

ADAPTIVE FINITE ELEMENTS FOR VISCOELASTIC DEFORMATION  
PROBLEMS

by

HARRY HILL

A thesis submitted for the degree of Doctor of Philosophy

School of Information Systems, Computing and Mathematics

Brunel University

February 2008

## Abstract

This thesis is concerned with the theoretical and computational aspects of generating solutions to problems involving materials with fading memory, known as viscoelastic materials. Viscoelastic materials can be loosely described as those whose current stress configuration depends on their recent past. Viscoelastic constitutive laws for stress typically take the form of a sum of an instantaneous response term and an integral over their past responses. Such laws are called hereditary integral constitutive laws.

The main purpose of this study is to analyse adaptive finite element algorithms for the numerical solution of the quasistatic equations governing the small displacement of a viscoelastic body subjected to prescribed body forces and tractions. Such algorithms for the hereditary integral formulation have appeared in the literature. However the approach here is to consider an equivalent formulation based on the introduction of a set of unobservable *internal variables*. In the linear viscoelastic case we exploit the structure of the quasistatic problem to remove the displacement from the equations governing the internal variables. This results in an elliptic problem with right hand side dependent on the internal variables, and a separate independent system of ordinary differential equations in a Hilbert space.

We consider a continuous in space and time Galerkin finite element approximation to the reformulated problem for which we derive optimal order *a priori* error estimates. We then apply the techniques of the theory of adaptive finite element methods for elliptic boundary value problems and ordinary differential equations, deriving reliable and efficient *a posteriori* error estimates and detailing adaptive algorithms. We consider the idea of splitting the error into space and time portions and present results regarding a splitting for space time projections. The ideas for splitting the error in projections is applied to the finite element approximation and a further set of *a posteriori* error estimates derived. Numerical studies confirm the theoretical properties of all of the estimators and we show how they can be used to drive adaptive in space and time solution algorithms.

We consider the extension of our results for the linear case to the constitutively nonlinear case. A model problem is formulated and the general techniques for dealing with *a posteriori* error estimation for nonlinear space time problems are considered.

## **Acknowledgements**

I would like to thank my supervisor Dr. Simon Shaw for his encouragement, guidance and direction, also the staff of the Maths department for their kind support and assistance. Furthermore I am grateful to EPSRC for providing financial support.

I would also like to acknowledge the Japanese Society for the Promotion of Science(JSPS) for providing the financial support for my Summer Fellowship in 2005, and thank Professor M. Yamamoto and Dr. S. Kim of University of Tokyo for being generous hosts and Professor Twizell for facilitating the visit.

# Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	The finite element method . . . . .	3
1.2	Adaptive finite element methods . . . . .	4
1.3	Convergence of AFEM . . . . .	7
1.4	Continuum mechanics and viscoelasticity . . . . .	10
1.5	Quasistatic linear viscoelasticity . . . . .	16
1.6	Preliminary material and notation . . . . .	17
1.7	Summary . . . . .	20
<b>2</b>	<b>Adaptive finite element methods</b>	<b>21</b>
2.1	The Galerkin method . . . . .	21
2.2	Approximation by finite elements . . . . .	24
2.3	Error analysis . . . . .	30
2.4	AFEM for linear elasticity . . . . .	34
2.5	AFEM for linear systems of ordinary differential equations . . . . .	46
2.6	Summary . . . . .	56
<b>3</b>	<b>Finite element approximation of quasistatic linear viscoelasticity</b>	<b>58</b>
3.1	Existence and uniqueness . . . . .	59
3.2	Internal variable formulation . . . . .	63
3.3	Finite element approximation . . . . .	68
3.4	Summary . . . . .	72

<b>4</b>	<b><i>A priori</i> error analysis</b>	<b>73</b>
4.1	Displacement . . . . .	73
4.2	Internal variables . . . . .	75
4.3	<i>A priori</i> estimates . . . . .	84
4.4	Numerical results . . . . .	86
4.5	Summary . . . . .	91
<b>5</b>	<b><i>A posteriori</i> error analysis</b>	<b>94</b>
5.1	Displacement . . . . .	95
5.2	Internal variables . . . . .	103
5.3	<i>A posteriori</i> estimates . . . . .	115
5.4	Adaptive algorithms . . . . .	116
5.5	Numerical results . . . . .	119
5.6	Summary . . . . .	123
<b>6</b>	<b>Exact <i>a posteriori</i> error estimators</b>	<b>125</b>
6.1	Space-time projections . . . . .	126
6.2	Error indicators . . . . .	128
6.3	Adaptive algorithms . . . . .	136
6.4	Numerical experiments . . . . .	138
6.5	Summary . . . . .	142
<b>7</b>	<b>Nonlinear viscoelasticity</b>	<b>145</b>
7.1	Schapery-Knauss-Emri constitutive model . . . . .	146
7.2	Finite element approximation . . . . .	150
7.3	Towards <i>a posteriori</i> error analysis . . . . .	151
7.4	Summary . . . . .	154
<b>8</b>	<b>Summary and recommendations for further work</b>	<b>155</b>
8.1	Summary . . . . .	155
8.2	Recommendations for future work . . . . .	156

<b>A</b>	<b>Adaptive mesh refinement in MATLAB</b>	<b>158</b>
A.1	Introduction . . . . .	158
A.2	Local mesh refinement . . . . .	160
A.3	Summary . . . . .	173
A.4	Main routine . . . . .	173

# Chapter 1

## Introduction

This thesis is concerned with the finite element solution of systems of equations modelling the behaviour of viscoelastic material bodies subjected to given forces and tractions. The theory of viscoelasticity provides an interesting set of problems in continuum mechanics, and is widely used in mechanical engineering for practical computations to model materials that display both fluid and solid like behaviour, such as polymers.

The general model for the displacement of a linear viscoelastic material is an integro-partial differential initial boundary value problem. Such models are analysed in the books of Christensen [21], Fabrizio and Morro [38] and Golden and Graham [40]. In this thesis two approximations that are common in engineering practice and make the problem more tractable are applied. First, only small strains are considered. This removes the geometric non-linearity associated with finite strain models. Second, the restriction to the quasistatic case. The quasistatic assumption means ignoring the inertia term in the balance of linear momentum. The resulting system is then an elliptic differential equation combined with a Volterra integral equation. The elliptic differential operator arises from the equilibrium equations of continuum mechanics and the Volterra integral due to the fading memory term in the hereditary integral formulation of linear viscoelasticity.

The aim of this thesis is to construct and analyse adaptive finite element methods (AFEMs) for calculating the solution of the quasistatic boundary value problem of linear viscoelasticity, in particular, where the constitutive law is expressed with *internal variables* [71], [38], [48]. By internal variables we mean a set of unobservable quantities, that allow

the constitutive equation to be rewritten in a simplified form, with a supplementary set of evolution equations governing the dynamics of the internal variables.

By exploiting the quasistatic structure of the problem, the internal variable equations can be separated from the displacement problem, leading to a system of ordinary differential equations (ODEs) in the energy space governing time dependent effects, together with an augmented linear elasticity problem. It is apparent that an implementation of a solution algorithm for the reformulated system modelling linear viscoelasticity would require little adjustment to existing implementations of algorithms for linear elasticity. Similarly, the theory for the analysis of the reformulated system follows from the well developed theory of AFEM for elliptic problems and ODEs. However, the ODEs are posed in Hilbert space and require a spatial approximation. This added complication provides a number of challenges, the majority of which arise from the requirements of space and time finite element approximations.

### **Overview of the thesis**

- Chapter 2.

The basics of finite element approximation are reviewed, leading towards AFEMs. Recent results on the convergence of AFEM for elliptic problems in the context of the linear elasticity problem are presented. Furthermore the theory of AFEM relating to ODEs is also covered. Adaptive algorithms for both problems are presented and numerical results presented.

- Chapter 3.

The formulation of the quasistatic linear viscoelastic boundary value problem in terms of internal variables and its finite element approximation are presented.

- Chapter 4.

Optimal order *A priori* error estimates for the finite element approximation of chapter 3 are derived.

- Chapter 5.

Reliable and efficient *A posteriori* error estimates for the finite element approximation of chapter 3 are presented. Adaptive space and time algorithms are discussed. Theoretical considerations indicate and numerical results confirm that the performance of the temporal indicator is in some sense sub-optimal, a state we try and improve in chapter 6.

- Chapter 6.

Taking a lead from the closing remarks of chapter 5, we present an idea about how errors in space and time approximations can be partitioned. More *a posteriori* error estimates are presented, and numerical evaluation is carried out.

- Chapter 7.

We look to extend our previous results to a constitutively nonlinear problem. The reformulation leaves us with a linear elliptic problem for which the results of the previous chapters apply, however the internal variable problem is now nonlinear. We focus on the nonlinearity in the internal variable problem and pose a model problem related to the original and consider the finite element approximation. We discuss ways forward for deriving *a posteriori* error estimates.

- Chapter 8.

A summary of the work carried out is given, with conclusions and recommendations for further work.

This chapter continues with a review of the literature on the theory of adaptive finite element methods together with the relevant details of continuum mechanics required for the problems we wish to consider and some mathematical preliminaries. The final section of the chapter concludes with a summary of the main boundary value problem that is considered in this study.

## 1.1 The finite element method

Finite element methods (FEMs) for elliptic boundary value problems have been the preferred method of the engineering community for the numerical solution of elliptic partial

differential equations (PDEs) since their introduction in the 1940s. Courant in [26] is widely acknowledged to have formulated the method, based on the earlier works of Galerkin, Rayleigh and Ritz. The development of the method grew in the succeeding years though it wasn't until the 1970s that a rigorous mathematical theory was put in place (see [23] and references within). Given the flexibility of the method, the boom in computational power over the last 25 years has led the adaptation of the FEM to treat an evergrowing class of problems, encompassing applications from medicine to finance. Furthermore the widening range of problems has led to an evergrowing class of methods and computational techniques. Various difficulties encountered by the classical FEM led to innovations such as mixed and discontinuous methods to name but a few.

Finite element methods have also been shown to be suitable for solving time dependent and mixed space and time problems. Space and time discretisations commonly involve the use of finite elements in space with finite difference based time stepping schemes, such as the Crank-Nicolson or Euler methods. However, a purely finite element approach is possible with FE approximations in both space and time. The resulting schemes are often related to classical finite difference based schemes. However the finite element approach has the advantages of weaker regularity requirements of the solution, and the variational formulation allows for a more general analysis and treatment of a wider class of problems than classical finite differences.

## 1.2 Adaptive finite element methods

In the 1980s a refined notion of computational scheme advanced (see [10] for references). The idea was that the error in a computed solution can be described by the action of some operator on the approximate solution and the initial problem data. This led to the idea of adaptive finite element methods (AFEMs). Generally speaking, AFEMs are a logical result of a successful *a posteriori* error analysis. While *a priori* error analysis of a given method looks to ascertain rates of convergence as the dimension of the approximating space increases, *a posteriori* error analysis aims to find computable representations, informative indicators and upper bounds for functionals and norms of the error in the approximation.

*A posteriori* estimators come in various forms (see [2] and [77] for reviews and comparisons), not only for norms of the approximation error but also for the error in arbitrary functionals of the solution. These are useful in applications where interest is not in the solution of the underlying PDE, but a functional of the solution. A further aim of *a posteriori* error analysis is to derive localised error indicators that are informative with respect to the error distribution. Local error indicators then inform where to adapt the underlying approximation space in a feedback loop. Hence the name adaptive finite elements. The dimension of the approximating space can be increased in two ways:

1. *h*-method. The polynomial order of the basis functions is fixed, and the dimension of the space is increased through the addition of more basis functions of the same order.
2. *p*-method. The dimension of the space is increased by increasing the order of the existing polynomial basis functions.

Another method of adapting the approximating space is the so-called *r*-method which does not increase the dimension but improves the current choice of basis functions by relocating the nodes of the discretisation. Combinations of methods have also proved popular with the *hp*-method topping the list. For an introduction to *p* and *hp* methods see the book by Schwab [63].

## Formal procedures of error analysis

In the papers [31] and [32] Eriksson and Johnson proposed a general scheme for deriving *a priori* and residual based *a posteriori* error estimates for finite element approximations to a wide class of problems. While they focused on norms of approximation errors, the formal procedures for arriving at error representations are applicable for deriving representations of errors in functionals of the solution, often called target functionals, since their calculation is the target of the computation. Summarised in [30] they propose the following general scheme for deriving *a priori* error estimates in the  $L^2$  norm:

### Scheme for deriving *a priori* estimates

1. Representation of the error in terms of the exact solution and the solution to a discrete linearised dual problem.
2. Use Galerkin orthogonality to introduce the interpolation error in the error representation.
3. Local estimates for the interpolation error.
4. Strong stability estimates for the discrete dual problem.

For linear problems, the dual problem is the formal adjoint of the primal problem. We follow this scheme in chapter 4. A similar scheme gives rise to *a posteriori* error estimates.

#### **Scheme for deriving *a posteriori* estimates**

1. Representation of the error in terms of the residual of the finite element problem the solution of a continuous linearised dual problem.
2. Use Galerkin orthogonality to introduce the interpolation error in the error representation.
3. Local interpolation estimates for the dual solution.
4. Strong stability estimates for the continuous dual problem.

