  denotes the orthogonal projection of  $[L_2(K)]^2$  into the subspace consisting of functions  $v$  such that  $v|_T$  is constant on each of the four subtriangles of  $K$ . By the definition of  $H_K$  we then have

$$\{\operatorname{div} \tau : \tau \in H_K\} = \pi_K(V_K).$$

Thus, we may set  $\operatorname{div} \tau = \pi_K \tilde{u}$  in (4.17) to obtain

$$\int_K \tilde{u} \pi_K \tilde{u} \, dx = \int_K (\pi_K \tilde{u})^2 \, dx = \int_K u \pi_K \tilde{u} \, dx.$$

Using here the easy-to-prove inequality

$$\|\pi_K v\|_{0,K}^2 \geq \frac{2}{3} \|v\|_{0,K}^2, \quad v \in V_K,$$

we see that  $u$  satisfies  $\|\tilde{u}\|_{0,K} \leq \frac{3}{2} \|u\|_{0,K}$ . So,  $\tilde{u}$  is uniquely determined. Moreover, since  $\tilde{u} = u$  if  $u$  is constant, we conclude from the Bramble-Hilbert lemma [4] that  $\|u - \tilde{u}\|_{0,K} \leq Ch|u|_{1,K}$ . This proves the first inequality in the lemma.

To prove the second inequality, note that if we write

$$\tilde{u}(x) = (\alpha_1, \alpha_2) + \alpha_3(x_2 - q_2, -x_1 + q_1) = \tilde{u}_0(x) + \tilde{u}_1(x), \quad x \in K,$$

where  $Q = (q_1, q_2)$  is the midpoint of  $K$ , then by (4.17)  $\int_K (u - \tilde{u}_0) dx = 0$ . Thus,  $\tilde{u}_0$  is the  $L_2$ -projection of  $u$  into  $[P_0(K)]^2$ . Since  $(u - \tilde{u})(Q) = (u - \tilde{u}_0)(Q) = 0$  if  $u \in [P_1(K)]^2$ , we obtain from the Bramble-Hilbert lemma the estimate

$$|(u - \tilde{u})(Q)| \leq Ch^2 |u|_{2,\infty,K},$$

and the proof is complete.  $\square$

We are now ready to state the main convergence result.

**THEOREM 4.1.** *Let  $(\sigma, u)$  be the solution of (2.2), and let  $(\sigma_h, u_h)$  be the solution of (2.3), where  $V_h$  is defined by (4.1a) and  $H_h \supset H_{h,0}$  is chosen so the stability condition (4.16) holds. Then we have the error estimates*

$$\|\sigma - \sigma_h\|_0 \leq Ch^2 |\sigma|_2,$$

$$|(u - u_h)(Q_K)| \leq Ch^2(1 + |\log h|)^{1/2} |\sigma|_2 + Ch^2 |u|_{2,\infty,K}, \quad K \in \mathcal{C}^h,$$

where  $Q_K$  is the midpoint of  $K$ .  $\square$

*Proof.* Let  $\tilde{\sigma} \in H_h$  be defined by (4.7) (or equivalently by (4.8)), and let  $\tilde{u} \in V_h$  be defined by (4.17). Then by (2.7) and (4.16) and since  $|\tau|_{0,h} \leq C\|\tau\|_0$  for  $\tau \in H_h$ , by Lemma 4.1, we conclude by the standard argument (cf. [1],[5]) that there exists  $(\tau, v) \in H_h \times V_h$  satisfying

$$(4.18a) \quad \|\tau\|_0 + |v|_{1,h} \leq C,$$

$$(4.18b) \quad \mathfrak{B}(\sigma_h - \tilde{\sigma}, u_h - \tilde{u}; \tau, v) \geq \|\sigma_h - \tilde{\sigma}\|_0 + |u_h - \tilde{u}|_{1,h},$$

where  $\mathfrak{B}$  is as in (3.2). Using (4.17) and (4.18b) in the identity (3.1), we obtain

$$\|\sigma_h - \tilde{\sigma}\|_0 + |u_h - \tilde{u}|_{1,h} \leq |(\sigma - \tilde{\sigma}, \tau)| + |(\operatorname{div}(\sigma - \tilde{\sigma}), v)|.$$

