

NASA Contractor Report 195293

1N-39
2419
127P

Analysis and Control of Hourglass Instabilities in Underintegrated Linear and Nonlinear Elasticity

Olivier P. Jacquotte and J. Tinsley Oden
The University of Texas at Austin
Austin, Texas

(NASA-CR-195293) ANALYSIS AND
CONTROL OF HOURGLASS INSTABILITIES
IN UNDERINTEGRATED LINEAR AND
NONLINEAR ELASTICITY Final Report
(Texas Univ.) 127 p

N94-28266

Unclas

March 1994

G3/39 0002419

Prepared for
Lewis Research Center
Under Grant NAG3-329

NASA
National Aeronautics and
Space Administration

ANALYSIS AND CONTROL OF HOURGLASS INSTABILITIES
IN UNDERINTEGRATED LINEAR AND NONLINEAR ELASTICITY*

Olivier P. Jacquotte and J. Tinsley Oden
Texas Institute for Computational Mechanics
The University of Texas at Austin
Austin, Texas, U.S.A.

July 1985

TICOM Report 85-10

*This work was supported by Grant NAG3-329 from the Lewis
Research Center, NASA.

PART I: INTRODUCTION

1.1 General

Among various choices involved in applying the finite element method for the resolution of problems of fluid or solid mechanics, the choice of the element type associated with a given mesh and that of the integration rule used on these elements often decide the properties of the approximation of the operator to be discretized. A bad approximation generally results in a kernel larger than the kernel of the continuous operator, and these spurious elements of this kernel may appear in the discrete solution. This solution then can exhibit undesirable oscillations and instabilities.

In particular, for the general class of problem:

$$Au - B*p = f$$

$$Bu = 0$$

oscillations and instabilities may arise from a bad approximation of the governing operator A or of the constraint operator B .

As far as the constraint is concerned, numerous theoretical and numerical studies of pressure instabilities have been done over the last decade [7, 10, 33, 36, 45]. In these studies, it has been proven that a key stability condition, the L.B.B. condition, must be satisfied in order to have existence and uniqueness of a solution [31, 9, 2], and that the constant appearing in this condition must be independent of the mesh size in order to have stability of the element considered [39, 40, 41]. Several stable and unstable elements have been studied with respect to the L.B.B. condition [19, 20, 42, 43, 44].

Our purpose in this report is to concentrate on a bad finite element approximation of the governing operator obtained when under-integration is used in numerical code for several model problems: the Poisson problem, the linear elasticity problem, and for problems in the nonlinear theory of elasticity. For each of these problems, the reasons for the occurrence of instabilities will be given, way to control or eliminate them will be presented, and theorems of existence, uniqueness and convergence for the given methods will be established. Finally, numerical results are given which illustrate the theory.

1.2 Major Results

Historical background as well as precise definitions and notations will be given in the introductions of Parts II and III. However, we list here the major results we have obtained:

- 1) In underintegrated finite element methods, a rank-deficiency counting for each element cannot work to predict the exact number of spurious modes of the global stiffness matrix.
- 2) When the spurious modes are known exactly, it is possible to eliminate them a-posteriori and obtain a unique, stable and accurate solution.
- 3) Even when the spurious zero-energy modes cannot be predicted from rank-deficiency observations, oscillations can be predicted from the eigenvalue analysis of the stiffness matrix.
- 4) Concentrated forces are the most effective source of oscillations; a way to handle such loads without exciting oscillations is theoretically obtained and its numerical implementation turns out to be efficient.

5) In linear elasticity, the spurious modes and their construction is described.

6) Stabilization methods proposed by several authors may not work.

7) For a general mesh, an element-by-element spurious mode control is presented, as well as its numerical results: this control is efficient and leads to an accurate solution.

8) This control is cheap, easy to implement, and can be used independently of the material.

9) In nonlinear incompressible elasticity, several numerical results are presented and confirm that the results previously listed can be extended to nonlinear problems.

1.3 Report Organization

Part II is devoted to the analysis and control of hourglass modes in underintegrated finite element methods. The framework of the theoretical study is the simple two-dimensional Poisson problem solved using bilinear elements; the comparison involves the 4- and 1-integration point rules. We prove that for some boundary conditions, the stiffness matrix is rank-deficient and spurious modes appear, but we also prove that the solution obtained from the underintegrated problem can be processed a-posteriori in order to obtain a stable and convergent solution in which the spurious modes have been eliminated. Moreover, the method of proof allows us to demonstrate precisely why spurious modes are excited, even though they cannot be predicted by a rank-deficiency of the stiffness matrix. Numerical experiments are discussed which not only illustrate the theory, but also generalize the results to various elements (8- and 9-node elements) and operators (linear elasticity operator).

In Part III, we numerically investigate the behavior of under-integrated 8- and 9-node elements associated with linear discontinuous pressures for the analysis of problems in finite elasticity. We observe that whereas the 9-node element is stable, the 8-node element exhibits pressure oscillations. We study the performance of the control presented in Part II to the control of oscillations appearing in the finite element solution of rubber elasticity when underintegration is practiced.

PART II: ANALYSIS OF INSTABILITIES IN
UNDERINTEGRATED FINITE ELEMENT METHODS

2.1 Introduction

For many years, a special type of numerical instability has been observed in finite difference approximations of flow fields, which has been referred to as "hourglassing", "keystoning", or "chickenwiring". These graphic terms refer to geometrical patterns which appear in computed flow fields (e.g. velocities) and which emerge as spurious oscillations superimposed on an otherwise smooth field, the spurious oscillations often taking a zig-zag form which resembles an hourglass or a chickenwire mesh. These spurious modes can be amplified upon refining the mesh, and to control such numerical instabilities, various schemes for incorporating "hourglass viscosity" or "hourglass damping" have been proposed by some authors.

It is now known that hourglass modes can arise from an incomplete (or poor) approximation of the kernel of the operators in the momentum equations in flow or solid mechanics problems (or, more generally, of the principal part of the operator in the governing differential equation of a given boundary-value problem). For example, in addition to the rigid body motions residing in the kernel of the standard operators appearing in the equilibrium (momentum) equations of solid and fluid mechanics, one finds hourglass modes in various crude discrete models of these operators.

In recent years, the occurrence of hourglass instabilities in

underintegrated finite element approximations has been observed. In the implementation of most finite element methods, integrals defining stiffness matrices are evaluated using numerical quadrature schemes. To improve computational efficiency, the practice of *underintegration* is often employed, by which is meant the use of a quadrature rule of an order lower than that required to integrate polynomial integrands exactly. This can produce rank-deficient stiffness matrices or, equivalently, an expanded kernel of the equilibrium operation which contains spurious hourglass modes, and the result is again a numerically unstable scheme.

In order to overcome this difficulty, artificial stiffness or viscosity methods, or other stabilization methods have been proposed by several authors (e.g. [3-6, 18, 30]). These methods involve computing an underintegrated matrix, and then adding a stabilization matrix which effectively eliminates the hourglass modes. They turn out to be fairly general and have been used for a long time in numerous codes. Whereas all of these methods based on intuitive feeling give good numerical results, their mathematical study remains often non-existent.

The most interesting challenge is to solve the problem using only the crude rank-deficient underintegrated stiffness matrix, the solution is obtained up to within an arbitrary spurious mode, and then to eliminate these modes from the solution in a post-processing operation.

Unfortunately, even when the stiffness matrix is rank-sufficient, similar oscillations are observed when underintegration is used. In that case, the process of the excitations of modes similar to the hourglass modes is not completely understood and these modes have never been

mathematically studied.

In this report, we give precise mathematical justifications and answers to the questions previously mentioned. The next Section (Section 2.2) is devoted to the proof that the Stabilization method is mathematically justified. Then, in Section 2.3, we present a method which involves solving an underintegrated and not well-posed problem, then in a-posteriori eliminating the unknown degree of freedom. The proof of the accuracy of the method is given in Section 2.4, and its numerical aspects and results are described in Section 2.5. In Section 2.6, we examine the case in which spurious oscillations cannot be predicted from the rank-deficiency of the stiffness matrix and we analyze why these modes may be excited. Finally, we apply the previous considerations to the linear elasticity problem.

It should be noted that the method and its results cannot be embedded in a classical elliptic theory: Strang's ellipticity condition [49] is here violated and this non-elliptic method cannot be studied by the classical theory of finite element methods and numerical integration [11, 12, 51, 52]. Note in these references that, both Ciarlet and Wahlbin crucially suppose the exactitude of the numerical scheme for the polynomials considered. This polynomial invariance plays a decisive role in their error estimations. We also refer to Girault [21, 22] for his approach to the same kind of problem, non-elliptic because of the use of partially underintegrated stiffness matrix, but where hourglass modes did not appear.

2.2 A-Priori Hourglass Control

2.2.1 Introduction. This section is devoted to mathematical preliminaries to several methods consisting of adding a stabilization matrix to the underintegrated matrix. For clarity, we shall confine our attention to a simple model problem. Let Ω be a regular domain in \mathbb{R}^2 with boundary $\partial\Omega$ and consider the model Neumann problem,

(P₀) Find $u = u(x,y)$ such that

$$\left. \begin{aligned} -\Delta u &= f && \text{on } \Omega \\ \frac{\partial u}{\partial n} &= 0 && \text{on } \partial\Omega \end{aligned} \right\} \quad (2.1)$$

where f is an $L^2(\Omega)$ -function satisfying

$$\int_{\Omega} f \, dx dy = 0 \quad (2.2)$$

The questions of the existence and uniqueness of solutions to (2.1) (which are well-known) are taken up in Part II.

We shall first consider a finite element approximation of (2.1) constructed using Q_1 -elements, i.e., four-node quadrilateral elements over which bilinear shape functions are used. Most of our notations and results are reproductions of those of Flanagan [18] and Belytschko [3, 5, 6]. Then we will attempt to extend our results to the Q_2 -elements (nine-node, biquadratic elements) and will indicate in which ways they differ from Belytschko's [4].

The construction of finite element approximations of (2.1) involves the calculation of the stiffness matrix K_e for a typical finite ele-

ment Ω_e , which is given by the formula,

$$\underline{K}_e = \int_{\Omega} \underline{\nabla N}^t \cdot \underline{\nabla N} \, dx dy \quad (2.3)$$

where \underline{N} is a vector representing the bilinear or biquadratic shape functions in each element Ω_e , $1 \leq e \leq E$.

When Q_1 - (respectively Q_2 -) elements are used to discretize the domain Ω , \underline{K}_e is a 4×4 matrix (resp. 9×9) and the \underline{N} 's contain four bilinear (resp. nine biquadratic) shape functions. We will distinguish exact-, full-, and under-integrations. The full integration is obtained using the number of Gauss integration points necessary to obtain the exact integration on regular square elements: 4 (resp. 9) points in our study. The underintegration will involve the Gauss rule of lower order: 1 (resp. 4) points. The stiffness matrix associated with a rule involving k points will be denoted $\underline{K}_e^{(k)}$, $k = 1, 4, 9$.

Several authors [3, 5, 6, 18, 30] proposed to add to the underintegrated stiffness matrix a stabilization matrix which exhibits several special properties. In this section, we will prove that these properties are indeed satisfied and that the exact stiffness matrix \underline{K}_e can be computed by this method. This will be accomplished by first carrying out the integration (2.3) exactly.

We first introduce some notations. Suppose that element Ω_e is defined by the coordinates of its nodes (x^I, y^I) , $1 \leq I \leq p$, $p = 4$ or 9 . We introduce the isoparametric mapping from a master element

$$\hat{\Omega} = \left[-\frac{1}{2}, +\frac{1}{2} \right] \times \left[-\frac{1}{2}, +\frac{1}{2} \right]$$

to Ω_e such that

$$\left. \begin{aligned} x &= \sum_{I=1}^P x^I N_I(\xi, \eta) \\ y &= \sum_{I=1}^P y^I N_I(\xi, \eta) \end{aligned} \right\} \quad (2.4)$$

where N_I , $1 \leq I \leq p$, are the shape functions for the quadrilateral element on the master element. The node numbering convention is shown in Figure 1.

The stiffness matrix \underline{K}_e is evaluated in (2.3) using the mapping (2.4) from $\hat{\Omega}$ to Ω_e :

$$\left. \begin{aligned} \underline{K}_e &= \int_{\Omega} \underline{\nabla N}^T \underline{\nabla N} \, dx dy \\ &= \int_{\hat{\Omega}} |J| \hat{\underline{\nabla N}}^T \begin{bmatrix} \frac{d\xi}{dx} \\ \frac{d\xi}{dx} \end{bmatrix}^T \begin{bmatrix} \frac{d\xi}{dx} \\ \frac{d\xi}{dx} \end{bmatrix} \hat{\underline{\nabla N}} \, d\xi d\eta \end{aligned} \right\} \quad (2.5)$$

where $\begin{bmatrix} \frac{d\xi}{dx} \\ \frac{d\xi}{dx} \end{bmatrix}$ is the Jacobian matrix of the mapping from $\hat{\Omega}$ to Ω_e , J is the inverse of its determinant, and where the gradients of the shape functions are derived with respect to the master element coordinates (ξ, η) . These matrices can be computed using (2.4):

$$\left. \begin{aligned} x &= \underline{x}^T \cdot \underline{N} \\ y &= \underline{y}^T \cdot \underline{N} \end{aligned} \right\} \quad (2.6)$$

$$\begin{bmatrix} \frac{dx}{d\xi} \\ \frac{dx}{d\xi} \end{bmatrix} = \begin{bmatrix} \underline{x}^T \frac{dN}{d\xi} & \underline{y}^T \frac{dN}{d\xi} \\ \underline{x}^T \frac{dN}{d\eta} & \underline{y}^T \frac{dN}{d\eta} \end{bmatrix} \quad (2.7)$$

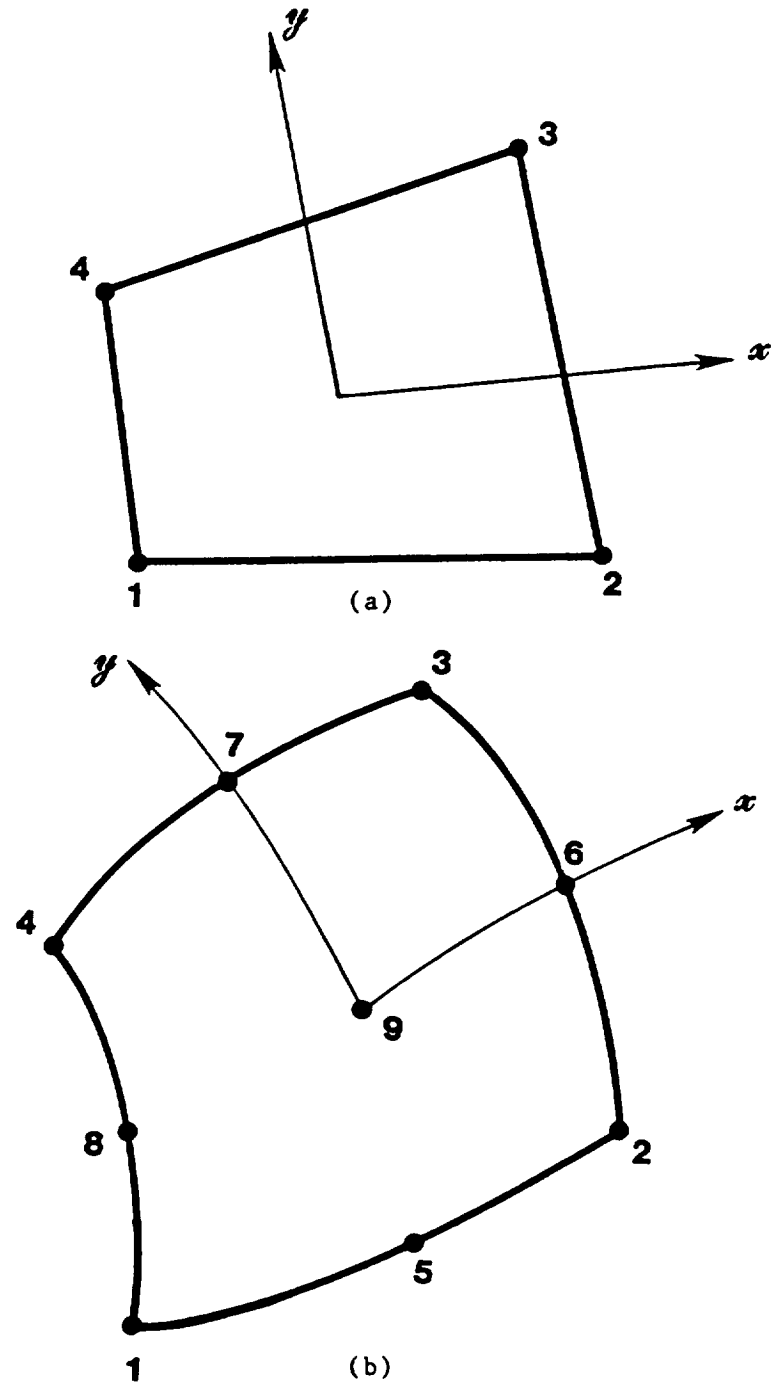


Figure 1. Node Numbering for a) 4-Node Element;
b) 9-Node Element.

$$\begin{bmatrix} \frac{d\xi}{dx} \end{bmatrix} = \begin{bmatrix} \frac{dx}{d\xi} \end{bmatrix}^{-1} = \frac{1}{J} \begin{bmatrix} y^T \frac{dN}{d\eta} & -y^T \frac{dN}{d\xi} \\ -x^T \frac{dN}{d\eta} & x^T \frac{dN}{d\xi} \end{bmatrix} \quad (2.8)$$

$$\hat{\nabla}N = \begin{bmatrix} \frac{dN}{d\xi} \\ \frac{dN}{d\eta} \end{bmatrix} \quad (2.9)$$

where J is the Jacobian of the mapping

$$J = \det \frac{dx}{d\xi} \quad (2.10)$$

Finally, we obtain the expression,

$$K_e = \int_{\hat{\Omega}} \left(\frac{\underline{\underline{Axx}}^T \underline{\underline{A}}^T}{\underline{\underline{y}}^T \underline{\underline{Ax}}} + \frac{\underline{\underline{Ayy}}^T \underline{\underline{A}}^T}{\underline{\underline{y}}^T \underline{\underline{Ax}}} \right) d\xi d\eta \quad (2.11)$$

where $\underline{\underline{A}}$ is the antisymmetric matrix

$$\underline{\underline{A}} = \frac{dN}{d\eta} \frac{dN}{d\xi}^T - \frac{dN}{d\xi} \frac{dN}{d\eta}^T \quad (2.12)$$

the Jacobian J can be expressed as $\underline{\underline{y}}^T \underline{\underline{Ax}}$.

A study of K_e expressed as in (2.11) and the properties of the matrix $\underline{\underline{A}}$ will then enable us to study the effect of the underintegration of the stiffness matrix. We will first concentrate on the 4 node element and derive the exact expression of the stabilization matrix. Then we will discuss what form this matrix may take for the 9-node element.

2.2.2. The stabilization matrix for the bilinear element. For the bilinear element, the shape function vector can be written as

$$\underline{\tilde{N}} = \frac{1}{4} \underline{\tilde{t}} - \frac{\xi}{2} \underline{\tilde{s}} + \frac{\eta}{2} \underline{\tilde{s}'} + \xi\eta \underline{\tilde{h}} \quad (2.13)$$

where

$$\begin{aligned} \underline{\tilde{g}}^T &= (1, -1, -1, 1) \\ \underline{\tilde{g}'}^T &= (1, 1, -1, -1) \\ \underline{\tilde{t}}^T &= (1, 1, 1, 1) \\ \underline{\tilde{h}}^T &= (1, -1, 1, -1) \end{aligned} \quad (2.14)$$

then the explicit form of the (4x4) matrix \underline{A} is

$$\underline{\tilde{A}} = \frac{1}{4}(\underline{\tilde{s}'}\underline{\tilde{s}}^T - \underline{\tilde{s}}\underline{\tilde{s}'}^T) + \frac{\xi}{2}(\underline{\tilde{s}}\underline{\tilde{h}}^T - \underline{\tilde{h}}\underline{\tilde{s}}^T) + \frac{\eta}{2}(\underline{\tilde{h}}\underline{\tilde{s}'}^T - \underline{\tilde{s}'}\underline{\tilde{h}}^T) \quad (2.15)$$

for $(\xi, \eta) = (0, 0)$, we obtain \underline{A}_0 which satisfies

$$\underline{\tilde{y}}^T \underline{\tilde{A}}_0 \underline{\tilde{x}} = \underline{\tilde{y}}^T \underline{\tilde{A}} \Big|_{\xi=\eta=0} \underline{\tilde{x}} = \frac{1}{2}(y_{24}x_{13} + y_{31}x_{24}) \quad (2.16)$$

which is merely the area of the element Ω_e , noted $|\Omega_e|$.

At this point, we can define precisely the matrix resulting from a 1-point rule; this underintegrated matrix, denoted by $\underline{K}_e^{(1)}$, is given by

$$\underline{K}_e^{(1)} = \frac{\underline{A}_{0xx}^T \underline{A}_0^T}{|\Omega_e|} + \frac{\underline{A}_{0yy}^T \underline{A}_0^T}{|\Omega_e|} \quad (2.17)$$

Also, if we denote $\underline{B} = (\underline{b}_1, \underline{b}_2)^T$, the discrete approximations of the gradient $\nabla \underline{\tilde{N}}$ evaluated at the integration point is given by

$$\begin{cases} \underline{b}_1 = -\underline{A}_0 \underline{\tilde{y}} \\ \underline{b}_2 = \underline{A}_0 \underline{\tilde{x}} \end{cases} \quad (2.18)$$

and, therefore, (2.17) takes the usual form

$$\underline{K}_e^{(1)} = \frac{1}{|\Omega_e|} (\underline{b}_{-1} \underline{b}_{-1}^T + \underline{b}_{-2} \underline{b}_{-2}^T) \quad (2.19)$$

The rank deficiency of $\underline{K}_e^{(1)}$ can now be verified.

Indeed, from (2.14) and (2.15) we note that

$$\left. \begin{array}{l} \underline{A}_{0\sim} h = 0 \\ \underline{A}_{0\sim} t = 0 \end{array} \right\} \quad (2.20)$$

and then

$$\left. \begin{array}{l} \underline{K}_e^1 h = 0 \\ \underline{K}_e^1 t = 0 \end{array} \right\} \quad (2.21)$$

Therefore, if we consider H and T the global hourglass and translation, and $\underline{K}^{(1)}$ the assembled underintegrated stiffness matrix, we have

$$\left. \begin{array}{l} \underline{K}^{(1)} \cdot H = \sum_e \underline{K}_e^{(1)} \cdot H = \sum_e \underline{K}_e^{(1)} \cdot h = 0 \\ \text{also } \underline{K}^{(1)} \cdot T = 0 \end{array} \right\} \quad (2.22)$$

and this proves the rank deficiency of $\underline{K}^{(1)}$. Note that this "+1" pattern is independent of the regularity of the mesh and that h will take alternating values +1 and -1 at neighbor nodes as shown in Fig. 2.

Our goal will now be to calculate a matrix $\underline{K}_e^{\text{stab}}$ such that, if added to $\underline{K}_e^{(1)}$, we obtain the exact stiffness matrix \underline{K}_e given by (2.11). This expression does not seem easy to integrate, but the image of certain vectors mapped by this matrix can be easily computed using

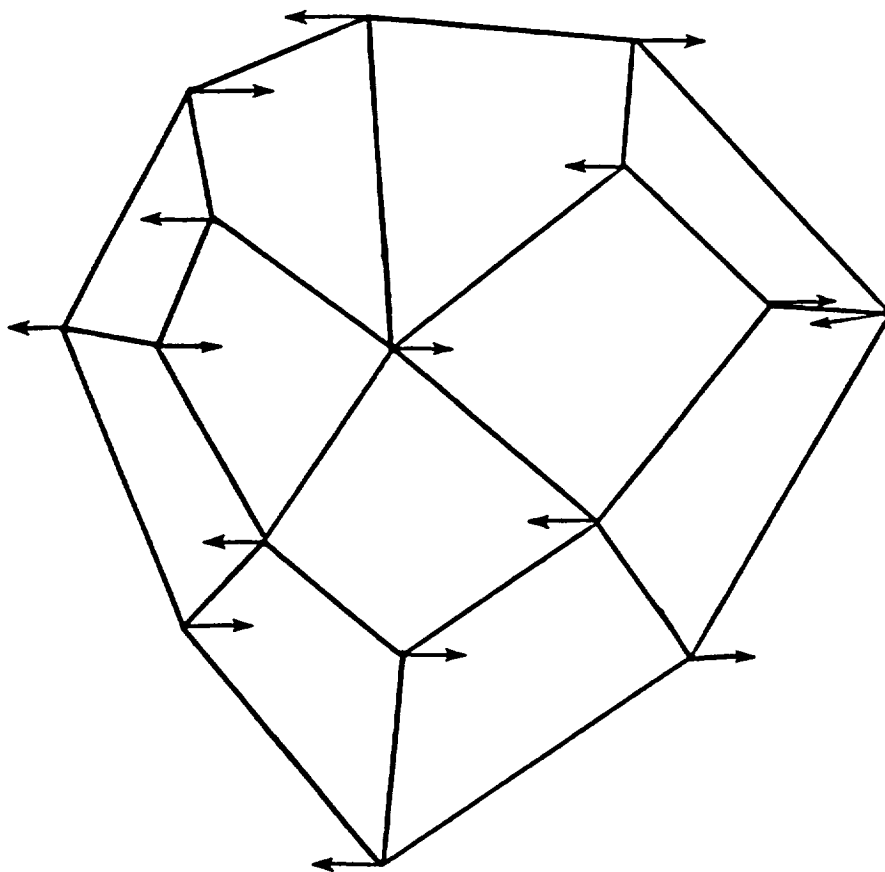


Figure 2. " ± 1 " Pattern of the Hourglass Mode in an Arbitrary Mesh.

orthogonality relations previously obtained (2.20) and the fact that

$$\begin{aligned} \tilde{b}_1^T \tilde{x} &= \tilde{b}_2^T \tilde{y} = |\Omega_e| \\ \tilde{b}_1^T \tilde{y} &= \tilde{b}_2^T \tilde{x} = 0 \end{aligned} \quad (2.23)$$

we obtain:

$$\left. \begin{aligned} \tilde{K}_e \tilde{t} &= 0 \\ \tilde{K}_e \tilde{x} &= \tilde{b}_1 \\ \tilde{K}_e \tilde{y} &= \tilde{b}_2 \end{aligned} \right\} \quad (2.24)$$

Equation (2.24) is not sufficient to compute \tilde{K}_e because it gives only 9 out of the 10 coefficients of \tilde{K}_e (4 x 4, symmetric). It is enough to know $\tilde{x}^t \tilde{K}_e \tilde{x}$, where \tilde{t} , \tilde{x} , \tilde{y} , and \tilde{h} form a set of independent vectors. That is the case for $\tilde{X} = \tilde{h}$ because

$$\det(\tilde{x}, \tilde{y}, \tilde{t}, \tilde{h}) = 4A \neq 0$$

provided the element is not singular. Then the knowledge of $\tilde{h}^T \tilde{K}_e \tilde{h}$ and the relations define uniquely \tilde{K}_e . If we set

$$\tilde{h}^T \tilde{K}_e \tilde{h} = 16 \bar{\epsilon} \quad (2.25)$$

then \tilde{K}_e is given by

$$\tilde{K}_e = \tilde{K}_e^{(1)} + \bar{\epsilon} \tilde{\gamma} \tilde{\gamma}^T \quad (2.26)$$

whereas, again, given by (2.25), $\bar{\epsilon}$ is a scalar, and

$$\tilde{\gamma} = \tilde{h} - \frac{\tilde{h}^T \tilde{x}}{|\Omega_e|} \tilde{b}_1 - \frac{\tilde{h}^T \tilde{y}}{|\Omega_e|} \tilde{b}_2 \quad (2.27)$$

While it is difficult to express $\bar{\epsilon}$ nicely as function of \underline{x} and \underline{y} , its exact value can be written

$$\bar{\epsilon} = \frac{1}{4} \int_{\hat{\Omega}} \frac{\left[\xi(\underline{s}^T \underline{x}) - \eta(\underline{s}^T \underline{x}) \right]^2 + \left[\xi(\underline{s}^T \underline{y}) - \eta(\underline{s}^T \underline{y}) \right]^2}{|\Omega_e| + \xi(y_{43}x_{12} + y_{12}x_{34}) + \eta(y_{32}x_{14} + y_{14}x_{23})} d\xi d\eta \quad (2.28)$$

We observe that for parallelogram elements, the denominator is constant and its value is the area of the domain $|\Omega_e|$. In this case

$$\bar{\epsilon} = \frac{1}{24|\Omega_e|} (x_{13}^2 + x_{24}^2 + y_{13}^2 + y_{24}^2) \quad (2.29)$$

or for rectangular elements

$$\bar{\epsilon} = \frac{\ell_x^2 + \ell_y^2}{12 \ell_x \ell_y} \quad (2.30)$$

where ℓ_x and ℓ_y are the lengths of the sides of the rectangular element. Also note that for such parallelogram elements $\underline{\gamma}$ reduces to \underline{h} .

The expression (2.26) is often used to eliminate a-priori spurious modes for the kernel of \underline{K} , but the determination of $\bar{\epsilon}$ remains a problem. The choice $\bar{\epsilon} = 0$ leads to the underintegrated matrix and to the method to be studied in the next section. On the other hand, a cheaper way than the full integration of the whole matrix would be to fully integrate $\bar{\epsilon}$ given by (2.28). This method would lead again to the full integration and is cheaper because it needs only one 4 x 4 integration per element instead of 10. A more common practice is to take for $\bar{\epsilon}$ a simple value independent of the geometry of the element, which

is often the value obtained for a square 1/6 or sometimes any arbitrary constant, as used in [5, 6].

2.2.3. The stabilization matrix for the biquadratic element. In this section, we will study the effect of the underintegration of the 9 x 9 stiffness matrix obtained in (2.11) with nine-node elements when a 4- Gauss integration point rule is used. Whereas Belytschko et al. [4] intuitively obtain another " $\underline{\gamma} \cdot \underline{\gamma}^T$ " stabilization, we prove that this decomposition is not even valid for regular meshes. We then propose a decomposition derived on regular mesh.

But first we exhibit the spurious modes out of $K_e^{(4)}$. For the biquadratic element, the shape function vector can be written as

$$\underline{N} = \underline{S} \underline{\xi} \quad (2.31)$$

where \underline{S} is the 9 x 9 matrix

$$\underline{S} = \begin{bmatrix} 0 & 0 & 0 & 1 & 0 & 0 & -2 & -2 & 4 \\ 0 & 0 & 0 & -1 & 0 & 0 & 2 & -2 & 4 \\ 0 & 0 & 0 & 1 & 0 & 0 & 2 & 2 & 4 \\ 0 & 0 & 0 & -1 & 0 & 0 & -2 & 2 & 4 \\ 0 & 0 & -1 & 0 & 0 & 2 & 0 & 4 & -8 \\ 0 & 1 & 0 & 0 & 2 & 0 & -4 & 0 & -8 \\ 0 & 0 & 1 & 0 & 0 & 2 & 0 & -4 & -8 \\ 0 & -1 & 0 & 0 & 2 & 0 & 4 & 0 & -8 \\ 1 & 0 & 0 & 0 & -4 & -4 & 0 & 0 & 16 \end{bmatrix} \quad (2.32)$$

and

$$\underline{\xi} = [1, \xi, \eta, \xi\eta, \xi^2, \eta^2, \xi\eta^2, \xi^2\eta, \xi^2\eta^2]^T \quad (2.33)$$

The integration rule we are interested in involves four integration points $(\xi^\alpha, \eta^\alpha)$, $\alpha = 1, 4$. Associated with each of them, we denote

by \underline{A}_α and J_α the corresponding matrix \underline{A} and Jacobian. The under-integrated matrix is then

$$\underline{K}_e^{(4)} = \sum_{\alpha=1}^4 \frac{\underline{A}_{\alpha xx}^T \underline{A}_\alpha + \underline{A}_{\alpha yy}^T \underline{A}_\alpha}{J_\alpha} \quad (2.34)$$

and can also be written

$$\underline{K}_e^{(4)} = \sum_{\alpha=1}^4 \frac{1}{J_\alpha} (b_{-1}^\alpha b_{-1}^{\alpha T} + b_{-2}^\alpha b_{-2}^{\alpha T}) \quad (2.35)$$

where

$$\underline{B}^\alpha = \begin{bmatrix} b_{-1}^{\alpha T} \\ b_{-2}^{\alpha T} \end{bmatrix} = \begin{bmatrix} (-A_\alpha y)^T \\ (A_\alpha x)^T \end{bmatrix} \quad (2.36)$$

generalizes (2.18) to a 4- point rule.