The derivation of residual based estimates for general functionals almost always uses steps 1 and 2. However, a common problem of the above scheme is that it requires strong stability of the dual problem. Depending on the choice of target functional, strong stability estimates may be hard to come by or unavailable. Therefore, step 2 is often the starting point for an alternative thread of analysis and computation. The Dual Weighted Residual (DWR) technique (see [11] and references within) attempts to compute the error in a target functional of the solution by evaluating the representation. This is no mean feat since the representation contains the unknown solution to a continuous dual problem. This leaves little alternative but to compute the dual solution via approximation. However there is evidence [75] that in deriving upper bounds for both norms and functionals, some information is lost through repeated application of the Cauchy-Schwarz inequality. It appears that the dual solution carries local information regarding the error distribution and maintaining the inner product structure of the error representation preserves this information.

An *a posteriori* error estimate gives a computable upper bound for the approximation error in norm or functional form. This is required for guaranteeing that the error is less than a given upper bound. When the estimator is an upper bound for the error, the estimator is said to be reliable. It is also important that the error estimator behaves like the error in the sense that it is of the same order as the discretisation becomes finer. When an estimator does so, it is called efficient. A theoretical measure of the performance of an *a posteriori* error estimate is its *effectivity index*, which is the ratio of the estimated error and the true error. A reliable and efficient estimator for which the effectivity index converges to one as the discretisation becomes finer is called asymptotically exact. However in practice asymptotic exactness is an optimistic target and typically the best we can expect is reliability and efficiency. Upper and lower bounds for the effectivity index of a residual based *a posteriori* error estimator for linear finite elements on triangles are analysed in [9].

### 1.3 Convergence of AFEM

For elliptic problems the theory of AFEM is mature ([77], [2]) and the results of the current theory are readily applicable ([28], [54]). Finite element methods for ODEs are well known ([30],[34]), and the work on AFEM beyond *a posteriori* error analysis is less complete.

Given a computable upper bound on the error, an adaptive process requires a user providing an error tolerance TOL, and a Solve - Estimate - Refine (SER) algorithm being applied until the *a posteriori* upper bound is less than the given tolerance. Locally adaptive methods are based on localising the error estimator so that those areas significant to the distribution of the error bound can be identified. The method of selecting areas to adapt is called a *marking strategy*. Until recently marking strategies were based on heuristic, common sense arguments with numerical experiments providing validation. However the work on the convergence of AFEM ([28], [54], [52], [72], [13]) shows that the marking strategy is an integral part of ensuring and determining a convergence rate for an adaptive algorithm.

The challenge to prove the convergence for an adaptive finite element method is to prove that there is an error reduction due to an adaptive step. The first proof of convergence of an adaptive finite element method is due to Dörfler [28] who constructed an adaptive algorithm

for the two dimensional Poisson problem. Dörfler’s proof establishes convergence of an SER algorithm under the restrictions that there is a sufficiently fine initial mesh and particular choices of marking and refinement strategies. The main result of [28] is that under specific conditions, the  $L^2$  norm of the gradient of the error is reduced by an SER step, or the specified tolerance has already been met.

The main point about proving the convergence of an adaptive algorithm is that whilst *a priori* error estimates tell us convergence rates of a method as the discretisation becomes finer in a uniform way, adaptive methods are based on discretisations that become finer only in localised portions of the computational domain. Proofs of convergence and furthermore convergence rates of adaptive algorithms are necessary to establish optimal choices in designing adaptive finite element algorithms. Convergence of an adaptive algorithm is the property that given an error tolerance target, the adaptive algorithm can achieve the tolerance in a finite number of steps. Morin, Nocketto and Siebert [54] (MNS) consider a sequence of FE approximations  $u_k$ ,  $k = 0, 1, \dots$ , and for constants  $C_0$ , and  $0 < \beta < 1$  depending on the problem data and initial mesh show that the sequence generated by their algorithm satisfies,

$$\|u - u_k\| \leq C_0 \beta^k, \quad (1.3.1)$$

where the norm is the energy norm. Furthermore the restriction upon the initial mesh of a degree of “fineness” required by Dörfler’s proof is lifted. However, it is not totally removed as it is hidden in a new term called the data oscillation, appearing in the lower bounds for the error. It is worth noting that there is a link between mesh fineness, data oscillation and the saturation assumption<sup>1</sup> as remarked in [28] and [54]. In fact it is shown in [56] that small data oscillation implies the saturation assumption.

The proof in [54] relies on an error and data oscillation reduction procedure, therefore retaining some degree of mesh fineness as described by Dörfler on a step by step basis and so not requiring such a fine initial triangulation. The MNS adaptive procedure has been generalised to deal with general second order elliptic problems [52], and more recently

---

<sup>1</sup>In its simplest form, the saturation assumption states that the solution to a PDE using finite elements can be approximated asymptotically better with quadratic finite elements than with linear ones.

provided the basis for a proof of the rate of convergence of AFEM in [13].

### Convergence rates for adaptive finite element methods

Given the proofs that AFEM converge, attention has turned to determining the rates of convergence, since in all of the works mentioned above there is nothing that guarantees an advantage in using locally adaptive methods over any others. In [13], a modified MNS algorithm is considered. Rather than only refining the underlying discretisation, a coarsening step is included in the algorithm. Using methods of nonlinear approximation theory, it is shown that the modified algorithm is optimal in the following sense: If the solution  $u : \mathbb{R}^2 \mapsto \mathbb{R}$  is such that it can be approximated by piecewise linear functions (using complete knowledge of  $u$ ) to accuracy  $O(n^{-s})$  in the energy norm, where  $s > 0$  and  $n$  is the number of subdomains in the triangulation, then the algorithm results in an approximation of the same type with the same asymptotic accuracy using only the information gained during the computation. See also [72] for extensions of [13] and [53] for results regarding convergence rates of the DWR method.

### AFEM for time dependent problems

Given the fairly complete picture for AFEMs for second order elliptic problems provided by the papers [77], [54], [55], [52], [13] and [72], AFEMs for time-dependent problems remain less fully understood. Space and time discretisations commonly involve the use of finite elements in space, and finite difference based stepping schemes such as the Crank-Nicolson or Euler methods, in time. The finite element approach is to consider finite element approximations in both space and time for which the result is again a stepping scheme (often related to a finite difference based scheme). However in the finite element instance, weaker regularity requirements of the solution to the variational formulation allows for a more general analysis, and for the treatment of a wider class of problems.

Historically, the discontinuous Galerkin method (DG) has received more attention than the continuous Galerkin (CG) in relation to time dependent problems. Adaptive finite element methods for ordinary differential equations (ODEs) using DG are presented in [49]. The extension of DG methods to space and time dependent problems follow in [31], [32],

where the approach is to take DG in time and CG in space for the approximation of some parabolic problems. Furthermore a general mechanism for deriving *a priori* and *a posteriori* error estimates is presented and the mechanism is summarised in [30].

Adaptive finite elements using the CG method are used for solving general ODEs in [34]. Both *a priori* and *a posteriori* error estimates are presented. The CG method is applied to the heat equation in [8], and similarly for the wave equation in [39]. However, in both cases only *a priori* error estimates are given. Therefore there are still gaps in the *a posteriori* error analysis for CG methods in space and time. In different directions, a review of AFEMs for general hyperbolic problems is given in [75] where both *p* and *hp* methods are considered.

The hereditary integral formulation of viscoelasticity involves a Volterra integral equation, and the quasistatic boundary value problem can be considered as an abstract Volterra problem. For Volterra problems, Bedivan and Fix [12] studied a Galerkin method and the implications of quadrature for non-coercive problems. For the quasistatic hereditary integral formulation of linear viscoelasticity, Shaw and Whiteman [69] provide *a priori* error estimates for a Galerkin finite element method that is discontinuous in time and continuous in space.

## 1.4 Continuum mechanics and viscoelasticity

In this section we present an overview of the continuum mechanics required for posing the main problem of study, the boundary value problem of quasistatic linear viscoelasticity. For details see, for example, the books of Antman [6], Fabrizio and Morro [38] and Golden and Graham [40].

Let  $I = [0, T] \subset \mathbb{R}$  be the time interval over which we will consider the deformation process to occur and let the open bounded Lipschitz domain  $\Omega \subset \mathbb{R}^d$  represent the interior of a continuous body. Consider a deformation of the set  $\Omega$ , given by  $\phi : \Omega \times I \mapsto \mathbb{R}^d$ . Written in components,  $\phi(x, t') = (\phi^1(x, t'), \dots, \phi^d(x, t'))$ , describes the position of a point originating from  $x$  at  $t = 0$ , at a later time  $t' > 0$ . Considering the difference between a point originally at  $x$ , in the deformed state at  $\phi(x, t)$  leads to the definition of the displacement

$u(x, t)$ , defined by,

$$u(x, t) := \phi(x, t) - x, \quad u^i(x, t) = \phi^i(x, t) - x^i, \quad i = 1, \dots, d, \quad x \in \Omega, t \in I. \quad (1.4.1)$$

### Balance Equations

Let  $\rho : \Omega \mapsto \mathbb{R}$  denote the density of the material at the point  $x$ . Under the exertion of known body forces  $f : \Omega \times I \mapsto \mathbb{R}^d$ , the balance of linear momentum [6] provides a relationship between the known forces  $f$ , the displacement  $u$ , and the second Piola-Kirchhoff stress tensor  $\Sigma : \Omega \times I \mapsto \mathbb{R}^{d \times d}$ . The balance of angular momentum implies that the second Piola-Kirchhoff stress is symmetric. These relationships are given by the equations,

$$\rho(x)u_{tt}(x, t) - \operatorname{div}((I + Du(x, t))\Sigma(x, t)) = f(x, t), \quad (x, t) \in \Omega \times I, \quad (1.4.2)$$

$$\Sigma(x, t) = \Sigma(x, t)^T, \quad (x, t) \in \Omega \times I, \quad (1.4.3)$$

where the divergence operator is defined by,

$$\operatorname{div} \cdot = \sum_{j=1}^d (\cdot)_{x_j}. \quad (1.4.4)$$

and where  $D \cdot$  is the derivative operator with respect to  $x$ , which with respect to Cartesian coordinates can be represented by a matrix with components,

$$(Du(x, t))_{ij} = u_{x_j}^i(x, t) \quad 1 \leq i, j \leq d. \quad (1.4.5)$$

The equations (1.4.2) and (1.4.3) are the local forms of the fundamental laws of continuum mechanics for an arbitrary continuous body. It is assumed that the initial state of the body is known, i.e., the displacement  $u$  and velocity  $u_t$  at  $t = 0$  are known functions, given by,

$$u(x, 0) = u_0(x), \quad x \in \Omega, \quad (1.4.6)$$

$$u_t(x, 0) = u_1(x), \quad x \in \Omega. \quad (1.4.7)$$

Often the situation where the body is in equilibrium is of interest. In this instance the acceleration term  $\rho(x)u_{tt}$  is neglected and the time dependence ignored. The equilibrium equations are then,

$$-\operatorname{div}((I + Du(x))\Sigma(x)) = f(x), \quad x \in \Omega, \quad (1.4.8)$$

$$\Sigma(x) = \Sigma(x)^T, \quad x \in \Omega. \quad (1.4.9)$$

Viscous effects are time dependent phenomena, therefore it makes little sense to consider an “equilibrium” problem. However, when the given forces or displacements are of small variation in time it is common in engineering practice to drop the acceleration term in the equations of motion. Furthermore, in materials with high internal friction losses such as rubbers and soft polymers, the inertial effects associated with the density  $\rho$  may be neglected compared to the viscous effects [29], [40]. The resulting situation is called the *quasistatic* case, and the equations are the same as the equilibrium equations (1.4.8) and (1.4.9) only with the dependence on  $t$  reinstated.

### Boundary Conditions

To complete the model specification, assume that on some portion of the body the displacement is zero, while on another part there are surface tractions being applied. More precisely, partition the boundary of  $\Omega$  into disjoint subsets  $\Gamma_D$  and  $\Gamma_N$ , where  $\Gamma_D$  is assumed to have positive  $(d - 1)$  Hausdorff measure, and let  $n(x)$  denote the outward unit normal of the boundary. Then suppose that the displacement,  $u(x, t) = (u^1(x, t), \dots, u^d(x, t))$ , satisfies  $u(x, t)|_{\Gamma_D} = 0$ ,  $\forall t \in I$  i.e., the restriction to  $\Gamma_D$  of the displacement is zero for all time. Furthermore, assume that on  $\Gamma_N$  there is a prescribed surface traction  $g : \Gamma_N \times I \mapsto \mathbb{R}^d$ . The boundary conditions for the problem defined by equations (1.4.2) and (1.4.3) are then,

$$u(x, t) = 0, \quad \forall (x, t) \in \Gamma_D \times I, \quad (1.4.10)$$

$$(I + Du(x, t))\Sigma(x, t) \cdot n(x) = g(x, t), \quad \forall (x, t) \in \Gamma_N \times I. \quad (1.4.11)$$

The same boundary conditions hold for the quasistatic problem, and similar time independent conditions hold for the equilibrium problem (1.4.8) and (1.4.9).