Using (4.18a) and (4.9), we have

$$|(\sigma - \tilde{\sigma}, \tau)| \leq \|\sigma - \tilde{\sigma}\|_0 \|\tau\|_0 \leq Ch^2 |\sigma|_2.$$

Similarly using (4.13) and (4.18a),

$$|(\operatorname{div}(\sigma - \tilde{\sigma}), v)| \leq |\sigma - \tilde{\sigma}|_{0,h} |v|_{1,h} \leq C |\sigma - \tilde{\sigma}|_{0,h} \leq C_1 h^2 |\sigma|_2.$$

Here the interpolation estimate in the norm  $|\cdot|_{0,h}$  is proved by standard techniques; cf. [2] for details of the argument.

Combining these estimates, using the triangle inequality and recalling (4.9), we obtain

$$(4.19) \quad \|\sigma - \sigma_h\|_0 + |u_h - \tilde{u}|_{1,h} \leq Ch^2 |\sigma|_2.$$

This proves the asserted estimate for  $\|\sigma - \sigma_h\|_0$ . Finally, by (4.19), Lemma 4.1 and Lemma 4.3, if  $Q_K$  is the midpoint of  $K \in \mathcal{C}^h$ , then

$$\begin{aligned} |(u - u_h)(Q_K)| &\leq |(u_h - \tilde{u})(Q_K)| + |(u - \tilde{u})(Q_K)| \\ &\leq C(1 + |\log h|)^{1/2} |u_h - \tilde{u}|_{1,h} + Ch^2 |u|_{2,\infty,K} \\ &\leq Ch^2(1 + |\log h|)^{1/2} |\sigma|_2 + Ch^2 |u|_{2,\infty,K}, \end{aligned}$$

which completes the proof.  $\square$

In the remaining part of this section we consider the practical problem of constructing  $H_h$  so that the stability condition (4.16) holds. We first formulate a general criterion which is sufficient (and probably also necessary) for the validity of (4.16). Let  $\mathcal{E}^h$  be a collection of "macroelements", i.e., a collection of open rectangles  $\tilde{K}$  such that if  $K \in \mathcal{C}^h$  and  $\tilde{K} \in \mathcal{E}^h$ , then either  $K \subset \tilde{K}$  or  $K \cap \tilde{K} = \emptyset$ . We associate to each  $\tilde{K} \in \mathcal{E}^h$  the following subspaces:

$$(4.20a) \quad U(\tilde{K}) = \left\{ \tau_{|\tilde{K}} : \tau \in H_h, \tau = 0 \text{ in } \Omega \setminus \tilde{K} \right\},$$

$$(4.20b) \quad N(\tilde{K}) = \left\{ v_{|\tilde{K}} : v \in V_h, \int_{\tilde{K}} \operatorname{div} \tau \cdot v \, dx = 0 \, \forall \tau \in U(\tilde{K}) \right\}.$$

Let us now assume that there exists an integer  $M \geq 1$  independent of  $h$  such that  $\mathcal{E}^h$  satisfies the following hypotheses:

(i) Each  $\tilde{K} \in \mathcal{E}^h$  contains at most  $M$  different rectangles  $K \in \mathcal{C}^h$ .

(ii) If  $l$  is a side of  $K \in \mathcal{C}^h$  and  $l$  is not on  $\Gamma$ , then  $l$  is contained in at least one and not more than  $M$  different rectangles  $\tilde{K} \in \mathcal{E}^h$ .

(iii) For each  $\tilde{K} \in \mathcal{E}^h$ ,  $N(\tilde{K}) = \{0\}$  if  $\partial\tilde{K} \cap \Gamma_1 \neq \emptyset$ , otherwise  $N(\tilde{K})$  is the space of rigid displacements of  $\tilde{K}$ , i.e.,  $\dim(N(\tilde{K})) = 3$ .