The rank deficiency of $\underline{K}_e^{(4)}$ can now be verified. Indeed if we call \underline{t} and \underline{h} the vectors defined by

$$\begin{aligned} \underline{t}^T &= [1, 1, 1, 1, 1, 1, 1, 1, 1] \\ \underline{h}^T &= [1, 1, 1, 1, -1, -1, -1, -1, 0] \end{aligned} \quad (2.37)$$

we easily obtain:

$$\underline{t}^T \cdot \underline{N} = \underline{t}^T \cdot \underline{S} \cdot \underline{\xi} = 1$$

and

$$\underline{h}^T \cdot \underline{N} = \underline{h}^T \cdot \underline{S} \cdot \underline{\xi} = -4(\xi^2 + \eta^2 - 12 \xi^2 \eta^2)$$

and then differentiating these expressions and using (2.12) we get:

$$\underline{A} \cdot \underline{t} = 0$$

$$\underline{A} \cdot \underline{h} = -8 \left[\eta(1-12\xi^2) \frac{dN}{d\xi} + \xi(1-12\eta^2) \frac{dN}{d\eta} \right]$$

the second expression vanishes when the point (ξ, η) is one of the four

Gauss integration points of $\hat{\Omega}$

$$(\xi^\alpha, \eta^\alpha) = \left(\pm \frac{1}{2\sqrt{3}}, \pm \frac{1}{2\sqrt{3}} \right); \alpha = 1, 4 \quad (2.38)$$

Therefore we have

$$\underline{A}_\alpha \cdot \underline{t} = \underline{A}_\alpha \cdot \underline{h} = 0 \quad \alpha = 1, 4 \quad (2.39)$$

and, consequently,

$$\underline{K}_e^{(4)} \cdot \underline{t} = \underline{K}_e^{(4)} \cdot \underline{h} = 0 \quad (2.40)$$

which proves the rank deficiency of $\underline{K}_e^{(4)}$. Once again, we note that the pattern of \underline{h} defined in (2.37) is independent of the geometry of the element and is therefore valid for a rectangular mesh as well as for irregular element meshes.

The analysis of the decomposition of underintegrated and stabilizing matrices for this element cannot be completed as completely as that for the 4-node element. However, Belytschko and co-workers [3] have intuitively arrived at a decomposition similar to (2.26) where $\underline{\gamma}$ and $\underline{\epsilon}$ are

$$\underline{\gamma} = \underline{h} - \frac{1}{4} \underline{h}^T \cdot \underline{x} \sum_{\alpha=1}^4 \frac{b_\alpha^\alpha}{J_\alpha} - \frac{1}{4} \underline{h}^T \cdot \underline{y} \sum_{\alpha=1}^4 \frac{b_\alpha^\alpha}{J_\alpha} \quad (2.41)$$

$$\underline{\epsilon} = \frac{1}{100} \sum_{\alpha=1}^4 \frac{1}{J_\alpha} \left[b_{\sim 1}^{\alpha T} \cdot b_{\sim 1}^\alpha + b_{\sim 2}^{\alpha T} \cdot b_{\sim 2}^\alpha \right] \quad (2.42)$$

This decomposition does in fact satisfy several properties also satisfied by the exact matrix,

$$\begin{aligned} \underline{K}_e \cdot \underline{t} &= 0 \\ \underline{K}_e \cdot \underline{x} &= \sum_{\alpha=1}^4 b_{\sim 1}^\alpha \end{aligned}$$

$$\underline{k}_e \cdot \underline{y} = \sum_{\alpha=1}^4 b_2^\alpha.$$

but for a simple square element*, \underline{k}_e and its decomposition (2.26) do not coincide. Indeed, for this simple geometry, the calculation of (2.11) can be carried out explicitly and the polynomial in (ξ, η) obtained can be split into one part exactly integrated with 4 Gauss points, and another part of higher order that requires 9 points. This calculation leads to the decomposition:

$$\left. \begin{aligned} K_{xx}^{(9)} &= K_{xx}^{(4)} + \Omega_e \left(\frac{1}{45} \underline{s}_7 \cdot \underline{s}_7^T + \frac{4}{135} \underline{s}_9 \cdot \underline{s}_9^T \right) \\ K_{yy}^{(9)} &= K_{xx}^{(4)} + \Omega_e \left(\frac{1}{45} \underline{s}_8 \cdot \underline{s}_8^T + \frac{4}{135} \underline{s}_9 \cdot \underline{s}_9^T \right) \end{aligned} \right\} \quad (2.43)$$

where

$$\left. \begin{aligned} K_{xx} &= \int_{\hat{\Omega}} \frac{\underline{A} \cdot \underline{y} \cdot \underline{y}^T \cdot \underline{A}^T}{\underline{y}^T \underline{A} \underline{x}} d\xi d\eta \\ K_{yy} &= \int_{\hat{\Omega}} \frac{\underline{A} \cdot \underline{x} \cdot \underline{x}^T \cdot \underline{A}^T}{\underline{y}^T \underline{A} \underline{x}} d\xi d\eta \end{aligned} \right\} \quad (2.44)$$

and $\underline{s}_7, \underline{s}_8, \underline{s}_9$ are the 7th, 8th and 9th column vectors of \underline{S} (2.32). These vectors correspond to the higher order of $\underline{\xi}$ (2.33) that cannot be exactly integrated by a 4-point rule. The form taken by the stabilization matrix involves now three matrices $(\underline{s}_i \cdot \underline{s}_i^T, i = 7, 8, 9)$, is exact for a square element and cannot coincide with the decomposition found in [4]. Finally, we note that both decompositions were used in our a-posteriori control described in Section 2.7 on a regular mesh,

* or also for a geometry for which the Jacobian is constant.

and optimal rates of convergence were only obtained with the decomposition (2.43).

2.2.4. The stabilization matrix for a general heat transfer equation. In this paragraph we give the stabilization matrix for a slightly more complicated operator. The case of the linear elasticity operator is discussed later.

Let us consider the case in which the operator is defined by

$$\underline{A} = \underline{\xi}^T \underline{C} \underline{\beta} \quad (2.45)$$

where

$$\underline{\beta} = \left(\frac{\partial}{\partial x} \quad \frac{\partial}{\partial y} \right)^T$$

and

$$\underline{C} = \begin{pmatrix} C_{11} & C_{21} \\ C_{12} & C_{22} \end{pmatrix}$$

Then the stiffness matrix associated with this operator is given by

$$\underline{K}_e = \int_{\Omega} \underline{vN}^T \cdot \underline{C} \cdot \underline{vN} \, dx dy \quad (2.46)$$

The generalization of the stabilization decomposition when Q_1 elements are used can then be written

$$\underline{K}_e = \underline{K}_e^{(1)} + \bar{\epsilon} \underline{\gamma} \cdot \underline{\gamma}^T \quad (2.47)$$

where

$$\underline{K}_e^{(1)} = \frac{1}{|\Omega_e|} \underline{B}^T \underline{C} \underline{B} \quad (2.48)$$

$$\bar{\epsilon} = C_{11} \bar{\epsilon}_{xx} + (C_{12} + C_{21}) \bar{\epsilon}_{xy} + C_{22} \bar{\epsilon}_{yy} \quad (2.49)$$

and

$$\left. \begin{aligned}
 \bar{\epsilon}_{xx} &= \frac{1}{4} \int_{\hat{\Omega}} \frac{1}{J} (\xi \underline{s}^T \cdot \underline{y} - \eta \underline{s}'^T \cdot \underline{y})^2 d\xi d\eta \\
 \bar{\epsilon}_{yy} &= \frac{1}{4} \int_{\hat{\Omega}} \frac{1}{J} (\xi \underline{s}^T \cdot \underline{x} - \eta \underline{s}'^T \cdot \underline{x})^2 d\xi d\eta \\
 \bar{\epsilon}_{xy} &= \frac{1}{4} \int_{\hat{\Omega}} \frac{1}{J} (\xi \underline{s}^T \cdot \underline{x} - \eta \underline{s}'^T \cdot \underline{x}) (\xi \underline{s}^T \cdot \underline{y} - \eta \underline{s}'^T \cdot \underline{y}) d\xi d\eta
 \end{aligned} \right\} (2.50)$$

The quantities $\underline{\gamma}$, \underline{B} and J are the ones defined previously. Expressions similar to those given in (2.29) and (2.30) can be used to simplify $\bar{\epsilon}$.

For a regular geometry, and corresponding to (2.29) and (2.30), we have

$$\left. \begin{aligned}
 \bar{\epsilon}_{xx} &= \frac{1}{24(\Omega_e)} (y_{13}^2 + y_{24}^2) \\
 \bar{\epsilon}_{yy} &= \frac{1}{24(\Omega_e)} (x_{13}^2 + x_{24}^2) \\
 \bar{\epsilon}_{xy} &= \frac{1}{24(\Omega_e)} (x_{13} y_{13} + x_{24} y_{24})
 \end{aligned} \right\} (2.51)$$

As far as the 9-node element is concerned, the decomposition can be obtained only for regular elements. First we note that

$$\underline{K}_{xy}^{(9)} = \underline{K}_{xy}^{(4)} \quad (2.52)$$

where the notations are similar to (2.43) and (2.44). Therefore, the decomposition can be written:

$$\begin{aligned}
 \underline{K}_e^{(9)} &= \underline{K}_e^{(4)} + \Omega_e C_{11} \left(\frac{1}{45} \underline{s}_7 \underline{s}_7^T + \frac{4}{135} \underline{s}_9 \underline{s}_9^T \right) \\
 &\quad + \Omega_e C_{22} \left(\frac{1}{45} \underline{s}_8 \underline{s}_8^T + \frac{4}{135} \underline{s}_9 \underline{s}_9^T \right)
 \end{aligned} \quad (2.53)$$

2.3 A-Posteriori Hourglass Control

2.3.1. Introduction and preliminaries. The basic ideas are more easily understood when demonstrated for the same simple model problem. We still focus on the model Neumann problem P_0 or its variational equivalent P .

Let Ω be a regular (e.g. Lipschitz) domain in \mathbb{R}^2 with boundary $\partial\Omega$ and let f be a given L^2 -function. Problem P_0 is then,

(P_0) Find u such that

$$\left. \begin{aligned} -\Delta u &= f \text{ in } \Omega \\ \frac{\partial u}{\partial n} &= 0 \text{ on } \partial\Omega \end{aligned} \right\} \quad (3.1)$$

where the data f satisfies the compatibility condition ,

$$\int_{\Omega} f dx = 0 \quad (3.2)$$

Later we shall put further restrictions on Ω and on f (e.g. we will need $f \in H(\Omega)$). The kernel of the governing operator $A = (-\Delta, \frac{\partial}{\partial n})$ in (3.1) is, of course, the space of constants. Thus, whenever (3.2) holds, there exists a solution to (3.1) which is unique up to an arbitrary constant.

To formulate a variational statement of problem P_0 , we introduce the spaces and inner products*,

* The elements of V (and $L^2(\Omega)/\mathbb{R}$) are cosets $[v]$ such that $u \in [v]$ implies that $u, v \in H^1(\Omega)$ (or $L^2(\Omega)$) and $v - u \in \mathbb{R}$. Throughout this paper we frequently refer to functions v in V , meaning, of course, that v is a representative function in the coset $[v]$.

$$V = H^1(\Omega) / \mathbf{R}$$

$(u, v)_1$ = an inner product on V

$$= \int_{\Omega} \nabla u \cdot \nabla v dx ; u, v \in V$$

$(f, g)_0$ = an inner product on $L^2(\Omega) / \mathbf{R}$

$$= \int_{\Omega} fg dx - \frac{1}{\text{meas } \Omega} \int_{\Omega} f dx \int_{\Omega} g dx \quad (3.3)$$

Three remarks are in order:

i) The norm $\|\cdot\|_0$ associated with the inner product $(\cdot, \cdot)_0$ is the canonical norm on the quotient space $L^2(\Omega) / \mathbf{R}$,

$$\|f\|_0 = \inf_{\lambda \in \mathbf{R}} \|f + \lambda\|_{L^2(\Omega)} \quad (3.4)$$

$$\lambda \in \mathbf{R}$$

ii) According to Temam [50], there exists a constant C_0 , depending only on Ω , such that

$$\|v\|_1 \geq C_0 \|v\|_0 \quad \forall v \in V \quad (3.5)$$

iii) For all f satisfying the compatibility condition (3.2) and any $v \in V$, we have

$$\begin{aligned} (f, v)_0 &= \int_{\Omega} f v dx \leq \|f\|_0 \|v\|_0 \\ &\leq \frac{1}{C_0} \|f\|_0 \|v\|_1 \end{aligned} \quad (3.6)$$

With these relations now established, we consider the variational statement of P_0 as problem P :

(P) Find $u \in V$ such that

$$(u, v)_1 = (f, v)_0 \quad \forall v \in V \quad (3.7)$$

We can easily verify that any solution of P_0 is a solution of P and, conversely, the solution of P satisfies the condition of P_0 in at least, a distributional sense. Moreover, since the bilinear form $(\cdot, \cdot)_1$ is continuous and coercive on V and since the linear form $(f, \cdot)_0$ is continuous on V if (and only if) f satisfies (3.2), the following result is an immediate consequence of the Lax-Milgram Theorem:

THEOREM I. *Let f satisfy (3.2).*

Then there exists one and only one solution $u \in V$ to problem P and this solution depends continuously on the data f . \square

We now consider a finite element approximation of the problem P . Let us now construct a finite element approximation of problem P . We begin by introducing a partition \mathcal{Q} of Ω into E finite elements so that

$$\Omega = \bigcup_{e=1}^E \Omega_e$$

We shall assume that Ω is such that it can be partitioned in this fashion into four-node quadrilateral elements over which bilinear shape functions are defined. Thus, if

$$Q_1(\Omega_e) = \text{space of bilinear functions defined on } \Omega_e$$

we can introduce the finite-dimensional space

$$V^h = \left\{ v^h \in C^0(\Omega) \text{ such that } v^h|_{\Omega_e} \in Q_1(\Omega_e), 1 \leq e \leq E / R \subset V \right\} \quad (3.8)$$

wherein, as usual, the label h is the mesh parameter (e.g. ,
 $h = \max_{1 \leq e \leq E} \text{dia}(\Omega_e)$). The functions in V^h are continuous and are still
 defined up to an arbitrary constant.

Our finite-element approximation of problem P is embodied in the
 discrete problem,

$$(P_h) \quad \text{Find } u^h \in V^h \text{ such that} \tag{3.9}$$

$$(u^h, v^h)_1 = (f, v^h)_0 \quad \forall v^h \in V^h$$

where, again, f satisfies condition (3.2).

In analogy with Theorem V, we have:

THEOREM II. *Let f satisfy (3.2). Then there is one and only
 one solution u^h to problem P_h in V^h and this solution depends
 continuously on the data f . \square*

In examining the convergence of such finite element approximations,
 we shall confine our attention throughout this section to regular mesh
 refinements. In such cases, we have the a priori asymptotic error
 estimates,

$$\|u - u^h\|_1 = O(h) , \quad \|u - u^h\|_0 = O(h^2) \tag{3.10}$$

2.3.2 The underintegrated problem. We now focus our attention
 on finite element approximations of problem P in which incomplete
 quadratures are used to evaluate the bilinear form $(\cdot, \cdot)_1$. To simpli-
 fy this study, we shall now introduce some additional assumptions:

- 1) Ω is the unit square,

$$\Omega = (0,1) \times (0,1)$$

ii) The finite elements are the squares,

$$\Omega_{ij} = \left(\frac{i-1}{N}, \frac{i}{N}\right) \times \left(\frac{j-1}{N}, \frac{j}{N}\right)$$

$$1 \leq i, j \leq N$$

$$\bar{\Omega} = \bigcup_{1 \leq i, j \leq N} \bar{\Omega}_{ij}$$

iii) The data f is L^2 -integrable; e.g.

$$f \in L^2(\Omega) \quad (3.11)$$

In this case, we take

$$h = \frac{1}{N}, \dim V^h = (N+1)^2 - 1 = O(h^{-2}) \quad (3.12)$$

In P_h we can replace f by f^h , its L^2 -projection on V^h is defined by

$$(f^h, v^h)_0 = (f, v^h)_0 \quad \forall v^h \in V^h \quad (3.13)$$

For further use, we note that the projection satisfies

$$\|f^h\|_0 \leq \|f\|_0 \quad (3.14)$$

and can be chosen such that

$$\int_{\Omega} f^h dx = 0 \quad (3.15)$$

Now we turn to the issue of numerical integration of the stiffnesses. Let $I(\cdot, \cdot)$ denote a discrete inner product on $C^0(\Omega)$ defined by a numerical quadrature rule as follows:

$$I_G(f, g) = \sum_{e=1}^E I_e^G(f, g) \quad (3.16)$$

$$I_e^G(f, g) = \sum_{j=1}^G W_j^e f(\xi_j^e) g(\xi_j^e)$$

Here W_j^e are the quadrature weights and ξ_j^e are the quadrature points for element e and G is the number of quadrature points used.

Assuming that Gaussian quadrature is used, the choice $G=4$ (2x2 - Gauss rule) leads to an exact integration of the stiffnesses for each element:

$$(u^h, v^h)_1 = I_4(u^h, v^h) = \underline{u}^T \underline{K} \underline{v} \quad (3.17)$$

for any $u^h, v^h \in V^h$. Here \underline{K} is the fully-integrated stiffness matrix and \underline{u} and \underline{v} are vectors of nodal degrees of freedom of u^h and v^h , respectively.

Instead of the correct bilinear form in (3.18), we wish to consider an underintegrated approximation to $(\cdot, \cdot)_1$ in which only one integration point per element is used:*

$$(u^h, v^h)_{1,h} = I_1(u^h, v^h) = \underline{u}^T \underline{K}^{(1)} \underline{v}$$

$$\forall u^h, v^h \in V^h \quad (3.18)$$

Here $\underline{K}^{(1)}$ is the underintegrated stiffness matrix. The difference between $(\cdot, \cdot)_1$ and $(\cdot, \cdot)_{1,h}$ (on V^h) is denoted $a'(\cdot, \cdot)$ and the corresponding stiffness matrix is $\underline{K}^{\text{stab}}$:**

* Recall Section (2.2).

** Recall that $\underline{K}^{\text{stab}} = \bar{\epsilon} \underline{\gamma} \underline{\gamma}^T$ where $\bar{\epsilon} = 1/6$ for a rectangular mesh and $\underline{\gamma}$ is given by (2.27).

$$\begin{aligned}
a'(u^h, v^h) &= (u^h, v^h)_1 - (u^h, v^h)_{1,h} \\
&= \underline{u}^T \underline{K}^{\text{stab}} \underline{v} \quad \forall u^h, v^h \in V^h
\end{aligned} \tag{3.19}$$

The "underintegrated problem",

$$\begin{aligned}
(P_h^*) \quad &\text{Find } u^h \in V^h \text{ such that} \\
&(u^h, v^h)_{1,h} = (f^h, v^h)_0 \quad \forall v^h \in V^h
\end{aligned} \tag{3.20}$$

is, in general, meaningless. This problem, in general, has no solution except for the special case in which f^h is orthogonal to the one-dimensional space of hourglass modes,

$$H = \{H \in V^h \mid (H, v^h)_{1,h} = 0 \quad \forall v^h \in V^h\} \tag{3.21}$$

A way to overcome this difficulty is to note that the underintegration of the righthand side also leads to a rank-deficient linear form

$(\cdot, \cdot)_{0,h}$:

$$\begin{aligned}
(f^h, H)_{0,h} = 0 \quad , \quad &\forall f^h \in V^h \\
&\forall H \in H
\end{aligned}$$

Note that if f^h satisfies (3.15) we also have

$$(f^h, 1)_{0,h} = 0 \tag{3.23}$$

Therefore we now consider the underintegrated problem \bar{P}_h :

$$\begin{aligned}
(\bar{P}_h) \quad &\text{Find } \bar{u}^h \in \bar{V}^h \text{ such that} \\
&(\bar{u}^h, v^h)_{1,h} = (f^h, v^h)_{0,h} \quad \forall v^h \in \bar{V}^h
\end{aligned} \tag{3.24}$$

where

$$\bar{v}^h : v^h/H \quad (3.25)$$

We can now state and prove

THEOREM III. *There exists one and only one solution \bar{u}^h to \bar{P}_h .*

Proof: This is an immediate consequence of the Lax-Milgram theorem. Since

$$(u^h, u^h)_{1,h} = 0 \iff u^h = \gamma_1 H + \gamma_2$$

we can consider $(\cdot, \cdot)_{1,h}^{1/2}$ as a norm on \bar{v}^h . It is therefore coercive and continuous on \bar{v}^h . As far as the continuity of the righthand side is concerned, a simple calculation shows that for any v^h in V^h we have

$$|(f^h, v^h)_{0,h}| \leq \|f^h\|_0 \|v^h\|_0$$

Also for any constants γ_1 and γ_2

$$|(f^h, v^h + \gamma_1 + \gamma_2 H)_{0,h}| = |(f^h, v^h)_{0,h}|$$

therefore

$$\begin{aligned} |(f^h, v^h)_{0,h}| &\leq \|f^h\|_0 \|v^h + \gamma_1 + \gamma_2 H\|_0 \quad \forall \gamma_1, \gamma_2 \\ &\leq \|f^h\|_0 \|v^h + \gamma_2 H\|_1 \quad \forall \gamma_2 \\ &\leq \alpha_h \|f^h\|_0 \|v^h\|_{1,h} \end{aligned}$$

Here we successively used (2.23), (2.22), (3.6) and the equivalence between the canonical norm of \bar{v}^h and the norm $\|\cdot\|_{1,h}$ \square

We have obtained a solution to the underintegrated problem \bar{P}_h .

This solution is unique in \bar{V}^h , from a computational point of view it is defined up to within an arbitrary hourglass mode. We now need a projection to obtain a reasonable solution from any representative \bar{u}^h chosen.

2.3.3. Projection of the underintegrated solution. In order to construct this projection, we remark that, since u^h is a solution of P_h and since $H \in V^h$, u^h satisfies

$$(u^h, H)_1 = (f^h, H)_0 \quad (3.26)$$

We wish to extend \bar{u}^h to all of V^h so that a new function $\tilde{u}^h \in V^h$ is obtained which contains an hourglass mode and which also satisfies (4.8). Thus, if π is an operator from \bar{V}^h into V^h , we define

$$\tilde{u}^h = \pi \bar{u}^h = \bar{u}^h + \lambda_0 H, \quad \lambda_0 \in \mathbb{R} \quad (3.27)$$

$$(\tilde{u}^h, H)_1 = (f^h, H)_0$$

This latter requirement determines λ_0 uniquely as

$$\lambda_0 = \frac{1}{\|H\|_1^2} \left[(f^h, H)_0 - (\bar{u}^h, H)_1 \right] \quad (3.29)$$

so that \tilde{u}^h is uniquely determined as the function

$$\tilde{u}^h = \bar{u}^h + \frac{(f^h, H)_0}{\|H\|_1^2} H - \frac{(\bar{u}^h, H)_1}{\|H\|_1^2} H$$

It is instructive to consider a geometrical interpretation of our projection defined in (4.9). Note that the "component" of the fully integrated solution u^h orthogonal (in V^h) to H^\perp is $(u^h, H)_1 = (f^h, H)_0$,

as indicated in Fig. 3. The solutions \bar{u}_h of \bar{P}_h constitute the vectors generating a line "parallel to" the space H in the figure. The projection \tilde{u}^h is then the vector defined by the orthogonal projection of u^h onto this line. Indeed, by construction,

$$(\tilde{u}^h - u^h, H)_1 = 0$$

At this point, we have established the following procedure for processing an underintegrated finite element approximation of problem P.

- i) Compute the underintegrated bilinear and linear forms $(\cdot, \cdot)_{1,h}$ and $(f^h, \cdot)_{0,h}$
- ii) Solve problem \bar{P}_h for \bar{u}^h
- iii) Compute $(\bar{u}^h, H)_1$
- iv) Construct the enhanced solution \tilde{u}^h using (3.30).

Thus, this procedure involves the computation of an underintegrated solution \bar{u}_h to a reduced problem \bar{P}_h and its enrichment via a post-processing operation to obtain a new approximation \tilde{u}^h . We shall now show that these post-processed solutions \tilde{u}^h converge to the exact solution u of problem P as the mesh is refined, and, remarkably, these approximations converge at precisely the same rate as the fully-integrated solution!

Indeed we have:

THEOREM IV: Let u , u^h and \bar{u}^h be the solutions of P, P_h and \bar{P}_h , let f be in $L^2(\Omega)$ and satisfy (3.2). Let \tilde{u}^h be ob-

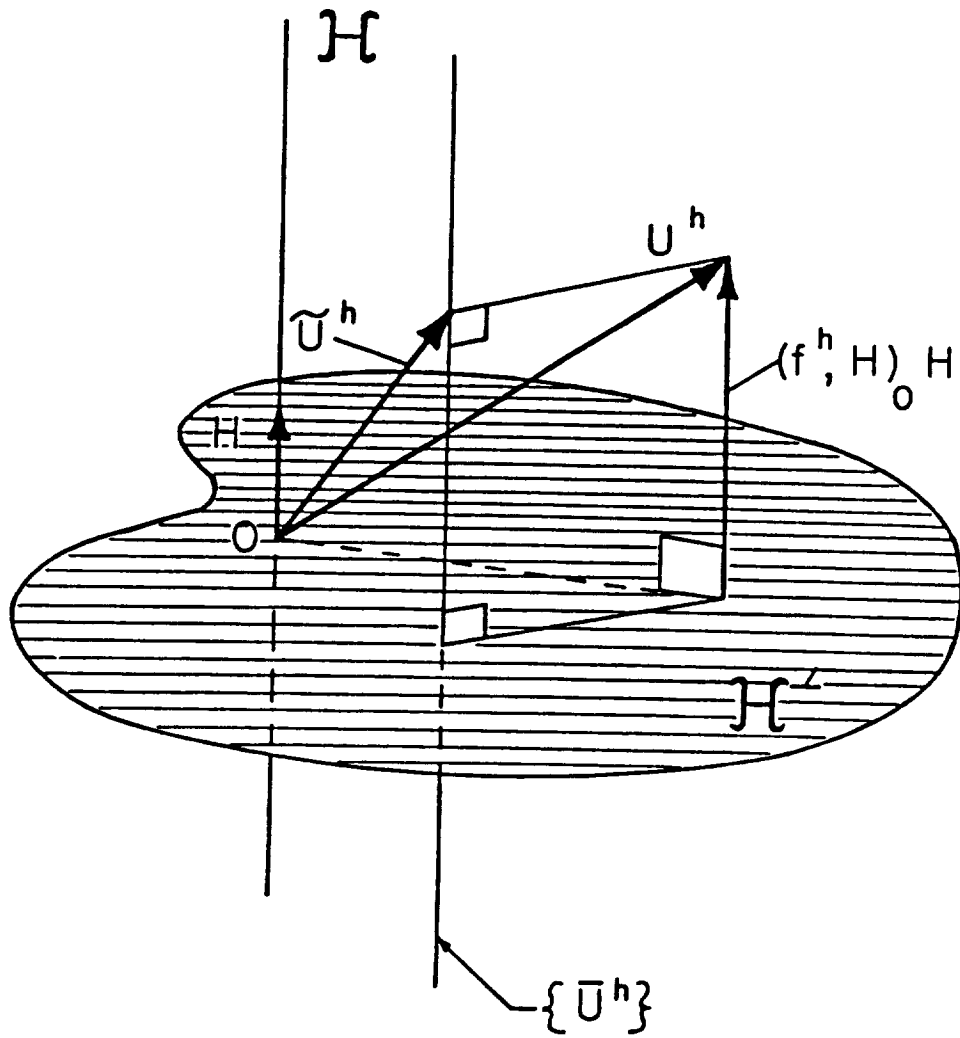


Figure 3. Geometrical Interpretation of the Projection
 $\tilde{u}^h = \pi u^h$.

tained by the projection of \bar{u}^h defined in (3.30). Then we have the following error estimates for $s = 0$ and 1

$$\|u^h - \tilde{u}^h\|_s \leq C_1 h^{2-s} \|f\|_0 \quad (3.31)$$

and

$$\|u - \tilde{u}^h\|_s \leq C_1 h^{2-s} \|f\|_0 \quad (3.32)$$

The next section will prove this theorem.

2.4. Convergence of the A-Posteriori Control

2.4.1 Introduction. This section is devoted to the proof of Theorem VIII. The method of proof relies on the tensor properties of the bilinear element and of the Gauss integration rules. The problems P_h and \bar{P}_h will be explicitly solved using an orthonormal basis of eigenvectors of $(\cdot, \cdot)_1$, $(\cdot, \cdot)_{1,h}$ and $(\cdot, \cdot)_{0,h}$. Then we note that for a regular domain and mesh, $f \in L^2(\Omega)$ implies $u \in H^2(\Omega)$ and that

$$\|u - u^h\|_1 < Ch \|f\|_0 \quad (4.1)$$

Likewise, the Aubin-Nitsche method provides also

$$\|u - u^h\|_0 \leq C'h^2 \|f\|_0 \quad (4.2)$$

By the triangle inequality,

$$\|u - \tilde{u}^h\|_1 \leq Ch \|f\|_0 + \|u^h - \tilde{u}^h\|_1 \quad (4.3)$$

with a similar estimate in the $\|\cdot\|_0$ -norm.

Thus, it suffices to estimate the relative error

$$e^h = \tilde{u}^h - u^h \quad (4.4)$$

The L^2 - and H^1 -norms of this relative error will be explicitly calculated and estimated.

2.4.2. Some one-dimensional results. For reasons to be made clear in the next subsection, it is convenient to review briefly some results on one-dimensional piecewise-linear approximations on a uniform mesh for $\Omega = (0,1)$. Our aim here is to establish concrete relationships between various bilinear forms $(\cdot, \cdot)_{0,h}$, $(\cdot, \cdot)_1$, and $(\cdot, \cdot)_{1,h}$ on spaces of piecewise-linear functions.