The general problem is to solve the equation (1.4.2), subject to (1.4.3), (1.4.6), (1.4.7), (1.4.10) and (1.4.11) (or the related quasistatic or equilibrium problems) for the displacement. However, the systems are underdetermined. In the case  $d = 3$ , there are 6 components of  $\Sigma$  (given the symmetry constraint) and 3 components of  $u$  to be determined from 3 equations. The constitutive equation simultaneously makes the system determined and provides a description of the material at hand by relating  $\Sigma$  to  $u$ , typically by relating the stress to the strain.

The strain tensor measures the local effect of the deformation and is defined in terms of the deformation by,

$$E := \frac{1}{2}(D\phi(x, t)^T D\phi(x, t) - I), \quad (1.4.12)$$

The strain can be written in terms of the displacement using  $\phi(x, t) = x + u(x, t)$  as,

$$E(u) = \frac{1}{2}(Du(x, t) + Du(x, t)^T + Du(x, t)^T Du(x, t)). \quad (1.4.13)$$

Models using (1.4.13) to model the strain are nonlinear and are called finite strain problems.

To make computations more amenable it is common to resort to a linear or small strain theory. Formally differentiating (1.4.13), for some displacements  $v$  and  $h$  we have,

$$DE(v; h) := \left. \frac{dE}{d\tau}(v + \tau h) \right|_{\tau=0} = \frac{1}{2}(Dh + Dh^T + Dv^T Dh + Dh^T Dv). \quad (1.4.14)$$

The small strain tensor is arrived at by looking at the linearisation in a neighbourhood of the zero displacement and is defined by,

$$\epsilon(u) := DE(0; u) = \frac{1}{2}(Du + Du^T), \quad \epsilon(u)_{ij} = \frac{1}{2}(u_{x_j}^i + u_{x_i}^j), \quad 1 \leq i, j \leq d. \quad (1.4.15)$$

## Constitutive Equation

Constitutive theory concerns itself with the search for relationships between fundamental quantities (stress, strain, temperature, etc...) which suitably model material behaviour. Empirical research can suggest functional dependencies between various quantities, however some ground can be gained in the search for such dependencies by considering certain universal requirements that a model should obey. The work of Coleman, Noll and Truesdell in the 1960s and 70s (see [25], [76], [57] for example) resulted in general representation theorems for constitutive laws, derived from theories based on purely hypothetical requirements. A full account of current constitutive theory and its applications is given in [81].

A fundamental requirement of a constitutive law is the *principle of material objectivity*, which states that the response of the material must be independent of the frame of reference. From material objectivity, consideration of the symmetry properties of the material and thermodynamic requirements lead to general forms and conditions that constitutive equations must obey. Our focus is on viscoelastic materials, which are contained in the class of simple materials as described by Day [27].

## Viscoelasticity

In [38] viscoelasticity is presented in the frame of materials with fading memory, similar results are achieved by considering conceptual rheological models based on spring and dashpots [71]. Either way the results are the same and in this section we present the basic forms of the most popular models. The theory of finite linear viscoelasticity is based on a representation of the stress of the form,

$$\Sigma(E(x, t), E(x, t - s)) = \Sigma_e(E(x, t)) + \Sigma_v(E(x, t), E(x, t - s)), \quad s \geq 0, \quad (1.4.16)$$

where the term  $\Sigma_e(E(x, t))$  describes the elastic response and  $\Sigma_v(E(x, t), E^t(x, s))$  describes the viscous response. The term  $\Sigma_e(E(x, t))$  is often taken from the theory of elasticity. A popular choice for both practical and theoretical [22] reasons is,

$$\Sigma_e(E) = \lambda \operatorname{tr} E + 2\mu E, \quad \operatorname{tr} A = \sum_{i=1}^d A_{ii}, \quad A \in \mathbb{R}^{d \times d}. \quad (1.4.17)$$

Materials with an elastic response given by (1.4.17) are referred to as St. Venant-Kirchhoff materials, and the relationship is a generalisation of Hooke's Law. The material constants  $\lambda$  and  $\mu$  describe the volumetric and shear behaviour of the material. Rather than determine them directly from physical experiments,  $\lambda$  and  $\mu$  are more commonly determined by their relationship to the engineering parameters, Young's modulus  $E$ , and Poisson's ratio  $\nu$  from the equations,

$$\lambda = \frac{E\nu}{(1+\nu)(1-2\nu)}, \quad \mu = \frac{E}{2(1+\nu)}. \quad (1.4.18)$$

The viscous component  $\Sigma_v(E(x, t), E(x, t - s))$  can be modelled using a convolution of the strain history with a stress relaxation function or by using a set of internal variables. Internal variable formulations of constitutive laws are becoming more popular as they provide added flexibility for theory [3] and computation.

The focus of this thesis is on small strain models, where the strain is given by (1.4.15). In this instance the second Piola-Kirchhoff stress tensor  $\Sigma$  is identified with the Cauchy stress tensor  $\boldsymbol{\sigma}$ , therefore we adopt the familiar convention of referring to the stress with the symbol  $\boldsymbol{\sigma}$ . The theory of linear viscoelasticity [38], is based on the constitutive law,

$$\boldsymbol{\sigma}(x, t) = \mathbf{C}(x, 0)\boldsymbol{\epsilon}(x, t) - \int_0^\infty \partial_s \mathbf{C}(x, s)\boldsymbol{\epsilon}(x, t - s)ds, \quad (1.4.19)$$

where  $\partial_s$  is the partial derivative with respect to the variable  $s$ , and the tensor  $\mathbf{C}(x, t) = (C_{ijkl}(x, t))_{1 \leq i, j, k, l \leq d}$  is a positive definite, fourth order tensor function representing the elastic response of the material, satisfying the symmetries,

$$C_{ijkl} = C_{klij}, \quad C_{ijkl} = C_{jikl}, \quad C_{ijkl} = C_{ijlk}. \quad (1.4.20)$$

Typically,  $\mathbf{C}(x, 0)\boldsymbol{\epsilon}(x, t)$  will be taken to be the of the form (1.4.17),

$$\mathbf{C}(x, 0)\boldsymbol{\epsilon}(x, t) = \lambda \text{tr}\boldsymbol{\epsilon}(x, t) + 2\mu\boldsymbol{\epsilon}(x, t). \quad (1.4.21)$$

Equation (1.4.19) is not useful in practical computations due to the convolution over the entire past history of the strain. Therefore it is commonly assumed that the material has a finite past, starting at the origin  $t = 0$ . Cutting off the past at the time origin, and changing variables allows us to write (1.4.19) as,

$$\boldsymbol{\sigma}(x, t) = \mathbf{C}(x, 0)\boldsymbol{\epsilon}(x, t) - \int_0^t \partial_s \mathbf{C}(x, t - s)\boldsymbol{\epsilon}(x, s) ds. \quad (1.4.22)$$

A common simplifying assumption is that of homogeneity of the temporal response. Under this assumption, the kernel in (1.4.22) can be decomposed as  $\mathbf{C}(x, t - s) = \mathbf{C}(x, 0)\varphi(t - s)$ . The function  $\varphi(t)$  is called the relaxation function. The relaxation encapsulates the fading memory of the material. Assumptions on the relaxation function are:

1. Fading memory hypothesis:

$$\varphi(t) > 0, \quad \forall t \in I, \quad (1.4.23)$$

$$\varphi'(t) < 0, \quad \forall t \in I. \quad (1.4.24)$$

2. Normalisation:  $\varphi(0) = 1$ . Note that combined with (1.4.23) and (1.4.24), normalisation implies that  $0 < \varphi(t) < 1, \forall t \in (0, T]$ .

A popular choice of relaxation function is a sum of decaying exponentials. This choice corresponds to a generalisation of the models derived from conceptual spring and dashpots and are termed spectrum models [40]. The relaxation function is given by,

$$\varphi(t) = \varphi_0 + \sum_{i=1}^{n_v} \varphi_i e^{-\alpha_i t}, \quad \sum_{i=0}^{n_v} \varphi_i = 1, \quad \varphi_0 > 0, \quad (1.4.25)$$

where  $\varphi_i, \alpha_i \geq 0$  for  $i = 1, \dots, n_v$ . The constitutive law (1.4.22) then becomes,

$$\boldsymbol{\sigma}(x, t) = \mathbf{C}(x, 0)\boldsymbol{\epsilon}(u(x, t)) - \sum_{i=1}^{n_v} \alpha_i \varphi_i \int_0^t e^{-\alpha_i(t-s)} \mathbf{C}(x, 0)\boldsymbol{\epsilon}(u(x, s)) ds. \quad (1.4.26)$$

## 1.5 Quasistatic linear viscoelasticity

In summary, the main problem we are studying is as follows. Identify a material body with a polygonal domain  $\Omega \subset \mathbb{R}^d$ . Partition the boundary  $\partial\Omega$  into two disjoint subsets  $\Gamma_N$  and  $\Gamma_D$ , where  $\Gamma_D$  has positive  $(d-1)$  surface measure. Denote by  $f : \Omega \times I \mapsto \mathbb{R}^d$  the forces acting throughout the body and by  $g : \Gamma_N \times I \mapsto \mathbb{R}^d$ , the tractions acting on  $\Gamma_N$ . Recall the equilibrium equations (1.4.8), then neglecting the nonlinear terms, the governing equations take the form,

$$-\operatorname{div}\boldsymbol{\sigma}(x, t) = f(x, t), \quad (x, t) \in \Omega \times I, \quad (1.5.1)$$

$$u(x, t) = 0, \quad (x, t) \in \Gamma_D \times I, \quad (1.5.2)$$

$$\boldsymbol{\sigma}(x, t)n(x) = g(x, t), \quad (x, t) \in \Gamma_N \times I, \quad (1.5.3)$$

where  $n$  is the outer unit normal to the surface  $\Gamma_N$ . For a synchronous, homogeneous linear viscoelastic solid the constitutive law relating the stress to the strain  $\boldsymbol{\epsilon}$  is,

$$\boldsymbol{\sigma}(x, t) = \mathbf{C}\boldsymbol{\epsilon}(u(x, t)) - \int_0^t \partial_s \varphi(t-s) \mathbf{C}\boldsymbol{\epsilon}(u(x, s)) ds, \quad (1.5.4)$$

where  $\varphi(t)$  is the relaxation function (1.4.25) and  $\mathbf{C} = (C_{ijkl})_{1 \leq i, j, k, l \leq d}$  is a positive definite fourth order tensor satisfying the symmetries (1.4.20). The setup for generating finite element solutions is a spatially weak form obtained in the usual way by integration by parts over the spatial domain. Define the the space  $V$ ,

$$V := \{v \in [H^1(\Omega)]^d \mid v|_{\Gamma_D} = 0\}. \quad (1.5.5)$$

We now consider functions as mappings from  $I \rightarrow V$ , so  $\forall t \in I$ ,  $u(\cdot, t) \in V$  and to focus on this aspect we suppress the  $x$  dependence of functions, but maintain the explicit  $t$  dependence. Define the symmetric bilinear form  $a : V \times V \rightarrow \mathbb{R}$  by,

$$a(u(t), v) := \int_{\Omega} \mathbf{C}\boldsymbol{\epsilon}(u(t)) : \boldsymbol{\epsilon}(v) dx, \quad (1.5.6)$$

which is the familiar bilinear form arising from the weak formulation of the linear elasticity problem, presented in section 2.4. Define the linear functional  $l : I \times V \mapsto \mathbb{R}$  by,

$$l(t; v) := \int_{\Omega} f(t) \cdot v dx + \int_{\Gamma_N} g(t) \cdot v d\gamma. \quad (1.5.7)$$

The Euclidean inner product of equation (1.5.1) with an arbitrary function  $v \in V$ , followed by integration by parts over  $\Omega$  and using (1.5.4) implies that,

$$a(u(t), v) - \int_0^t \partial_s \varphi(t-s) a(u(s), v) ds = l(t; v) \quad \forall v \in V. \quad (1.5.8)$$

We now present a brief summary of some notation, basic definitions and frequently used results that are required for the sequel.

## 1.6 Preliminary material and notation

As always it is necessary to set out some notation, conventions and basic material. Most notation is introduced as it appears, however, the purpose of this section is to introduce the various function spaces, related notation and frequently used inequalities that will be used throughout, for full details see [1], [36].