We can now prove

**LEMMA 4.4.** *If there is a collection  $\mathcal{E}^h$  which satisfies the above assumptions, then (4.16) holds.  $\square$*

*Proof.* For  $\tilde{K} \in \mathcal{E}^h$  let

$$W(\tilde{K}) = \left\{ w_{|\tilde{K}} : w \in V_h, \int_{\tilde{K}} w \cdot v \, dx = 0 \, \forall v \in N(\tilde{K}) \right\}.$$

and let

$$|v|_{1,h,\tilde{K}}^2 = h^{-1} \int_{\Gamma_h \cap \tilde{K}} |v^+ - v^-|^2 \, ds + h^{-1} \int_{\Gamma_1 \cap \partial\tilde{K}} |v|^2 \, ds, \quad v \in V_h.$$

It follows from assumption (iii) that  $|\cdot|_{1,h,\tilde{K}}$  is a norm on  $W(\tilde{K})$ . For, if  $v \in W(\tilde{K})$  and  $|v|_{1,h,\tilde{K}} = 0$ , then  $v = 0$  if  $\tilde{K}$  has a side on  $\Gamma_1$ , and otherwise  $v$  is a rigid displacement of  $\tilde{K}$ . By assumption (iii), this implies that  $v \in N(\tilde{K})$ , so  $v \in W(\tilde{K}) \cap N(\tilde{K})$ , which is possible only if  $v = 0$ .

On the other hand, if we define

$$|v|_{h,\tilde{K}} = \sup_{\tau \in U(\tilde{K})} \frac{\int_{\tilde{K}} \operatorname{div} \tau \cdot v \, dx}{\|\tau\|_{0,\tilde{K}}},$$

then  $|\cdot|_{h,\tilde{K}}$  is another norm on  $W(\tilde{K})$ , by the definition of  $W(\tilde{K})$ . Using now a scaling argument and the equivalence of norms in a finite-dimensional space, recalling the assumption (i), we conclude the existence of a constant which depends only on  $M$  such that, for all  $v \in W(\tilde{K})$ ,

$$(4.21) \quad |v|_{h,\tilde{K}} \geq C |v|_{1,h,\tilde{K}}.$$

Since  $|v|_{h,\tilde{K}} = |v|_{1,h,\tilde{K}} = 0$  if  $v \in N(\tilde{K})$ , we further conclude that (4.21) actually holds for all  $v \in V_{h|\tilde{K}}$ . In view of the definition of  $|\cdot|_{h,\tilde{K}}$ , this may be interpreted as: For each  $v \in V_{h|\tilde{K}}$ , there exists  $\tau_{\tilde{K}} \in U(\tilde{K})$  such that

$$(4.22a) \quad \|\tau_{\tilde{K}}\|_0 \leq C |v|_{1,h,\tilde{K}},$$

$$(4.22b) \quad \int_{\tilde{K}} \operatorname{div} \tau_{\tilde{K}} \cdot v \, dx \geq |v|_{1,h,\tilde{K}}^2,$$

where  $C$  depends only on  $M$ .

Now let  $v \in V_h$  be given, define  $\tau_{\tilde{K}} \in H_h$  for each  $\tilde{K} \in \mathcal{S}^h$  so that (4.22) holds and  $\tau_{\tilde{K}} = 0$  on  $\Omega \setminus \tilde{K}$ , and let  $\tau \in H_h$  be defined by

$$\tau = \sum_{\tilde{K} \in \mathcal{S}^h} \tau_{\tilde{K}}.$$

Then, by (4.22a) and by assumptions (i) and (ii),

$$\begin{aligned} \|\tau\|_0^2 &= \sum_{K \in \mathcal{C}^h} \left\| \sum_{\substack{\tilde{K} \in \mathcal{S}^h \\ \tilde{K} \supset K}} \tau_{\tilde{K}} \right\|_{0,K}^2 \leq M \sum_{\tilde{K} \in \mathcal{S}^h} \|\tau_{\tilde{K}}\|_{0,\tilde{K}}^2 \\ &\leq CM \sum_{\tilde{K} \in \mathcal{S}^h} |v|_{1,h,\tilde{K}}^2 \leq 2CM^2 |v|_{1,h}^2. \end{aligned}$$