Let $\underline{D}(k, \alpha)$ and \underline{I} denote the $N+1$ -order matrices

$$\underline{D}(k, \alpha) = \begin{bmatrix} k & \alpha & \cdot & \cdot & \cdot & 0 & 0 \\ \alpha & 2k & \cdot & \cdot & \cdot & 0 & 0 \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ 0 & 0 & \cdot & \cdot & \cdot & 2k & \alpha \\ 0 & 0 & \cdot & \cdot & \cdot & \alpha & k \end{bmatrix}, \underline{I}' = \begin{bmatrix} 1 & 0 & \cdot & \cdot & \cdot & 0 & 0 \\ 0 & 2 & \cdot & \cdot & \cdot & 0 & 0 \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ 0 & 0 & \cdot & \cdot & \cdot & 2 & 0 \\ 0 & 0 & \cdot & \cdot & \cdot & 0 & 1 \end{bmatrix} \quad (4.5)$$

(i.e. $\underline{I}' = \underline{D}(1,0)$). Then, for $\alpha \neq 0$, one can show that

$$\begin{aligned} \det \underline{D}(k, \alpha) &= (-\alpha)^{N+1} \det \underline{D}\left(-\frac{k}{\alpha}, -1\right) \\ &= (-\alpha)^{N+1} \det \underline{D}\left(-\frac{k}{\alpha}\right) \end{aligned} \quad (4.6)$$

where

$$\underline{D}(k) \stackrel{\text{def}}{=} \underline{D}(k, -1) \quad (4.7)$$

The values of k for which $\det \underline{D}(k)$ vanishes are

$$k_i = \cos \frac{i\pi}{N}, \quad 0 \leq i \leq N \quad (4.8)$$

and the corresponding vectors $(D(k_1)v_1 = 0)$ are

$$v_1 = \left\{ \cos \frac{j\pi}{N} \right\}, \quad 0 \leq j \leq N \quad (4.9)$$

The significance of the above matrices is that in one dimension, the discrete $H^1(0,1)$ - , $L^2(0,1)$ - and underintegrated $L^2(0,1)$ -norms, on the space V_1^h of piecewise linear C_N^0 -functions on a uniform mesh of N elements on $(0,1)$,

$$V_1^h = \{v^h \in C^0(0,1) \mid v^h \text{ is linear on } [eh, (e+1)h], e=0, \dots, N-1\} \quad (4.10)$$

are associated with the matrices

$$A_{\sim 0} = \frac{h}{3} D(1, \frac{1}{2}), \quad A_{\sim 1} = \frac{1}{h} D(1, -1) \text{ and } A_{\sim 0, h} = \frac{h}{4} D(1, 1) \quad (4.11)$$

respectively. In other words,

$$\|v^h\|_s^2 = v_{\sim s}^T A v_{\sim s} \quad s = 0, 1, (0, h) \quad (4.12)$$

where v_{\sim} is the vector of nodal values of v^h .

By using (4.6) through (4.8), one can verify that the numbers α_i and β_i which render $A_{\sim 0, h} - \alpha_i A_{\sim 0}$ and $A_{\sim 1} - \beta_i A_{\sim 0}$ singular are

$$\alpha_i = \frac{3(1 + \cos \frac{i\pi}{N})}{2(2 + \cos \frac{i\pi}{N})} \quad (4.13)$$

$$\beta_i = \frac{6}{h^2} \frac{1 - \cos \frac{i\pi}{N}}{2 + \cos \frac{i\pi}{N}} \quad (4.14)$$

In particular, let $\phi^1 = \phi^1(x)$, $x \in [0, 1]$ denote the piecewise

linear functions associated with the vectors v_i :

$$\left. \begin{aligned} \phi^i(jh) &= \cos \frac{ij\pi}{N}, \quad 0 \leq i, j \leq N \\ \text{span} \{ \phi^i \}_{0 \leq i \leq N} &= V_1^h \end{aligned} \right\} \quad (4.15)$$

Then,

$$(v^h, \phi^i)_{0,h} = \alpha_i (v^h, \phi^i)_0 \quad (4.16)$$

$$(v^h, \phi^i)_1 = \beta_i (v^h, \phi^i)_0 \quad (4.17)$$

$$v^h \in V_1^h$$

Notice that the base functions ϕ^i are orthogonal for each of the scalar products under consideration.

The following remarks are in order:

- i) The denominators in (4.13) and (4.14) are non-zero.
- ii) For $i = N$, $\alpha_i = 0$ and the corresponding eigenfunction is the one-dimensional hourglass mode:

$$(1, -1, 1, -1, \dots)$$
- iii) For $i = 0$, $\beta_i = 0$ and the corresponding eigenfunction is constant. Then we have the condition $(v^h, 1)_1 = 0$ as expected.

2.4.3. Discrete norms for two dimensional meshes. The extension of the above results to two-dimensional rectangular meshes is straightforward. Since the bilinear basis functions for V^h are tensor products of piecewise linear functions of one variable, we can define

$$\begin{aligned} \phi^{ij}(x,y) &= \phi^i(x)\phi^j(y) \\ 0 \leq i, j &\leq N \end{aligned} \quad (4.18)$$

Further, let us normalize these basis functions so that

$$\|\phi^{ij}\|_0 = 1$$

We can then establish the following:

Lemma 1.1. For $v^h \in V^h$, we have

$$(v^h, \phi^{ij})_{0,h} = \alpha_i \alpha_j (v^h, \phi^{ij})_0 \quad (4.19)$$

$$(v^h, \phi^{ij})_1 = (\beta_i + \beta_j) (v^h, \phi^{ij})_0 \quad (4.20)$$

$$(v^h, \phi^{ij})_{1,h} = (\alpha_j \beta_i + \alpha_i \beta_j) (v^h, \phi^{ij})_0 \quad (4.21)$$

Moreover, if arbitrary $v^h \in V^h$ is expressed in the form,

$$\left. \begin{aligned} v^h &= \sum_{0 \leq i, j \leq N} v_{ij} \phi^{ij} \\ v_{ij} &= (v^h, \phi^{ij})_0 \end{aligned} \right\} \quad (4.22)$$

then

$$\|v^h\|_0^2 = \sum_{0 \leq i, j \leq N} v_{ij}^2 \quad (4.23)$$

$$\|v^h\|_1^2 = \sum_{0 \leq i, j \leq N} (\beta_i + \beta_j) v_{ij}^2 \quad (4.24)$$

Proof: First note that

$$\begin{aligned} (v^h, \phi^{ij})_{0,h} &= I(v^h \phi^i(x) \phi^j(y)) \\ &= \alpha_i \alpha_j \int_{\Omega} v^h \phi^i(x) \phi^j(y) \, dx dy \\ &= \alpha_i \alpha_j (v^h, \phi^{ij})_0 \end{aligned}$$

We also have

$$\begin{aligned}
 (v^h, \phi^{ij})_1 &= \int_{\Omega} \frac{\partial v^h}{\partial x} \phi^{i'}(x) \phi^j(y) \\
 &\quad + \frac{\partial v^h}{\partial y} \phi^{j'}(y) \phi^i(x) \, dx dy \\
 &= \beta_i \int_{\Omega} v^h \phi^i(x) \phi^j(y) \, dx dy \\
 &\quad + \beta_j \int_{\Omega} v^h \phi^i(x) \phi^j(y) \, dx dy \\
 &= (\beta_i + \beta_j) (v^h, \phi^{ij})_0
 \end{aligned}$$

Finally

$$\begin{aligned}
 (v^h, \phi^{ij})_{1,h} &= I_1 \left(\frac{\partial v^h}{\partial x} \phi^{i'} \phi^j + \frac{\partial v^h}{\partial y} \phi^i \phi^{j'} \right) \\
 &= \beta_i \alpha_j (v^h, \phi^{ij})_0 + \beta_j \alpha_i (v^h, \phi^{ij})_0
 \end{aligned}$$

The norms (4.23) and (4.24) are then directly obtained \square

In analogy with our remarks on the one-dimensional case, we observe that for $i = j = N$, $\phi^{ij} = H$, the two dimensional hourglass mode. Then

$$(v^h, H)_1 = \frac{24}{h^2} (v^h, H)_0 \quad (4.25)$$

and

$$(v^h, H)_{0,h} = 0 \quad (4.26)$$

Also, for $i = j = 0$, $\phi^{ij} = 1$ and the equilibrium condition (3.2) can be written

$$f_{00} = 0 \quad (4.27)$$

with

$$f_{ij} = (f^h, \phi^{ij})_0 \quad (4.28)$$

2.4.4. Explicit resolution of P_h and $(\bar{P}_h + \pi)$. With the above results in hand, let us now return to the fully-integrated finite-element approximate problem P_h given in (3.9). The solution u^h to that problem can be written

$$\left. \begin{aligned} u^h &= \sum_{0 \leq i, j \leq N} u_{ij} \phi^{ij} \\ u_{ij} &= (u^h, \phi^{ij})_0 \end{aligned} \right\} \quad (4.29)$$

and since for $[v^h = \phi^{ij} \text{ in (3.9)}]$,

$$\begin{aligned} (u^h, \phi^{ij})_1 &= (\beta_i + \beta_j) (u^h, \phi^{ij})_0 \\ &= (f^h, \phi^{ij})_0 = f_{ij} \end{aligned}$$

we have

$$u_{ij} = \frac{1}{\beta_i + \beta_j} f_{ij} ; (i, j) \neq (0, 0) \quad (4.30)$$

Using constructions similar to those in (4.29) for the fully-integrated problem, we easily verify that the solution \bar{u}^h to the underintegrated problem \bar{P}_h is representable in the form,

$$\bar{u}^h = \sum_{\substack{(i, j) \neq (N, N) \\ (i, j) \neq (0, 0)}} \bar{u}_{ij} \phi^{ij} \quad (4.31)$$

with

$$\bar{u}_{ij} = \frac{\alpha_i \alpha_j}{\alpha_i \beta_j + \alpha_j \beta_i} f_{ij} ; (i, j) \neq (0, 0) \text{ and } (N, N) \quad (4.32)$$

The cases $(i,j) = (N,N)$ and $(i,j) = (0,0)$ correspond to the arbitrary hourglass mode and arbitrary constant, respectively.

The projected approximation \tilde{u}^h defined by $\tilde{u}^h = \bar{\pi}u^h$ is constructed so that projections of \tilde{u}^h and u^h coincide; i.e.

$$\left. \begin{aligned} \tilde{u}_{ij} &= \bar{u}_{ij} & (i,j) \neq (0,0) \text{ and } (N,N) \\ \tilde{u}_{N,N} &= \bar{u}_{N,N} \end{aligned} \right\} \quad (4.33)$$

2.4.5. Proof of theorem IV. Since the error function $e^h = u^h - \tilde{u}^h$ is in V^h , we use (4.29) and (4.31) to obtain

$$e^h = \sum_{\substack{(i,j) \neq (N,N) \\ (i,j) \neq (0,0)}} e_{ij} \phi^{ij} \quad (4.34)$$

where

$$\begin{aligned} e_{ij} &= (e^h, \phi^{ij})_0 \\ &= (\tilde{u}^h, \phi^{ij})_0 - (u^h, \phi^{ij})_0 \\ &= \tilde{u}_{ij} - u_{ij} \end{aligned}$$

Thus, from (4.30) and (4.33),

$$e_{ij} = \left(\frac{\alpha_i \alpha_j}{\alpha_i \beta_j + \alpha_j \beta_i} - \frac{1}{\beta_i + \beta_j} \right) f_{ij} \quad (4.35)$$

Then, using (4.13) and (4.14), e_{ij} can be written as

$$e_{ij} = h^2 K_{ij} f_{ij} \quad (4.36)$$

where

$$K_{ij} = K\left(\cos \frac{i\pi}{N}, \cos \frac{j\pi}{N}\right) \quad (4.37)$$

and

$$K(x,y) = \frac{1}{4} \frac{(1+x)(1+y)}{(1+x)(1-y)+(1+y)(1-x)} - \frac{1}{6} \frac{(2+x)(2+y)}{(2+x)(1-y)+(2+y)(1-x)} \quad (4.38)$$

On the square $S = [-1,+1] \times [-1,+1] / \{(-1,-1), (1,1)\}$, $K(\cdot, \cdot)$ is bounded and there exists a positive constant K such that

$$|K(x,y)| \leq K \quad \forall (x,y) \in S \quad (4.39)$$

Therefore we have

$$|K_{ij}| \leq K \quad \forall (i,j) \neq (0,0) \text{ and } (N,N) \quad (4.40)$$

and we can obtain using (3.13) and (4.23)

$$\begin{aligned} \|e^h\|_0^2 &= h^4 \sum_{\substack{(i,j) \neq (0,0) \\ (i,j) \neq (N,N)}} K_{ij}^2 f_{ij}^2 \\ &\leq h^4 K^2 \|f^h\|_0^2 \leq h^4 K^2 \|f\|_0^2 \end{aligned}$$

Also, after calculation and use of (4.24), (4.14) and (3.14), we have

$$\begin{aligned} \|e^h\|_1^2 &= h^4 \sum_{\substack{(i,j) \neq (0,0) \\ (i,j) \neq (N,N)}} (\beta_i + \beta_j) K_{ij}^2 f_{ij}^2 \\ &\leq 12 h^2 K^2 \|f\|_0^2 \quad \blacksquare \end{aligned}$$

2.5. Implementation and Numerical Results of the A-Posteriori Control For the Laplace Equation.

In this section we first would like to indicate how the a-posteriori

control method is implemented, and how its time efficiency compares to the a-priori method. Then several numerical results will be given, illustrating the accuracy of the method and confirming the results obtained in the previous sections.

2.5.1. Implementation of the a-posteriori method. First let us indicate that from a mathematical point of view the problem \bar{P}^h is well-posed but computationally, the matrix obtained from this formulation is singular and the dimension of its kernel is 2. Consequently, we must pick two nodes, fix them a value, and solve. The first value fixes the constant mode, and the second one fixes the hourglass mode to be eliminated later. Let us fix \bar{u}^h equal to zero at the origin and at the next point on the boundary (coordinates : $h,0$) (Figure 4.a). According to the error estimates (3.30) and (3.31), we may write

$$u^h = \bar{u}^h + \lambda H + O(h^{2-s}) \quad (5.1)$$

and therefore, if we normalize H such that its nodal values are 0 or 1, λ measures precisely the value of u^h at $(h,0)$ (Figure 4.b), and approaches $u(h,0)$

$$\lambda = u(h,0) + O(h^\sigma) \quad (5.2)$$

But $u(h,0)$ is $O(h^2)$ for a smooth enough solution ($u(0,0) = 0$, $\partial u / \partial \eta(0,0) = 0$) and using L^∞ -estimates [11], σ can be evaluated to $2-\epsilon$, ϵ arbitrary. Finally, we have the estimate

$$\lambda = O(h^{2-\epsilon}), \quad \epsilon \text{ arbitrary} \quad (5.3)$$

Also, the choice of H leads to

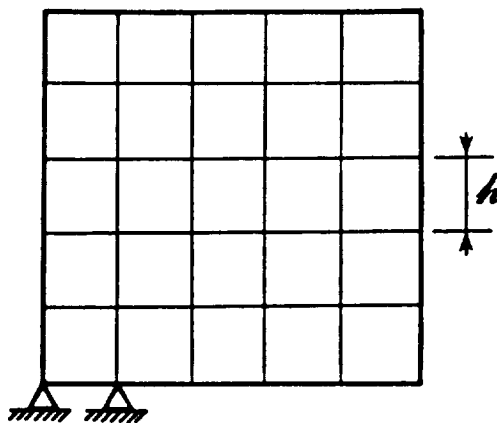


Figure 4.a. Two Fixed Nodes

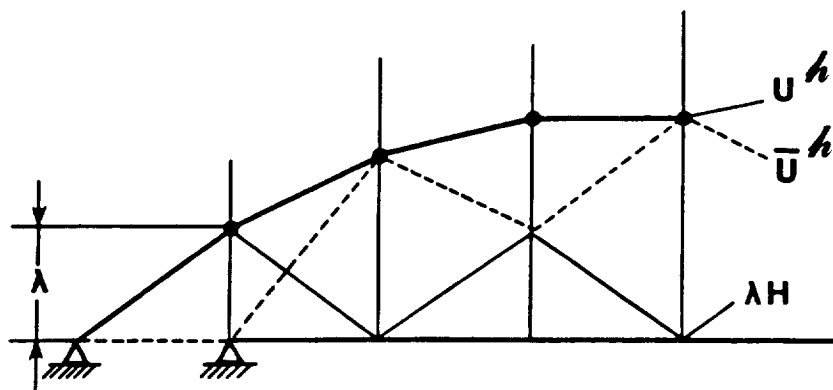
Figure 4.b. $\underline{u}^h = \underline{\bar{u}}^h + \underline{\lambda H} + \underline{O}(h^{2-s})$

Figure 4. Justification of the Omission of the Projection.

$$\begin{aligned} \|H\|_0 &= 1/3 \\ \|H\|_1 &= \frac{2\sqrt{2}}{\sqrt{3}h} \end{aligned} \quad (5.4)$$

and therefore we obtain

$$\|\lambda H\|_s = O(h^{2-s-\epsilon}) \quad s = 0,1 \quad (5.5)$$

and that proves that the post processor contribution λH can be neglected if the fixed nodes are chosen as indicated for this type of boundary condition. The error estimates of Theorem IV still hold up to within $h^{-\epsilon}$.

Unfortunately this remark has two major drawbacks: it supposes that u is smooth ($u \in H^2(\Omega)$) and it is not valid to 9-node elements that will later be discussed.

Before discussing the implementation of (3.30), we indicate that this projection can be simplified. Indeed, taking $v^h = H$ in (4.25) we obtain

$$\left\| \frac{(f_1^h H)_0}{\|H\|_1^2} H \right\| \leq \|f^h\|_0 \frac{\|H\|_0^2}{\|H\|_1^2} = \frac{h^2}{24} \|f^h\|_0 \quad (5.6)$$

$$\left\| \frac{(f_1^h H)_0}{\|H\|_1^2} H \right\|_1 \leq \|f^h\|_0 \frac{\|H\|_0}{\|H\|_1} = \frac{h}{2\sqrt{6}} \|f^h\|_0$$

Therefore we have

$$\left\| \frac{(f_1^h H)_0}{\|H\|_1^2} H \right\|_s \leq C \|f^h\|_0 h^{2-s}, \quad s = 0,1 \quad (5.8)$$

and this term can be neglected without affecting the estimate of Theorem

VIII. The formula used in the post processor is then

$$\tilde{u}^h = \bar{u}^h - \lambda H \quad (5.9)$$

$$\lambda = (\bar{u}^h, H)_1 \|H\|_1^{-2} \quad (5.10)$$

In order to preserve the efficiency of the method provided by underintegration, one must find an efficient way to compute the parameter λ in (5.9). One way that suggests itself is to calculate the H^1 inner products of (5.10) using numerical integration. The use of a one point rule would be absurd and would lead to a ratio 0/0. The use of a 4 Gauss point rule has been numerically implemented and gives good results (similar to those to be presented next) but cost of this integration is expensive, as shown in Table 3. This method is therefore rejected.

We shall now describe a more efficient method with related numerical results shown in the next subsection. This method relies on the fact that, for the bilinear element, the stiffness matrix can be decomposed into two parts, one of which contains H in its kernel. The other part is such that the image of H is cheap to calculate.

This decomposition proved in Section 2.2.2 can be written as

$$K_{\approx e}^{\text{exact}} = K_{\approx e}^{\text{under}} + \bar{\epsilon}_e \gamma_{\approx e} \cdot \gamma_{\approx e}^T \quad (5.11)$$

where $K_{\approx e}^{\text{exact}}$ and $K_{\approx e}^{\text{under}}$ respectively are the exact element stiffness matrix and its under-integrated form, $\bar{\epsilon}_e$ and $\gamma_{\approx e}$ are obtained from (2.27) and (2.28). In particular

$$\gamma_{\approx e} = \bar{h} - \frac{h \cdot x}{|\Omega_e|} b_{\approx 1} - \frac{h \cdot y}{|\Omega_e|} b_{\approx 2} \quad (5.12)$$

Even though ϵ_e given by (2.27) is still difficult to calculate, note that, during the element calculations, the vectors b_1 and b_2 and the Jacobian $|\Omega_e|$ are necessarily computed. Therefore γ_e is very easy to calculate.

The inner product $(\bar{u}^h, H)_1$ may, therefore, be calculated by using the decomposition (2.1) of $K_{\approx}^{\text{exact}}$. Introducing the nodal vectors \bar{U} and \bar{H} associated with the function \bar{u}^h and H , we have

$$\begin{aligned} (\bar{u}^h, H)_1 &= \bar{U}^T K_{\approx}^{\text{exact}} \bar{H} = \bar{U}^T \cdot \sum_e \bar{\epsilon}_e \gamma_e \gamma_e^T \cdot \bar{H} \\ &= \sum_e \bar{\epsilon}_e (\gamma_e^T \cdot \bar{U}) (\gamma_e^T \cdot \bar{H}) \end{aligned}$$

where \bar{U}_e and \bar{H}_e are the values of \bar{U} and \bar{H} at the nodes of the element e . We note that if the values of \bar{H} are ± 1 or -1 , the scalar vector product $\gamma_e^T \cdot \bar{H}_e$ is always ± 4 . Therefore

$$(\bar{u}^h, H)_1 = 4 \sum_e \pm \bar{\epsilon}_e \sum_{i=1}^4 \gamma_e^i u_e^i \quad (5.13)$$

and

$$(H, H)_1 = 16 \sum_e \bar{\epsilon}_e \quad (5.14)$$

These expressions are still exact since no approximation has been made on $\bar{\epsilon}_e$. If we suppose that the Jacobian of the element is approximately constant (true for parallelogram element), $\bar{\epsilon}_e$ is simply expressed as

$$\bar{\epsilon}_e = \frac{1}{12|\Omega_e|} (b_1^T b_1 + b_2^T b_2) \quad (5.15)$$

The calculation of the approximate projection can be summarized in

the following algorithm:

- Loop on Elements

↑ Calculate $\underline{\gamma}_e, \epsilon_e$ using (2.2) and (2.6)

Calculate $\pm \epsilon_e (\underline{\gamma}_e^T \cdot \bar{U}_e)$

└ Add

$$\begin{cases} \lambda_2 = \lambda_1 \pm \epsilon_e (\underline{\gamma}_e^T \cdot \bar{U}_e) \\ \lambda_2 = \lambda_2 + \epsilon_e \end{cases}$$

- $\lambda = \lambda_1 / 4 \lambda_2$

- Loop on Nodes

↑ $\bar{u}^h|_{\text{Node}} = \bar{u}^h|_{\text{Node}} \pm \lambda$

Remark: The notations previously used are essentially those found in the work of Belytschko and co-workers [5, 6] on stabilization methods. These methods rely on the decomposition (5.11) but the stabilization term $\epsilon \underline{\gamma} \cdot \underline{\gamma}^T$ is a-priori added to the under-integrated matrix to prevent the spurious modes from the kernel of the stiffness matrix; whereas our control method uses the very same term a-posteriori, after solving with the underintegrated matrix. Therefore, our method seems to be cheaper than the stabilization methods as summarized in Table 1.

2.5.2. Numerical results.

2.5.2.a. Regular mesh of 4-node elements. In order to illustrate what has been stated, we have considered the Laplacian problem solved on a square domain partitioned into $N^2 (= h^{-2})$ subdomains, for various values of N and we have studied the norms of the difference between the solution obtained with a full integration u^h (4 point rule

Table 1. Operations Cost per Element for Both
Stabilization and A-posteriori Methods

Stabilization Method	Operations Cost		A-posteriori Method
	20x ; 21+	20x ; 21+	
Computations of $\bar{\epsilon}_e, \gamma_e$	20x ; 21+	20x ; 21+	Computations of $\bar{\epsilon}_e, \gamma_e$
Multiply $\gamma_e \cdot \gamma_e^T$	16x	4x	Multiply $U_e^T \cdot \gamma_e$
Multiply $\epsilon_e \cdot \gamma_e \gamma_e^T$	16x	1x	Multiply $\epsilon_e \cdot U_e^T \gamma_e$
Add $K_e + \epsilon_e \gamma_e \gamma_e^T$	16+	2+	Add $\gamma_1 \pm \epsilon_e U_e^T \gamma_e$ $\gamma_2 + \epsilon_e$
			Then: (4 nodes/element) Add $\bar{u} \pm \lambda$
TOTAL	52x ; 37+	29x ; 27+	TOTAL

and with underintegration (1 point rule) \bar{u}^h and \hat{u}^h (before and after post processing). The results are shown as plot of $\text{Log} \|u^h - \bar{u}^h\|$ or $\text{Log} \|u - \hat{u}^h\|_s$ in function of $|\text{Log } h|$, for $s = 0, 1$. Data of various regularities have been used:

i) f_1 is a C^0 -function, but not C^1 :

$$\begin{cases} f_1(x,y) = \frac{3}{2}(1-x) - \frac{14}{9}y & \text{if } Y_1(x,y) \geq 0 \\ f_1(x,y) = y(\frac{3}{2}(1-x) - y - \frac{5}{9}) & \text{if } Y_1(x,y) \leq 0 \end{cases}$$

where the C^1 -discontinuity line is

$$Y_1(x,y) = \frac{3}{2}(1-x) - y$$

ii) f_2 is a non-continuous function

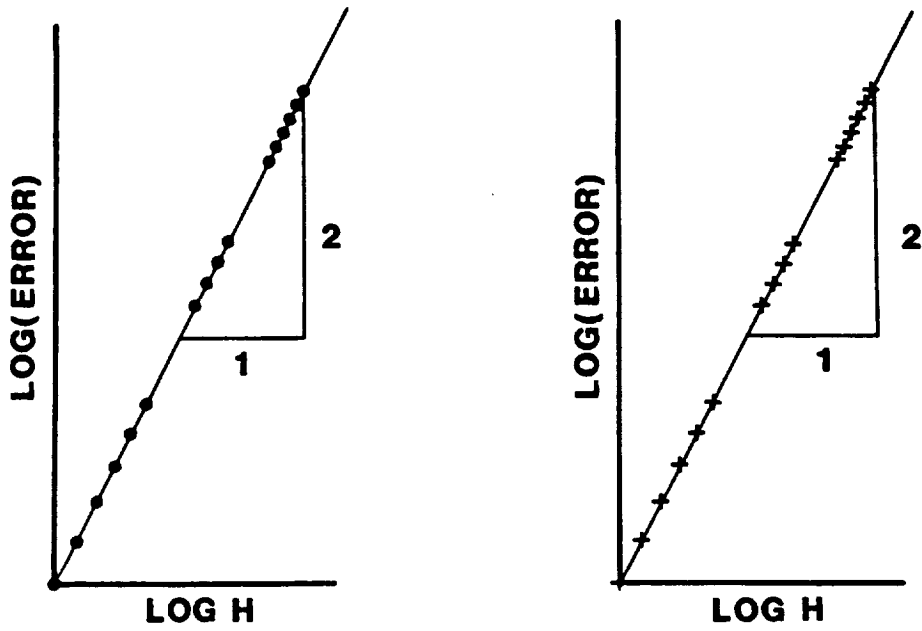
$$\begin{cases} f_2(x,y) = 1 & \text{if } Y_1(x,y) > 0 \\ f_2(x,y) = -2 & \text{if } Y_1(x,y) < 0 \end{cases}$$

where Y_1 is the same as in i).

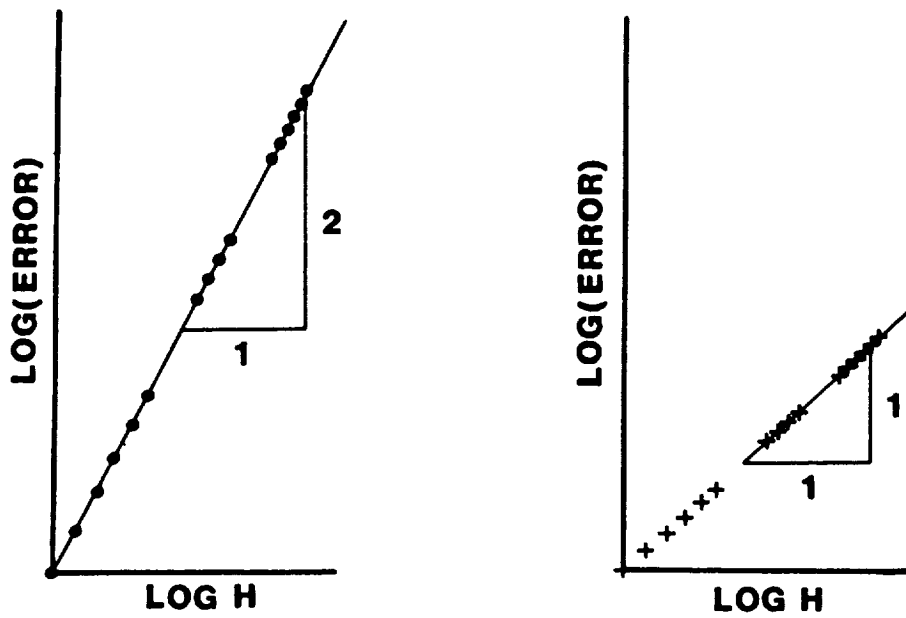
Remark: Both of these functions satisfy the compatibility condition (3.2).

Results obtained with the continuous function f_1 are shown in Figure 5. When the solution has been treated by the post processor (Fig. 5a.), for both L^2 and H^1 norms, the representing points lie on lines of slope 2. This proves that whereas the estimate (1-13) is optimal for the L^2 norm ($s = 0$), it is not in the H^1 norm and seems to be in fact better than what was expected in our study. This does not affect in any case the comparison with the exact solution (1.14).

Figure 5.b shows the comparison with the crude solution \bar{u}^h , obtained with two fixed nodes, and not treated by the post-processor.



(a)



(b)

Figure 5. Results Obtained with a Continuous Data Function:
 Comparison between : a) u^h and \bar{u}^h ; b) u^h and \bar{u}^h .

Slopes 2 and 1 are observed and the loss (1. instead of 2. for the H^1 norm) corroborates the final remark of Section 3.32.

When the function f_2 is used (Figure 6), the points show oscillations around two lines of slope 2. (for the L^2 -norm) and 1.65 (for the H^1 -norm) proving that the estimate (3.31) still holds (Fig. 6.a). When the solution has not been treated by the post-processor (Fig. 6.b), the slope 1.65 becomes 1. as expected.