### Function Spaces

Let  $\Omega$  be an open subset of  $\mathbb{R}^d$ . The space denoted by  $L^p(\Omega)$  is the linear space of functions with norm defined by,

$$\|u\|_{L^p(\Omega)} := \begin{cases} \left( \int_{\Omega} |u|^p dx \right)^{1/p} & 1 \leq p < \infty, \\ \text{ess sup}_{\Omega} |u| & p = \infty. \end{cases}$$

Let  $1 \leq p < \infty$  and let  $q$  be such that  $\frac{1}{p} + \frac{1}{q} = 1$ . Let  $\phi \in L^p(\Omega)^*$ , then there exists a unique  $g \in L^q(\Omega)$  such that,

$$\phi(f) = \int_{\Omega} f g dx, \quad \forall f \in L^p(\Omega),$$

and the map  $\phi \mapsto g$  is an isometric isomorphism of  $L^p(\Omega)^*$  and  $L^q(\Omega)$ . A vector  $\alpha := (\alpha_1, \dots, \alpha_d)$ ,  $\alpha_i \geq 0$  is called a multi-index of order,

$$|\alpha| := \sum_{i=1}^d \alpha_i.$$

Given a multi-index  $\alpha$  define,

$$D^{\alpha} u(x) := \frac{\partial^{|\alpha|} u(x)}{\partial x_1^{\alpha_1} \dots \partial x_n^{\alpha_n}}$$

$$x^\alpha := x_1^{\alpha_1} \dots x_d^{\alpha_d}, \quad \alpha! := \alpha_1! \alpha_2! \dots \alpha_d!.$$

The space denoted by  $W^{k,p}(\Omega)$  is the linear space of all locally summable functions  $u : \Omega \mapsto \mathbb{R}$  such that for each multi-index  $\alpha$  with  $\alpha \leq k$ ,  $D^\alpha u$  exists in the weak sense and belongs to  $L^p(\Omega)$ . The space  $W^{k,p}(\Omega)$  has norm defined by,

$$\|u\|_{W^{k,p}(\Omega)} := \begin{cases} \left( \sum_{|\alpha| \leq k} \int_{\Omega} |D^\alpha u|^p dx \right)^{1/p} & 1 \leq p < \infty, \\ \max_{|\alpha| \leq k} \|D^\alpha u\|_{L^\infty(\Omega)} & p = \infty. \end{cases}$$

The spaces  $L^p(\Omega)$  are called Lebesgue spaces, and the spaces  $W^{m,p}(\Omega)$  are called Sobolev spaces. Lebesgue and Sobolev spaces are both Banach spaces, and in the case  $p = 2$  it is usual to write  $H^m(\Omega) = W^{m,2}(\Omega)$ . Both  $L^2(\Omega)$  and  $H^m(\Omega)$  are Hilbert spaces with inner products,

$$(u, v)_\Omega = \int_{\Omega} uv dx, \quad (u, v)_{H^m(\Omega)} = \sum_{|\alpha| \leq m} (D^\alpha u, D^\alpha v)_\Omega. \quad (1.6.1)$$

Let  $E \subset \Omega$  be a  $(d-l)$  hyperplane passing through  $\Omega$ . Let  $n_E$  be a normal vector to the hyperplane  $E$ , then define the jump of a function  $u : \Omega \rightarrow \mathbb{R}$  across  $E$  by,

$$[[u]]_E(x) := \lim_{t \rightarrow 0^+} u(x + tn_E) - u(x - tn_E), \quad \forall x \in E. \quad (1.6.2)$$

To be explicit in the one dimensional case,

$$[[u]]_i := u(t_i^+) - u(t_i^-), \quad w(t^\pm) := \lim_{s \rightarrow 0^+} w(t \pm s). \quad (1.6.3)$$

## Useful Inequalities

**Proposition 1.6.1** (Young's inequality). *For  $\alpha, \beta, p, q \in \mathbb{R}$  and  $p, q > 1$  with  $\frac{1}{p} + \frac{1}{q} = 1$  then,*

$$\alpha\beta \leq \frac{\alpha^p}{p} + \frac{\beta^q}{q}.$$

**Proposition 1.6.2** (Cauchy's inequality). *Let  $(X, (\cdot, \cdot))$  be a real inner product space with norm defined by  $\|\cdot\| = \sqrt{(\cdot, \cdot)}$ . Then,*

$$|(u, v)| \leq \|u\| \|v\|.$$

**Proposition 1.6.3** (Hölder's inequality). *Let  $u \in L^p(\Omega)$  and  $v \in L^q(\Omega)$ ,  $1 \leq p, q \leq \infty$  with  $\frac{1}{p} + \frac{1}{q} = 1$  then  $uv \in L^1(\Omega)$  and,*

$$\|uv\|_{L^1(\Omega)} \leq \|u\|_{L^p(\Omega)} \|v\|_{L^q(\Omega)}.$$

**Proposition 1.6.4** (Gronwall's lemma). *Let  $a \in \mathbb{R}$ ,  $u \in C^1([0, T])$  and  $f \in C^0([0, T])$  be such that  $\partial_t u \leq au + f$ , then,*

$$\forall t \in [0, T], \quad u(t) \leq e^{at}u(0) + \int_0^t e^{a(t-s)} f(s) ds. \quad (1.6.4)$$

*Proof.* See [33]. □

**Lemma 1.6.5.** *Suppose that  $s_i \geq 0$ ,  $\forall i \in \{1, \dots, N\}$  and  $C > 0$ . Then,*

$$\prod_{i=1}^N (1 + Cs_i) \leq \exp\left(C \sum_{i=1}^N s_i\right)$$

*Proof.* See [60]. □

**Definition 1.6.1** (Convolution). The convolution of  $v, w \in L^1(\mathbb{R}_+)$  is defined as,

$$(v * w)(t) := \int_0^t w(t-s)v(s)ds. \quad (1.6.5)$$

**Lemma 1.6.6** (Young's inequality for convolutions). *Let  $c \in L^p(a, b)$  and  $w \in L^q(a, b)$  and let  $p, q, r \geq 1$ . If  $\frac{1}{p} + \frac{1}{q} = 1 + \frac{1}{r}$  then,*

$$\|u * v\|_{L^r(a,b)} \leq \|u\|_{L^p(a,b)} \|v\|_{L^q(a,b)}.$$

*Proof.* See [1]. □

**Definition 1.6.2.** Let  $X$  and  $Y$  be Banach spaces and let  $u \in U \subseteq X$ . Let  $F : U \mapsto Y$  be a mapping from the open subset  $U$  into  $Y$ . Then  $F$  is differentiable at  $u$  if there exists a linear operator  $DF(u) : X \mapsto Y$  such that,

$$F(u+h) = F(u) + DF(u)h + o(\|h\|), \quad h \rightarrow 0.$$

For notational simplicity of derivatives of operators, define  $F'(u; h)$  as follows

$$F'(u; h) := \lim_{t \rightarrow 0} \frac{F(u+th) - F(u)}{t} = DF(u)h, \quad (1.6.6)$$

then higher order derivatives are denoted similarly with the order indicated by the number of primes, and we have observed that when the derivative exists in the sense outlined above (Fréchet) it can be formally computed via the Gateaux derivative.

**Theorem 1.6.7.** *Let  $F : U \subseteq X \mapsto Y$  be  $C^n$  on the open set  $U$ , then,*

$$F(u + h) = F(u) + \sum_{k=1}^{n-1} \frac{1}{k!} F^{(k)}(u; h^k) + R_n, \quad \forall u + h \in U \quad h \in X, \quad (1.6.7)$$

where  $F^{(k)}(u; h^k) = F^{(k)}(u; h, \dots, h)$  is a  $k$ -linear mapping and,

$$R_n := \int_0^1 \frac{(1 - \tau)^{n-1}}{(n-1)!} F^{(n)}(u + \tau h) h^n d\tau. \quad (1.6.8)$$

*Proof.* The proof follows from setting  $\phi(t) := \langle v^*, F(u + th) \rangle$ ,  $v^* \in Y^*$  and  $0 \leq t \leq 1$  and using the classical Taylor formula with continuity of the mapping  $t \mapsto \langle v^*, u(t) \rangle$ .  $\square$

## 1.7 Summary

This chapter presented a review of the relevant literature from the field of AFEM together with the basics of continuum mechanics that we required to specify our model problem of quasistatic linear viscoelasticity. In the next chapter, we review the basics of finite element approximations, leading towards AFEMs. We present the latest results on convergence of AFEM for elliptic problems in the context of the linear elasticity problem. The linear elasticity problem is important as it is the basis of the analysis of the elliptic component of the viscoelastic problem and the results are easily applied to the reformulated system, with a small modification for time dependency. Chapter 2 also contains the theory of AFEM relating to ODEs. This theory is less readily applicable to the system of ODEs we are dealing with due to the complication of spatial variation, however it provides sufficient guidance and structure for our later work. In chapter 3 the reformulation of the quasistatic linear viscoelastic boundary value problem and its finite element approximation are presented. Chapter 4 contains an *a priori* error analysis of the finite element approximation and chapter 5 gives an *a posteriori* error analysis. Chapter 6 provides a summary and conclusions together with a discussion of potential further work.

## Chapter 2

# Adaptive finite element methods

### 2.1 The Galerkin method

In this chapter the Galerkin finite element method is presented, together with a summary of the current work on adaptive methods. Most of the early sections of this chapter are taken from the books of Brenner and Scott [16], Ciarlet [23] and Ern and Guermond [33]. The FEM rests on a variational formulation, or weak form of the original PDE achieved by multiplying the original problem by an arbitrary smooth function and integrating. The resulting problem is generally of the following form. Let  $U$  be a Banach space and  $V$  be a reflexive Banach space with topological dual spaces  $U^*$  and  $V^*$ . Let  $b : U \times V \mapsto \mathbb{R}$  be a continuous bilinear form, let  $f \in V^*$ : find  $u \in U$  such that,

$$b(u, v) = f(v), \quad \forall v \in V. \quad (2.1.1)$$

The following theorem shows under what conditions the above problem is known to have a unique solution.

**Theorem 2.1.1.** *Let  $U$  be a Banach space and let  $V$  be a reflexive Banach space. Let  $b : U \times V \mapsto \mathbb{R}$  be a continuous bilinear form and  $f \in V^*$ . Suppose that the bilinear form  $b$  satisfies,*

$$\exists c_b > 0, \quad \inf_{w \in U} \sup_{v \in V} \frac{b(w, v)}{\|w\|_U \|v\|_V} \geq c_b, \quad (2.1.2)$$

and if for all  $v \in V$ ,

$$\forall w \in U, \quad b(w, v) = 0 \Rightarrow v = 0. \quad (2.1.3)$$

Then problem (2.1.1) has a unique solution, with a priori stability estimate,

$$\|u\|_U \leq \frac{1}{c_b} \|f\|_{V^*}. \quad (2.1.4)$$

The solution  $u$  of (2.1.1) is called the weak solution of the original PDE.

Problem (2.1.1) can equivalently be treated by defining the operator  $B : U \mapsto V^*$ , by,

$$\langle Bu, v \rangle_{V^*, V} := b(u, v). \quad (2.1.5)$$

Equation (2.1.1) can then be written as,

$$\langle Bu, v \rangle_{V^*, V} = \langle f, v \rangle_{V^*, V}, \quad \text{or} \quad Bu = f, \quad \text{in } V^*. \quad (2.1.6)$$

The above form will be more convenient in later sections where the concept of adjoint operators is required.

### Consistent and conforming Galerkin approximations

A Galerkin approximation to (2.1.1) is made by selecting finite dimensional subspaces  $U^h$  and  $V^h$ , and looking for an approximate solution in  $U^h$ , the trial space, by sampling (2.1.1) on  $V^h$ , the test space. When  $U^h \subset U$  and  $V^h \subset V$  the approximation is called *conforming*. An approximation is said to be *consistent* if the exact solution satisfies the approximate problem. For more details see [33] (Chap 2, p 89). A consistent and conformal Galerkin approximation results in the finite dimensional problem: Find  $u_h \in U^h \subset U$ , such that,

$$b(u_h, v_h) = f(v_h), \quad \forall v_h \in V^h \subset V. \quad (2.1.7)$$

The existence and uniqueness of the solution to (2.1.7) rests on the following discrete equivalent of theorem (2.1.1).

**Theorem 2.1.2.** *Let  $U^h$  and  $V^h$  be two finite dimensional spaces with  $\dim U^h = \dim V^h$ . Let  $b_h : U^h \times V^h \mapsto \mathbb{R}$  be a continuous bilinear form and let  $f_h$  be continuous on  $V^h$ . Suppose that the bilinear form  $b_h$  satisfies,*

$$\exists c_{b_h} > 0, \quad \inf_{w_h \in U^h} \sup_{v_h \in V^h} \frac{b_h(w_h, v_h)}{\|w_h\|_U \|v_h\|_V} \geq c_{b_h}. \quad (2.1.8)$$

Then problem (2.1.7) has a unique solution, with a priori stability estimate,

$$\|u_h\|_U \leq \frac{1}{c_{b_h}} \|f_h\|_{V^*}. \quad (2.1.9)$$

It is important to note that in the consistent and conforming case,  $b_h = b$ . Of fundamental importance to the analysis of the finite element method, and to Galerkin methods in general is the following orthogonality property.

**Lemma 2.1.3** (Galerkin Orthogonality). *Let  $u$  be the solution of (2.1.1) and  $u_h$  be the solution to (2.1.7), then,*

$$b(u - u_h, v_h) = 0, \quad \forall v_h \in V^h. \quad (2.1.10)$$

*Proof.* Since  $V^h \subset V$ , choose  $v = v_h$  in (2.1.1) and subtract (2.1.7) from it.  $\square$

The Galerkin orthogonality property states that the error is orthogonal with respect to the bilinear form  $b(\cdot, \cdot)$  to the test space. Using (2.1.10), a basic error estimate is available showing that the energy norm of the error due to the Galerkin approximation is quasi-optimal in the sense that it is proportional to the best approximation error using the space  $U^h$ . The most general forms of the following result are the famous lemmas of Strang [73], however in the conforming consistent case, the earlier lemma of C ea [16] can be applied.