Here we also used the fact that if  $l$  is a side of  $K \in \mathcal{C}^h$  on  $\Gamma_1$ , then  $l$  is contained in  $\partial\tilde{K}$  for at most  $2M$  different rectangles  $\tilde{K} \in \mathcal{S}^h$ . This is a consequence of assumption (ii). From (ii) it also follows that each side of each  $K \in \mathcal{C}^h$  is contained in at least one  $\tilde{K} \in \mathcal{S}^h$  if  $l$  is in the interior of  $\Omega$  or in at least one  $\partial\tilde{K}$  if  $l \subset \Gamma$ . Therefore, and by (4.22b),

$$(\operatorname{div} \tau, v) = \sum_{\tilde{K} \in \mathcal{S}^h} \int_{\tilde{K}} \operatorname{div} \tau_{\tilde{K}} \cdot v \, dx \geq \sum_{\tilde{K} \in \mathcal{S}^h} |v|_{1,h,\tilde{K}}^2 \geq |v|_{1,h}^2.$$

Combining these inequalities, (4.16) follows.  $\square$

We consider now two examples where the above hypotheses can be verified.

*Example 1.* Let  $H_h$  be defined by

$$H_h = \left\{ \tau = \sum_{\substack{P \in \mathcal{R} \\ P \notin \Gamma}} \sum_{i=1}^6 \alpha_i^P \tau_i^P + \sum_{\substack{P \in \mathcal{R} \\ P \in \Gamma}} \sum_{i=1}^3 \alpha_i^P \tau_i^P, \alpha_i^P \in \mathbf{R}, \tau \cdot n = 0 \text{ on } \Gamma_2 \right\},$$

where the functions  $\tau_i^P$  are as above.  $\square$

Assume that  $m_1, m_2 \geq 3$  in (2.5), and let  $\mathcal{S}^h$  be a collection of rectangles  $\tilde{K}$  such that assumption (ii) is satisfied for some finite  $M$  and such that each  $\tilde{K} \in \mathcal{S}^h$  is of the form

$$\tilde{K} = \{x \in \Omega: ih_1 < x_1 < (i + \nu_1)h_1, jh_2 < x_2 < (j + \nu_2)h_2\},$$

where  $\nu_i = 3$  if  $\tilde{K}$  has a side  $l$  on  $\Gamma$  parallel with the  $x_i$ -axis and otherwise  $\nu_1 = \nu_2 = 2$ . For such  $\mathcal{E}^h$ , assumption (i) holds with  $M = 9$ , so it remains to verify (iii).

Assume first that  $\partial\tilde{K} \cap \Gamma = \emptyset$ . Then if  $\tau \in H_h$  and  $\tau = 0$  on  $\tilde{K} \setminus \Omega$ , we have (see Figure 5)

$$(4.23) \quad \tau = \sum_{i=1}^6 \alpha_i \tau_i^A + \alpha_7 (\tau_1^B - \tau_4^B) + \alpha_8 (\tau_1^C + \tau_4^C) + \alpha_9 (\tau_2^D - \tau_5^D) + \alpha_{10} (\tau_2^E + \tau_5^E), \quad \alpha_i \in \mathbf{R}.$$

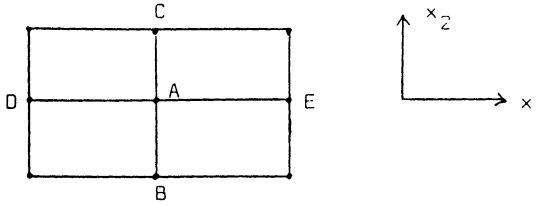


FIGURE 5

Hence if  $U(\tilde{K})$  is defined by (4.20a) we have  $\dim(U(\tilde{K})) = 10$ . If  $N(\tilde{K})$  is now defined as in (4.20b), then  $v \in N(\tilde{K})$  if and only if

$$\int_K \operatorname{div} \tau \cdot v \, dx = \int_{\Gamma_h \cap K} (\tau \cdot n) \cdot (v^+ - v^-) \, ds = 0$$

for all  $\tau$  of the form (4.23). It is straightforward to verify that this is possible only if  $v^+ = v^-$  on  $\Gamma_h$ , i.e.,  $v$  is a rigid displacement of  $\tilde{K}$ .