The next series of examples was intended to study the influence of a singularity (at the origin) for a unit square domain regularly partitioned. The data functions are of the form

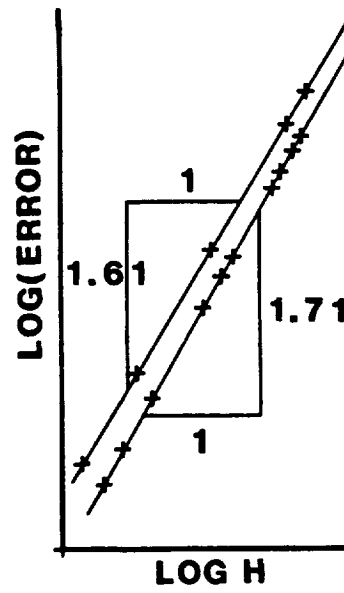
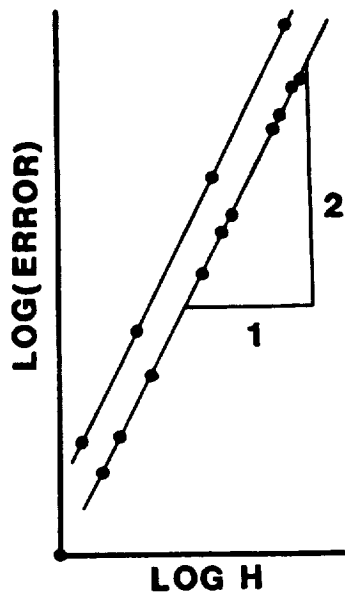
$$f_\alpha(x,y) = r^\alpha - C, \quad \alpha > -2 \quad (5.16)$$

where C is a real number chosen such that the equilibrium condition (2.2) is satisfied. The family $\{f_\alpha\}$ represents various regularities of data:

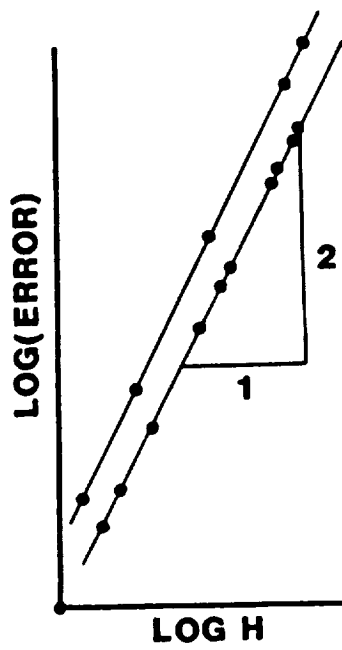
$$f_\alpha \in H^s(\Omega) \quad \alpha > s - 1 \quad (5.17)$$

The result shown in Fig. 7.a is a plot of α (regularity) versus σ (rate of convergence of $\|\tilde{u}^h - u^h\|_{s=0,1}$). The pattern of the (α, σ) plot seems to show a linear increase of slope 1 towards the maximum value 2 reached for $f \in L^2$ ($\alpha = -1$) for the L^2 -norm ($s=0$). As far as the H^1 -norm ($s=1$) of the error is concerned, the linear increase of slope 1 reached 1 for $f \in L^2$ but keeps increasing towards 2. This shows that the expected error estimate

$$\|u^h - \tilde{u}^h\|_s \leq C h^k \|f\|_m, \quad s = 0,1 \quad (5.18)$$



(a)



(b)

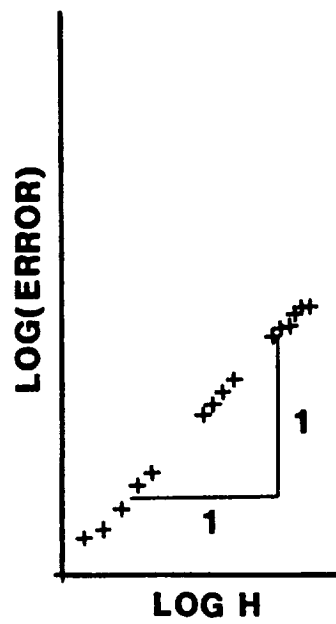


Figure 6. Results Obtained with the Continuous Data f_2 :
 Comparison between : a) u^h and \bar{u}^h ; b) u^h and \bar{u}^h .

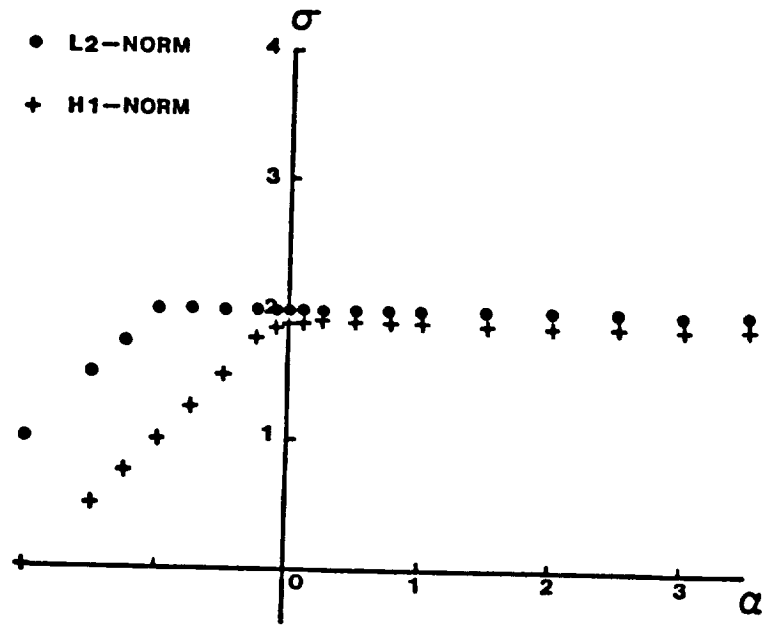


Figure 7.a. Bilinear Elements

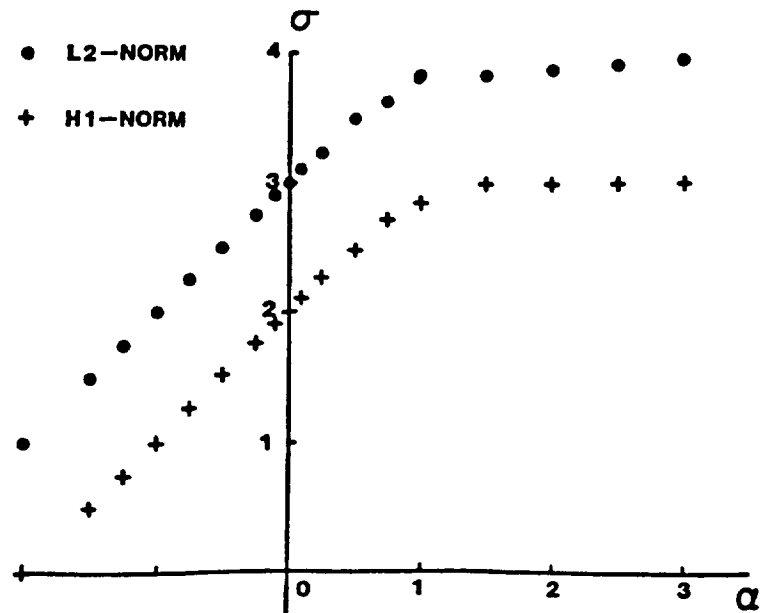


Figure 7.b. Biquadratic Elements

Figure 7. (α, σ) Plot for a Square Domain

where

$$k = 1 + \min(1, m) - s \quad (5.19)$$

is not optimal for $s = 1$ and $m > 0$. The estimate (3.32) remains optimal however.

In conclusion, these numerical results suggest that the method is accurate for regular meshes with a convergence rate equal to the fully integrated case.

2.5.2.b. Regular mesh of 8- and 9-node elements. Since the beginning of Part II we have not discussed the underintegration of the stiffness matrix of the 8-node elements. It is well known that this matrix is not rank-deficient, and the practice of the underintegration has been widely used with good results when the mesh is regular. Since there is not any spurious mode, the a-posteriori control previously described is not needed.

Unfortunately, the method of proof presented in Section 2.4 cannot be used because this element does not possess the nice tensor product properties on which the method relies. The only hope for a proof of convergence would be to obtain the result as a by-product of a result for the 9-node elements.

As far as this element is concerned (9-node element), we have proved (Section 2.2.3) that the underintegration of this element leads to a rank-deficient matrix; in fact, the procedure described in Section 2.3, for the resolution of the underintegrated problem and the projection of its solution is completely applicable to a mesh of 9-node elements. Thus, Theorem VII is valid and the projection defined in (3.3)

a convergence theorem is concerned, one can establish generalizations of (4.16) and (4.17) to 3-node, one-dimensional elements: there exist α_1 , β_1 , ϕ^1 and ψ^1 such that

$$\begin{aligned} (v^h, \phi^1)_{0,h} &= \alpha_1 (v^h, \phi^1)_0 \\ (v^h, \psi^1)_1 &= \beta_1 (v^h, \psi^1)_0 \end{aligned} \quad \forall v^h$$

Unfortunately, the basis functions ϕ^1 and ψ^1 are different for the L^2 -underintegrated and H^1 -norms and a lemma as Lemma 1.1 cannot be obtained.

However, in this subsection we will show numerical results obtained by use of the projection (3.30) for regular meshes of 9-node elements. Note that two types of control have been tested with similar results: the control only involving the term in $\gamma \cdot \gamma^T$ predicted by Belytschko [4] and the complete control calculated with $s_1^T s_1^T$, $i = 7, 8, 9$. (See Section 2.2.3). The results obtained with either of them are similar for this operator $(-\Delta)$.

For 8 and 9-node elements, the optimal rates of convergence are given by

$$\|u - u^h\|_s \leq C h^k \|f\|_m \quad s = 0, 1 \quad (5.20)$$

where

$$k = 2 + \min(1, m) - s \quad (5.21)$$

and the best rates of convergence $O(h^{3-s})$ are obtained when $f \in H^1(\Omega)$. The results obtained with functions presenting a singularity line (such as the functions f_1 and f_2 previously defined and others) are pre-

sented in Table 2 (first and second lines). We obtained 1.99 and 1.74 for a discontinuous function ($f \in L^2$), then 2.43 and 1.97 for a continuous, not C^1 , function ($f \in H^1$), 2.95 and 1.95 for a C^1 , not C^2 function ($f \in H^2$) and finally 4 and 3 for a C^∞ data. Therefore, the rates 3 and 2 are reached when f is at least H^2 or equivalent when the solution u is in H^4 . In this case, the convergence rate (5.1) does not seem to be reached.

The second series of data involving the singularity at the origin (5.16) has been tested and results are shown in Fig. 6.b. The pattern of the (α, σ) plot shows linear increases of slope 1, the predicted values 3 and 2 are reached for $f \in H^1(r)$ according to (5.21), but the maximum values 4 and 3 are reached for $f \in H^2(\Omega)$.

2.5.2.c. Irregular mesh of 4- and 9- node elements. Finally, the method has been tested on the quarter unit disk shown in Fig. 8 with

$$f = r^\alpha - \frac{2}{\alpha+2} \quad \alpha > -2$$

The plot (α, σ) is shown in Fig. 9 and we can point out:

- The general pattern is respected (linear increase towards a maximum value)
- The maximum values 4 and 3 (9-node elements) are reduced to values slightly lower than 3 and 2.

2.6 Excitation of Spurious Modes

The previous sections were devoted to the study of the Laplace equations with Neumann boundary conditions. The choice of these boundary conditions is convenient for the analysis of the hourglass instabil-

Table 2. Rate of Convergence $\|\log \|e\|_s\| \text{ v.s. } \|\log h\| (0=0,1)$
for 8- and 9-node Elements

REGULARITY BOUNDARY CONDITIONS AND ELEMENT	$f \in L^2$ $e \in H^1$	$f \in H^1$ $e \in H^2$	$f \in H^2$ $e \in H^3$	$f \in C^\infty$
NEUMANN, 9-NODE +SPECIAL PROJECTION	1.99 1.74	2.43 1.97	2.95 1.94	4.00 3.00
NEUMANN, 8-NODE EL.	1.99 1.79	2.00 1.93	2.97 1.95	4.00 3.00
DIRICHLET, 9-NODE EL.	2.35 1.47	2.85 1.99	2.99 1.99	3.00 2.00
DIRICHLET, 8-NODE EL.	2.30 1.46	2.71 2.12	2.99 1.99	3.00 2.00
MIXED, 9-NODE EL.	2.00 1.67	2.00 2.28	3.83 2.84	4.00 3.00
MIXED, 8-NODE EL.	2.00 1.74	2.00 2.27	3.84 2.84	4.00 3.00

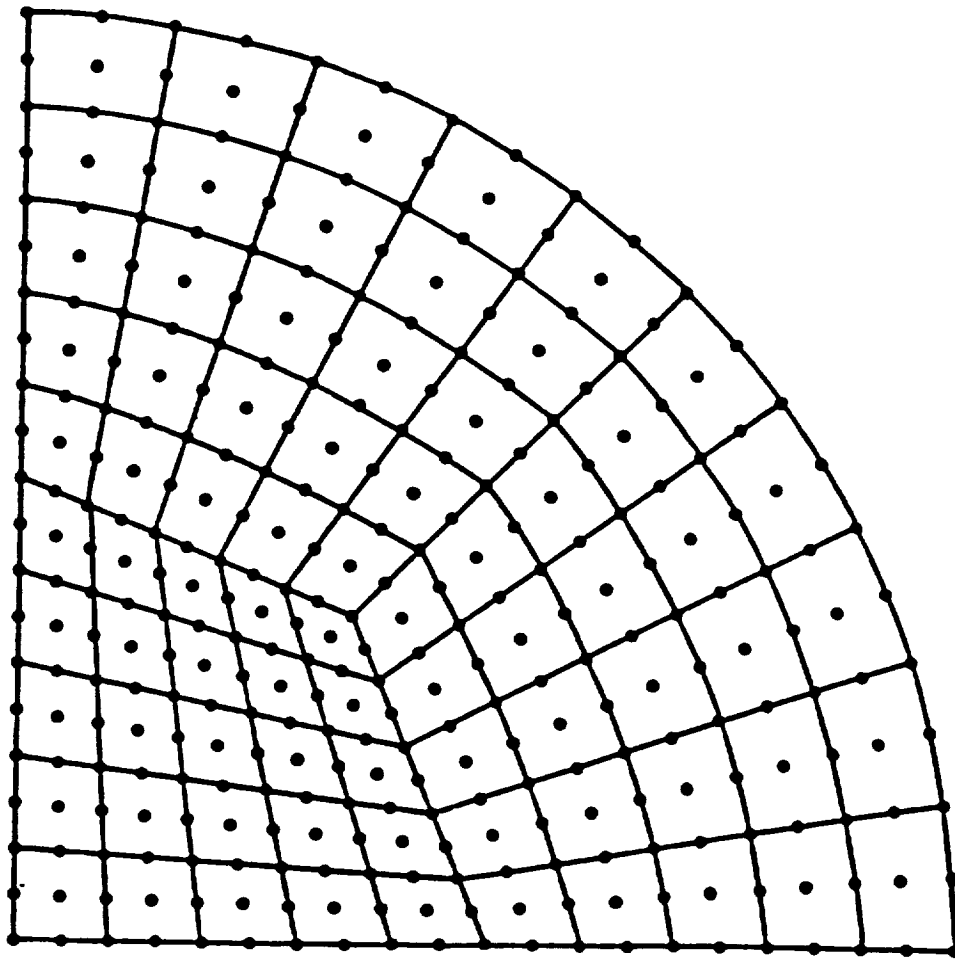


Figure 8. Typical Mesh on a Quarter Circle Domain.

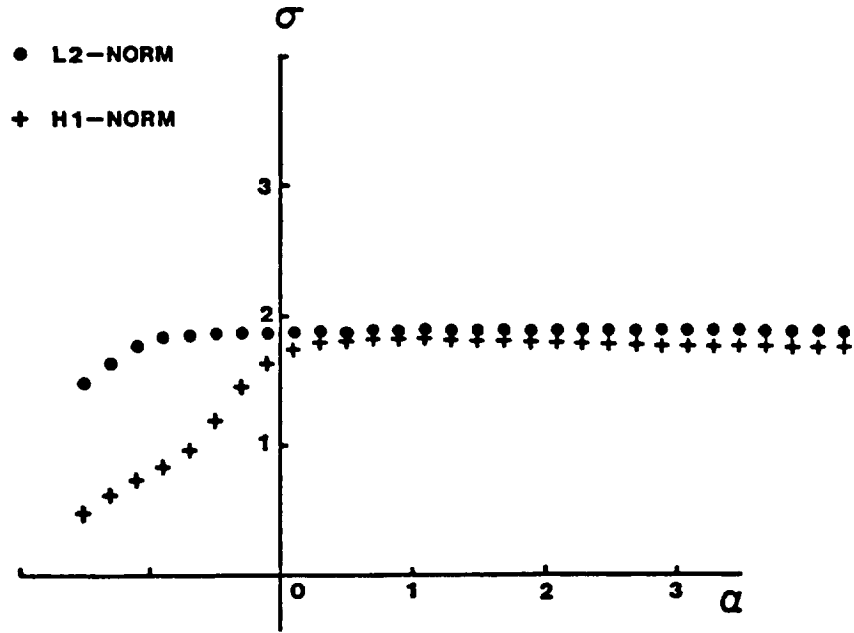


Figure 9.a. Bilinear Elements.

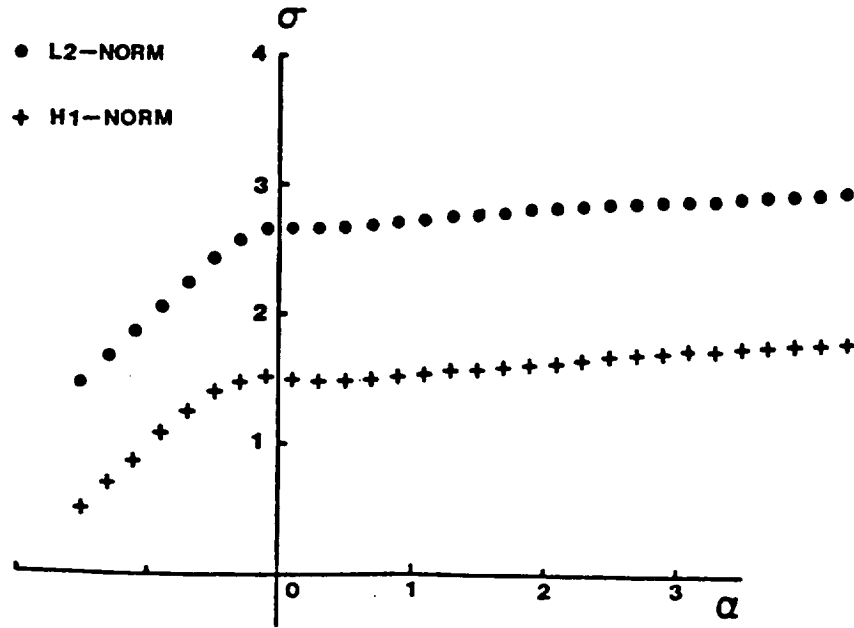


Figure 9.b. Biquadratic Elements.

Figure 9. (α, σ) Plot for a Quarter Circle Domain.

ities because these modes appear explicitly in the kernel of the underintegrated discrete operator. When Dirichlet conditions are applied on a part of the boundary, even though the kernel of the underintegrated stiffness matrix is not rank-deficient, instabilities may appear.

In this section we would like to study the influence the boundary conditions have on the solution of the underintegrated problem, and obtain results analogous to Theorem VIII. Also we would like to explain how the oscillations may be excited in certain problems. The method of proof is similar to that presented in Section 2.4. For various boundary conditions, we are able to exhibit the exact eigenvalues and eigenfunctions of the various linear and bilinear form involved. The explanation of the excitation of oscillations will result from the comparison of these eigenvalues. The procedure also allows us to study the underintegration of the operator $-\Delta+1$ and the control of resulting spurious modes. Numerical results will illustrate the theory.

2.6.1. The underintegrated problem with Dirichlet or mixed boundary conditions. This section is devoted to a generalization of the results obtained in Section 2.4 to the Laplacian equation with Dirichlet or mixed Dirichlet-Neumann boundary conditions. Only proofs for the Dirichlet case will be given in this section, but their equivalent for the mixed case can be found in Appendix A.

The Dirichlet case is simpler than the Neumann case because the hourglass mode does not belong to the new approximation space defined to handle the boundary condition. Therefore the stiffness matrix is no longer singular and can be inverted directly. In the variational formulation, similar to (3.7), the projection of the data function is not

necessary and the problem \bar{P}^h is written:

$$\begin{aligned}
 (\bar{P}^h) \quad & \text{Find } \bar{u}^h \in V_0^h \text{ such that} \\
 & (\bar{u}^h, v^h)_{1,h} = (f^h, v^h)_{0,h} \quad v^h \in V_0^h \quad (6.1)
 \end{aligned}$$

where

$$\begin{aligned}
 V_0^h = \{ & v^h/v^h \in C^0(\bar{\Omega}), v^h|_{\Omega_{ij}} \in Q_1(\Omega_{ij}), 1 \leq i, j \leq N, \\
 & v^h|_{\partial\Omega} = 0 \}
 \end{aligned}$$

Remarks

i) The fact that $(\cdot, \cdot)_{1,h}$ is not singular on V_0^h does not make the problem classically elliptic in the sense that the constant in the Lax Milgram Theorem is h-dependent.

ii) When Dirichlet or mixed boundary conditions are applied,

$$\text{Ker } A_1 = \text{Ker } A_{1,h} = \{0\}$$

Thus the post processor is not justified anymore and we may compare directly \bar{u}^h and u^h .

This comparison is carried the same way as in Section 2.4 and a basis of the approximation-space V_0^h can be obtained. One useful basis is the common eigen-basis of the matrices of the H^1 -, L^2 -, and underintegrated H^1 - or L^2 - inner product. Let us consider the $(N-1) \times (N-1)$ matrix

$$D(k) = \begin{bmatrix} 2k & & & & 0 \\ & -1 & & & \\ & & & & \\ & & & & -1 \\ 0 & & & -1 & 2k \end{bmatrix}$$

The values for which $\det D(k)$ vanishes are

$$k_i = \cos \frac{i\pi}{N} \quad 1 \leq i \leq N-1 \quad (6.2)$$

and the corresponding vectors $(D(k_i)v_i = 0)$ are

$$v_i = \left\{ \sin \frac{j\pi}{N} \right\}_{j=1, N-1} \quad (6.3)$$

Let $\phi^i = \phi^i(x)$, $x \in [0,1]$, denote the piecewise linear function associated with the vector v_i :

$$\phi^i(jh) = \sin \frac{j\pi}{N} \quad 1 \leq i, j \leq N-1 \quad (6.4)$$

$$\text{span}\{\phi^i\}_{1 \leq i \leq N-1} = V_{1,0}^h \quad (6.5)$$

where

$$V_{1,0}^h = \{v^h \in C^0(0,1), v^h(0) = v^h(1) = 0\}$$

$$v^h \text{ is linear on } [eh, (e+1)h], \quad 0 \leq e \leq N-1$$

From this point, the remainder of the proof goes as in Section 2.4 and the variational problem and its underintegrated formulation can be explicitly solved and the decomposition (4.29), (4.30) and (4.31), (4.32) are obtained for $1 \leq i, j \leq N-1$, and we finally obtain the result for Dirichlet boundary conditions:

THEOREM V: Let f be a function in $L^2(\Omega)$. Let u be the solution of P :

$$P : \text{Find } u \in H_0^1 / (u, v)_1 = (f, v)_0 \quad \forall v \in H_0^1 \quad (6.6)$$

Let f^h be the L^2 -projection of f onto V^h and let \bar{u}^h be the solution of \bar{P}^h :

$$\bar{P}^h: \text{Find } \bar{u}^h \in V_0^h / (\bar{u}^h, v^h)_{1,h} = (f^h, v^h)_{0,h} \quad \forall v^h \in V_0^h \quad (6.7)$$

Then we have the following error estimate:

$$\|u - \bar{u}^h\|_s \leq C h^{2-s} \|f\|_0 \quad s = 0, 1 \quad \square \quad (6.8)$$

This theorem proves that the use of the underintegrated matrix does not affect the rate of convergence of the solution. The method is therefore accurate and efficient.

Data with various regularity have been tested for meshes of 4-, 8-, and 9-node elements, with various boundary conditions. Results are summarized in Tables 2 and 3. They indicate that the optimal rates of convergence are obtained for $f \in L^2(Q)$ for the 4-node case and $f \in H^2(Q)$ for the 8- and 9-node case.

2.6.2. The underintegration of the operator $-\Delta+1$. In this subsection we consider the underintegration of the operator associated with the problem

$$P_0 : \text{Find } u \in H^1(\Omega) \text{ such that}$$

$$\begin{cases} -\Delta u + u = f & \text{in } \Omega \\ \frac{\partial u}{\partial n} = 0 & \text{on } \partial\Omega \end{cases} \quad (6.9)$$

The usual variational formulation of P_0 is

$$P : \text{Find } u \in H^1(\Omega) \text{ such that}$$

$$(u, v)_1 + (u, v)_0 = (f, v)_0, \quad v \in H^1(\Omega) \quad (6.10)$$

The results of existence, uniqueness of solutions of P are well-known

Table 3: Rate of Convergence $|\text{Log } \|e\|_s|$ v.s. $|\text{Log } h|$ for
4-node Elements ($s=0,1$)

REGULARITY		$f \in L^2$	$f \in H^1$	$f \in C^\infty$
OPERATOR AND BOUNDARY CONDITIONS		$e \in H^1$	$e \in H^2$	
- Δ	NEUMANN	1.99	2.00	2.00
	+PROJECTION	1.61	2.00	2.00
	DIRICHLET	1.99	2.00	2.00
		1.50	1.85	2.00
	MIXED	2.00	2.00	2.00
		1.50	1.99	1.99
- $\Delta + 1$	NEUMANN	2.00	2.00	2.00
	+PROJECTION	1.50	2.00	2.00
	DIRICHLET	2.00	2.00	2.00
		1.50	1.85	2.00
	MIXED	2.00	2.00	2.00
		1.50	1.85	2.00

as are those for its discrete formulation:

P_h : Find $u^h \in V^h$ such that

$$(u^h, v^h)_1 + (u^h, v^h)_0 = (f, v^h)_0, \quad \forall v^h \in V^h \quad (6.11)$$

where V^h is an approximation of $H^1(\Omega)$ using bilinear elements.

The underintegration of $(\cdot, \cdot)_1 + (\cdot, \cdot)_0$ leads to the following under-integrated problem:

\bar{P}_h : Find $\bar{u}^h \in \bar{V}^h$ such that

$$(\bar{u}^h, v^h)_{1,h} + (\bar{u}^h, v^h)_{0,h} = (f, v^h)_{0,h}, \quad \forall v^h \in \bar{V}^h \quad (6.12)$$

where the choice of approximation space

$$\bar{V}^h = V^h/H \quad (6.13)$$

is justified by

$$(v^h, H)_{1,h} + (v^h, H)_{0,h} = 0, \quad \forall v^h \in V^h \quad (6.14)$$

Then, the method of proof used in Section 2.3 allows us to obtain the existence and uniqueness of \bar{u}^h . A projection similar to (3.30) can be obtained by analogy: we have

$$(u^h, H)_1 + (u^h, H)_0 = (f, H)_0 \quad (6.15)$$

We therefore construct the projection as:

$$\bar{u}^h = \pi u^h = \bar{u} + \lambda_0 H \quad (6.16)$$

$$(\bar{u}^h, H)_1 + (\bar{u}^h, H)_0 = (f, H)_0 \quad (6.17)$$

This defines uniquely λ_0 as

$$\lambda_0 = \frac{(f, H)_0 - (\bar{u}, H)_1 - (\bar{u}, H)_0}{\|H\|_1^2 + \|H\|_0^2} \quad (6.18)$$

Similarly to what was done in Section 2.5.1, we can use (4.25), (5.6) through (5.8), simplify λ_0 without any loss of accuracy and still use (5.9), (5.10) for the projection

$$\tilde{u}^h = \bar{u}^h - \lambda H \quad (5.9)(\text{repeated})$$

$$\lambda = (\bar{u}^h, H)_1 \|H\|_1^{-2} \quad (5.10)(\text{repeated})$$

The proof of the convergence of \tilde{u}^h towards u is again done by direct calculation of u^h and \tilde{u}^h : the explicit resolution of P_h and $(\bar{P}_h + \pi)$ leads to :

$$u_{i,j} = \frac{1}{1 + \beta_i + \beta_j} f_{ij} \quad 0 \leq i, j \leq N \quad (6.19)$$

and

$$\left\{ \begin{array}{l} \tilde{u}_{i,j} = \frac{\alpha_i \alpha_j}{\alpha_i \alpha_j + \alpha_i \beta_j + \alpha_j \beta_i} f_{ij} \quad \begin{array}{l} 0 \leq i, j \leq N \\ (i, j) \neq (N, N) \end{array} \\ \tilde{u}_{NN} = u_{NN} \end{array} \right. \quad (6.20)$$

These decompositions allow us to obtain $u^h - \tilde{u}^h$ as done in Section (2.4). Provided that $f \in L^2(\Omega)$, we can obtain

$$\|u^h - \tilde{u}^h\|_s \leq C h^{2-s} \|f\|_0 \quad s = 0, 1 \quad (6.21)$$

Once again, the underintegration does not seem to affect the rate of convergence. The result can also be obtained with various boundary

conditions. Numerical results are summarized in Tables 5 and 6 and confirm the theory.

2.6.3. Excitation of oscillations. The existence of spurious oscillations when underintegration is used is not only encountered when Neumann boundary conditions are applied on the whole boundary. In this subsection, we analyze precisely how modes that oscillate with wavelength of order h are excited when underintegration is used, whereas they are damped when the integration is exact.

For this discussion we consider the unit square

$$\Omega =]0,1[\times]0,1[\quad (6.22)$$

discretized into $N \times N$ elements. We consider the Laplace equation on Ω

$$\left. \begin{aligned} -\Delta u &= f && \text{in } \Omega \\ u &= 0 && \text{on } \partial\Omega \cap \{x = 0\} \\ \frac{\partial u}{\partial n} &= g && \text{on } \partial\Omega / \{x = 0\} \end{aligned} \right\} \quad (6.23)$$

For the first time we include two kinds of load: body forces and surface loads, and we will observe separately the effects of each of them.

The eigenfunctions associated with these particular mixed boundary conditions are constructed as in Section 2.4.

$$\chi^{ij}(x,y) = \psi^i(x)\phi^j(y); \quad \begin{array}{l} 1 < i < N \\ 0 < \underline{j} < \underline{N} \end{array} \quad (6.24)$$

where ϕ^j is defined in (4.15) (associated with Neumann boundary conditions at both ends) and ψ^i is similarly defined (see Appendix A). These functions are defined through sine and cosine functions and therefore oscillate. Among them we will distinguish "smooth" modes with

longer wavelengths ($O(1)$) from "irregular" modes with shorter wavelengths ($O(h)$). Smooth (respectively irregular) modes correspond to smaller (respectively larger) values of i or j . Examples of each extreme are shown in Fig. 10 for $N=10$.

The resolution of the fully integrated problem leads to the search for coefficients u_{ij} such that

$$u^h = \sum_{\substack{1 < i < N \\ 0 < j < N}} u_{ij} \chi^{ij} \quad (6.25)$$

The basis $\{\chi^{ij}\}$ is an eigenbasis for $(\cdot, \cdot)_1$ and therefore we have

$$u_{ij} = A_{ij} \left[(f, \chi^{ij})_0 + (g, \chi^{ij})_{0, \partial\Omega} \right] \quad (6.26)$$

where

$$A_{ij} = \frac{1}{\beta'_i + \beta_j} \quad (6.27)$$

with

$$\left. \begin{aligned} \beta'_i &= \frac{6}{h^2} \frac{1 - \cos\left(\frac{i\pi}{N} - \frac{\pi}{2N}\right)}{2 + \cos\left(\frac{i\pi}{N} - \frac{\pi}{2N}\right)} \\ \beta_j &= \frac{6}{h^2} \frac{1 - \cos\left(\frac{j\pi}{N}\right)}{2 + \cos\left(\frac{j\pi}{N}\right)} \end{aligned} \right\} \quad (6.28)$$

The values A_{ij} have been calculated exactly with these formulae and their values are reported in Table 4.a for $N=10$. The 20 highest values are in the shaded zone. We clearly can observe that

- i) these values range from the highest value to 1% of this value,
- ii) these values are associated with smooth modes (tensor products of smooth modes).