**Lemma 2.1.4** (C ea's Lemma). *Let  $u$  be the solution of (2.1.1),  $u_h$  the solution to (2.1.7). then,*

$$\|u - u_h\|_U \leq C \inf_{w_h \in U^h} \|u - w_h\|_U. \quad (2.1.11)$$

*Proof.* From Galerkin orthogonality (2.1.10), it follows that  $\forall v_h \in V^h$

$$b(u_h - w_h, v_h) = b(u - w_h, v_h). \quad (2.1.12)$$

Using (2.1.8) and (2.1.12),

$$c_{b_h} \|u_h - w_h\|_U \leq \sup_{v_h \in V^h} \frac{b(u_h - w_h, v_h)}{\|v_h\|_V} = \sup_{v_h \in V^h} \frac{b(u - w_h, v_h)}{\|v_h\|_V} \leq \|b\|_{U,V} \|u - w_h\|_U. \quad (2.1.13)$$

Then since  $\|u - u_h\|_U \leq \|u - w_h\|_U + \|u_h - w_h\|_U$ , it follows that,

$$\|u - u_h\|_U \leq \left(1 + \frac{\|b\|_{U,V}}{c_{b_h}}\right) \inf_{w_h \in U^h} \|u - w_h\|_U. \quad (2.1.14)$$

$\square$

### Construction of the approximate solution

Let  $\{\psi_U^i\}_{i=1}^n$  be a basis for  $U^h$ , and let  $\{\psi_V^i\}_{i=1}^n$  be a basis for  $V^h$ . The approximate solution  $u_h$  can be expanded in terms of the basis for  $U^h$ ,

$$u_h = \sum_{j=1}^n u_j \psi_U^j. \quad (2.1.15)$$

Substituting for  $u_h$  its expansion given in (2.1.15) and sampling (2.1.7) at each basis function of  $V^h$  results in a square linear system with matrix  $\mathcal{B}$  given by,

$$\mathcal{B} \in \mathbb{R}^{n \times n}, \quad \mathcal{B}_{ij} = b(\psi_U^j, \psi_V^i). \quad (2.1.16)$$

The right hand side of the linear system is,

$$\mathbf{f} \in \mathbb{R}^n, \quad f_i = f(\psi_V^i). \quad (2.1.17)$$

The resulting finite dimensional problem is then: Given  $\mathcal{B} \in \mathbb{R}^{n \times n}$ ,  $\mathbf{f} \in \mathbb{R}^n$ , find  $\mathbf{u} = (u_1, \dots, u_n)^T \in \mathbb{R}^n$ , such that,

$$\mathcal{B}\mathbf{u} = \mathbf{f}. \quad (2.1.18)$$

Theorem (2.1.2) ensures that the matrix  $\mathcal{B}$  is invertible so the problem now is to choose an algorithm to solve the linear system (2.1.18). Many choices exist [41], however for symmetric positive definite systems the Cholesky method is typically used. Furthermore we mention that, by design, finite element methods give rise to sparse matrices, so ideally implementations utilising sparse matrix data structures should be used. For the problems under consideration in this thesis, the systems are typically symmetric and positive definite so a natural choice of solution algorithm would be the Cholesky algorithm.

## 2.2 Approximation by finite elements

### Lagrange Finite Elements

By finite element we mean the triple  $(K, P, \Sigma)$  as defined by Ciarlet [23] of a non-empty, compact, connected Lipschitz domain  $K \subset \mathbb{R}^d$ , a finite dimensional space of functions  $P$  and a set of linear functionals  $\Sigma$  forming a basis for the algebraic dual of  $P$ . In this thesis we restrict ourselves to simplicial Lagrange finite elements which are now described.

The set  $K$  is called the element domain. Let  $\{z_i\}_{i=0}^d$  be a set of points in  $\mathbb{R}^d$  such that the set of vectors  $\{z_1 - z_0, \dots, z_d - z_0\}$  are linearly independent. Set  $K$  to be the convex hull of those points,  $K = \text{conv}\{z_0, \dots, z_d\}$ . Set  $P = \mathbb{P}_k(K)$ , the space of polynomials in  $d$  variables on  $K$  of global degree at most  $k$  from which the local shape functions will be taken. The functionals in  $\Sigma$  are called the local degrees of freedom. The local degrees of freedom are taken to be nodal evaluations associated to the set of points  $\{a_i\}_{i=1}^n$ ,  $a_i \in K$ . Therefore for each  $\sigma_i \in \Sigma$  and all  $v \in \mathbb{P}_k(K)$ ,  $\sigma_i(v) = v(a_i)$ . The local shape functions determining a basis for  $\mathbb{P}_k(K)$  are then determined by solving the equations  $\psi_i(a_j) = \delta_{ij}$  for  $1 \leq i, j \leq n$  where  $\psi_i \in \mathbb{P}_k(K)$ .

The local interpolation operator  $\Pi_K : C^0(K) \rightarrow \mathbb{P}_k(K)$  is defined by,

$$\Pi_K^k v(x) = \sum_{i=1}^n \sigma_i(v) \psi_i(x) = \sum_{i=1}^n v(a_i) \psi_i(x). \quad (2.2.1)$$

The extension of the above interpolation operator to vector valued functions follows by considering the interpolant componentwise, with related polynomial space,

$$\mathbb{P}_k(K; \mathbb{R}^d) := \{v : \mathbb{R}^d \mapsto \mathbb{R}^d, v = (v^1, \dots, v^d), v^i|_K \in \mathbb{P}_k(K)\}. \quad (2.2.2)$$

### Triangulation

Let  $\Omega$  be an open, bounded, polyhedral region of  $\mathbb{R}^d$ ,  $d = 1, 2$ , or  $3$ , with Lipschitz boundary  $\partial\Omega$ . A geometrically conforming affine mesh is a partition of  $\Omega$  into a family of element domains  $\mathcal{T} = \{K\}$ , with edges  $\mathcal{E}$  and nodes  $\mathcal{N}$  such that,

1. The partition covers the closure of  $\Omega$ ,

$$\bar{\Omega} = \bigcup_{K \in \mathcal{T}} K. \quad (2.2.3)$$

2. The intersection of the interiors of distinct elements is empty, so that if  $K \neq K'$  then  $\overset{\circ}{K} \cap \overset{\circ}{K}' = \emptyset$ .
3. For any  $K$  and  $K'$  there exists an affine mapping  $F : K \rightarrow K'$  such that  $F(K) = K'$ .
4. The intersection of any two distinct elements results in an edge, a node or is empty.

That is, if  $K \neq K'$  then,

$$K \cap K' = \begin{cases} E \in \mathcal{E}, & \text{or,} \\ z \in \mathcal{N}, & \text{or,} \\ \emptyset. \end{cases} \quad (2.2.4)$$

Define the piecewise constant gridsize function  $h : \Omega \rightarrow \mathbb{R}^d$ , measuring the size of elements and edges of the mesh by,

$$h_S(x) = \begin{cases} \text{diam}(S), & x \in S, \\ 0, & \text{otherwise,} \end{cases} \quad \text{with } h := \max_{K \in \mathcal{T}} h_K. \quad (2.2.5)$$

For  $K \in \mathcal{T}$  and  $E \in \mathcal{E}$  define the following neighbourhoods, illustrated in figure 2.1,

$$\begin{aligned} \omega_K &:= \bigcup_{\mathcal{E}(K) \cap \mathcal{E}(K') \neq \emptyset} K', & \tilde{\omega}_K &:= \bigcup_{\mathcal{N}(K) \cap \mathcal{N}(K') \neq \emptyset} K', \\ \omega_E &:= \bigcup_{E \in \mathcal{E}(K')} K', & \tilde{\omega}_E &:= \bigcup_{\mathcal{N}(E) \cap \mathcal{N}(K') \neq \emptyset} K', \end{aligned}$$

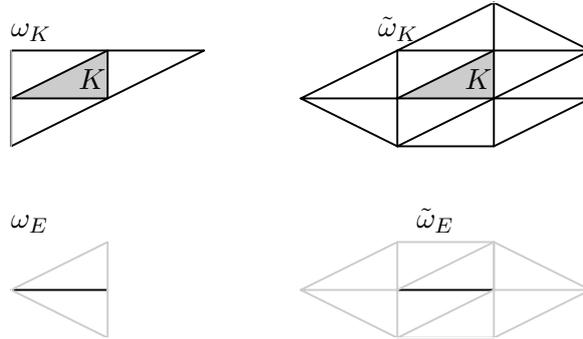


Figure 2.1: Illustration of the neighbourhoods  $\omega_K$ ,  $\tilde{\omega}_K$ ,  $\omega_E$  and  $\tilde{\omega}_E$ .

Approximation theory, from which a number of results will be required later on, requires stricter conditions on the form of the element domains. Define the element shape parameters,

$$\gamma_K := \frac{h_K}{\rho_K}, \quad \rho_K := \sup\{r \mid B_r(x_0) \subset K, \forall x_0 \in K\}, \quad (2.2.6)$$

where  $B_r(x_0)$  represents the ball of radius  $r$  centred a  $x_0$ . A family of meshes  $\{\mathcal{T}_h\}_{h \geq 0}$  is then said to be shape regular if,

$$\forall h, \forall K \in \mathcal{T}, \quad \gamma_K \leq \bar{\gamma} < \infty. \quad (2.2.7)$$

Also to characterise the degree with which neighbourhoods overlap in the mesh, let  $M_K := \text{card}\{\omega_{K'} | K \neq K', K \in \omega_{K'}\}$  then set  $M = \max_{K \in \mathcal{T}_h} M_K$ .

### Lagrange finite element space

Define the reference domain  $K_{\text{ref}}$  to be the unit simplex,

$$K_{\text{ref}} := \{x \in \mathbb{R}^d, x_i > 0, \quad 1 \leq i \leq d+1, \sum_{i=1}^{d+1} x_i \leq 1\}, \quad (2.2.8)$$

then each element domain of the triangulation  $\mathcal{T}$  is characterised by the affine mapping,

$$F_K : K_{\text{ref}} \ni x \rightarrow J_K x + b_k \in K. \quad (2.2.9)$$

The element associated with the domain (2.2.8) is called the reference element. An affine equivalent family of Lagrange finite elements can then be generated using (2.2.9) (see [16] for details). Affine equivalent families of elements are required for the interpolation error bounds that will be presented below. For a given affine mesh  $\mathcal{T}$  and the family of finite elements  $(K, \mathbb{P}_k(K), \Sigma_K)_{K \in \mathcal{T}}$  define the global interpolation operator  $\Pi^k$  by,

$$\Pi^k v(x)|_{K \in \mathcal{T}} = \Pi_K^k v(x). \quad (2.2.10)$$

Let  $S_k(\mathcal{T})$  denote the space of simplicial Lagrange finite elements based on a geometrically conforming mesh, where nodes on interfaces between elements are matched up, then it is well known that for  $k \geq 1$ ,  $S_k(\mathcal{T}) \subset C^0(\Omega)$ . Then we have,

$$S_k(\mathcal{T}) = C^0(\bar{\Omega}) \cap \mathbb{P}_k(\mathcal{T}), \quad \mathbb{P}_k(\mathcal{T}) := \{v : \Omega \rightarrow \mathbb{R}, \quad v|_{K \in \mathcal{T}} \in \mathbb{P}_k(K)\}. \quad (2.2.11)$$

Since  $S_1(\mathcal{T}) \subset C^0(\bar{\Omega})$  the global interpolant in the piecewise linear case has the alternative representation,

$$\Pi^1 v(x) = \sum_{z \in \mathcal{N}} \sigma_z(v) \eta_z = \sum_{z \in \mathcal{N}} v(z) \eta_z, \quad (2.2.12)$$

where  $\eta_z$  is the global basis function formed by local basis functions meeting at  $z$ . Such a function for the case  $k = 1$  is illustrated in figure 2.2. The extension of the above interpolation operator to vector valued functions follows by considering the interpolant componentwise.

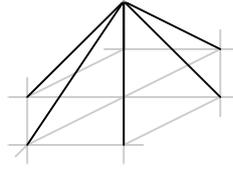


Figure 2.2: A hat function made up of basis functions from different elements associated to a central node.

### Interpolation error

Approximation theory plays a vital role in the analysis of the FEM. Céa's lemma (lemma 2.1.4) shows that the the Galerkin error is bounded above by a constant depending on the given bilinear form and the error from approximating the function  $u$  with functions from the space  $U^h$ . In this section we take the required results from chapter 1 in Ern and Guermond [33] in the context of affine families of Lagrange finite elements.

**Theorem 2.2.1** (Local Interpolation). *Let  $(K, \mathbb{P}_k, \Sigma)$  be a Lagrange finite element. Let  $1 \leq p \leq \infty$  and assume that  $\frac{d}{p} - 1 < l \leq k$  so that  $W_p^{k+1}(K) \subset C^0(K)$  with continuous embedding. Let  $\Pi_K^k$  be the local interpolant defined in (2.2.1), then there exists  $c > 0$  such that, for all  $m \in 0, \dots, l + 1$ ,*

$$\forall K, \forall v \in W^{l+1,p}(K), \quad |v - \Pi_K^k v|_{W^{m,p}(K)} \leq c \gamma_K^m h_K^{l+1-m} |v|_{W^{l+1,p}(K)}. \quad (2.2.13)$$

To extend this local estimate to a global one the hypothesis of shape regularity for the mesh is required.