Assume next that  $\tilde{K}$  has a side on  $\Gamma_1$  which contains the nodes  $A, B \in \mathcal{N}$  (see Figure 6). It follows from the above consideration that if  $v \in N(\tilde{K})$ , then  $v$  can be at most a rigid displacement of  $K$ . But by (4.20) and by the definition of  $H_h$  we also have

$$(4.24) \quad \int_{\tilde{K}} \operatorname{div} \tau_i^P \cdot v \, dx = \int_{\partial\tilde{K} \cap \Gamma_1} (\tau_i^P \cdot n) \cdot v \, ds = 0, \quad P = A, B, i = 1, 2, 3.$$

It is easy to see that if  $v$  is a rigid displacement on  $\tilde{K}$ , then (4.24) implies that  $v = 0$ . Thus, assumption (iii) is verified and (4.16) now follows from Lemma 4.4.

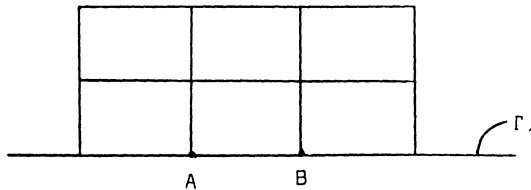


FIGURE 6

*Example 2.* We assume here that  $m_1$  and  $m_2$  in (2.5) are even numbers, with  $m_1, m_2 \geq 4$ . Denote by  $\tilde{\mathcal{C}}^h$  a coarser rectangular subdivision of  $\Omega$  such that each  $\tilde{K} \in \tilde{\mathcal{C}}^h$  is of the type

$$\tilde{K} = \{x \in \Omega: 2ih_1 < x_1 < 2(i + 1)h_1, 2jh_2 < x_2 < 2(j + 1)h_2\}$$

for some  $i, j, 0 \leq i \leq m_1/2, 0 \leq j \leq m_2/2$ . Let  $\mathfrak{N}_i, i = 0, \dots, 3$ , be subsets of the nodal set  $\mathfrak{N}$  such that  $\mathfrak{N}_0$  is the set of midpoints of  $\tilde{K} \in \tilde{\mathcal{C}}^h$ ,  $\mathfrak{N}_i, i = 1, 2$ , is the set of midpoints of the sides of  $\tilde{K} \in \tilde{\mathcal{C}}^h$  that are parallel with the  $x_i$ -axis, and  $\mathfrak{N}_3$  is the set of vertices of  $\tilde{K} \in \tilde{\mathcal{C}}^h$ . Then define  $H_h$  as

$$H_h = \left\{ \tau = \sum_{P \in \mathfrak{N}_0} \sum_{i=1}^3 \alpha_i^P \tau_i^P + \sum_{P \in \mathfrak{N}_0} \sum_{i=4,6} \alpha_i^P \tau_i^P + \sum_{P \in \mathfrak{N}_1} \alpha_4^P \tau_4^P + \sum_{P \in \mathfrak{N}_2} \alpha_5^P \tau_5^P, \alpha_i^P \in \mathbf{R}, \tau \cdot n = 0 \text{ on } \Gamma_2 \right\}. \quad \square$$

Note that

$$\dim(H_h) = 4m_1m_2 + O(h^{-1}),$$

so there are only about four degrees of freedom per node in this case.

To verify assumptions (i) through (iii) for the above choice of  $H_h$ , let  $\mathcal{E}^h$  be a collection of rectangles  $\tilde{K}$  such that (ii) is satisfied for some finite  $M$  and such that each  $\tilde{K} \in \mathcal{E}^h$  is of the type

$$\tilde{K} = \{x \in \Omega: 2ih_1 < x_1 < (2i + 4)h_1, 2jh_2 < x_2 < (2j + 4)h_2\},$$

where  $i, j \geq 0, (2i + 4)h_1 \leq a_1, (2j + 4)h_2 \leq a_2$ . Then assumption (i) holds with  $M = 16$ .