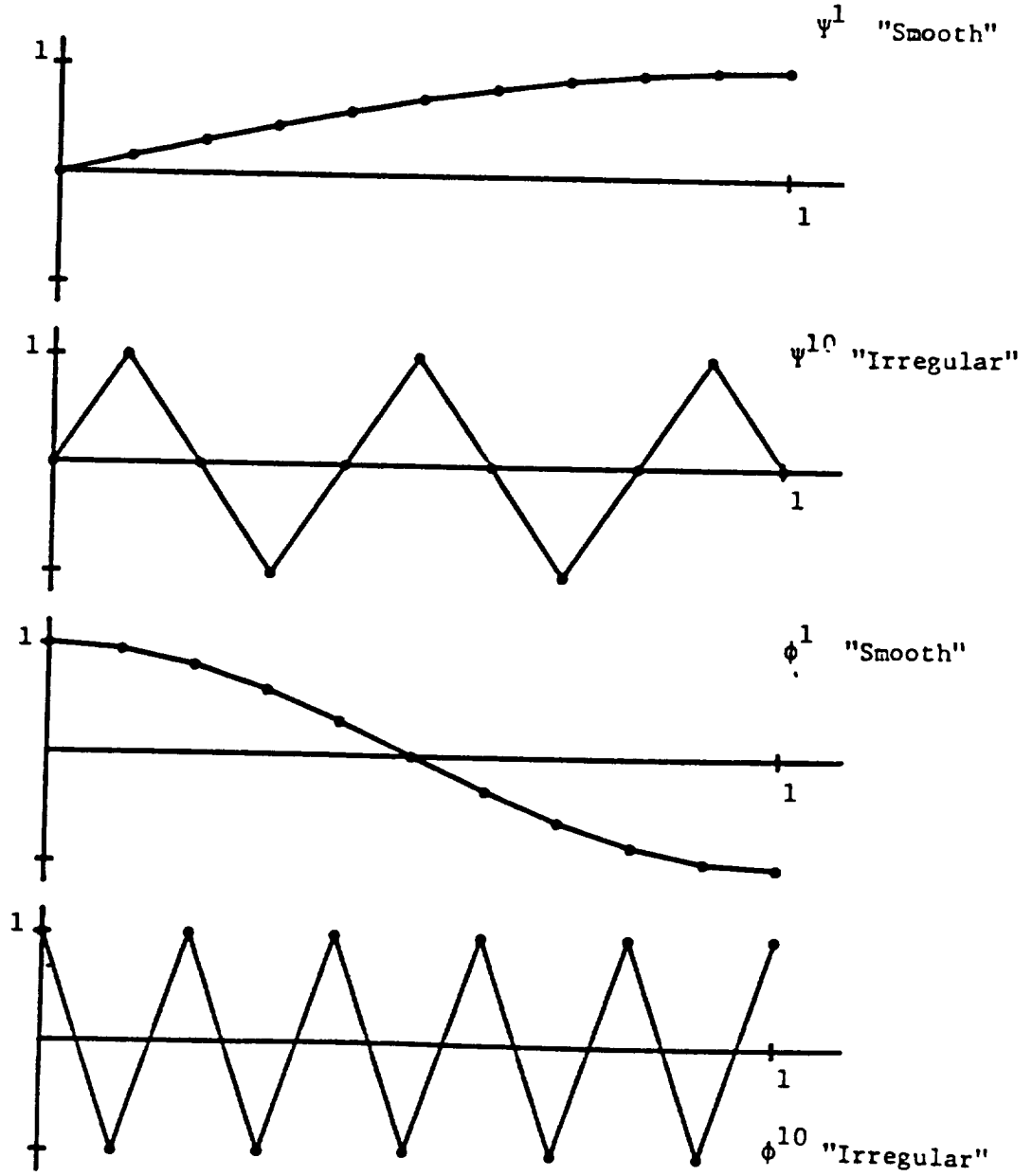


Figure 10. Examples of "Smooth" or "Irregular" Eigenfunctions.

Table 4: Arrays of Eigen Values A_{ij} , \bar{A}_{ij} and $\bar{\bar{A}}_{ij}$ Table 4a. Array A_{ij}

i	1	2	3	4	5	6	7	8	9	10	
j	0	40.45	4.42	1.54	0.75	0.43	0.27	0.18	0.13	0.10	0.08
1	8.05	3.07	1.34	0.70	0.41	0.26	0.17	0.12	0.10	0.08	
2	2.31	1.58	0.95	0.57	0.36	0.24	0.17	0.12	0.09	0.08	
3	1.02	0.85	0.62	0.44	0.30	0.21	0.15	0.11	0.09	0.08	
4	0.55	0.49	0.41	0.32	0.24	0.18	0.13	0.10	0.08	0.07	
5	0.33	0.31	0.27	0.23	0.19	0.15	0.12	0.09	0.08	0.07	
6	0.21	0.21	0.19	0.17	0.14	0.12	0.10	0.08	0.07	0.06	
7	0.15	0.14	0.14	0.12	0.11	0.10	0.08	0.07	0.06	0.05	
8	0.11	0.11	0.10	0.10	0.09	0.08	0.07	0.06	0.05	0.05	
9	0.09	0.09	0.08	0.08	0.07	0.07	0.06	0.05	0.05	0.04	
10	0.08	0.08	0.08	0.07	0.07	0.06	0.06	0.05	0.04	0.04	

Table 4b. Array \bar{A}_{ij}

i	1	2	3	4	5	6	7	8	9	10	
j	0	40.36	4.34	1.46	0.67	0.34	0.18	0.09	0.04	0.01	0.00
1	7.99	3.02	1.27	0.62	0.33	0.18	0.09	0.04	0.01	0.00	
2	2.24	1.53	0.90	0.52	0.30	0.17	0.09	0.04	0.01	0.00	
3	0.94	0.79	0.58	0.39	0.25	0.15	0.09	0.04	0.01	0.00	
4	0.47	0.43	0.36	0.28	0.20	0.13	0.08	0.04	0.01	0.00	
5	0.25	0.24	0.21	0.18	0.14	0.11	0.07	0.04	0.01	0.00	
6	0.13	0.13	0.12	0.11	0.10	0.08	0.05	0.03	0.01	0.00	
7	0.06	0.06	0.06	0.06	0.05	0.05	0.04	0.03	0.01	0.00	
8	0.03	0.03	0.03	0.03	0.02	0.02	0.02	0.02	0.01	0.00	
9	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.00	0.00	
10	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	

Table 4c. Array $\bar{\bar{A}}_{ij}$

i	1	2	3	4	5	6	7	8	9	10	
j	0	40.45	4.42	1.54	0.75	0.43	0.27	0.18	0.13	0.10	0.08
1	8.08	3.11	1.36	0.71	0.42	0.26	0.18	0.13	0.10	0.09	
2	2.32	1.62	0.99	0.61	0.39	0.26	0.18	0.13	0.10	0.09	
3	1.02	0.87	0.67	0.48	0.34	0.24	0.18	0.13	0.10	0.09	
4	0.55	0.51	0.44	0.37	0.29	0.23	0.17	0.14	0.11	0.10	
5	0.33	0.32	0.30	0.27	0.24	0.20	0.17	0.14	0.12	0.11	
6	0.22	0.21	0.21	0.20	0.19	0.18	0.17	0.16	0.14	0.14	
7	0.15	0.15	0.15	0.15	0.15	0.16	0.17	0.17	0.18	0.19	
8	0.11	0.11	0.11	0.12	0.13	0.14	0.16	0.20	0.26	0.33	
9	0.09	0.09	0.09	0.10	0.11	0.13	0.16	0.23	0.42	0.97	
10	0.08	0.08	0.09	0.09	0.10	0.12	0.16	0.25	0.57	4.57	

On the other hand, the eigenvalues of irregular modes are smaller and because of this, these modes will be damped; only smooth modes will contribute in (6.25).

When the underintegration is used, and when g is zero, the solution \tilde{u}^h is

$$\tilde{u}^h = \sum_{\substack{1 < i < N \\ 0 < j < N}} \tilde{u}_{ij} \chi^{ij} \quad (6.29)$$

with

$$\tilde{u}_{ij} = \tilde{A}_{ij}(f, \chi^{ij})_0 \quad (6.30)$$

where

$$\tilde{A}_{ij} = \frac{\alpha'_i \alpha_j}{\alpha_j \beta'_i + \alpha'_i \beta_j} \quad (6.31)$$

and

$$\left. \begin{aligned} \alpha'_i &= \frac{3(1 + \cos(\frac{i\pi}{N} - \frac{\pi}{2N}))}{2(2 + \cos(\frac{i\pi}{N} - \frac{\pi}{2N}))} \\ \alpha_j &= \frac{3(1 + \cos(\frac{j\pi}{N}))}{2(2 + \cos(\frac{j\pi}{N}))} \end{aligned} \right\} \quad (6.32)$$

Again, the values of \tilde{A}_{ij} have been calculated exactly and they are reported in Table 4b. The 20 highest values are in the shaded zone. The comparison between Tables 7a and b shows that these 20 values are approximately the same and they are associated with the same smooth modes. In this case, irregular modes will still be damped, and one can predict that no oscillation will occur.

When a load is only applied on the boundary ($f = 0, g \neq 0$),

\tilde{u}_{ij} is now

$$\tilde{u}_{ij} = \bar{A}_{ij} (g, \chi^{ij})_{0, \partial\Omega} \quad (6.33)$$

where

$$\bar{A}_{ij} = \frac{1}{\alpha_j \beta_i + \alpha_i \beta_j} \quad (6.34)$$

Again the values of \bar{A}_{ij} are reported in Table 4.c and the 20 highest values are in the shaded zone. Among these 20 values, three correspond to very irregular modes. In particular, the third value is associated with $\chi^{10,10}$. Therefore, we can predict a strong contribution of irregular modes within the solution \tilde{u}^h , which will show oscillations.

Finally, one could wonder if the calculation of the boundary integral can be calculated such that $(g, \chi^{ij})_{0, \partial\Omega}$ is damped for large i and j . Unfortunately, no precise method has been obtained. In particular, if the load g is a concentrated load at (x_0, y_0) , then

$$(g, \chi^{ij})_{0, \partial\Omega} = \chi^{ij}(x_0, y_0) \quad (6.35)$$

and this value is not necessarily zero. A procedure, consisting in forcing $(g, \chi^{ij})_{0, \partial\Omega}$ to zero in (6.35), can be obtained in the case in which the load is concentrated. This procedure is discussed in a more general context in the next subsection.

2.6.4. A-priori Orthogonalization of the Data. Our purpose in this subsection is to illustrate how simple orthogonalization consideration can drastically prevent the excitation of spurious modes. In particular, one interprets the existence of artificial anti-hourglass

forces that can be applied to damp spurious oscillations.

In order to solve the underintegrated system

$$\bar{A} \bar{U} = \bar{F} \quad (6.36)$$

we must have the orthogonality condition

$$\bar{F} \in (\ker \bar{A})^\perp \quad (6.37)$$

From a practical point of view, the applied load must be orthogonal to the spurious mode (compare with 6.35). If this is not the case, concentrated spurious forces may be induced at the nodes that have been fixed to solve the system *modulo* spurious modes.

This is clearly demonstrated in the following example where Ω is a square partitioned into 4 or 9 4-node elements. For the Laplace equation, we consider two loads concentrated in two opposite points and satisfying the equilibrium condition as shown in Fig. 11. When the number of elements is even (Fig. 11.a), the system of forces are orthogonal to the hourglass mode $H(\pm 1$ -pattern indicated at the nodes) and, therefore, the reactions at the fixed points are zero as desired. Conversely, when the number of elements is odd (Fig. 11.b), the system of forces is not orthogonal to H and two reactions R_1 and R_2 appear at the fixed points such that the system F_A, F_B, R_1 and R_2 is orthogonal to both translation and hourglass mode. Such considerations can explain the peculiar rates of convergence obtained on the unit square discretized with $N \times N$ elements with the data

$$f = \begin{cases} -6x(y^2-1)^2 + 4x^3(1-3y^2) & \text{if } x < \frac{1}{2} \\ \frac{3}{2}(y^2-1)^2 + \frac{1}{2}(-12x^2 + 24x - 7)(1-3y^2) & \text{if } x > \frac{1}{2} \end{cases} \quad (6.38)$$

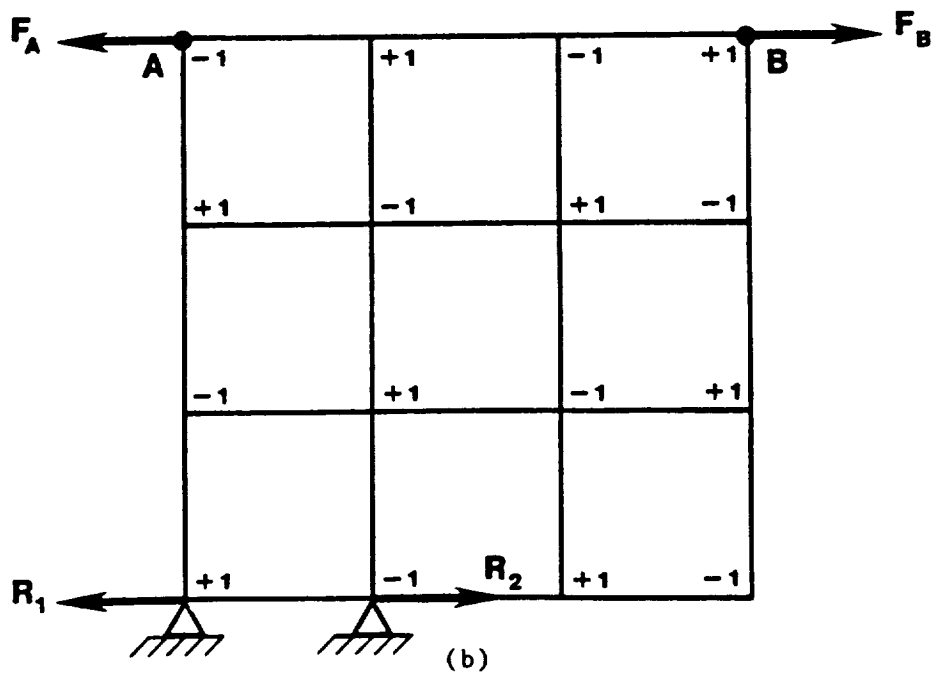
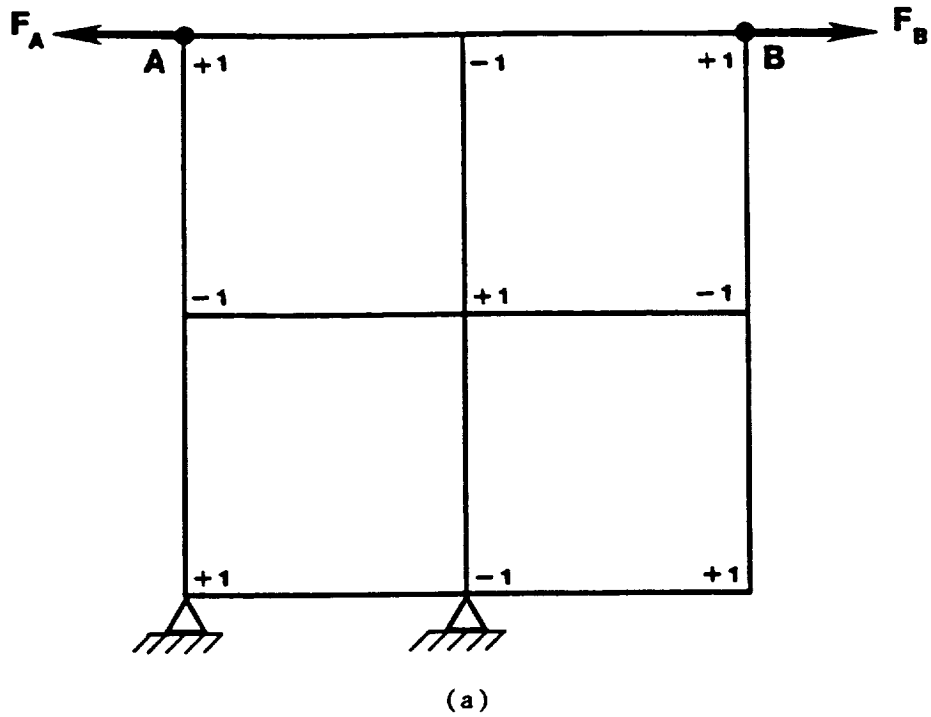


Figure 11: Spurious Reactions: a) Even Mesh and b) Odd Mesh

Figure 12.a shows the rates of convergence of \tilde{u}^h towards u^h in the L^2/\mathbb{R} - and H^1 norms, where \tilde{u}^h is obtained from the underintegrated problem and the projection, and u^h is the solution of the fully integrated problem. We can clearly see the good (more than optimal) behavior of the rates when the mesh is even; in this case, the discontinuity line ($x = \frac{1}{2}$) corresponds to a mesh line. However, the quality of the solution deteriorates when the mesh is odd or when the discontinuity line is across the mesh and coincides with the integration points. Note that when f_1 is averaged on $x = \frac{1}{2}$:

$$f(\frac{1}{2}, y) = \frac{1}{2}(f(\frac{1}{2}^+, y) + f(\frac{1}{2}^-, y)) \quad (6.39)$$

then all the points on both plots lay on one straight line. Also note that the knowledge of the exact solution

$$u = \begin{cases} x^3 (y^2 - 1)^2 & \text{if } x < \frac{1}{2} \\ \frac{1}{16}(-12x^2 + 24x - 7)(y^2 - 1) & \text{if } x > \frac{1}{2} \end{cases} \quad (6.40)$$

The interpretation of (6.37) has been possible when Neumann boundary conditions are applied. When Dirichlet conditions are applied, any right-hand side vector can produce a solution, but essentially the same type of behavior is observed (Figure 12.b). These results, the interpretations of (6.35) and (6.37) suggest that for any problem, one must pay attention to the data and make sure that it is orthogonal to the spurious modes. In particular, the procedure consisting in splitting a concentrated load between neighboring elements seems to give good results and will be discussed for elasticity problems in the next section.

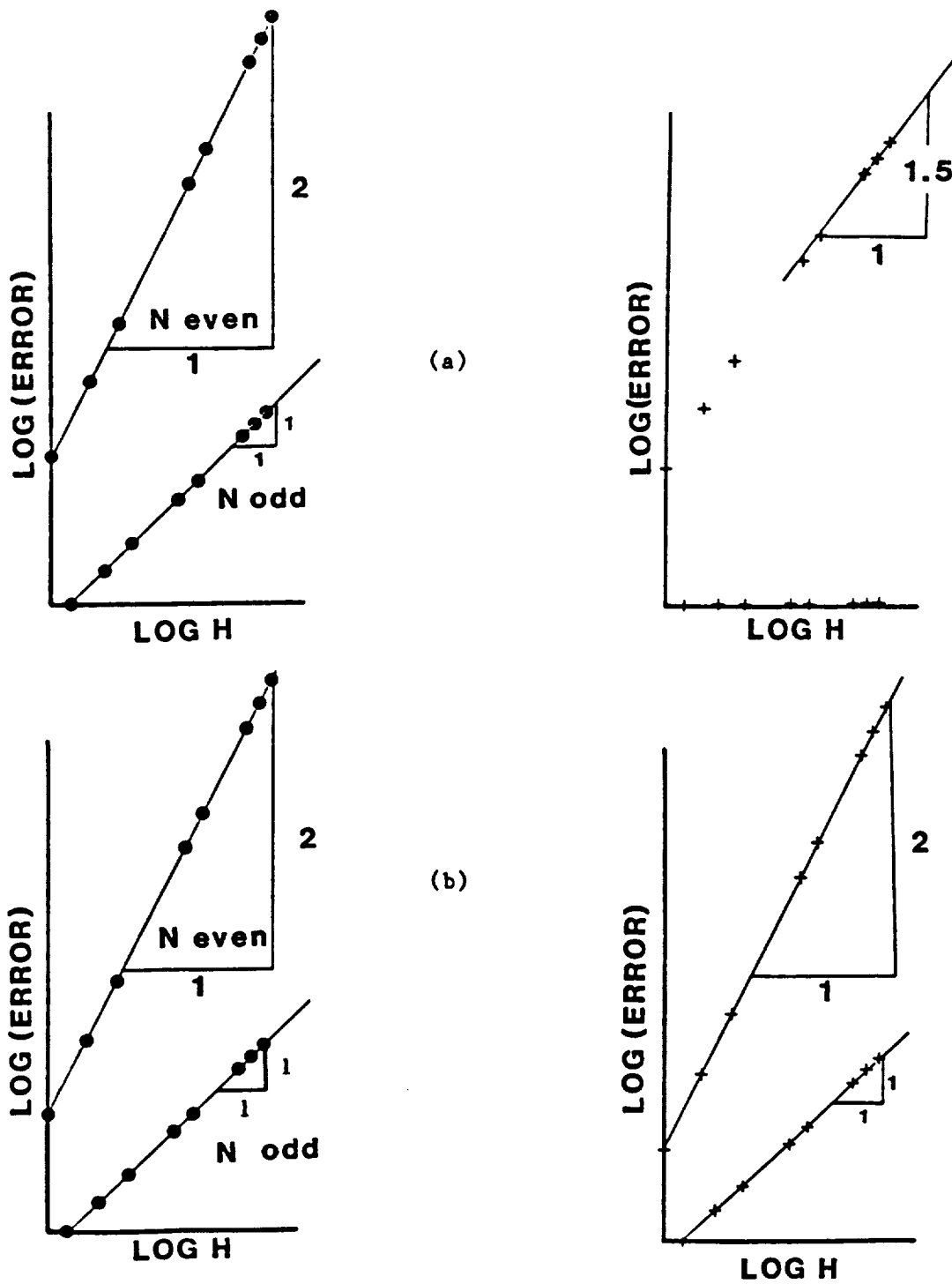


Figure 12 Rates of Convergence Obtained with a Discontinuous Data.

2.7 The Practice of Underintegration in Linear Elasticity

We devote this section to the discussion of the effects of underintegration in linear elasticity. We first exhibit the kernel of the underintegrated operator and obtain a post-processor formula similar to (3.30) to a-posteriori control the spurious modes. This global control can only be used in a very limited number of problems, but it suggests a local control that gives satisfactory results for all the examples considered. We discuss this control, its easy implementation, and several numerical results.

This study is entirely qualitative -- the basis functions obtained in Section 2.4 cannot be used to obtain eigen-functions for the elasticity operator. Both 4- and 9-node elements are discussed with an emphasis on the 9-node element in the numerical studies.

2.7.1. The kernel of the discrete underintegrated linear elasticity operator. We consider the linear elasticity operator defined by

$$\underline{A} = \underline{\beta}^T \underline{C} \underline{\beta} \quad (7.1)$$

where

$$\underline{\beta} = \begin{pmatrix} \frac{\partial}{\partial x} & 0 & \frac{\partial}{\partial y} \\ 0 & \frac{\partial}{\partial y} & \frac{\partial}{\partial x} \end{pmatrix}^T \quad (7.2)$$

and \underline{C} is a 3x3 symmetric matrix. In the plane strain case we may particularize \underline{C} :

$$\underline{C} = \begin{pmatrix} \lambda+2\mu & \lambda & 0 \\ \lambda & \lambda+2\mu & 0 \\ 0 & 0 & \mu \end{pmatrix} \quad (7.3)$$

In order to exhibit the spurious modes we consider the operator A associated with Neumann boundary conditions. In that case, the kernel of A consists of the usual 2-dimensional rigid body modes denoted by \underline{t}_x , \underline{t}_y and \underline{r}

$$\text{RBM} = \text{span} \left\{ \underline{t}_x = \begin{pmatrix} 1 \\ 0 \end{pmatrix}; \underline{t}_y = \begin{pmatrix} 0 \\ 1 \end{pmatrix}; \underline{r} = \begin{pmatrix} -y \\ x \end{pmatrix} \right\} \quad (7.4)$$

We consider the problem

P : Find $\underline{u} = \begin{pmatrix} u_1 \\ u_2 \end{pmatrix} \in [H^1(\Omega)]^2 / \text{RBM}$ such that

$$\int_{\Omega} \underline{u}^T \underline{\beta}^T \underline{C} \underline{\beta} \underline{v} \, dx dy = \int_{\Omega} \underline{f}^T \cdot \underline{v} \, dx dy \quad (7.5)$$

where $\underline{f} = (f_1, f_2)^T$ is a force satisfying the equilibrium conditions

$$\underline{f} \in \text{RBM}^T \quad (7.6)$$

or equivalently

$$\left. \begin{aligned} \int_{\Omega} f_1 \, dx dy &= \int_{\Omega} f_2 \, dx dy = 0 \\ \int_{\Omega} (f_1 y - f_2 x) \, dx dy &= 0 \end{aligned} \right\} \quad (7.7)$$

The existence and uniqueness of a solution for P are well known. The construction of finite element approximations of (7.5) involves the calculation of the $(2N \times 2N)$ stiffness matrix \underline{K}_e for a typical element Ω_e , which is given by the formula

$$\underline{K}_e = \int_{\Omega} \underline{N}^T \underline{\beta}^T \underline{C} \underline{\beta} \underline{N} \, dx dy \quad (7.8)$$

where \underline{N} is a vector representing the bilinear ($N=4$) or biquadratic ($N=9$) shape functions in each element Ω_e , $1 \leq e \leq E$. In computational applications \underline{K}_e is evaluated using an integration rule:

$$\underline{K}_e = \sum_{\alpha=1}^L w_{\alpha} \underline{B}^{\alpha T} \underline{C} \underline{B}^{\alpha} \quad (7.9)$$

where, similar to (2.36),

$$\underline{B}^{\alpha} = \begin{pmatrix} b_1^{\alpha} & 0 & b_2^{\alpha} \\ b_{-1} & & b_{-2} \\ 0 & b_2^{\alpha} & b_{-1} \end{pmatrix}^T \quad (7.10)$$

and w_{α} is the weight at the integration point α . Simple rank considerations allow us to predict the rank of \underline{K}_e . Indeed, since

$$\left. \begin{aligned} \text{rk}(A \cdot B) &\leq \max(\text{rk } A, \text{rk } B) \\ \text{and } \text{rk}(A + B) &\leq \text{rk } A + \text{rk } B \end{aligned} \right\} \quad (7.11)$$

we have

$$\text{rk } \underline{K}_e \leq 3L \quad (7.12)$$

When the full integration is used, (7.12) provides no information, but we know that $\underline{K}_e^{\text{full}}$ has the correct kernel containing only rigid body modes. However, when underintegration ($L=1$) is used on 4-node elements, we have

$$\text{rk } \underline{K}_e \leq 3 \quad (7.13)$$

Therefore, the 8x8 matrix \underline{K}_e possesses at least two spurious modes.

In fact, only two are present. Similarly, when underintegration ($L=4$) is used and 8- or 9-node elements, we have

$$\text{rk } \underline{K}_e \leq 12 \quad (7.14)$$

This inequality predicts one spurious mode for the 16x16 matrix associ-

ated with 8-node elements, but when the procedure is repeated for two neighboring 8-node elements, the spurious modes can no longer exist in the global matrix [53]. We can also interpret this elimination of the spurious mode by noticing that neighboring element cannot share the mode [14]. Also note that, in spite of failing an element stability test [24], there are no extraneous zero energy modes in the kernel of the underintegrated stiffness matrix and this element is stable.

As far as 9-node elements are concerned, the inequality (7.14) indicates that the 18x18 stiffness matrix has at least three spurious modes. In fact, there are exactly three such modes and they can be shared by adjacent elements. Next the modes will be explicitly described.

2.7.1.a. The Spurious Modes for 4-node Elements. Let \underline{H}_x and

\underline{H}_y be the two hourglass vectors defined as

$$\underline{H}_x = \begin{pmatrix} h \\ 0 \end{pmatrix} \quad \underline{H}_y = \begin{pmatrix} 0 \\ h \end{pmatrix} \quad (7.15)$$

where h^* is the hourglass nodal displacement defined in (2.14). Then when $L=1$, we obtain from (2.18) and (2.20)

$$\underline{B}_x^1 \cdot \underline{H}_x = \begin{pmatrix} b_1^T \cdot h \\ -1 & 0 \\ b_2^T \cdot h \\ -2 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ 0 \end{pmatrix}$$

and similarly

$$\underline{B}_y^1 \cdot \underline{H}_y = 0$$

* In this section, nodal values and associated functions will be denoted by the same letter, the underlining "-" differentiating them. The nodal values are expressed component by component.

Therefore,

$$K_{\sim e}^{(1)} \cdot H_{\sim x} = K_{\sim e}^{(1)} \cdot H_{\sim y} = 0 \quad (7.16)$$

These element displacements can be put together to obtain two global spurious modes, also denoted by $H_{\sim x}$ and $H_{\sim y}$ and we have :

$$\text{Ker } K^{\text{under}} = \{\text{span } t_{\sim x}, t_{\sim y}, r, H_{\sim x}, H_{\sim y}\} \quad (7.17)$$

This defines entirely the kernel of the underintegrated matrix and the spurious modes for 4-node elements.

Remark: In problems where symmetry is used for simplifications, the kernel of K^{under} must respect the symmetry. If one axis of symmetry (say, the x-axis) exists, then

$$\text{Ker } K^{\text{under}} = \text{span } \{t_{\sim x}, H_{\sim x}\}$$

If the problem has two axes of symmetry (x- and y-axis)

$$\text{Ker } K^{\text{under}} = \{0\}$$

The spurious modes are eliminated by the symmetry conditions.

2.7.1.b. The spurious modes for 9-node elements. Let $H_{\sim x}$ and $H_{\sim y}$ be the two vectors defined as

$$H_{\sim x} = \begin{pmatrix} h \\ 0 \end{pmatrix} \quad H_{\sim y} = \begin{pmatrix} 0 \\ h \end{pmatrix} \quad (7.18)$$

where h is the spurious mode of 9-node elements defined in (2.37).

Using (2.36), (2.39) and (7.9), we easily get

$$K_{\sim e}^{(4)} \cdot H_{\sim x} = K_{\sim e}^{(4)} \cdot H_{\sim y} = 0 \quad (7.19)$$

Therefore, $H_{\sim x}$ and $H_{\sim y}$ are two out of the three spurious modes of

$K_e^{(4)}$. We remark that the pattern (2.37) defining them does not depend on the geometry of the mesh. As far as the third spurious mode, denoted by

$$\underline{w} = \begin{pmatrix} w_1 \\ w_2 \end{pmatrix} \quad (7.20)$$

is concerned, one can show that the equations defining it are

$$\underline{b}_1^{\alpha T} \cdot \underline{w}_1 = \underline{b}_2^{\alpha T} \cdot \underline{w}_2 = \underline{b}_1^{\alpha T} \cdot \underline{w}_2 + \underline{b}_2^{\alpha T} \cdot \underline{w}_1 = 0 ; \alpha=1,4 \quad (7.21)$$

or equivalently :

$$\left. \begin{aligned} \underline{y}^T \underline{A}^{\alpha} \underline{w}_1 &= 0 \\ \underline{x}^T \underline{A}^{\alpha} \underline{w}_2 &= 0 \\ \underline{x}^T \underline{A}^{\alpha} \underline{w}_1 &= \underline{y}^T \underline{A}^{\alpha} \underline{w}_2 \end{aligned} \right\} \alpha = 1,4 \quad (7.22)$$

Note that for this system of 12 equations, we have 18 unknowns. If we add 5 orthogonality equations between \underline{w} and $\underline{t}_x, \underline{t}_y, \underline{r}, \underline{H}_x$ and \underline{H}_y , the system will define only one \underline{w} (up to within a multiplicative factor). For a general geometry of Ω_e , one cannot exhibit an explicit form for \underline{w} ; however, when Ω_e is a quadrilateral (strait-sided), and when \underline{x} and \underline{y} are of the form

$$\underline{x} = (x_1, x_2, x_3, x_4, \frac{1}{2}(x_1 + x_2), \frac{1}{2}(x_2 + x_3), \frac{1}{2}(x_3 + x_4), \frac{1}{2}(x_4 + x_1), \frac{1}{2}(x_1 + x_2 + x_3 + x_4)) \quad (7.23)$$

we can prove that one candidate for \underline{w} can be written as

$$\left. \begin{aligned} \underline{w}_1 &= \underline{T} \underline{y}' \\ \underline{w}_2 &= -\underline{T} \underline{x}' \end{aligned} \right\} \quad (7.24)$$

where

$$\tilde{T} = \begin{pmatrix} 4 & 2 & 0 & 2 & -1 & 0 & 0 & -1 & 0 \\ -2 & -4 & -2 & 0 & 1 & 1 & 0 & 0 & 0 \\ 0 & 2 & 4 & 2 & 0 & -1 & -1 & 0 & 0 \\ -2 & 0 & -2 & -4 & 0 & 0 & 1 & 1 & 0 \end{pmatrix}^T \quad (7.25)$$

and \underline{x}' , \underline{y}' are the vectors constructed with the first four components of \underline{x} and \underline{y} . An example of \underline{W} for a geometry satisfying (7.23) is shown in Figure 13 and can be constructed as follows:

- i) the displacement of a mid-side node is normal to the side, alternatively inwardly and outwardly oriented, with magnitude proportional to the length of the side.
- ii) the displacement of a corner is obtained by multiplication by -2 of the sum of the two displacements of the closest mid-side nodes.
- iii) the displacement of the centroid is zero.