**Theorem 2.2.2** (Global Interpolation). *Let  $p, k$  and  $l$  be as in theorem 2.2.1. Let  $\{\mathcal{T}_h\}$  be a shape regular family of affine meshes on the polyhedral domain  $\Omega$  and let  $\Pi^k$  denote the global interpolant defined in (2.2.10), then there exists  $c > 0$  such that  $\forall v \in W_p^{k+1}(\Omega)$ ,*

$$\begin{aligned} \|v - \Pi^k v\|_{L^p(\Omega)} + \sum_{m=1}^{k+1} h^m \left( \sum_{K \in \mathcal{T}_h} |v - \Pi^k v|_{W^{m,p}(K)}^p \right)^{1/p} &\leq c h^{k+1} |v|_{W^{k+1,p}(\Omega)}, \quad 1 \leq p < \infty \\ \|v - \Pi^k v\|_{L^\infty(\Omega)} + \sum_{m=1}^{k+1} h^m \max_{K \in \mathcal{T}_h} |v - \Pi^k v|_{W^{m,\infty}(K)} &\leq c h^{k+1} |v|_{W^{k+1,\infty}(\Omega)}. \end{aligned}$$

and for  $p < \infty$ ,

$$\lim_{h \rightarrow 0} \left( \inf_{v_h \in S_k(\mathcal{T})} \|v - v_h\|_{L^p(\Omega)} \right) = 0. \quad (2.2.14)$$

Furthermore, if  $S_k(\mathcal{T}_h) \subset W^{1,p}(\Omega)$  there holds,

$$|v - \Pi^k v|_{W^{1,p}(\Omega)} \leq ch^l |v|_{W^{l+1,p}(\Omega)}, \quad l \geq 0. \quad (2.2.15)$$

The interpolation operators defined above rely on pointwise values for their definition. This is reflected in their error estimates by the regularity needed for boundedness of the error. Since the FEM approximates weak solutions it often does not make sense to take pointwise values of functions. More general interpolation operators with degrees of freedom that are well defined even for non - smooth functions were introduced by Clément [24]. Scott and Zhang [64] introduced an alternative which preserved homogeneous boundary conditions and is a projection. In either case the basic error estimate does not change.

**Theorem 2.2.3** (Properties of the Quasi-Interpolant). *Let  $\mathcal{T}_h$  denote a shape regular mesh and  $\mathbb{P}_k(\mathcal{T}_h)$  the associated polynomial approximating space. Let  $v \in W_p^k(\Omega)$  for  $0 \leq k \leq m$  and  $1 \leq p \leq \infty$ , then there exists a mapping  $\mathcal{I}_h^k : W^{k,p}(\Omega) \rightarrow S_k(\mathcal{T}_h)$  such that,*

$$\left( \sum_{K \in \mathcal{T}} h_K^{-p(k-s)} \|v - \mathcal{I}_h^k v\|_{W^{s,p}(K)}^p \right)^{1/p} \leq C |v|_{W^{k,p}(\Omega)}, \quad 0 \leq s \leq k, \quad (2.2.16)$$

$$\|\mathcal{I}v\|_{W^{k,p}(\Omega)} \leq C_{\mathcal{I}} \|v\|_{W^{k,p}(\Omega)}. \quad (2.2.17)$$

It will often prove more useful to consider the localised forms of the above estimate [33],

$$\|v - \mathcal{I}_h^k v\|_{W^{s,p}(K)} \leq c_K h_K^{k-s} |v|_{W^{k,p}(\tilde{\omega}_K)}, \quad 0 \leq s \leq k, \quad (2.2.18)$$

$$\|v - \mathcal{I}_h^k v\|_{W^{s,p}(E)} \leq c_E h_E^{k-s-\frac{1}{p}} |v|_{W^{k,p}(\tilde{\omega}_E)}, \quad 0 \leq s \leq k. \quad (2.2.19)$$

An important concept and result from approximation theory is that of orthogonal projections. We briefly state the definition and required result, taken from [7]. Define the mapping  $P : V \rightarrow V_h$  by the equation,

$$(Pv, w_h) = (v, w_h), \quad \forall w_h \in V_h. \quad (2.2.20)$$

**Theorem 2.2.4** (Orthogonal Projection). *Let  $(\cdot, \cdot)$  be an inner product with associated norm  $\|\cdot\|$ , on the Hilbert space  $V$ . Let  $V_h$  be a non-empty, closed, convex subset of  $V$ . Then  $\forall v \in V$ , there exists a  $v_h \in V_h$  such that,*

$$\|v - v_h\| = \min_{w \in V_h} \|v - w\|. \quad (2.2.21)$$

Of particular interest for time discretisations is the projection onto piecewise constant functions in 1-D. For  $I \subset \mathbb{R}$ , define  $\pi_I^0 : L^2(I) \rightarrow \mathbb{P}_0(I)$  by,

$$\int_I \pi_I^0 v \cdot w \, dt = \int_I v \cdot w \, dt, \quad \forall w \in \mathbb{P}_0(I). \quad (2.2.22)$$

For such a simple projection, we can solve this equation explicitly,

$$\pi_I^0 v = \frac{1}{|I|} \int_I v \, dt, \quad (2.2.23)$$

and we have the following error estimate.

**Theorem 2.2.5.** *The  $L^2(I)$  projection onto the space of constant functions defined in (2.2.22) satisfies the error estimate,*

$$\|v - \pi_I^0 v\|_{L^p(I)} \leq |I| \|v\|_{W^{1,p}(I)}, \quad 1 \leq p \leq \infty. \quad (2.2.24)$$

*Proof.* Follows by Taylor's theorem and the Cauchy-Schwarz inequality.  $\square$

## 2.3 Error analysis

Let  $u$  be the solution to problem (2.1.1) and  $u_h$  be the solution to problem (2.1.7). *A priori* error analysis aims to determine rates of convergence of the approximation scheme. Recall from Céa's lemma (lemma 2.1.4) the characterisation of the Galerkin error as proportional to the best approximation error from the space  $U^h$ ,

$$\|u - u_h\|_U \leq C \inf_{w_h \in U^h} \|u - w_h\|_U. \quad (2.3.1)$$

Then in conjunction with the interpolation error estimates, Céa's lemma leads to an *a priori* error estimate. An example will be given in a later section. An alternative approach is that proposed by Eriksson and Johnson, ([31], [32], [30]) that utilises properties of suitably designed dual problems which can be used to derive  $L^2$  norm bounds.

*A posteriori* error analysis aims to find computable representations and upper bounds for functionals and norms of the error in the approximation, together with estimators that are informative with respect to the error distribution.

**Definition 2.3.1** (Residual). Define the residual of approximation (2.1.7) to problem (2.1.1),  $R(u_h) \in V^*$  by,

$$\langle R(u_h), v \rangle_{V^*, V} := b(e, v) = f(v) - b(u_h, v). \quad (2.3.2)$$

For problems that fit into the framework presented so far, the residual plays an important part in characterising the error, as can be seen in the following lemma.

**Lemma 2.3.1.** *Let  $b : U \times V \rightarrow \mathbb{R}$  satisfy the hypotheses of theorem 2.1.1 and let  $e = u - u_h$  denote the error in the finite element approximation. Then there holds,*

$$c_b \|e\|_U \leq \|R(u_h)\|_{V^*} \leq C_b \|e\|_U, \quad (2.3.3)$$

where  $C_b$  is the constant implied by the continuity hypothesis of theorem 2.1.1, and  $c_b$  is the constant appearing in (2.1.2).

*Proof.* Since  $b : U \times V \rightarrow \mathbb{R}$  is continuous, we have from the definition of the residual (2.3.2),

$$\langle R(u_h), v \rangle_{V^*, V} = b(e, v) \leq C_b \|e\|_U \|v\|_V. \quad (2.3.4)$$

Dividing by  $\|v\|_V$  and taking the supremum over all  $v \neq 0$  implies that,

$$\|R(u_h)\|_{V^*} \leq C_b \|e\|_U. \quad (2.3.5)$$

To prove the other way, we have from condition (2.1.2),

$$c_b \|e\|_U \leq \sup_{v \in V} \frac{b(e, v)}{\|v\|_V} = \sup_{v \in V} \frac{\langle R(u_h), v \rangle_{V^*, V}}{\|v\|_V} = \|R(u_h)\|_{V^*}. \quad (2.3.6)$$

□

**Lemma 2.3.2.** *The kernel of the residual is the discrete test space,  $\ker R(u_h) = V^h$ , that is,*

$$\langle R(u_h), v_h \rangle_{V^*, V} = 0, \quad \forall v_h \in V^h. \quad (2.3.7)$$

*Proof.* Taking  $v \in V^h$  in (2.3.2) leads to (2.1.7). □

In the Hilbert space context the above property is called ‘‘Galerkin orthogonality’’.

**Definition 2.3.2.** Let  $B^* : V \mapsto U^*$ , be the adjoint of the operator  $B : U \mapsto V^*$  satisfying the relationship,

$$\langle w, B^*v \rangle_{U, U^*} = \langle Bw, v \rangle_{V^*, V}, \quad \forall w \in U, v \in V. \quad (2.3.8)$$

Then for  $\phi \in U^*$  given, define the dual problem as: Find  $z \in V$ , such that,

$$B^*z = \phi, \quad \text{in } U^*. \quad (2.3.9)$$

Suppose that the dual problem has a unique solution  $z$ , and that when  $z \in Z \subset V$ ,  $Z$  a subspace of  $V$ , it satisfies the estimate,

$$\|z\|_V \leq C_{stab} \|\phi\|_{U^*}. \quad (2.3.10)$$

**Lemma 2.3.3** (Error Representation). *Let  $z$  be the solution of the dual problem (2.3.9), then the following representation of a linear functional of the error holds,*

$$\langle e, \phi \rangle_{U, U^*} = \langle R(u_h), z - z_h \rangle_{V^*, V}, \quad \forall z_h \in V^h. \quad (2.3.11)$$

*Proof.* Using the definition of the dual problem (2.3.9), the residual (2.3.2) and the lemma (2.3.2), it follows that,

$$\begin{aligned} \langle e, \phi \rangle_{U, U^*} &= \langle e, B^*z \rangle_{U, U^*}, \\ &= \langle Be, z \rangle_{V^*, V}, \\ &= \langle f - Bu_h, z \rangle_{V^*, V}, \\ &= \langle R(u_h), z \rangle_{V^*, V}, \\ &= \langle R(u_h), z - z_h \rangle_{V^*, V}. \end{aligned}$$

□

As mentioned in section 1.2 there are two ways of proceeding from this point. The DWR approach is to compute a localised form of the right-hand side of the above equality. The alternative is to use strong stability of the dual solution (2.3.10) and we now sketch that approach. First, we must have an  $L^2$  representation of the residual, so that we in fact have a representation of the error as,

$$\langle e, \phi \rangle_{U, U^*} = (R(u_h), z - z_h)_{L^2(\Omega)}, \quad \forall z_h \in V^h. \quad (2.3.12)$$

Suppose that there exists an operator (e.g., an interpolation operator as given in section 2.2)  $\Pi : V \mapsto V^h$ , satisfying the estimate,

$$\|h^{-r}(I - \Pi)z\|_{L^2(\Omega)} \leq C_{int}\|z\|_V. \quad (2.3.13)$$

Then,

$$\begin{aligned} \langle e, \phi \rangle_{U, U^*} &= (R(u_h), (I - \mathcal{I})z)_{L^2(\Omega)} \\ &\leq \|h^r R(u_h)\|_{L^2(\Omega)} \|h^{-r}(I - \mathcal{I})z\|_{L^2(\Omega)} \\ &\leq C_{int} \|h^r R(u_h)\|_{L^2(\Omega)} \|z\|_V. \end{aligned}$$

Then using (2.3.10) we can arrive at the abstract error estimate,

$$\|e\|_U = \sup_{\phi \in U^*} \frac{\langle e, \phi \rangle_{U, U^*}}{\|\phi\|_{U^*}} \leq C_{int} C_{stab} \|h^r R(u_h)\|_{L^2(\Omega)}. \quad (2.3.14)$$

This gives us a computable upper bound on the error. The determination of the constant  $C_{int}$  is a problem of approximation theory and there are several works dedicated to determining optimal constants for various quasi-interpolation operators, see for example [19] and [79] for those with a specific slant towards FEMs. The constant  $C_{stab}$  can be determined by a stability analysis of the dual problem.

Given the computable upper bound given in (2.3.14), the problem is now to derive informative error indicators. Let  $\eta$  denote a given error indicator, two concepts that can be used to determine the usefulness of an error estimator are reliability and efficiency [15].