To verify assumption (iii), consider first a given macroelement  $\tilde{K}$  which contains four rectangles of  $\mathcal{C}^h$ . With the notation of Figure 5, if  $\tau \in H_h$  vanishes outside  $\tilde{K}$ , then

$$\tau = \sum_{i=1}^4 \alpha_i \tau_i^A + \alpha_5 \tau_6^A + \alpha_6 (\tau_1^B - \tau_4^B) + \alpha_7 (\tau_1^C + \tau_4^C) + \alpha_8 (\tau_2^D - \tau_5^D) + \alpha_9 (\tau_2^E + \tau_5^E), \quad \alpha_i \in \mathbf{R}.$$

We omit the details of showing that if  $v \in V_h$  and  $\int_{\tilde{K}} \operatorname{div} \tau \cdot v \, dx = 0$  for all  $\tau$  of the above form, then  $v|_{\tilde{K}}$  is a rigid displacement.

Consider now a rectangle  $\tilde{K} \in \tilde{\mathcal{E}}^h$  consisting of four subrectangles  $\tilde{K}_i$ , as in Figure 7, and let  $U(\tilde{K})$  and  $N(\tilde{K})$  be defined by (4.20).

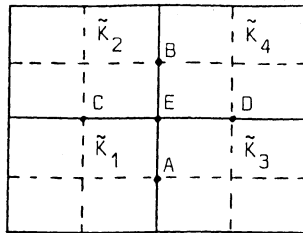


FIGURE 7

We noted above that if  $v \in N(\tilde{K})$ , then  $v|_{\tilde{K}_i}$  is a rigid displacement for  $i = 1, \dots, 4$ . Therefore, if  $v \in N(\tilde{K})$ , we may write

$$(4.25) \quad v(x) = \frac{h_2}{h} \alpha_i(1, 0) + \frac{h_1}{h} \beta_i(0, 1) + \frac{1}{2h} \gamma_i(x_2 - q_{2i}, q_{1i} - x_1), \quad x \in \tilde{K}_i, i = 1, \dots, 4,$$

where  $\alpha_i, \beta_i, \gamma_i \in \mathbf{R}$  and  $(q_{1i}, q_{2i})$  is the midpoint of  $\tilde{K}_i$ . Now, since  $\tau_{i\tilde{K}}^P \in U(\tilde{K})$  for  $i = 1, 2, 3$  and  $P = A, B, C, D, E$  (see Figure 7), any  $v \in N(\tilde{K})$  satisfies, for example,

$$\int_{\tilde{K}} \operatorname{div} \tau_i^P \cdot v \, dx = 0 \quad \text{for} \quad \begin{cases} P = A, B, & i = 1, 3, \\ P = C, D, & i = 2, 3, \\ P = E, & i = 1. \end{cases}$$

With  $v$  given by (4.25), these equations are written equivalently as

$$\begin{cases} \alpha_1 - \alpha_3 & = 0 \\ \beta_1 - \beta_3 - \gamma_1 - \gamma_3 & = 0 \\ \alpha_2 - \alpha_4 & = 0 \\ \beta_2 - \beta_4 - \gamma_2 - \gamma_4 & = 0 \\ \beta_1 - \beta_2 & = 0 \\ \alpha_1 - \alpha_2 + \gamma_1 + \gamma_2 & = 0 \\ \beta_3 - \beta_4 & = 0 \\ \alpha_3 - \alpha_4 + \gamma_3 + \gamma_4 & = 0 \\ \alpha_1 - \alpha_2 - \alpha_3 - \alpha_4 + \frac{2}{3}(\gamma_1 - \gamma_2 - \gamma_3 + \gamma_4) & = 0 \end{cases}$$

From this it is easy to see that the equations are all linearly independent. Thus, if  $\partial\tilde{K} \cap \Gamma_1 = \emptyset$ ,  $\dim(N(\tilde{K})) = 12 - 9 = 3$ , and so  $N(\tilde{K})$  consists of the rigid displacements only.

Finally if  $\partial\tilde{K} \cap \Gamma_1 \neq \emptyset$ , one can further show that  $N(\tilde{K}) = \{0\}$  using the same argument as in Example 1 above. Thus, assumption (iii) is verified, and so (4.16) follows.