On a square, the pattern of \underline{W} is well-known:

$$\tilde{W} = \begin{pmatrix} -2, 2, 2, -2, 0, -1, 0, 1, 0 \\ 2, 2, -2, -2, -1, 0, 1, 0, 0 \end{pmatrix} \quad (7.26)$$

Contrary to 8-node elements, and because of the presence of H_x and H_y , this mode can "propagate" from one element to another. For example, on a square mesh, if the nodal displacement vector is \underline{W} given by (7.26) on an element Ω_0 , then the displacement vectors $\underline{W} + 3H_x + t_x$ and $\underline{W} - 3H_y - t_y$ on the elements to the right of Ω_0 and above Ω_0 allow us to construct a continuous global displacement also denoted \underline{W} , on the mesh as shown in Figure 14.

We finally have

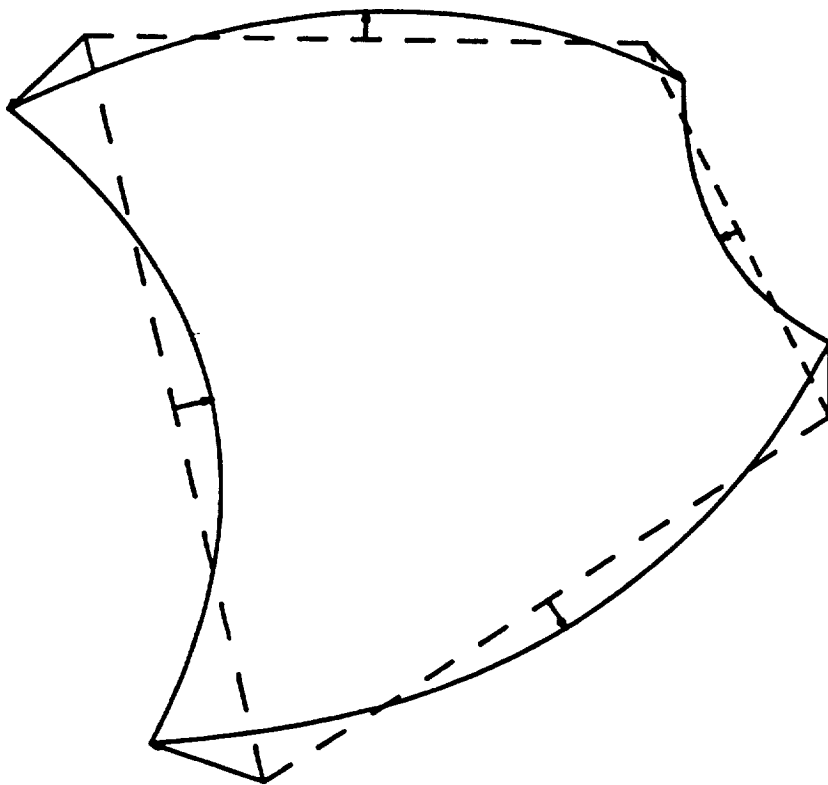


Figure 13. Spurious Mode W for a General Quadrilateral Element.

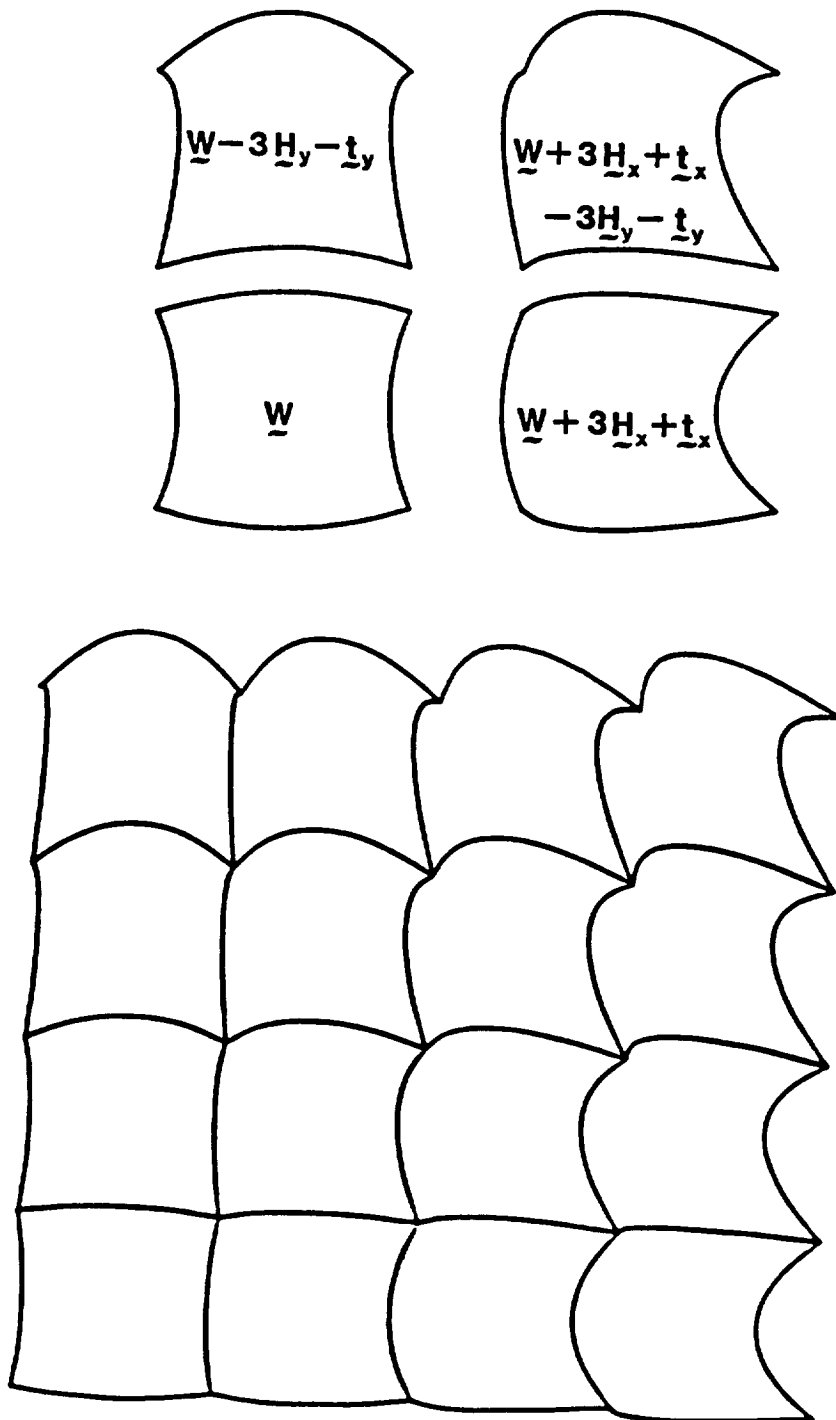


Figure 14. The Spurious Mode \underline{W} : (a) Construction,
 (b) The Mode on a 16 Element Mesh.

$$\text{Ker } \underline{K}^{\text{under}} = \text{span} \{ \underline{t}_x, \underline{t}_y, r, \underline{H}_x, \underline{H}_y, \underline{W} \} \quad (7.27)$$

Remark: Similar to what we have with 4-node elements, the existence of one axis of symmetry (say, the x-axis) reduces the kernel of the underintegrated stiffness matrix:

$$\text{Ker } \underline{K}^{\text{under}} = \text{span} \{ \underline{t}_x, \underline{H}_1, \underline{H}_2 + \underline{H}_3 \} \quad (7.28)$$

where

$$\left. \begin{aligned} \underline{H}_1 &= 3/2(\underline{H}_x - \underline{t}_x) \\ \underline{H}_2 &= -3/2(\underline{H}_y - \underline{t}_y) \\ \underline{H}_3 &= \underline{W} + 2(\underline{t}_x - \underline{t}_y) \end{aligned} \right\} \quad (7.29)$$

have been chosen such that the displacements of these modes are zero at the intersection of both axes for a square mesh. Contrary to the 4-node case, we still have a spurious mode when two axes of symmetry exist:

$$\text{Ker } \underline{K}^{\text{under}} = \text{span} \{ \underline{H}_1 + \underline{H}_2 + \underline{H}_3 \} \quad (7.30)$$

This mode is shown in Figure 15.

It is also important to point out that whereas the pattern of the spurious modes \underline{H}_x and \underline{H}_y are independent of both the geometry and the element, the mode \underline{W} depends upon both of them. Moreover, we can see by construction on a square mesh that the amplitude varies strongly when we consider successive elements. In fact, the pattern we may observe is a succession of pattern \underline{H}_x and \underline{H}_y with increasing amplitude.

2.7.2. Global A-Posteriori Control in Linear Elasticity. In this subsection we wish to generalize (3.30) with regard to the discrete operators, using various kernels discussed in the previous subsection.

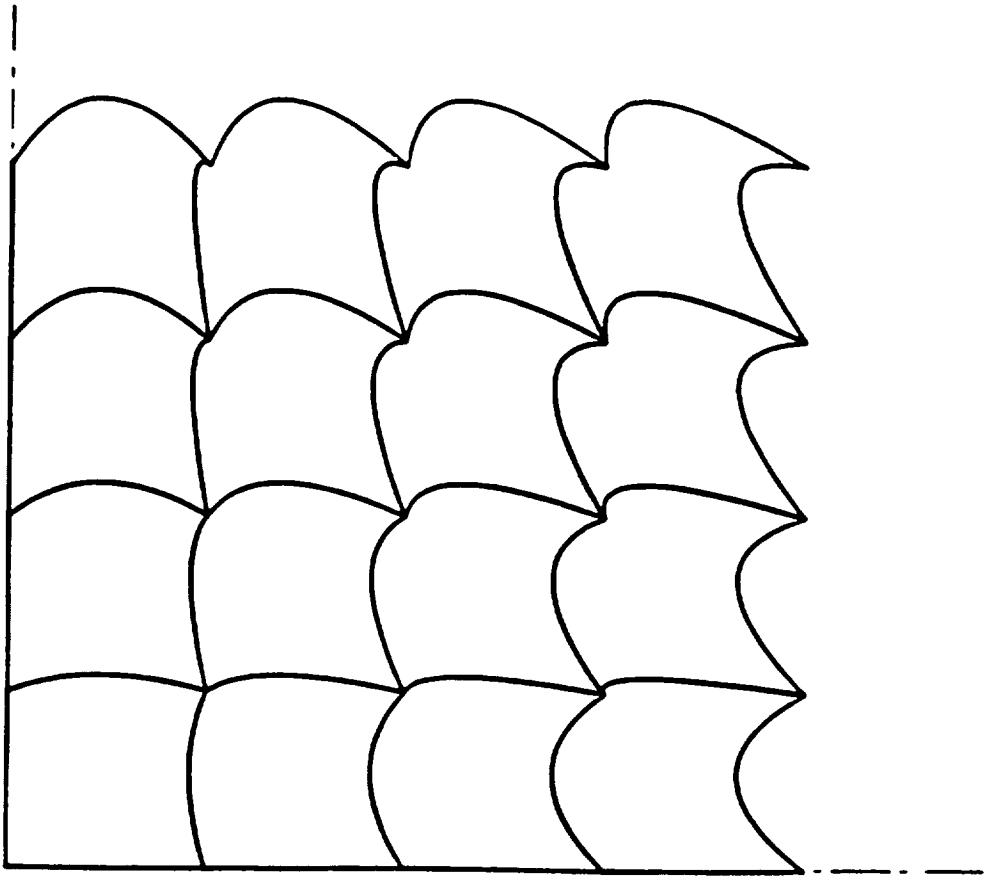


Figure 15. The Spurious Mode $H_1 + H_2 + H_3$.

We consider the general case where

$$\text{Ker } \tilde{K}^{\text{under}} = \text{RBM} \oplus \text{span}\{H_1, i = 1, I\} \quad (7.31)$$

where I may have the values 1, 2 or 3. We recall that for $I = 1$, we obtained a control formula similar to

$$\tilde{u}^h = \bar{u}^h - \frac{a(\bar{u}^h, H_1)}{a(H_1, H_1)} H_1 \quad (7.32)$$

where the bilinear form $a(\cdot, \cdot)$ was obtained in the variational formulation of the initial problem. This projection satisfies:

$$a(\tilde{u}^h, H_1) = 0 \quad (7.33)$$

or, in other words, \tilde{u}^h is orthogonal to the spurious mode. We generalize this property to the elasticity problem by supposing the projection to be orthogonal to all the spurious modes. Therefore the control will consist of looking for I constants λ_i ($i = 1, I$) such that

$$\begin{cases} \tilde{u}^h = \bar{u}^h - \sum_{i=1, I} \lambda_i H_i \\ a(\tilde{u}^h, H_i) = 0 \quad \text{for } i=1, I \end{cases} \quad (7.34)$$

This leads to the system of I equations with I unknowns :

Find λ_i , $i = 1, I$ such that

$$\sum_{j=1, I} \lambda_j a(H_j, H_i) = a(\bar{u}^h, H_i), \quad i=1, I \quad (7.35)$$

The computations involved in the control are computations of products of \bar{u}^h and the spurious modes by themselves. The implementation of these computations is discussed in the next section.

2.7.2.a Implementation of the Spurious Modes Control. For the computation of the coefficients in (7.35), we again use the decomposition

$$\underline{\underline{K}}^{\text{full}} = \underline{\underline{K}}^{\text{under}} + \underline{\underline{K}}^{\text{h}} \quad (7.36)$$

where $\underline{\underline{K}}^{\text{under}}$ satisfies

$$\underline{\underline{K}}^{\text{under}} \cdot \underline{\underline{H}}_1 = 0 \quad (7.37)$$

Then

$$a(\underline{\underline{u}}^{\text{h}}, \underline{\underline{H}}_1) = \sum_{e=1, E} \underline{\underline{U}}^{\text{T}} \cdot \underline{\underline{K}}^{\text{h}} \cdot \underline{\underline{H}}_1 \quad (7.39)$$

The expressions used for $\underline{\underline{K}}^{\text{h}}$ are next given for 4- or 9-node elements.

2.7.2.b. Control for 4-node Elements. For the operator defined in (7.1) and (7.2) with

$$\underline{\underline{C}} = \begin{pmatrix} C_{11} & C_{12} & C_{13} \\ C_{21} & C_{22} & C_{23} \\ C_{31} & C_{32} & C_{33} \end{pmatrix} \quad (7.40)$$

we have the *exact* decomposition for *any* geometry of Ω_e :

$$\underline{\underline{K}}_{\text{e}}^{\text{exact}} = \underline{\underline{K}}^{\text{under}} + \begin{pmatrix} \alpha_{11} \underline{\underline{Y}} \cdot \underline{\underline{Y}}^{\text{T}} & \alpha_{12} \underline{\underline{Y}} \cdot \underline{\underline{Y}}^{\text{T}} \\ \alpha_{21} \underline{\underline{Y}} \cdot \underline{\underline{Y}}^{\text{T}} & \alpha_{22} \underline{\underline{Y}} \cdot \underline{\underline{Y}}^{\text{T}} \end{pmatrix} \quad (7.41)$$

where

$$\left[\underline{\underline{\alpha}} \right] = \begin{pmatrix} C_{11} \bar{\epsilon}_{xx} + (C_{13} + C_{31}) \bar{\epsilon}_{xy} + C_{33} \bar{\epsilon}_{yy}; C_{13} \bar{\epsilon}_{xx} + (C_{12} + C_{33}) \bar{\epsilon}_{xy} + C_{32} \bar{\epsilon}_{yy} \\ C_{31} \bar{\epsilon}_{xx} + (C_{21} + C_{33}) \bar{\epsilon}_{xy} + C_{23} \bar{\epsilon}_{yy}; C_{33} \bar{\epsilon}_{xx} + (C_{32} + C_{23}) \bar{\epsilon}_{xy} + C_{22} \bar{\epsilon}_{yy} \end{pmatrix} \quad (7.42)$$

The vector $\underline{\gamma}$ and the $\bar{\epsilon}$'s are defined in Section 2 ((2.41) and (2.50)). For practical use, the expressions (2.51) are used for $\bar{\epsilon}$. For linear isotropic linear material, \underline{C} is given by (7.3) and

$$\begin{bmatrix} \alpha \\ \sim \end{bmatrix} = \begin{pmatrix} (\lambda+2\mu)\bar{\epsilon}_{xx} + \mu\bar{\epsilon}_{yy} & \mu\bar{\epsilon}_{xy} \\ \mu\bar{\epsilon}_{xy} & \mu\bar{\epsilon}_{xx} + (\lambda+2\mu)\bar{\epsilon}_{yy} \end{pmatrix} \quad (7.43)$$

This expression of $\begin{bmatrix} \alpha \\ \sim \end{bmatrix}$ can be compared to the general strain-stress relationship:

$$\begin{bmatrix} \sigma \\ \sim \end{bmatrix} = \begin{pmatrix} (\lambda+2\mu)\epsilon_x + \mu\epsilon_y & \mu\epsilon_{xy} \\ \mu\epsilon_{xy} & \mu\epsilon_x + (\lambda+2\mu)\epsilon_y \end{pmatrix} \quad (7.44)$$

An algorithm similar to the one presented in Section 2.5.1 can be constructed. It involves the computation of $\underline{\gamma}$, ϵ and $\underline{\alpha}$, then the computation of $a(H_i, H_j)$ and $a(\underline{u}^h, H_i)$, and finally the coefficients λ_i are obtained by resolution of a $N \times N$ system, N measuring the rank deficiency of K^{under} ($N=1$ or 2).

2.7.2.c. Control for 9-node Elements. In this subsection, devoted to 9-node elements, we first show why the results obtained by Belytschko are not sufficient to obtain a generalization of the linear elasticity problem, and then we propose an implementation of the control that leads to a stable solution converging to the exact solution with the optimal rate of convergence. However, for 9-node elements, we have not yet been able to obtain a computationally easy way to exhibit the third spurious mode, and the proposed results are only applicable to regular discretizations of a domain.

As far as the stabilization method proposed in [4] is concerned, algebra similar to that in Subsection 2.7.3.a leads to (7.41) where $\underline{\gamma}$ was defined in 2.41. But, whereas the stabilization matrix constructed with the submatrix $\underline{\gamma} \cdot \underline{\gamma}^T$ eliminates \underline{H}_x and \underline{H}_y from the kernel of the stiffness matrix, it does not take \underline{W} into account. Indeed, we have

$$\begin{pmatrix} \alpha_{11} \underline{\gamma} \cdot \underline{\gamma}^T & \alpha_{21} \underline{\gamma} \cdot \underline{\gamma}^T \\ \alpha_{12} \underline{\gamma} \cdot \underline{\gamma}^T & \alpha_{11} \underline{\gamma} \cdot \underline{\gamma}^T \end{pmatrix} \cdot \underline{W} = \underline{0}$$

Therefore this procedure cannot be used to control \underline{W} .

In order to obtain an accurate control, we have to consider a generalization of (2.43). Now we have

$$\begin{pmatrix} \alpha_{11} \underline{s}_i \cdot \underline{s}_i & \alpha_{21} \underline{s}_i \underline{s}_i \\ \alpha_{12} \underline{s}_i \cdot \underline{s}_i & \alpha_{22} \underline{s}_i \underline{s}_i \end{pmatrix} \cdot \underline{H}_j \neq \underline{0}$$

for $7 < i < 9$
 $1 < j < 3$

where the vectors \underline{s}_i are defined in Section 2.2.3. Finally, using (2.43) and (2.52), we have

$$\begin{aligned} \underline{K}_e^{(9)} = & \underline{K}_e^{(4)} + \frac{4\Omega_e}{135} \begin{bmatrix} (C_{11}+C_{33}) \underline{s}_9 \underline{s}_9^T & (C_{13}+C_{32}) \underline{s}_9 \underline{s}_9^T \\ (C_{31}+C_{23}) \underline{s}_9 \underline{s}_9^T & (C_{22}+C_{33}) \underline{s}_9 \underline{s}_9^T \end{bmatrix} \\ & + \frac{\Omega_e}{45} \left(\begin{bmatrix} C_{11} \underline{s}_7 \underline{s}_7^T & C_{13} \underline{s}_7 \underline{s}_7^T \\ C_{31} \underline{s}_7 \underline{s}_7^T & C_{33} \underline{s}_7 \underline{s}_7^T \end{bmatrix} + \begin{bmatrix} C_{33} \underline{s}_8 \underline{s}_8^T & C_{32} \underline{s}_8 \underline{s}_8^T \\ C_{23} \underline{s}_8 \underline{s}_8^T & C_{22} \underline{s}_8 \underline{s}_8^T \end{bmatrix} \right) \end{aligned} \quad (7.45)$$

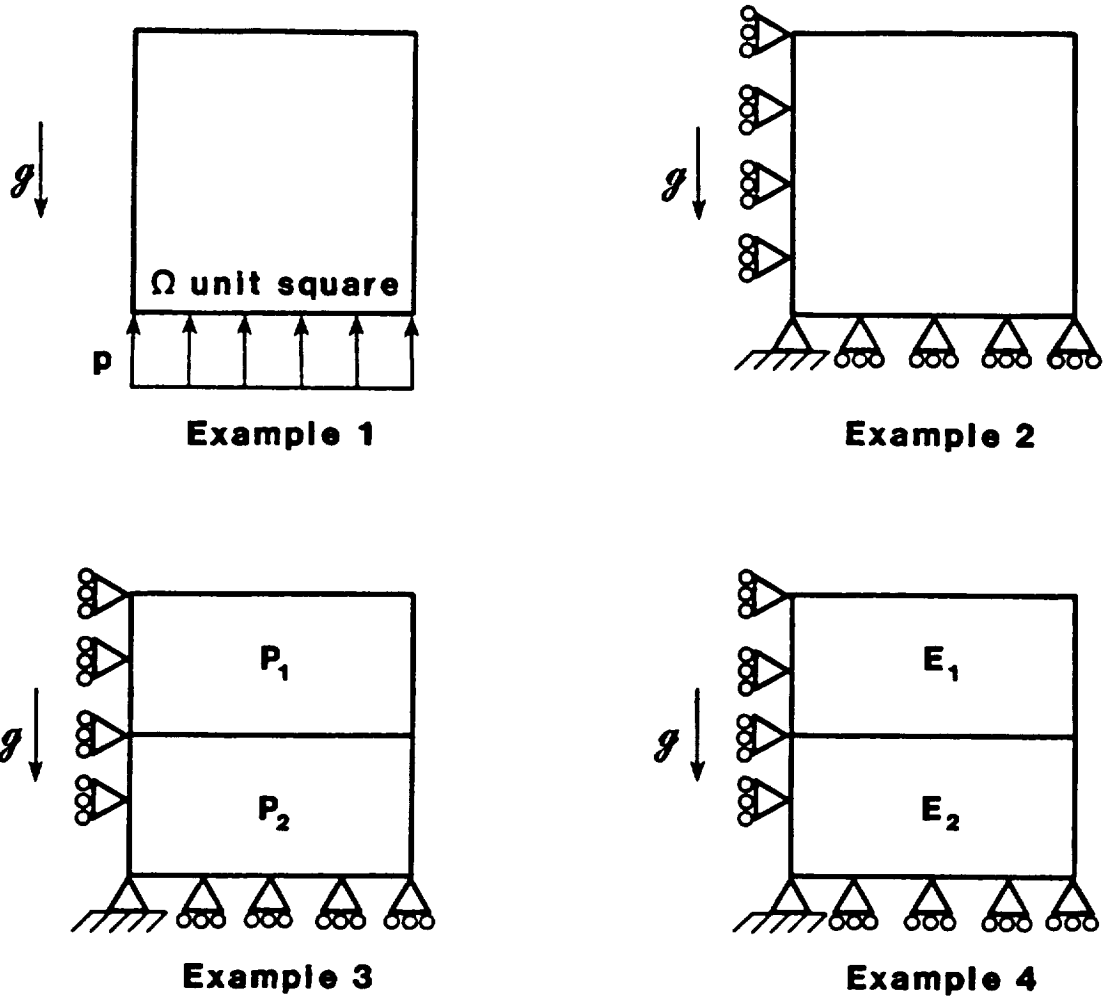
Similarly, for the 4-node case, the algorithm for the computations

of the coefficients in (7.35) has been obtained and implemented. Numerical results agree with our presumptions concerning a " $\underline{\gamma} \cdot \underline{\gamma}^T$ "-type of control and incline in favor of the decomposition (7.45). On a square domain discretized with $N \times N$ elements, we have calculated and compared the solutions obtained with full (exact) and underintegration for various boundary and symmetry conditions. Various examples considered are described in Figure 16. The rates of convergence were calculated by comparing the error norms ($\delta=0$: L^2 /RBM norm; $\delta=1$: energy norm) obtained with $N=5,6$ and 7 . We consistently obtained the rate $O(h^{2-\delta})$ using a $\underline{\gamma} \cdot \underline{\gamma}^T$ decomposition and $O(h^{3-\delta})$ with (7.45) for homogeneous materials under the action of gravity ($\underline{f} \in C^\infty$); the order $3-\delta$ being optimal, we may conclude that the method presented below is accurate. It is also efficient: for one second taken for the fully integrated stiffness matrix, only .61 are taken when the underintegration is used and only .05 seconds are taken for the control. Also note that the Example 4 in Figure 30 is also optimal ($u \in H^1/H^2$).

Unfortunately, this control is far from general. In particular, it can only be used when the exact shape of the spurious modes is known, which is the case only when Neumann (traction) boundary conditions are applied on a domain discretized with regular square elements.

2.7.3. A Local Control of the Spurious Modes

2.7.3.a. Introduction. The major drawbacks of the global control are overcome by considering the procedure consisting of eliminating, element by element, the components of \underline{H}_x , \underline{H}_y and \underline{W} and then of averaging the nodal values obtained in neighboring elements. For any element, we choose to do the following simplifications:



Example	Number of Spurious Modes	Rate of Convergence	
		L^2/\mathbb{R}	H^1
1	3	3	2
2	1	3	2
3	2	3	2
4	1	2	1

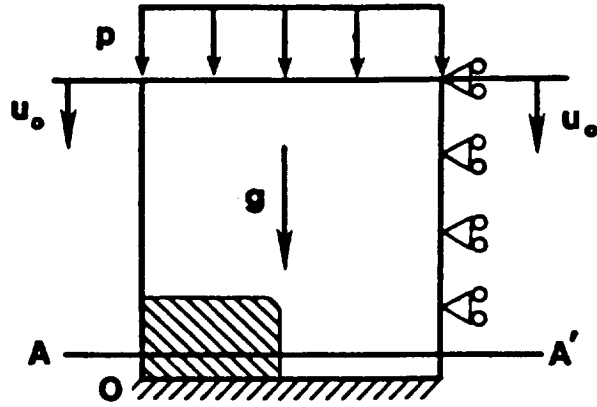
Figure 16. Results Obtained with the Global Control.

- i) the nodal values of \underline{W} are taken as if the element was strait-sided quadrilateral
- ii) \underline{C} is diagonal.
- iii) \underline{K}^{stab} is given by (7.41) or (7.45).

These simplifications lead to simple calculations detailed in Appendix B for the 9-node element and that are easily implemented in the subroutine listed in Appendix C. Note that the expression obtained for the control is uniquely geometric, does not depend upon the material properties and can be used in any linear or nonlinear problem. It only requires that the shape of the element not deviate too much from a quadrilateral.

2.7.3.b. Numerical Results. The examples displayed in this section illustrate the efficiency and the accuracy of the local control previously described. However, we only consider linear elastic material in plane strain, on domains discretized with biquadratic (9-node) elements. Three examples are described.

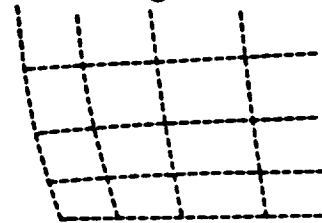
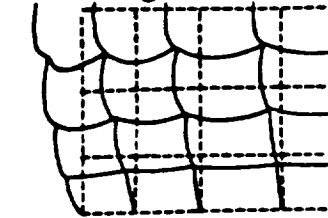
The first example is defined in Figure 17.a. We consider a square domain with one fixed side, under the action of pressure, gravity or a prescribed compressive displacement. Under any of these loads, a singularity appears at the neighborhood of the origin. Figure 17.b shows how the underintegrated solution behaves in the singularity region and how the control affects the results. Whereas the underintegrated solution shows oscillations, the displacements obtained after control are smooth and similar to those obtained with the full integration of the stiffness matrix. The shear along a line AA' across the singularity region is also shown in Figure 18. Whereas undesirable oscillations



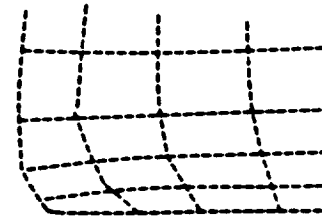
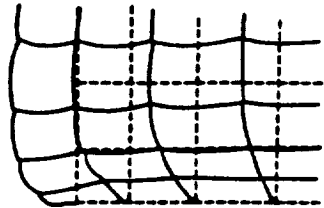
(a)

**Undeformed Configuration
and
Underintegrated Solution**

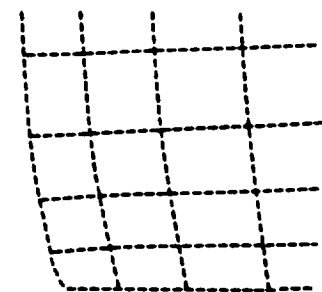
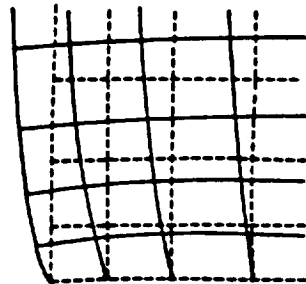
**Fully and Controlled-
Underintegrated Solution**



$p \neq 0$



$g \neq 0$



$u_0 \neq 0$

(b)

Figure 17. Displacement of a Fixed Square :
a) Definition of the Problem; b) Results.