**Definition 2.3.3** (Reliability). An estimator  $\eta$  is called reliable if there is a constant,  $C_R > 0$  and a bound such that,

$$\|u - u_h\| \leq C_R \eta + o(\|u - u_h\|). \quad (2.3.15)$$

**Definition 2.3.4** (Efficiency). An estimator is called efficient if there is a constant,  $C_E > 0$  and a bound such that,

$$\eta \leq C_E \|u - u_h\| + o(\|u - u_h\|). \quad (2.3.16)$$

Reliability provides insurance that the error is bounded above by the estimator and terms that decay faster than the error as the approximation improves. Efficiency then provides reassurance that the estimator is of the same order as the error as the error decays.

**Definition 2.3.5** (Asymptotic Exactness). An estimator is called asymptotically exact if it is reliable and efficient with  $C_R = C_E^{-1}$ .

The main technique for proving efficiency of residual based error estimators for stationary problems is due to Verfürth [77]. It is usual that for more general problems an explicit form of the inequality (2.3.16) is difficult to prove. Therefore it is typical in such cases to prove an *a priori* upper bound for the *a posteriori* error estimate, and show that the estimator and the error converge at the same rate as the discretisation is refined.

## 2.4 AFEM for linear elasticity

In this section the adaptive finite element algorithm of Morin, Nochetto and Siebert (see [54] and [55]) is presented in the context of linear elasticity. A residual based error estimator ([44],[78]) is used in conjunction with the data oscillation measure introduced in [54] to drive the adaptive process. The proof of convergence is based around constructing a procedure which is a contraction mapping of the error.

### The boundary value problem of linear elasticity

Identify a material body with a polygonal domain  $\Omega \subset \mathbb{R}^d$ . Partition the boundary  $\partial\Omega$  into two disjoint subsets  $\Gamma_N$  and  $\Gamma_D$ , where  $\Gamma_D$  has positive  $(d-1)$  surface measure. Denote by  $f : \Omega \mapsto \mathbb{R}^d$  the forces acting throughout the body and by  $g : \Gamma_N \mapsto \mathbb{R}^d$ , the tractions acting on  $\Gamma_N$ . The equilibrium equations are (1.5.1) together with the boundary conditions (1.5.2) and (1.5.3). For isotropic linear elasticity, the stress tensor is related to the strain tensor by Hooke's law, which describes the stress at a given point by the action of the tensor  $\mathbf{C} = (C_{ijkl})_{i,j,k,l=1}^d$  on the strain  $\epsilon$ ,

$$\mathbf{C}\epsilon = \lambda \text{tr}\epsilon + 2\mu\epsilon, \quad \text{tr}A = \sum_{i=1}^d A_{ii}, \quad A \in \mathbb{R}^{d \times d}. \quad (2.4.1)$$

The material constants  $\lambda$  and  $\mu$  are described in section 1.4. However, for the existence result, we only assume that  $\mathbf{C} = (C_{ijkl})_{1 \leq i,j,k,l \leq d}$  is a positive definite fourth order tensor satisfying the symmetries (1.4.20), with  $C_{ijkl} \in L^\infty(\Omega)$ ,  $1 \leq i, j, k, l \leq d$ .

To construct the weak form of the problem formed by equations (1.5.1) and (2.4.1) together with the boundary conditions (1.5.2) and (1.5.3), let  $V$  be the space introduced in (1.5.5). Taking the Euclidean inner product of equation (1.5.1) with a test function  $v \in V$ , integrating by parts and using symmetry of the stress  $\boldsymbol{\sigma}$ , results in,

$$(\mathbf{C}\boldsymbol{\epsilon}(u), \boldsymbol{\epsilon}(v))_{\Omega} = (f, v)_{\Omega} + (g, v)_{\Gamma_N}, \quad \forall v \in V. \quad (2.4.2)$$

Define the bilinear form  $a(\cdot, \cdot)$  by,

$$a(u, v) := (\mathbf{C}\boldsymbol{\epsilon}(u), \boldsymbol{\epsilon}(v))_{\Omega}, \quad (2.4.3)$$

and the linear functional  $l$  by,

$$\langle l, v \rangle := (f, v)_{\Omega} + (g, v)_{\Gamma_N}. \quad (2.4.4)$$

The weak problem is to find  $u \in V$  such that,

$$a(u, v) = \langle l, v \rangle, \quad \forall v \in V. \quad (2.4.5)$$

The existence of a unique solution can be established by theorem 2.1.1 and we briefly quote the results that verify the hypotheses of theorem 2.1.1 for (2.4.5).

**Lemma 2.4.1.** *Let  $\text{meas}(\Gamma_D) \neq 0$ , and assume that  $\mathbf{C} = (C_{ijkl})_{1 \leq i, j, k, l \leq d}$  is positive definite with  $C_{ijkl} \in L^{\infty}(\Omega)$ , then there exist positive constants  $c_a, C_a$  such that the bilinear form defined in equation (2.4.3) satisfies,*

$$c_a \|v\|_{H^1(\Omega)} \leq a(v, v)^{1/2} \quad \text{and} \quad |a(v, w)| \leq C_a \|v\|_{H^1(\Omega)} \|w\|_{H^1(\Omega)}, \quad \forall v, w \in H^1(\Omega).$$

*Proof.* The existence of the constant  $c_a$  is a corollary of Korn's Inequality ([29],[46]), together with the positive definiteness of  $\mathbf{C}$ . The existence of the constant  $C_a$  does not require the assumption on  $\Gamma_D$  and follows from  $C_{ijkl} \in L^{\infty}(\Omega)$  and the Cauchy-Schwarz and Youngs' inequalities.  $\square$

Lemma 2.4.1 implies that  $a(u, u)^{1/2}$  is a norm on  $V$  so define the "energy" norm,

$$\|\cdot\| = a(\cdot, \cdot)^{1/2}. \quad (2.4.6)$$

To show continuity of the linear functional (2.4.4), an upper bound on the boundary values of functions from  $H^1(\Omega)$  is needed. Such a bound follows from the Trace theorem.

**Theorem 2.4.2** (Trace Theorem). *Let  $\Omega$  be open bounded with Lipschitz boundary, let  $1 \leq p \leq \infty$ . Then there exists a unique continuous linear map  $\gamma : W^{1,p}(\Omega) \mapsto L^p(\partial\Omega)$ , called the trace operator, such that, if  $v \in C^0(\bar{\Omega}) \cap W^{1,p}(\Omega)$ , then  $\gamma v = v|_{\partial\Omega}$ .*

*Proof.* See [42]. □

Theorem 2.4.2 implies the existence of a constant such that,

$$\|\gamma v\|_{L^p(\partial\Omega)} \leq C_{\partial\Omega} \|v\|_{W^{1,p}(\Omega)}, \quad \forall v \in W^{1,p}(\Omega). \quad (2.4.7)$$

**Lemma 2.4.3.** *Suppose that  $f \in L^2(\Omega)$  and  $g \in L^2(\Gamma_N)$ , then the linear functional  $l : V \mapsto \mathbb{R}$  defined by (2.4.4) is continuous on  $V$ .*

*Proof.* By Cauchy-Schwarz and theorem 2.4.2,

$$\begin{aligned} |\langle l, v \rangle| &\leq \|f\|_{L^2(\Omega)} \|v\|_{L^2(\Omega)} + \|g\|_{L^2(\Gamma_N)} \|v\|_{L^2(\Gamma_N)} \\ &\leq \|f\|_{L^2(\Omega)} \|v\|_{H^1(\Omega)} + C_\Gamma \|g\|_{L^2(\Gamma_N)} \|v\|_{H^1(\Omega)} \\ &\leq c_a^{-1} \|f\|_{L^2(\Omega)} \|v\| + c_a^{-1} C_\Gamma \|g\|_{\Gamma_N} \|v\| \\ \|l\|_{V^*} &\leq c_a^{-1} \|f\|_{L^2(\Omega)} + c_a^{-1} C_\Gamma \|g\|_{L^2(\Gamma_N)}. \end{aligned}$$

□

Therefore, by lemmas 2.4.1, 2.4.3 together with theorem 2.4.2, theorem 2.1.1 implies the existence and uniqueness of a solution to problem (2.4.2) satisfying the bound,

$$\|u\| \leq c_a^{-1} \|f\|_{L^2(\Omega)} + c_a^{-1} C_{\Gamma_N} \|g\|_{L^2(\Gamma_N)}. \quad (2.4.8)$$

### Finite element approximation

We make an approximation using piecewise linear Lagrange finite element methods as described in section 2.2. Let  $V_h = V \cap S_1(\mathcal{T}_h; \mathbb{R}^d)$ , then the finite element problem is: Find,  $u_h \in V_h$  such that,

$$a(u_h, v) = \langle l, v \rangle \quad \forall v \in V_h. \quad (2.4.9)$$

Problem (2.4.9) results in the symmetric system of linear equations: Find  $\mathbf{u} \in \mathbb{R}^d$ , such that,

$$A\mathbf{u} = \mathbf{f}, \quad A \in \mathbb{R}^{n \times n}, \quad A_{ij} = a(\psi_j, \psi_i), \quad f_i = \langle l, \psi_i \rangle. \quad (2.4.10)$$

Since the bilinear form is symmetric positive definite on  $V$  and the approximation is conforming, the matrix  $A$  in the system (2.4.10) is symmetric positive definite, and so there exists a unique solution to the finite element problem (2.4.9). Céa's lemma (2.1.4) and the interpolation error estimates (2.2.2) give the *a priori* error estimate,

$$\|u - u_h\|_{H^1(\Omega)} \leq ch|u|_{H^2(\Omega)}. \quad (2.4.11)$$

### Residual based *a posteriori* estimation

Let  $e = u - u_h$  denote the error between the finite element solution and the true solution. It is shown in [78] that the energy norm of the error is bounded above and below by a dual norm of the residual of the finite element solution.

**Lemma 2.4.4.** *The residual of the approximation of problem (2.4.9) to problem (2.4.2) has the following localised representation,*

$$\langle R(u_h), v \rangle = \sum_{K \in \mathcal{T}} \left\{ (f, v)_K + \sum_{E \in \mathcal{E}(K)} (R_E, v)_E \right\}, \quad (2.4.12)$$

where,

$$R_E = \begin{cases} \frac{1}{2} \llbracket \boldsymbol{\sigma}(u_h) \rrbracket, & \text{on } E \in \mathcal{E}(\bar{\Omega}) \setminus \mathcal{E}(\Gamma_N), \\ g - \boldsymbol{\sigma}(u_h)n_E, & \text{on } E \in \mathcal{E}(\Gamma_N). \end{cases} \quad (2.4.13)$$

*Proof.* From the definition of the residual, integration by parts over  $\Omega$  gives,

$$\begin{aligned} \langle R(u_h), v \rangle &= \langle l, v \rangle - a(u_h, v), \\ &= \sum_{K \in \mathcal{T}} (f, v)_K + \sum_{E \in \mathcal{E}(\Gamma_N)} (g, v)_{\Gamma_N}, \\ &\quad + \sum_{K \in \mathcal{T}} \left\{ (\operatorname{div} \boldsymbol{\sigma}(u_h), v)_K - (\boldsymbol{\sigma}(u_h)n_{\partial K}, v)_{\partial K} \right\}. \end{aligned}$$

Since  $u_h$  is piecewise linear and  $\mathbf{C}$  is constant over the domain, the divergence of the stress term is zero on each element. Collection of the boundary integrals to form jumps and portioning half to each element sharing that edge leads to the definition of  $R_E$  in (2.4.13).  $\square$

**Theorem 2.4.5** (Upper Bound). *There exists a constant  $C_{\text{rel}}$  depending on the domain  $\Omega$ , the coercivity constant  $c_a$ , the minimum angle in the domain through  $c_E$ ,  $c_K$  and the*

maximum number of overlapping element neighbourhoods  $c_M$  such that the residual (2.4.12) satisfies the bound,

$$\|R(u_h)\|_{V^*} \leq C_{\text{rel}} \left( \sum_{K \in \mathcal{T}} \eta_K^2 \right)^{1/2},$$

where,

$$\eta_K^2 := h_K^2 \|f\|_{L^2(K)}^2 + \sum_{E \in \mathcal{E}(K)} h_E \|R_E\|_{L^2(E)}^2. \quad (2.4.14)$$

*Proof.* By Galerkin orthogonality,

$$\begin{aligned} \langle R(u_h), v \rangle &= \langle R(u_h), v - \mathcal{I}_h^1 v \rangle \\ &= \sum_{K \in \mathcal{T}} \left\{ (f, v - \mathcal{I}_h^1 v)_K + \sum_{E \in \mathcal{E}(K)} (R_E, v - \mathcal{I}_h^1 v)_E \right\}. \end{aligned}$$