**5. A Numerical Example.** We give here the results of numerical computations using the methods presented in the preceding sections. We consider a simple model problem where  $\Omega$  is the unit square,  $\lambda = -0.3$  and  $\mu = 1.3$  in (2.1), and  $f = (f_1, f_2)$  is chosen so that the exact displacements under the boundary condition  $u = 0$  on  $\Gamma$  are

$$\begin{aligned} u_1 &= 16x_1(1 - x_1)x_2(1 - x_2)e^{(x_1 - x_2)}, \\ u_2 &= \sin \pi x_1 \sin \pi x_2, \quad (x_1, x_2) \in \Omega. \end{aligned}$$

In solving the discrete equations, the following iterative version of the penalty method (2.4) is used (cf. [11, p. 321])

$$\begin{aligned} (5.1) \quad a(\sigma_h^k, \tau) &+ \frac{1}{\epsilon} (\pi_h \operatorname{div} \sigma_h^k, \pi_h \tau) \\ &= -\frac{1}{\epsilon} (\pi_h f, \pi_h \operatorname{div} \tau) - (u_h^{k-1}, \operatorname{div} \tau), \quad \tau \in H_h, \\ u_h^k &= u_h^{k-1} + \frac{1}{\epsilon} \pi_h (\operatorname{div} \sigma_h^k + f), \quad k = 1, 2, \dots \end{aligned}$$

Using  $u_h^0 = 0$  as the starting guess, we see that the first step in (5.1) is equivalent to the penalty method (2.3a, b').

The main benefit of using (5.1) instead of (2.3a, b') is that one has more freedom in choosing the parameter  $\epsilon$ . For, if the scheme (2.3a, b') is used, one has to take  $\epsilon = O(h^2)$  in order to obtain the convergence rate  $\|\sigma - \sigma_h\|_0 = O(h^2)$  (cf. [10]). This

makes the system (2.4) rather ill-conditioned (the condition number is of the order of  $h^{-4}$ ). Also, one has to compute  $\pi_h f$  with high precision in (2.4) since the error is multiplied by  $1/\varepsilon$ .

In practice the iteration (5.1) seems to converge quite fast: not more than six iterations were required in the computations. Since only one Cholesky decomposition is required in (5.1), the additional cost due to iteration is relatively small.

Table I shows the computed  $L_2$ -errors  $\|\sigma - \sigma_h\|_0$  in the above model problem. Method IIA corresponds to the choice  $H_h = H_{h,1}$ , in IIB  $H_h$  is chosen as in Example 1 and in IIC as in Example 2 of Section 4. In parentheses is shown the relative number of operations required in the Cholesky reduction of the matrix in (5.1).

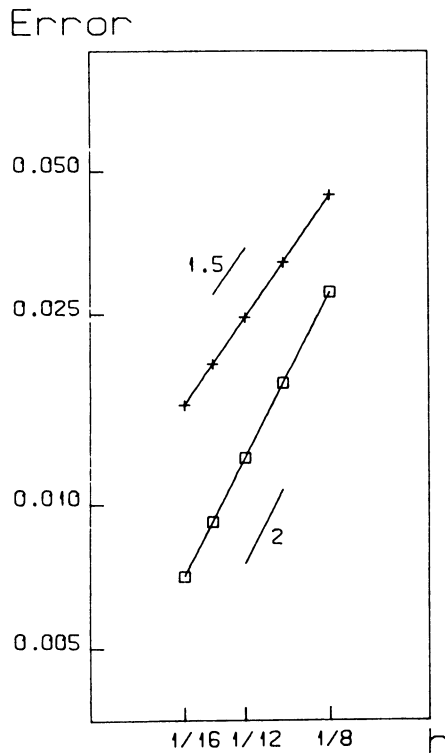


FIGURE 8

$L_2$ -errors  $\|\sigma - \sigma_h\|_0$  in a model problem for Method I (+) and Method IIC (□). Two line segments with slopes  $3/2$  and  $2$  are drawn for comparison.

As shown in Figure 8, the rate of convergence is  $O(h^{3/2})$  for Method I and  $O(h^2)$  for Method II, as expected theoretically. Notice also that the constant in the error estimate is practically the same in the three variants of Method II. Thus, Method IIC should be preferred as it involves the least amount of computational work. As compared with Method I, Method IIC is superior roughly below the error level  $\|\sigma - \sigma_h\|_0 \approx 0.03$ .



TABLE I

$L_2$ -errors  $\|\sigma - \sigma_h\|_0$  for various methods in a model problem. Numbers in parentheses indicate the relative cost of Cholesky reduction.