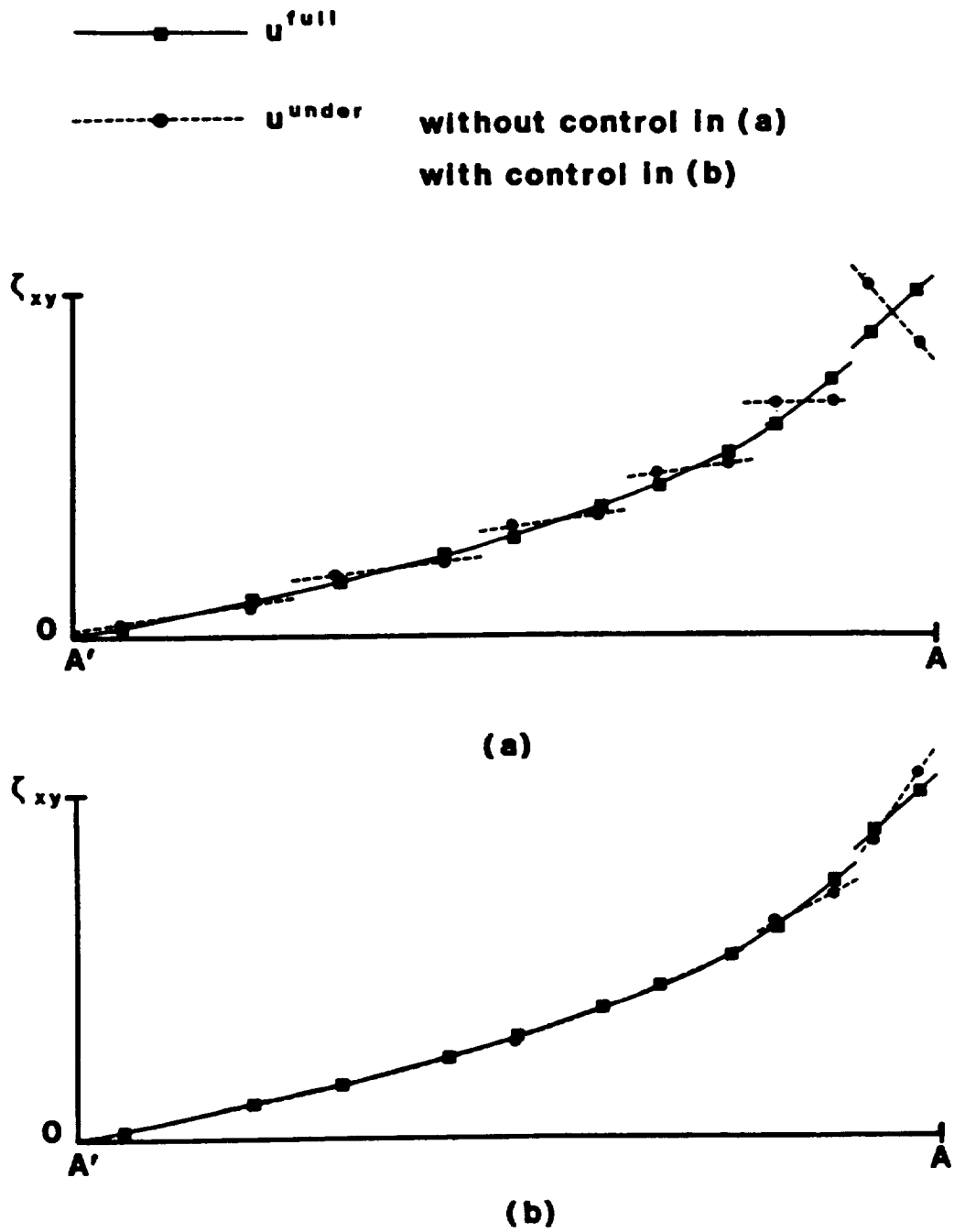


Figure 18. Shear along a Line AA'.

are observed before control, the shear behaves properly after control.

In the second example, we consider a ring under the action of an external pressure (Figure 19). Here again the oscillations generated by the underintegration are damped when the control is applied. Only very slight oscillations remain, not exceeding 5%, and these can be easily interpreted: in the control, the expression taken for the mode \underline{W} was obtained for quadrilateral elements. For the mesh considered, the elements are slightly bent and this difference explains these slight oscillations. The same domain (quarter ring) has also been discretized using quadrilateral elements and the control of the underintegrated solution has led to a displacement field without any oscillations and similar to the underintegrated displacements. Calculations of the stress along a radius show behavior identical to the previous example.

The third example involves a concentrated force and illustrates our discussion concerning the excitation of spurious modes and the orthogonalization of the data. A point force is applied at a corner of a fixed side square discretized with a mesh refining in the neighborhood of the singularity. Strong oscillations appear in this region when underintegration is used, whereas the full integration solution is smooth (Figure 20). These oscillations show a pattern similar to the one used to construct the mode \underline{W} (Figure 14): amplification of the mode $3\underline{H}_x + \underline{t}_x$ (respectively $3\underline{H}_y + \underline{t}_y$) along the x- (resp. y-) direction. Therefore, according to the previous interpretations of (6.35) and (6.37), a way to prevent oscillations is to consider a system of loads similar to the load concentrated in a point A but ortho-

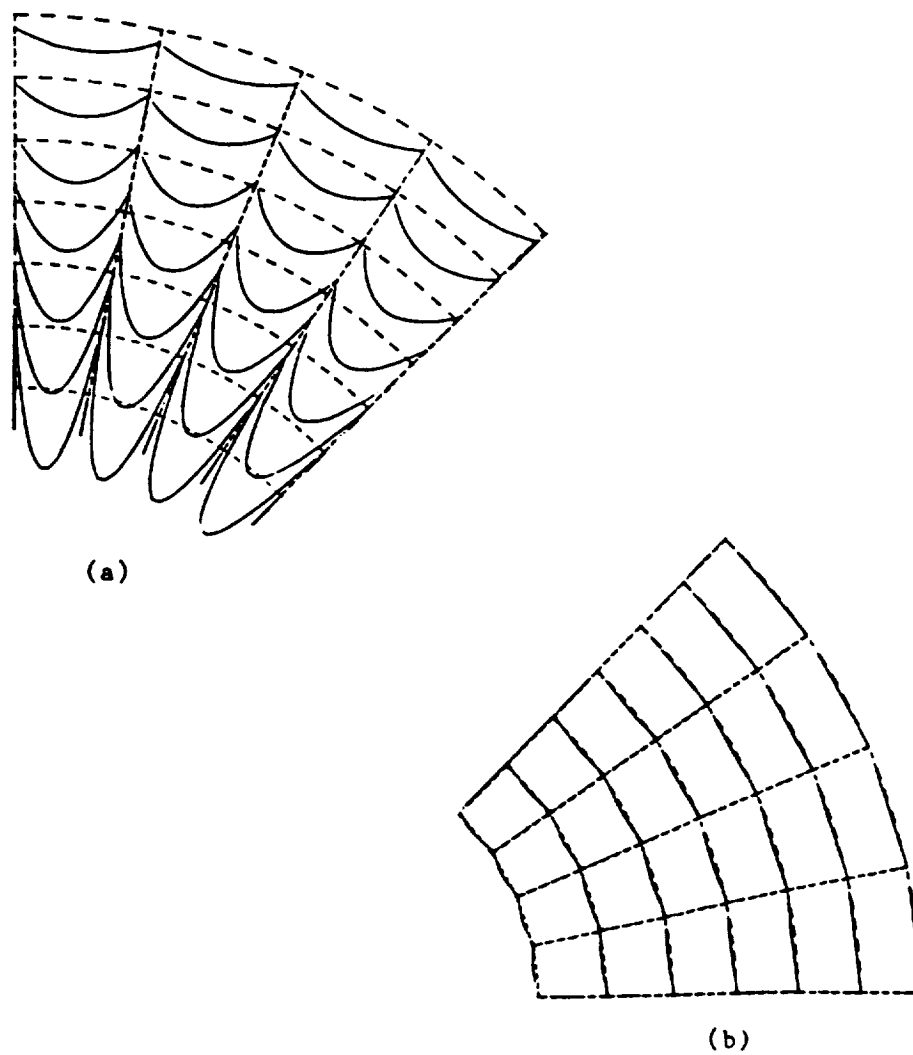
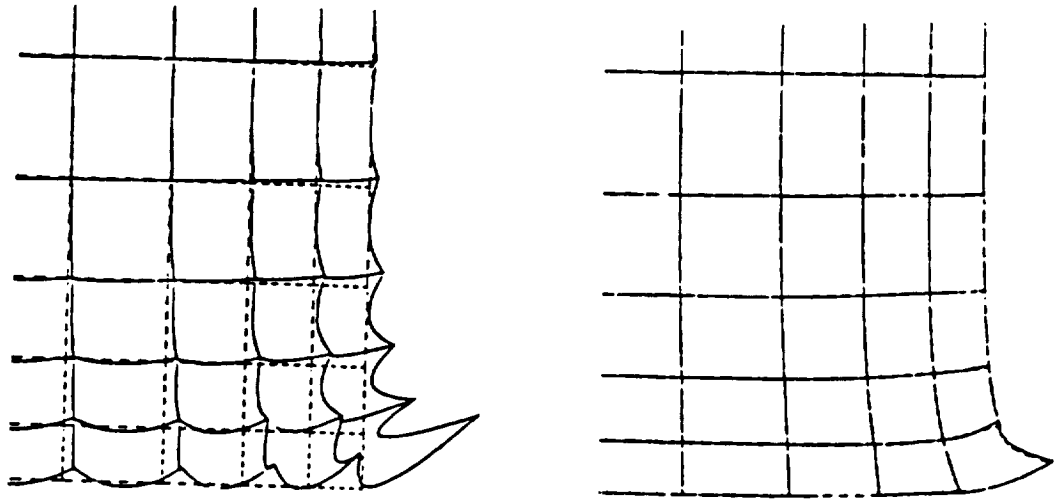


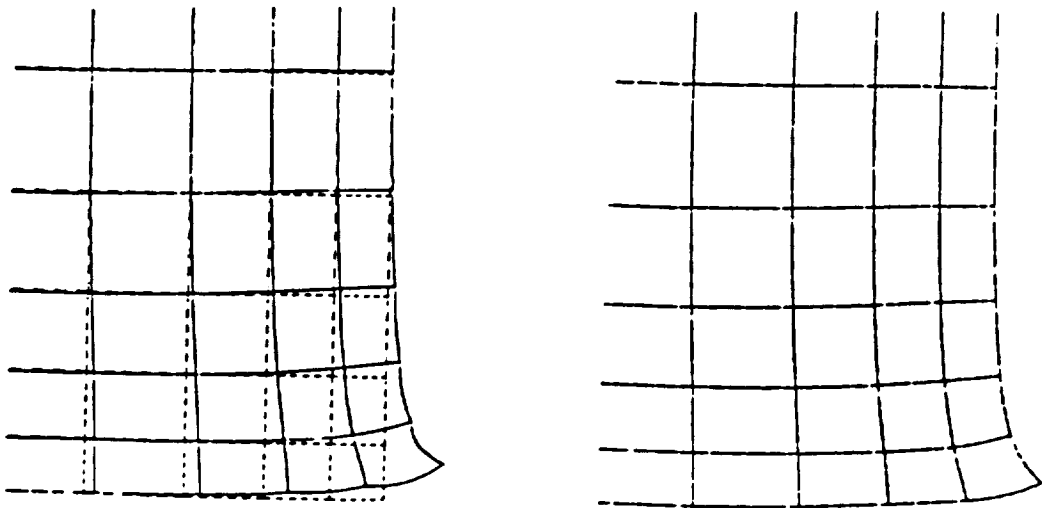
Figure 19. Quarter Ring Under External Pressure

(a) Undeformed Configuration and Underintegrated Solution, and

(b) Fully- and Controlled Underintegrated Solution.



(a)



(b)

Figure 20. Displacements Due to a Concentrated Load: Undeformed, Underintegrated, Controlled and Fully Integrated Solution Due to a) the Force, and b) its Orthogonalized Equivalent.

gonal to $3\mathbf{H} + \underline{\underline{t}}$. This is obtained by splitting the force into 3 equal forces applied at A and its two closest nodes. Indeed, the displacements obtained with this system of load only show slight oscillations. Finally, note that this control produces displacement fields similar to the fields obtained using full integration.

PART III: APPLICATION TO NONLINEAR
INCOMPRESSIBLE ELASTICITY

3.1 Introduction

In this concluding part, we analyze instabilities observed in discrete solutions of nonlinear problems in finite elasticity involving incompressible materials. We compare the behavior of the biquadratic (9-node) and isoparametric (8-node) elements associated with linear (P_1) discontinuous pressures. Then we focus on 9-node elements to discuss the efficiency and accuracy of control of the underintegrated solution introduced in Part II for linear operators.

Only Mooney-Rivlin materials are considered. They are characterized by the strain energy function

$$\sigma = C_1(I_1 - 3) + C_2(I_2 - 3) \quad (3.1)$$

where I_i , $i=1,2,3$ are the principal invariants of the Cauchy-Green deformation tensor

$$\underline{C} = (\underline{1} + \underline{\nabla u})^T (\underline{1} + \underline{\nabla u}) \quad (3.2)$$

where $\underline{1}$ is the unit tensor and $\underline{\nabla u}$ is the displacement gradient.

They are:

$$\begin{aligned} I_1 &= \text{tr } \underline{C} \\ I_2 &= \frac{1}{2}(\text{tr } \underline{C}^2 - (\text{tr } \underline{C})^2) \\ I_3 &= \det \underline{C} \end{aligned} \quad (3.3)$$

The condition of incompressibility can be expressed as

$$I_3 = 1 \quad (3.4)$$

and is taken into account in a mixed formulation of the equilibrium problem by introducing a Lagrange multiplier P . The energy function σ can be replaced by σ^{Lag} :

$$\sigma^{\text{Lag}} = C_1(I_1 - 3) + C_2(I_2 - 3) + P(\sqrt{I_3} - 1) \quad (3.5)$$

For this choice of energy function, P is the hydro-static pressure.

We consider the usual virtual work equilibrium equation:

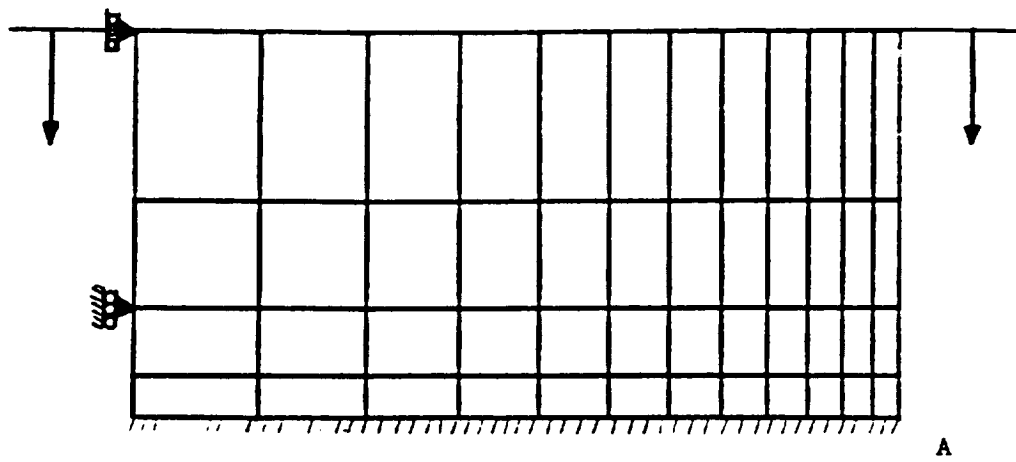
$$\begin{aligned} \delta \int_{\Omega_0} \sigma^{\text{Lag}} dv_0 - \int_{\Omega_0} \rho_0 \underline{f} \cdot \delta \underline{u} dv_0 \\ - \int_{\partial\Omega_2} \underline{t} \cdot \delta \underline{u} ds = 0 \end{aligned} \quad (3.6)$$

The solution of this highly nonlinear problem is accomplished in this work using Newton's method. Details of the finite element method applied to this particular class of problems are discussed at length in the book of Oden [38]. A more recent account is given in Aly [1].

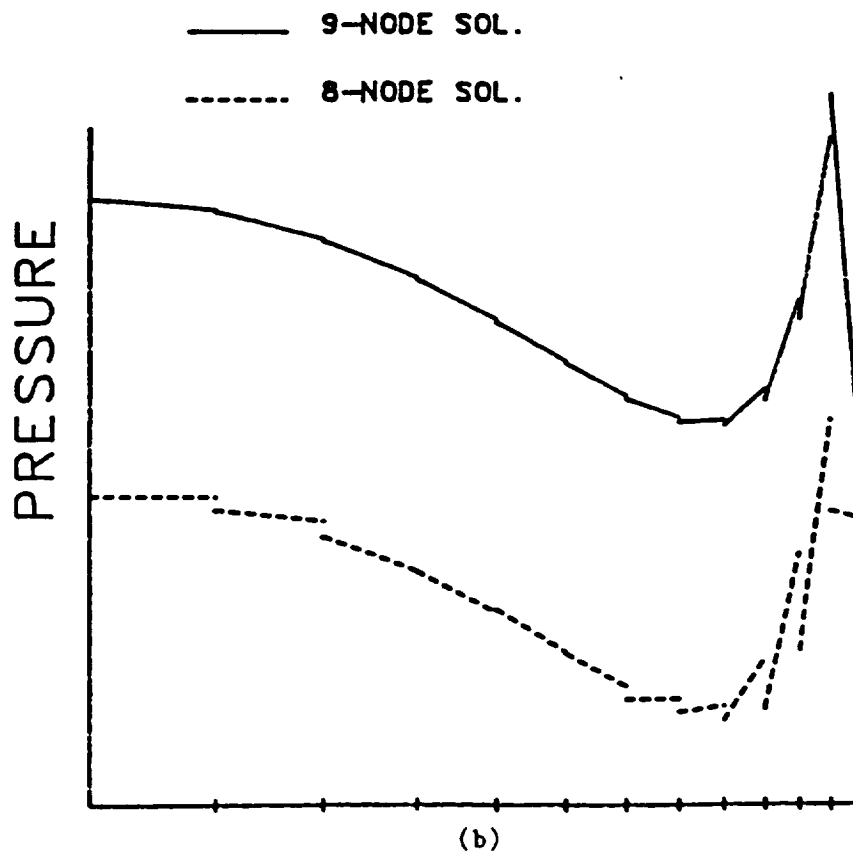
3.2 Behavior of 8- and 9-Node Elements (Full Ingegration)

In this section, we briefly review some observations made by Miller [37] that are now clearly understood with the results obtained by Oden and Jacquotte [42, 43].

We consider a fixed side rectangular domain discretized with the refining mesh shown in Figure 21.a. We compress this domain imposing a displacement on the top side. The displacement increments are 5, 10, 15 and 17.5% of compression with respect to the original shape.



(a)



(b)

Figure 21. Compression of a Square a) Mesh at Definition in Problem; b) Pressure Obtained with 8- and 9-Node Elements ($u_0 = 5\%$).

The solution of this problem is complicated by a stress singularity that occurs in the neighborhood of the corner A . Both 8- and 9-node displacement elements were tested associated with a linear discontinuous Lagrange multiplier. The displacements show similar behaviors. As far as pressure is concerned, oscillations similar to the pattern of $\ker B_h^*$ in [44] are observed (Figure 21.b) when 8-node elements are used, whereas the pressure distribution is smoother with 9-node elements. Finally, note that the nodal averaging technique described in Section 1.8 has been tested by Miller [37] for a problem with similar singularity and his conclusions corroborate ours from Part I.

3.3 Control of the 9-Node Underintegrated Element

In this section, we analyze how the displacements and pressures behave when underintegration is used. We noticed in Section 2.7 how the control obtained was only geometric. This is particularly useful for rubber-like materials and makes it very efficient when constitutive relations are numerically expensive to obtain. Also note that the calculation of the projected element solution (Appendix B) involves the knowledge of the mode \underline{w} which is computed using the *current* deformed geometry of the element. We may foresee what will be one of the major drawbacks of the method: \underline{w} is only exactly known for quadrilateral elements. When the element is too distorted, the approximation we do assuming it to be quadrilateral is too poor and the control is not accurate.

Finally, we point out that the control is applied at each load increment. We present three examples.

3.3.1 Stretch of a Rubber Material Domain. The first example involves the stretch of a rubber square (Fig. 22). The domain is partitioned in 49 elements (450 degrees of freedom). The material is a Mooney Rivlin material with

$$\begin{aligned} C_1 &= 100 \\ C_2 &= 20 \end{aligned} \tag{3.7}$$

and the incremental stretches are 25, 50, 75 and 100% of its undeformed configuration. Whereas the underintegrated solution shows slight oscillations of the displacement in the singularity region, the controlled solution is smooth and similar to a fully integrated solution. As far as pressure is concerned, similar observations to the ones in Section 2.7.3.b for stresses can be made. However, the convergence is obtained slower: for the various displacement increment, 6, 4, 4 and 4 (respectively 6, 5, 5 and 4) Newton iterations have been needed to obtain convergence with the full (resp. under) integration. Nevertheless, the gain in time is almost 40% (673 sec. versus 413 sec.).

3.3.2 Behavior in the Neighborhood of a Concentrated Force. This second example illustrates the ability of the orthogonalization of the data to obtain an underintegrated solution. We consider the same Mooney Rivlin material (1.1, 1.7) and the problem defined in Figure 23.a. The application of the force at only one point leads to the behaviors:

- i) Slight oscillations are observed in the singularity region when full integration is performed (Fig. 23.b)
- ii) Uncontrollable oscillations appear when underintegration is

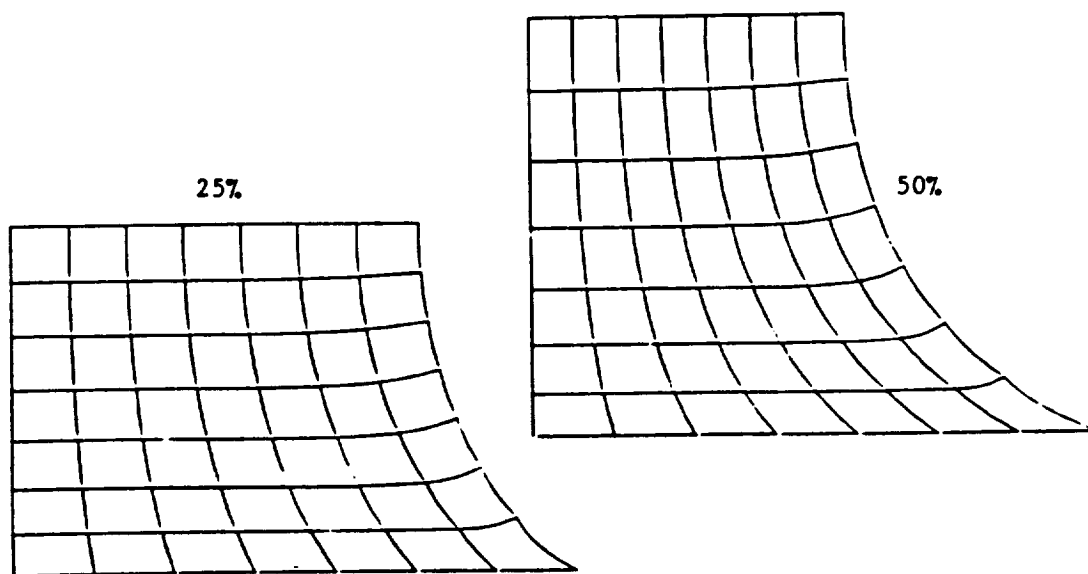


Figure 22.a. Underintegrated Solution.

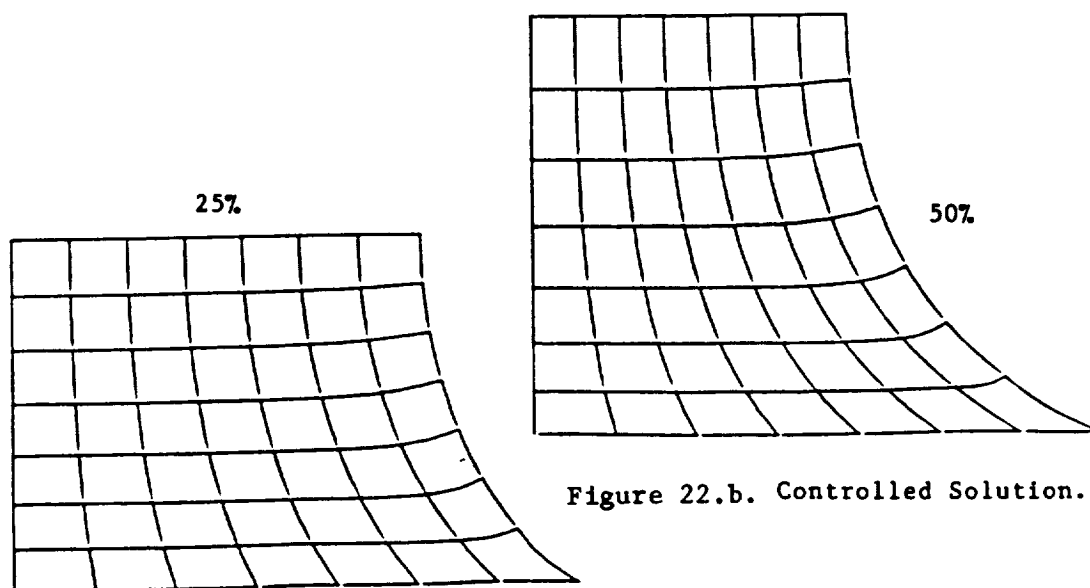


Figure 22.b. Controlled Solution.

Figure 22. Stretch (25, 50%) of a Rubber Material Domain.

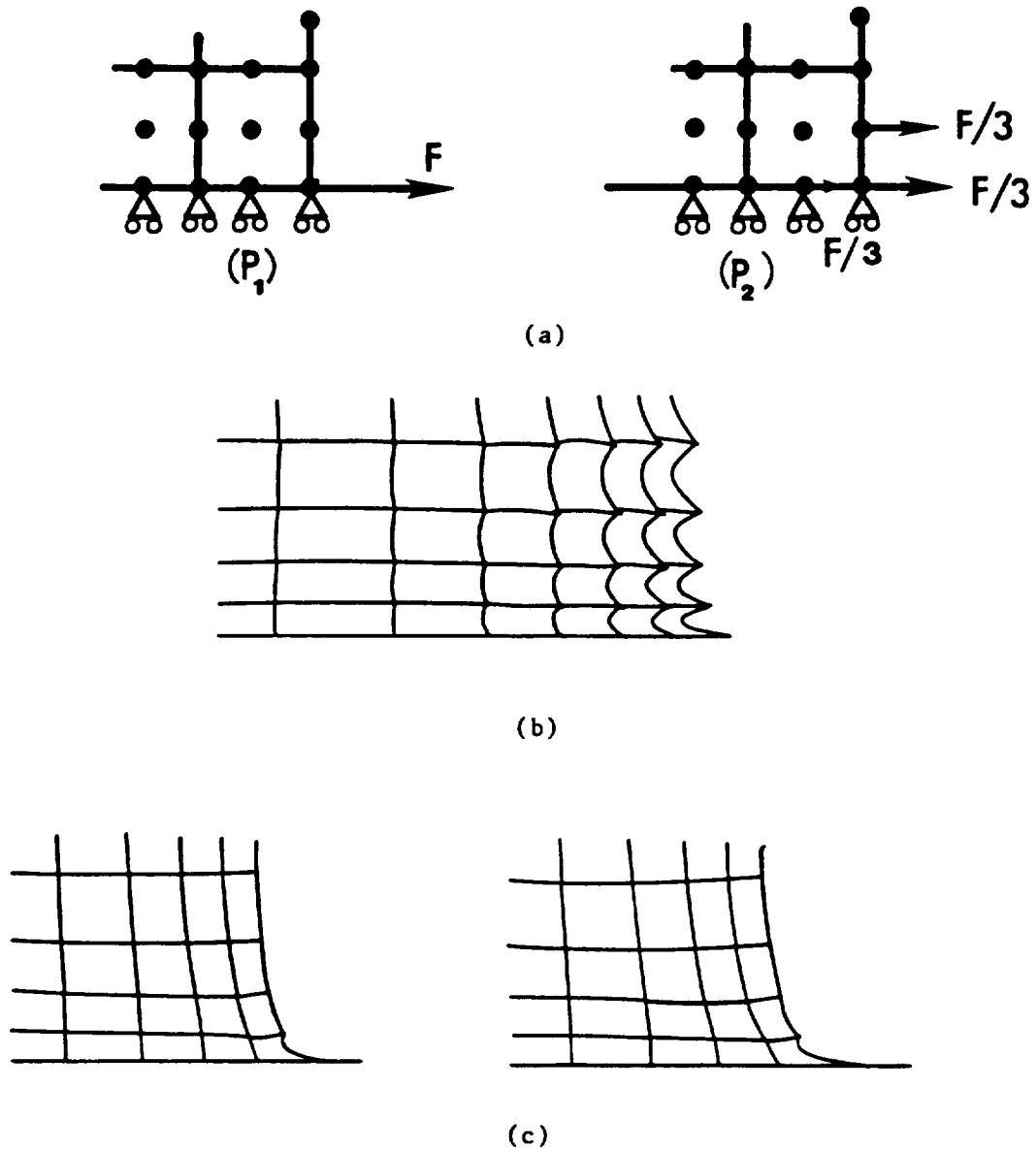


Figure 23. Behavior in the Neighborhood of a Concentrated Force: a) Definition of the Problems; b) Underintegration of P_1 ; c) Controlled Solutions at two Increments.

used: Figure 23.c shows the displacements after only 3 iterations for the first load increment. Later the solution diverges.

As in the linear case discussed in Section 2.7.3.b, we split the concentrated force into three equal forces at the closest nodes; we observed that (Fig. 24)

iii) When underintegration is performed, without control, the solution converges, but oscillations still develop (Fig. 24.a).

iv) The control of this solution is smooth and similar to the one obtained with full integration (Figs. 24.b and c).

For this example, one supplemental iteration was needed in the third load increment and the gain in time also approaches 40%.

3.3.3 Compression of a Fixed Rubber Material Domain. This final example reconsiders the problem used to compare the performances of the 8- and 9-node elements (Section 3.2). We consider two discretizations of the domain with 25 and 49 elements, and displacement increments of 5, 10, 15 and 17.5% with respect to the original shape. For the crude mesh, oscillations appear very soon when underintegration is performed, but the control easily corrects the solution and a displacement field close to the fully integrated field is obtained (Fig. 25.a). But when the mesh is refined, the oscillations become more important and deform the element to a degree such that the control is not able to restore the shape of the element corner (Fig. 26.b). We interpret this lack of performance to the fact that the control has been exactly obtained for quadrilateral elements. In this case, the element sides curve and the element is too deformed. Also this lack

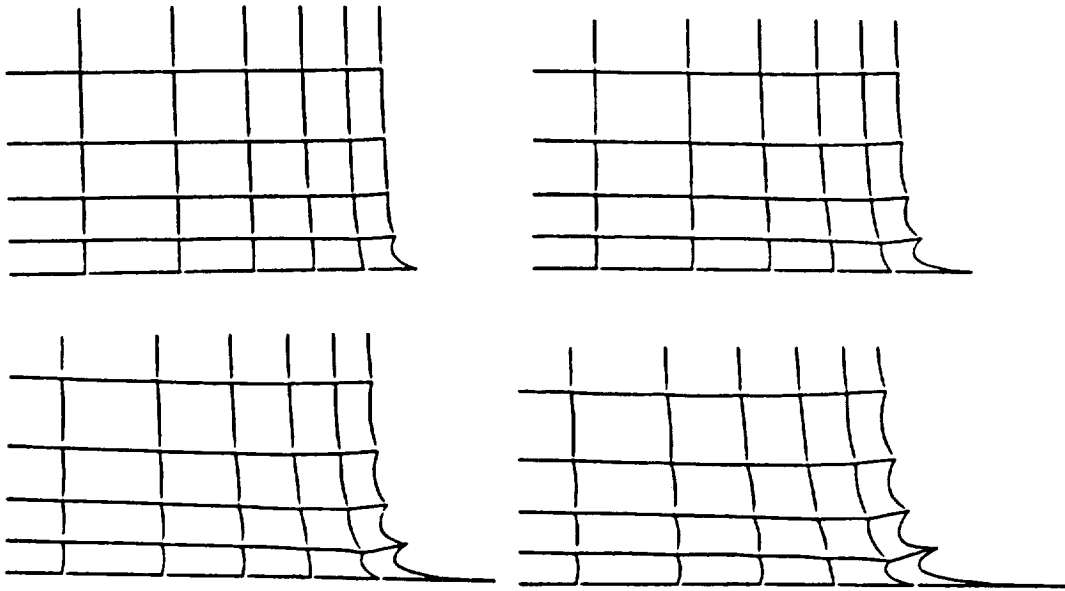


Figure 24.a. Underintegration of P_2 .

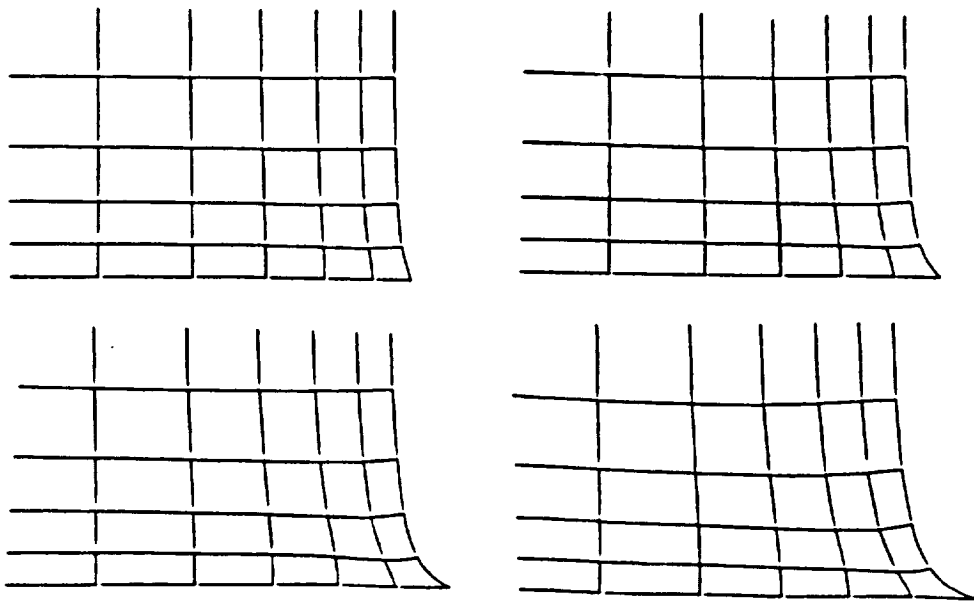


Figure 24.b. Controlled Underintegrated Solution of P_2 .

Figure 24. Behavior in the Neighborhood of a Concentrated Force (Problem P_2).

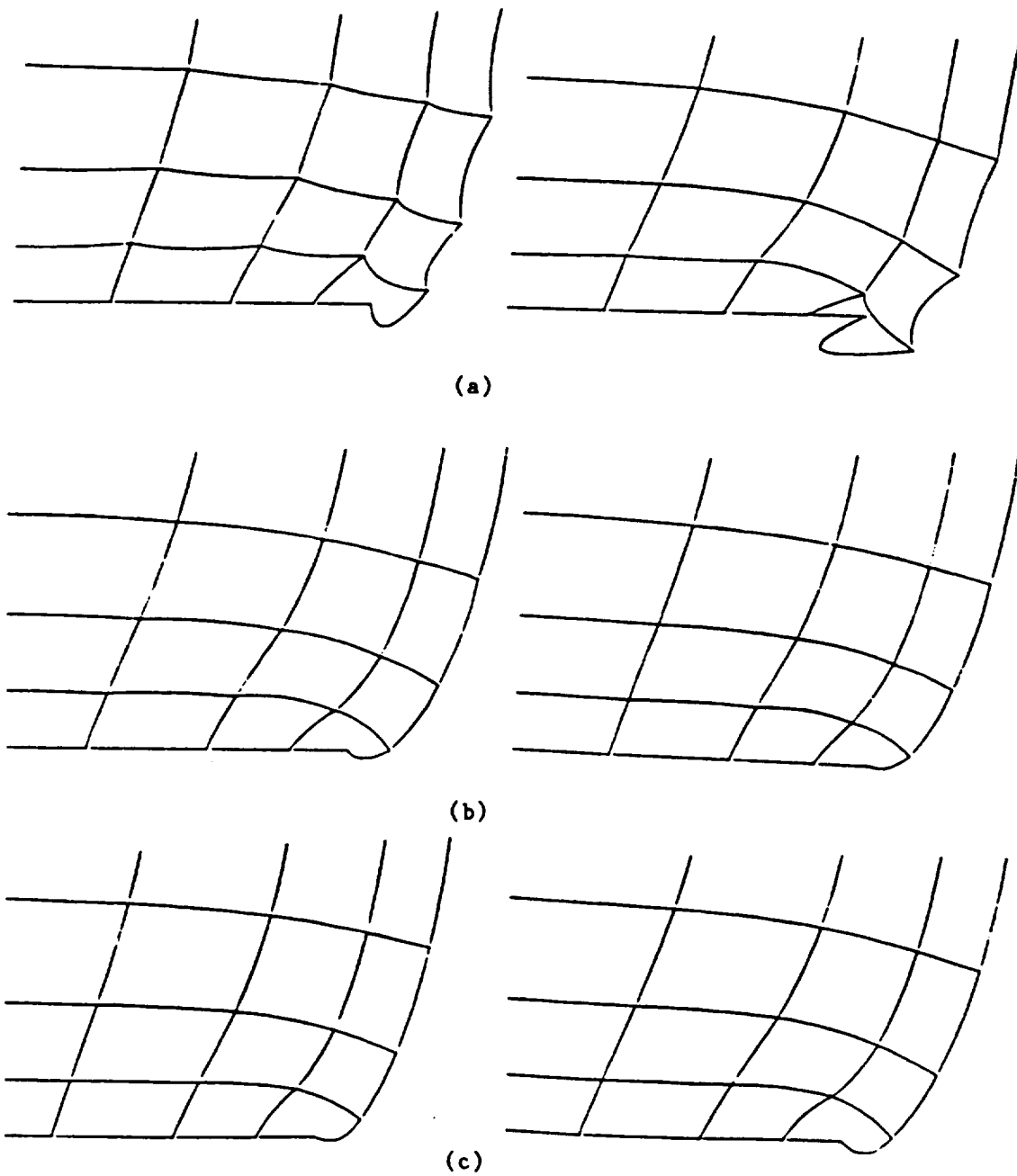


Figure 25. Compression (15%, 17.5%) of a Rubber Material Domain (25 Element Mesh): a) Underintegration, b) Controlled Underintegration, and c) Full Integration.

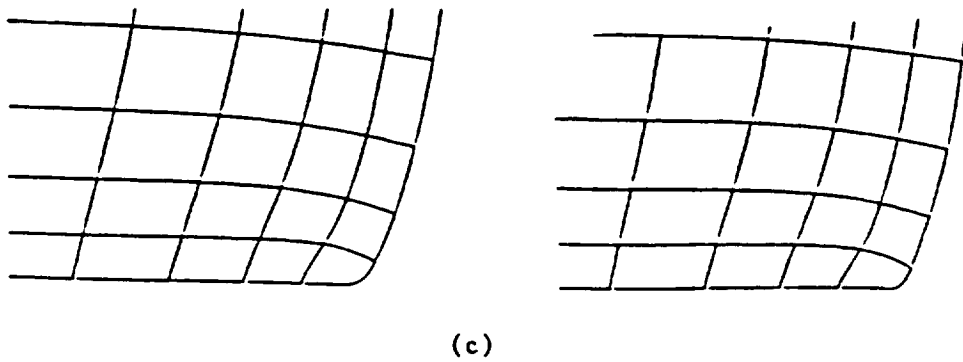
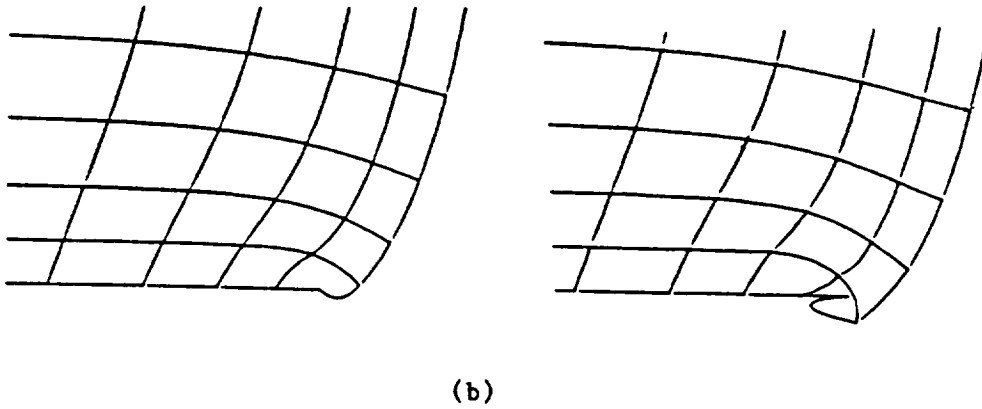
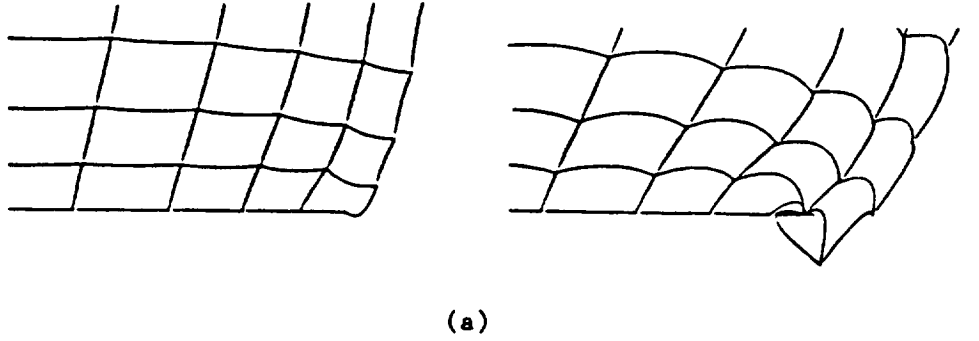


Figure 26. Compression (10, 15%) of a Rubber Material Domain (49 Element Mesh): a) Underintegration, b) Controlled Underintegration, and c) Full Integration.

of performance is observed when looking at the number of iterations required for the convergence of the Newton algorithm: they are 5, 5, 5 and 5 (respectively 5, 5, 6, 7) for the fully (resp. under-) integration in the 25 element mesh and 6, 6, and 6 (resp. 5, 6 and 13) in the 49 element mesh.

PART IV: CONCLUSIONS

In this work, we study instabilities appearing in the finite element resolution of linear and nonlinear, compressible and incompressible elasticity. The study is carried both mathematically and numerically. Some of the principal conclusions are listed as follows:

1) The use of underintegration in the stiffness matrix calculations results in rank-deficient stiffness matrix. These rank-deficiencies correspond to additional modes supplied to the rigid body modes that appropriately belong to the kernel of the operator.

2) There is a significant class of problems in which, with appropriate filtering, it can be shown that an underintegrated solution with hourglass control can yield very satisfactory answers, and produce a finite element method which has the same rate of convergence as the fully integrated method. The fact that this does indeed hold has been rigorously proved in this dissertation for a class of scalar elliptic boundary value problems, and numerically verified for a class of linear and nonlinear elasticity problems.

The method developed in this work seems to give satisfactory results in a broad class of problems. Several questions have, however, arisen:

• When an element is too distorted, the control cannot restore a reasonable shape. The accuracy of the control relies on the

approximation made on the mode \tilde{W} , supposing that the element remains quadrilateral. Does an exact computation of \tilde{W} (7.21) for very distorted elements give a better answer?

- The lack of stiffness seems to slow the speed of convergence of the Newton's scheme. Can the method be modified in order to increase the speed of convergence?

- Finally, do the results generalize to three dimensional elasticity?

The answer to these questions may provide a tremendous gain in computation time.

APPENDIX A

As far as mixed boundary conditions are concerned, we suppose that a Dirichlet boundary condition is applied at 0 and a Neumann boundary condition at 1. For the interval $[0,1]$, we consider the $N \times N$ matrix

$$D(k) = \begin{pmatrix} 2k & -1 & 0 \\ -1 & 2k & -1 \\ 0 & -1 & k \end{pmatrix}$$

The values for which $\det K(k)$ vanishes are:

$$k_i = \cos\left(\frac{-\pi}{2N} + \frac{i\pi}{N}\right), \quad 1 \leq i \leq N$$

and the corresponding vectors ($D(k_i)v_i = 0$) are :

$$v_i = \left\{ \sin \frac{ij\pi}{2N} \right\} \quad 1 \leq j \leq N$$

The corresponding approximation space $V_{1, \frac{1}{2}}^h$ with basis $\{\phi^j\}$ is constructed as in [25] or in Section 2.4. Then, depending upon the sides where the various boundary conditions (D or N) are applied, tensor product of V_1^h , $V_{1,0}^h$ or $V_{1, \frac{1}{2}}^h$ are to be considered. The results of Theorem II hold for the Mixed Problem.

APPENDIX B

According to the simplification introduced in Section 2.7.3.a we have

$$\underline{\tilde{K}}^{\text{stab}} = \frac{Q_e}{135} \begin{pmatrix} 3(\delta_7 \cdot \delta_7^T + \delta_8 \cdot \delta_8^T) + \delta_9 \cdot \delta_9^T & ; & 0 \\ 0 & ; & 3(\delta_7 \cdot \delta_7^T + \delta_8 \cdot \delta_8^T) + 4\delta_9 \cdot \delta_9^T \end{pmatrix}$$

Note that

$$\delta_7 \cdot h^T = \delta_8 \cdot h^T = 0$$

$$\delta_9 \cdot h^T = 12$$

and that we can choose $\underline{w} = (w_1, w_2)^T$ such that

$$\delta_9 \cdot w_1^T = \delta_9 \cdot w_2^T = 0$$

Then the control is

$$\begin{pmatrix} \tilde{u}_1 \\ \tilde{u}_2 \end{pmatrix} = \begin{pmatrix} u_1 \\ u_2 \end{pmatrix} - \lambda_1 \begin{pmatrix} h \\ 0 \end{pmatrix} - \lambda_2 \begin{pmatrix} 0 \\ h \end{pmatrix} - \lambda_3 \begin{pmatrix} w_1 \\ w_2 \end{pmatrix}$$

with

$$\lambda_i = \delta_9 \cdot u_i^T / 12 \quad i=1,2$$

$$\lambda_3 = \frac{(\delta_7 \cdot w_1^T)(\delta_7 \cdot u_1^T) + (\delta_8 \cdot w_1^T)(\delta_8 \cdot u_1^T) + (\delta_7 \cdot w_2^T)(\delta_7 \cdot u_2^T) + (\delta_8 \cdot w_2^T)(\delta_8 \cdot u_2^T)}{(\delta_7 \cdot w_1^T)^2 + (\delta_8 \cdot w_1^T)^2 + (\delta_7 \cdot w_2^T)^2 + (\delta_8 \cdot w_2^T)^2}$$

APPENDIX C

```

SUBROUTINE PROJ(U,XY)
C
C THIS SUBROUTINE PROJECTS THE ELEMENT SOLUTION
C ORTHOGONALLY W.R.T. HX, HY AND W
C INPUT :U SOLUTION
C        XY NODAL COORDINATES(CURRENT CONFIGURATION)
C OUTPUT :U PROJECTED SOLUTION
C
DIMENSION U(2,9),XY(2,9),W(2,9),S7(9),S8(9),S9(9),H(9)
DIMENSION S7U(2),S8U(2),S9U(2),S7W(2),S8W(2)
INTEGER SIGN
DATA S7U,S8U,S9U,S7W,S8W/10*0./
DATA S7/-1., 1., 1.,-1., 0.,-2., 0., 2., 0./
DATA S8/-1.,-1., 1., 1., 2., 0.,-2., 0., 0./
DATA S9/ 1., 1., 1., 1.,-2.,-2.,-2.,-2., 4./
DATA H/ 1., 1., 1., 1.,-1.,-1.,-1.,-1., 0./
SIGN=-1
DO 1 K=1,2
K1=3-K
IF(K.EQ.2) SIGN=1
W(K,1)=SIGN*(+3.*XY(K1,1)-XY(K1,2)-XY(K1,3)-XY(K1,4))
W(K,2)=SIGN*(+XY(K1,1)-3.*XY(K1,2)+XY(K1,3)+XY(K1,4))
W(K,3)=SIGN*(-XY(K1,1)-XY(K1,2)+3.*XY(K1,3)-XY(K1,4))
W(K,4)=SIGN*(+XY(K1,1)+XY(K1,2)+XY(K1,3)-3.*XY(K1,4))
W(K,5)=SIGN*(XY(K1,3)-XY(K1,4))
W(K,6)=SIGN*(XY(K1,1)-XY(K1,4))
W(K,7)=SIGN*(XY(K1,1)-XY(K1,2))
W(K,8)=SIGN*(XY(K1,3)-XY(K1,2))
1 W(K,9)=0.
DO 2 K=1,9
DO 2 J=1,2
S7U(J)=S7U(J)+S7(K)*U(J,K)
S8U(J)=S8U(J)+S8(K)*U(J,K)
S9U(J)=S9U(J)+S9(K)*U(J,K)
S7W(J)=S7W(J)+S7(K)*W(J,K)
2 S8W(J)=S8W(J)+S8(K)*W(J,K)
S9U(1)=S9U(1)/12.
S9U(2)=S9U(2)/12.
W1=(S7W(1)*S7U(1)+S8W(2)*S8U(2)+S7W(2)*S7U(2)+S8W(1)*S8U(1))/
. (S7W(1)*S7W(1)+S8W(2)*S8W(2)+S7W(2)*S7W(2)+S8W(1)*S8W(1))
DO 3 K=1,9
DO 3 J=1,2
3 U(J,K)=U(J,K)-S9U(J)*(H(K)+1./3.)-W1*W(J,K)
RETURN
END

```

REFERENCES

1. ALY, A., S., "A Finite Element Analysis for Problem of Large Strain and Large Displacement", TICOM Report, 81-14, 1981.
2. BABUSKA, I., "The Finite Element Method with Lagrange Multipliers", Num. Math., Vol.20, 1973.
3. BELYTSCHKO, T. and KENNEDY, J.M., "Computer Models for Sub-assembly Simulation", Nuclear Engineering & Design, Vol.49, pp.17-38, July 1978.
4. BELYTSCHKO, T., ONG, J., S.-J. and LIU, W.K., "A Consistent Control of Spurious Singular Modes in the 9-Node Lagrange Element for the Laplace and Mindlin Plate Equation", Comp. Meth. in Appl. Mech. & Engrg., Vol.44, pp.269-295, 1984.
5. BELYTSCHKO, T. and TSAY, C.S., "A Stabilization Procedure for the Quadrilateral Plate Element with One-Point Quadrature", Comp. Meth. in Appl. Mech. & Engrg., Vol.19, pp.409-419, 1983.
6. BELYTSCHKO, T., TSAY, C.S. and LIU, W.K., "A Stabilization Matrix for the Bilinear Mindlin Plate Element", Comp. Meth. in Appl. Mech. & Engrg., Vol.29, pp.313-327, 1981.
7. BERCOVIER, M., "Perturbation of Mixed Variational Problems - Applications to Mixed Finite Element Methods", R.A.I.R.O., Vol.12., No.3, 1978.
8. BICANIC, N. and HINTON, E., "Spurious Modes in Two-Dimensional Isoparametric Elements", Int. J. Num. Methods in Engrg., Vol.14, pp.1545-1557, 1979.
9. BREZZI, R., "On the Existence, Uniqueness and Approximation of Saddle-Point Problems Arising from Lagrangian Multipliers", RAIRO, Numerical Analysis, 8, 1974.
10. CAREY, G.F. and KRISHNAN, R., "Penalty Approximation of Stokes Flow, Part I: Stability Analysis, Part II: Error Estimates and Numerical Results", TICOM Report, 82-5, June 1982.
11. CIARLET, P., Numerical Analysis of the Finite Element Method for Elliptic Boundary Problems, North Holland, Amsterdam, 1978.
12. CIARLET, P. and RAVIART, P.A., "The Combined Effect of Curved Boundaries and Numerical Integration in Isoparametric Finite Element Methods", The Mathematical Foundations of the Finite Element Method with Application to Partial Differential Equations

13. COOK, R.D., Concepts and Applications of Finite Element Analysis, New-York, Willey, 1981.
14. COOK, R.D. and ZHAO-HUA, F., "Control of Spurious Modes in the Nine-Node Quadrilateral Element", Int. J. Num. Methods in Engrg., Vol. 18, pp.1576-1580, 1982.
15. ENGLEMAN, M., FIDAP User's Manual, Boulder, Colorado, August 1981.
16. ENGLEMAN, M. AND SANI, R., The Mathematics of Finite Elements, Edited by J.R.Whiteman, Academic Press, Ltd., London, 1982.
17. FALK, R.S., "An Analysis of the Penalty Method and Extrapolation for the Stationary Stokes Problem", Advances in Computer Methods for Partial Differential Equations, Edited by R.Vichnevetsky, AICA Publication, pp.66-69, 1975.
18. FLANAGAN, D.P. and BELYTSCHKO, T., "A Uniform Strait Hexahedron and Quadrilateral with Orthogonal Hourglass Control", Int. J. Num. Methods in Engrg., Vol.17, pp.676-706, 1981.
19. FORTIN, M., "A Analysis of Convergence of Mixed Finite Element Methods", R.A.I.R.O., Vol.II, No.4, 1977.
20. FORTIN, M., "Old and New Finite Elements for Incompressible Flows", Int. J. Num. Methods in Fluids, Vol.1, pp.347-364, 1979.
21. GIRAULT, V., "Theory of a Finite Difference Method on Irregular Networks", SIAM J. Numer. Anal., Vol.II, pp.409-474, 1974.
22. GIRAULT, V., "Nonelliptic pproximation of a Class of Partial Differential Equations with Neumann Boundary Conditions", Mathematics of Computation, Vol.30, No.133, pp.68-91, January 1976.
23. GIRAULT, V. and RAVIART, P.A., Lecture Notes in Mathematics 749, Finite Element Approximation of the Navier-Stokes Equations, Springer-Verlag, Berlin, 1979.
24. HAYES, L.J., "Practical Stability Test for Finite Elements with Reduced Integration", Int. J. Num. Methods in Engrg., Vol. 17, pp. 1689-1695, 1981.
25. HUGHES, T.J.R., "Equivalence of Finite Elements for Nearly Incompressible Elasticity", J. Appl.Mech., March 1977.
26. JACQUOTTE, O.-P., "Stability, Accuracy and Efficiency of Some Underintegrated Methods in Finite Element Computations", Comp. Meth. in Appl. Mech. & Engrg., (to appear).

27. JACQUOTTE, O.-P. and ODEN, J.T., "Analysis of Hourglass Instabilities and Control in Underintegrated Finite Element Methods", Comp. Meth. in Appl. Mech. & Engrg., Vol.44, pp.339-363, 1984.
28. JACQUOTTE, O.-P. and ODEN, J.T., "Analysis and Treatment of Hourglass Instabilities in Underintegrated Finite Element Methods", Innovative Methods in Nonlinear Computational Mechanics", Edited by T.Belytschko, K.C.Park and W.K.Liu, Pineridge Press, Swansea, 1984 [presented at the ASEM Winter Annual Meeting, New Orleans, Louisiana, December 9-14, 1984]
29. JACQUOTTE, O.-P., ODEN, J.T. and BECKER, E.B., "Numerical Control of the Hourglass Instabilities", Computers and Structures, (to appear).
30. KOSLOFF, D. and FRAZIER, G.A., "Treatment of Hourglass Patterns in Low Order Finite Element Codes", Int. J. Num. Anal. Meth. in Geo-Mech., Vol.2, pp.57-72, 1978.
31. LADYSZHENSKAYA, O.A., The Mathematical Theory of Viscous Incompressible Flows, Gordon Breach, New York, 1969.
32. LE TALLEC, P., "Numerical Analysis of Equilibrium Problems in Incompressible Nonlinear Elasticity", Ph.D. Dissertation, The University of Texas at Austin, 1980.
33. MALKUS, D.S., "Finite Element Analysis of Incompressible Solids", Dissertation, Department of Mathematics, Boston University, 1972.
34. MALKUS, D.S., "A Finite Element Displacement Model Valid for Any Value of the Compressibility", International Journal of Solids and Structures Vol.12, pp.731-738, 1976.
35. MALKUS, D.S., "Finite Elements with Penalties for Incompressible Elasticity: A Progress Report", Illinois Inst. of Technology, Chicago, 1981.
36. MALKUS, D.S., HUGHES, T.J.R., "Mixed Finite Element Methods - Reduced and Selective Integration Technique - A Unification of Concepts", Comp. Meth. in Appl. Mech. & Engrg., Vol.15, pp.63-81, 1978.
37. MILLER, T., "A Finite Element Study of Instabilities in Rubber Elasticity", Ph.D. Dissertation, The University of Texas at Austin, 1983.
38. ODEN, J.T., Finite Elements in Nonlinear Continua, McGraw-Hill, New York, 1970.

39. ODEN, J.T., "R.I.P. Methods for Stokesian Flows", Finite Elements in Fluids, Vol.IV, Edited by Galaagher et al., John Wiley and Sons, Ltd, London, 1982 [presented at the Third Symposium of Finite Elements in Flow Problems, Bnaff, Canada, June 1980].
40. ODEN, J.T., "Penalty Methods for Constrained Problems in Nonlinear Elasticity", Proceedings, IUTAM Symposium on Finite Elasticity, Edited by D.E. Carlson and R.T. Shield, (Lehigh 1980), Martenus Nijhoff, The Hague, pp.281-300, 1982.
41. ODEN, J.T., "Penalty Method and Reduced Integration for the Analysis of Fluids", Proceedings, Symposium on Penalty Finite Element Methods, ASME Winter Annual Meeting, Phoenix, Arizona, November 14-19, 1982.
42. ODEN, J.T. and JACQUOTTE, O.-P., "Stable and Unstable RIP/Perturbed Lagrangian Methods for Two-Dimensional Viscous Flow Problems", Finite Elements in Fluids, Vol.V, ed. R.H. Gallagher et al, John Wiley and Sons, Ltd., London, 1982.
43. ODEN, J.T. and JACQUOTTE, O.-P., "A Stable Second-Order Accurate, Finite Element Scheme for Analysis of Two-Dimensional Incompressible Viscous Flows", Proceedings of The Fourth International Symposium Finite Element in Flow Problems, Tokyo, Japan, December 19-25, 1982.
44. ODEN, J.T. and JACQUOTTE, O.-P., "Stability of Some Mixed Finite Element Methods for Stokesian Flows", Comp. Meth. in Appl. Mech. & Engrg., Vol. 43, pp. 231-247, 1984.
45. ODEN, J.T. and KIKUCHI, N., "Penalty Methods for Constrained Problems in Elasticity", International Journal for Numerical Methods in Engineering, Vol.18, pp. 701-725, 1982.
46. ODEN, J.T., KIKUCHI, N. and SONG, Y.J., "Penalty-Finite Elements Methods for the Analysis of Stokesian Flows", Comp. Meth. in Appl. Mech. & Engrg., Vol.31, pp.297-329, 1982.
47. ODEN, J.T. and WELLFORD, L.C., Jr., "Analysis of Flow of Viscous Fluids by the Finite Element Method", AIAA Journal, Vol. 10, No 12, pp. 1591-1599, 1972.
48. REDDY, J.N., "On the Accuracy and Existence of Solutions to Primitive Variable Models of Viscous Incompressible Fluids", International Journal of Engineering Science, 16-21, 921, 1978.
49. STRANG, G. and FIX, G., Analysis of the Finite Element Method, Prentice Hall, Englewood Cliffs, N.J. 1973.
50. TEMAN, R., Navier-Stokes Equations, North-Holland, Amsterdam, New-York, Oxford, 1979.

51. WAHLBIN, L.B., "Maximum Norm Error Estimates in the Finite Element Method with Quadratic Elements and Numerical Integration", R.A.I.R.O. Anal. Numer., Vol. 12, pp 173-202, 1978.
52. WAHLBIN, L.B., "A Remark on Parabolic Smoothing and the Finite Element Method", SIAM J. Numer. Anal., Vol. 17, No. 1, 1980.
53. ZIENKIEWICZ, O.C., Finite Element in Engineering, McGraw Hill, London, 1980.
54. ZIENKIEWICZ, O.C., The Finite Element Methods, McGraw Hill, London, 1977.

REPORT DOCUMENTATION PAGE

Form Approved
OMB No. 0704-0188

Public reporting burden for this collection of information is estimated to average 1 hour per response, including the time for reviewing instructions, searching existing data sources, gathering and maintaining the data needed, and completing and reviewing the collection of information. Send comments regarding this burden estimate or any other aspect of this collection of information, including suggestions for reducing this burden, to Washington Headquarters Services, Directorate for Information Operations and Reports, 1215 Jefferson Davis Highway, Suite 1204, Arlington, VA 22202-4302, and to the Office of Management and Budget, Paperwork Reduction Project (0704-0188), Washington, DC 20503.

1. AGENCY USE ONLY (<i>Leave blank</i>)	2. REPORT DATE March 1994	3. REPORT TYPE AND DATES COVERED Final Contractor Report	
4. TITLE AND SUBTITLE Analysis and Control of Hourglass Instabilities in Underintegrated Linear and Nonlinear Elasticity		5. FUNDING NUMBERS WU-509-10-11 G-NAG3-329	
6. AUTHOR(S) Olivier P. Jacquotte and J. Tinsley Oden		8. PERFORMING ORGANIZATION REPORT NUMBER E-8659	
7. PERFORMING ORGANIZATION NAME(S) AND ADDRESS(ES) The University of Texas at Austin Texas Institute for Computational Mechanics Austin, Texas 78712		10. SPONSORING/MONITORING AGENCY REPORT NUMBER NASA CR-195293	
9. SPONSORING/MONITORING AGENCY NAME(S) AND ADDRESS(ES) National Aeronautics and Space Administration Lewis Research Center Cleveland, Ohio 44135-3191		11. SUPPLEMENTARY NOTES Project Manager, Christos C. Chamis, Structures Division, organization code 5200, NASA Lewis Research Center, (216) 433-3252.	
12a. DISTRIBUTION/AVAILABILITY STATEMENT Unclassified - Unlimited Subject Category 39		12b. DISTRIBUTION CODE	
13. ABSTRACT (<i>Maximum 200 words</i>) Methods are described to identify and correct a bad finite element approximation of the governing operator obtained when under-integration is used in numerical code for several model problems: the Poisson problem, the linear elasticity problem, and for problems in the nonlinear theory of elasticity. For each of these problems, the reason for the occurrence of instabilities are given, a way to control or eliminate them are presented, and theorems of existence, uniqueness and convergence for the given methods are established. Finally, numerical results are included which illustrate the theory.			
14. SUBJECT TERMS Error estimates; Poisson problems; Elasticity problems; Nonlinear problems; Existence; Uniqueness; Convergence; Numerical results		15. NUMBER OF PAGES 127	
		16. PRICE CODE A07	
17. SECURITY CLASSIFICATION OF REPORT Unclassified	18. SECURITY CLASSIFICATION OF THIS PAGE Unclassified	19. SECURITY CLASSIFICATION OF ABSTRACT Unclassified	20. LIMITATION OF ABSTRACT