1/h	Method I		Method II					
			A		B		C	
4	0.123	(1.0)	0.090	(19)	0.101	(8.0)	0.102	(3.0)
6	0.068	(5.1)	0.042	(97)	0.047	(41)	0.048	(17)
8	0.045	(16)	0.024	(300)	0.026	(130)	0.027	(48)
10	0.032	(39)	0.016	(740)	0.017	(310)	0.018	(120)
12	0.025	(81)	0.011	(1500)	0.012	(650)	0.012	(240)

Institute of Mathematics  
Helsinki University of Technology  
SF-02150 Espoo 15, Finland

1. I. BABUŠKA, "Error bounds for finite element method," *Numer. Math.*, v. 16, 1971, pp. 322–333.
2. I. BABUŠKA, J. OSBORN & J. PITKÄRANTA, "Analysis of mixed methods using mesh dependent norms," *Math. Comp.*, v. 35, 1980, pp. 1039–1062.
3. M. BERCOVIER, "Perturbation of mixed variational problems. Application to mixed finite element methods," *RAIRO Anal. Numér.*, v. 12, 1978, pp. 211–236.
4. J. BRAMBLE & S. HILBERT, "Estimation of linear functionals on Sobolev spaces with application to Fourier transforms and spline interpolation," *SIAM J. Numer. Anal.*, v. 7, 1970, pp. 112–124.
5. F. BREZZI, "On the existence, uniqueness and approximation of saddle-point problems arising from Lagrange multipliers," *RAIRO Ser. Rouge*, v. 8, 1974, pp. 129–151.
6. P. G. CIARLET, *The Finite Element Method for Elliptic Problems*, North-Holland, Amsterdam, 1978.
7. G. DUVAUT & J. L. LIONS, *Inequalities in Mechanics and Physics*, Springer-Verlag, Berlin and New York, 1976.
8. I. HLAVAČEK, "Convergence of an equilibrium finite element method for plane elastostatics," *Appl. Math.*, v. 24, 1979, pp. 427–456.
9. C. JOHNSON & B. MERCIER, "Some equilibrium finite element methods for two-dimensional elasticity problems," *Numer. Math.*, v. 30, 1978, pp. 103–116.
10. C. JOHNSON & J. PITKÄRANTA, "Analysis of some mixed finite element methods related to reduced integration," *Math. Comp.*, v. 38, 1982, pp. 375–400.
11. D. G. LUENBERGER, *Introduction to Linear and Nonlinear Programming*, Addison-Wesley, Reading, Mass., 1973.
12. D. MALKUS & T. HUGHES, "Mixed finite element methods—reduced and selective integration techniques: A unification of concepts," *Comp. Methods Appl. Mech. Engrg.*, v. 15, 1978, pp. 63–81.
13. J. NITSCHKE, "On Korn's second inequality," *RAIRO Anal. Numér.*, v. 15, 1981, pp. 237–248.
14. G. SANDER, *Application of the Dual Analysis Principle*, Proc. of IUTAM Symp. on High Speed Computing of Elastic Structures, Univ. de Liège, 1971, pp. 167–207.
15. R. TAYLOR & O. C. ZIENKIEWICZ, "Complementary energy with penalty functions in finite element analysis," in *Energy Methods in Finite Element Analysis* (R. Glowinski, E. Y. Rodin and O. C. Zienkiewicz, eds.), Wiley, New York, 1979, pp. 143–174.
16. V. B. WATWOOD, JR. & B. J. HARTZ, "An equilibrium stress field model for finite element solutions of two-dimensional elastostatic problems," *Internat. J. Solids and Structures*, v. 4, 1968, pp. 857–873.
17. B. FRAEJIS DE VEUBEKE, "Displacement and equilibrium models in the finite element method," in *Stress Analysis* (O. C. Zienkiewicz and G. S. Holister, eds.), Wiley, New York, 1965, pp. 145–197.
18. B. FRAEJIS DE VEUBEKE, *Finite Element Method in Aerospace Engineering Problems*, Proc. Internat. Sympos. Computing Methods in Appl. Sci. and Eng., Versailles, 1973, Part 1, pp. 224–258.