The interpolation estimates (2.2.18), (2.2.19), followed by repeated application of the Cauchy-Schwarz inequality give,

$$\begin{aligned} \langle R(u_h), v \rangle &\leq \sum_{K \in \mathcal{T}} \left\{ c_K h_K \|f\|_{L^2(K)} |v|_{H^1(\tilde{\omega}_K)} + \sum_{E \in \mathcal{E}(K)} c_E h_E^{1/2} \|R_E\|_{L^2(K)} |v|_{H^1(\tilde{\omega}_E)} \right\}, \\ &\leq \max\{c_K, c_E\} \left\{ \sum_{K \in \mathcal{T}} h_K^2 \|f\|_{L^2(K)}^2 + \sum_{E \in \mathcal{E}} h_E \|R_E\|_{L^2(E)}^2 \right\}^{1/2} \\ &\quad \times \left\{ \sum_{K \in \mathcal{T}} |v|_{H^1(\tilde{\omega}_K)}^2 + \sum_{E \in \mathcal{E}} |v|_{H^1(\tilde{\omega}_E)}^2 \right\}^{1/2}, \\ &\leq c_M \max\{c_K, c_E\} \left( \sum_{K \in \mathcal{T}} \eta_K^2 \right)^{1/2} |v|_{H^1(\Omega)}. \end{aligned}$$

Then using coercivity,  $|v|_{H^1(\Omega)} \leq \|v\|_{H^1(\Omega)} \leq \frac{1}{c_a} \|v\|$ , the result follows with,

$$C_{\text{rel}} = \frac{M \max\{c_K, c_E\}}{c_a}. \quad (2.4.15)$$

□

The term  $\eta_K$  defined in (2.4.14) is called the local error indicator of the element  $K$ . We generalise the notion of local error indicators to deal with error indicators for collections of elements. The indicator for the sub-domain  $\omega \subset \Omega$  is defined as,

$$\eta_\omega^2 := \sum_{K \subseteq \omega} \eta_K^2.$$

To prove lower bounds on the error estimator, the construction of Verfürth [77] will be used, which is based on the properties of the bubble functions of elements and edges defined as,

$$b_K := (d+1)^{d+1} \prod_{z \in \mathcal{N}(K)} \phi_z, \quad (2.4.16)$$

$$b_E := d^d \prod_{z \in \mathcal{N}(E)} \phi_z. \quad (2.4.17)$$

Given the definitions (2.4.16) and (2.4.17) then for all  $v \in \mathbb{P}_k(T)$  and  $w \in \mathbb{P}_k(E)$  the following inverse estimates hold [77], [80],

$$\|v\|_{L^2(K)} \leq \varepsilon_1 \|b_K^{1/2} v\|_{L^2(K)}, \quad (2.4.18)$$

$$\|\nabla(b_K v)\|_{L^2(K)} \leq \varepsilon_2 h_K^{-1} \|v\|_{L^2(K)}, \quad (2.4.19)$$

$$\|w\|_{L^2(E)} \leq \varepsilon_3 \|b_E^{1/2} w\|_{L^2(E)}, \quad (2.4.20)$$

$$\|\nabla(b_E w)\|_{L^2(\omega_E)} \leq \varepsilon_4 h_E^{-1/2} \|w\|_{L^2(E)}, \quad (2.4.21)$$

$$\|b_E w\|_{L^2(\omega_E)} \leq \varepsilon_5 h_E^{1/2} \|w\|_{L^2(E)}. \quad (2.4.22)$$

For further details on the values of the constants see [80].

**Theorem 2.4.6** (Lower Bound). *(Verfürth) There exists a constant  $C_{\text{eff}}$  depending on  $C_a$  and the constants  $\{\varepsilon_i\}_{i=1}^5$  such that the local error indicator (2.4.14) satisfies,*

$$\eta_K^2 \leq C_{\text{eff}}^2 \left\{ \|e\|_{\omega_K}^2 + \sum_{K' \in \omega_K} h_{K'}^2 \|f - f_{K'}\|_{L^2(K')}^2 + \sum_{E \in \mathcal{E}(K) \cap \mathcal{E}(\Gamma_N)} h_E \|g - g_K\|_{L^2(E)}^2 \right\}. \quad (2.4.23)$$

*Proof.* The proof proceeds in three stages, one stage for each term in the error indicator. Step 1: Element terms. Let  $f_K$  denote the  $L^2$ -projection of  $f$  onto  $\mathbb{P}_0(K)$ . From (2.4.18) it holds that,

$$\|f_K\|_{L^2(K)}^2 \leq \varepsilon_1^2 (f_K, f_K b_K)_K = \varepsilon_1^2 (f, f_K b_K)_K + \varepsilon_1^2 (f_K - f, f_K b_K)_K. \quad (2.4.24)$$

From the representation of the residual (2.4.12) and the fact that  $b_K$  vanishes on  $\partial K$ , it follows that,

$$a(e, f_K b_K) = \begin{cases} (f, f_K b_K)_K, & \text{on } K, \\ 0, & \text{otherwise.} \end{cases} \quad (2.4.25)$$

Combining (2.4.24) and (2.4.25), together with the estimate (2.4.19) results in,

$$\begin{aligned}
\|f_K\|_{L^2(K)}^2 &\leq \varepsilon_1^2 a(e, f_K b_K) + \varepsilon_1^2 (f_K - f, f_K b_K)_K, \\
&\leq \varepsilon_1^2 \|e\|_K \|f_K b_K\|_K + \varepsilon_1^2 \|f - f_K\|_{L^2(K)} \|f_K b_K\|_{L^2(K)}, \\
&\leq C_a \varepsilon_1^2 \|e\|_K |f_K b_K|_{H^1(K)} + \varepsilon_1^2 \|f - f_K\|_{L^2(K)} \|f_K b_K\|_{L^2(K)}, \\
&\leq C_a \varepsilon_1^2 \varepsilon_2 h_K^{-1} \|e\|_K \|f_K\|_{L^2(K)} + \varepsilon_1^2 \|f - f_K\|_{L^2(K)} \|f_K\|_{L^2(K)}, \\
\Rightarrow h_K \|f_K\|_{L^2(K)} &\leq C_a \varepsilon_1^2 \varepsilon_2 \|e\|_K + \varepsilon_1^2 h_K \|f - f_K\|_{L^2(K)},
\end{aligned}$$

after multiplication by  $h_K$ . Using the triangle inequality,

$$h_K \|f\|_{L^2(K)} \leq h_K \|f_K\|_{L^2(K)} + h_K \|f - f_K\|_{L^2(K)}, \quad (2.4.26)$$

and the previous bound on  $h_K \|f_K\|_{L^2(K)}$ , the following holds for the volume term in the error indicator,

$$h_K \|f\|_{L^2(K)} \leq C_a \varepsilon_1^2 \varepsilon_2 \|e\|_K + (1 + \varepsilon_1^2) h_K \|f - f_K\|_{L^2(K)}. \quad (2.4.27)$$

Step 2: Internal edges. Terms corresponding to internal edges are half the jump in the stress across that edge. Since the approximation is by piecewise linear functions and  $\mathbf{C}$  is constant throughout the domain, this term is constant. From (2.4.20) and the representation of the residual (2.4.12) tested on  $R_E b_E$  we get,

$$\begin{aligned}
\|R_E\|_E^2 &\leq \varepsilon_3^2 (R_E, R_E b_E)_E, \\
&\leq \varepsilon_3^2 a(e, R_E b_E) - \varepsilon_3^2 \sum_{K \in \omega_E} (f, R_E b_E)_K.
\end{aligned}$$

Applying the Cauchy-Schwarz inequality and using the estimates (2.4.20) and (2.4.21) results in,

$$\begin{aligned}
\|R_E\|_E^2 &\leq \varepsilon_3^2 \|e\|_{\omega_E} \|R_E b_E\|_{\omega_E} + \varepsilon_3^2 \|f\|_{L^2(\omega_E)} \|R_E b_E\|_{L^2(\omega_E)}, \\
&\leq \varepsilon_3^2 \|e\|_{\omega_E} |R_E b_E|_{H^1(\omega_E)} + \varepsilon_3^2 \|f\|_{L^2(\omega_E)} \|R_E b_E\|_{L^2(\omega_E)}, \\
&\leq C_a \varepsilon_3^2 \varepsilon_4 h_E^{-1/2} \|e\|_{\omega_E} \|R_E\|_{L^2(E)} + \varepsilon_3^2 \varepsilon_5 h_E^{1/2} \|f\|_{L^2(\omega_E)} \|R_E\|_{L^2(E)}, \\
\Rightarrow h_E^{1/2} \|R_E\|_E &\leq C_a \varepsilon_3^2 \varepsilon_4 \|e\|_{\omega_E} + \varepsilon_3^2 \varepsilon_5 h_E \|f\|_{L^2(\omega_E)}.
\end{aligned}$$

Since  $h_E \leq h_K$ , (2.4.27) implies that,

$$h_E^{1/2} \|R_E\|_E \leq C_a \varepsilon_3^2 (\varepsilon_1^2 \varepsilon_2 \varepsilon_5 + \varepsilon_4) \|e\|_{\omega_E} + 2(1 + \varepsilon_1^2) \varepsilon_3^2 \varepsilon_5 \sum_{K \in \omega_K} h_K \|f - f_K\|_{L^2(K)}. \quad (2.4.28)$$

Step 3. Boundary Edges. In this instance,  $R_E = g - \boldsymbol{\sigma}(u_h)n_E$ . Let  $\bar{R}_E := g_E - \boldsymbol{\sigma}(u_h)n_E$  where  $g_E$  is the  $L^2$ -projection of  $g$  onto piecewise constants. From (2.4.20),

$$\|\bar{R}_E\|_E^2 \leq \varepsilon_3^2 (\bar{R}_E, \bar{R}_E b_E)_E = \varepsilon_3^2 (R_E, \bar{R}_E b_E)_E + \varepsilon_3^2 (g_E - g, \bar{R}_E b_E)_E. \quad (2.4.29)$$

Combining the representation of the residual (2.4.12) tested against  $\bar{R}_E b_E$  of the residual with (2.4.29) results in,

$$\|\bar{R}_E\|_E^2 \leq \varepsilon_3^2 a(e, \bar{R}_E b_E) - \varepsilon_3^2 (f, \bar{R}_E b_E)_K + \varepsilon_3^2 (g_E - g, \bar{R}_E b_E)_E,$$

where  $\omega_E = K$  since the edge is on the boundary. Estimating as before gives,

$$\begin{aligned} \|\bar{R}_E\|_E^2 &\leq \varepsilon_3^2 \|e\|_K \|\bar{R}_E b_E\|_K + \varepsilon_3^2 \|f\|_{L^2(K)} \|\bar{R}_E b_E\|_{L^2(K)} + \varepsilon_3^2 \|g_E - g\|_{L^2(E)} \|\bar{R}_E b_E\|_E, \\ &\leq C_a \varepsilon_3^2 \|e\|_K \|\bar{R}_E b_E\|_{H^1(K)} + \varepsilon_3^2 \varepsilon_5 h_E^{1/2} \|f\|_{L^2(K)} \|\bar{R}_E\|_{L^2(E)} \\ &\quad + \varepsilon_3^2 \|g_E - g\|_{L^2(E)} \|\bar{R}_E\|_{L^2(E)}, \\ &\leq C_a \varepsilon_3^2 \varepsilon_4 h_E^{-1/2} \|e\|_K \|\bar{R}_E\|_{L^2(E)} + \varepsilon_3^2 \varepsilon_5 h_E^{1/2} \|f\|_{L^2(K)} \|\bar{R}_E\|_{L^2(E)} \\ &\quad + \varepsilon_3^2 \|g_E - g\|_{L^2(E)} \|\bar{R}_E\|_{L^2(E)}. \end{aligned}$$

Multiplying by  $h_E^{1/2}$  and dividing through by  $\|\bar{R}_E\|_{L^2(E)}$  implies,

$$h_E^{1/2} \|\bar{R}_E\|_E^2 \leq C_a \varepsilon_3^2 \varepsilon_4 h_E^{-1/2} \|e\|_K + \varepsilon_3^2 \varepsilon_5 h_K \|f\|_{L^2(K)} + \varepsilon_3^2 h_E^{1/2} \|g_E - g\|_{L^2(E)}, \quad (2.4.30)$$

then using (2.4.27) it follows that there exists a constant independent of the mesh width such that,

$$h_E^{1/2} \|\bar{R}_E\|_E^2 \leq C \left\{ \|e\|_K + h_K \|f - f_K\|_{L^2(K)} + h_E^{1/2} \|g_E - g\|_{L^2(E)} \right\}. \quad (2.4.31)$$

Using the triangle inequality,  $\|R_E\|_E \leq \|\bar{R}_E\|_E + \|g - g_E\|_E$ , the following estimate holds for the boundary edges,

$$h_E^{1/2} \|R_E\|_E^2 \leq C \left\{ \|e\|_K + h_K \|f - f_K\|_{L^2(K)} + h_E^{1/2} \|g_E - g\|_{L^2(E)} \right\}. \quad (2.4.32)$$

Squaring each of (2.4.27), (2.4.28) and (2.4.32) and using Young's inequality on the mixed product terms it follows that there is a constant  $C_{\text{eff}}$  depending on the constants  $\{\varepsilon_i\}_{i=1}^5$  such that the conclusion holds.  $\square$

Using Verfürth's construction of a lower bound for the residual based error estimator of theorem 2.4.5 we have shown the following.

**Corollary 2.4.7.** *The error estimator of theorem 2.4.5 is reliable and efficient.*

Therefore, the residual based estimator of theorem 2.4.5 describes, up to higher order terms, global upper and lower bounds on the error of the finite element solution. In recognition of the higher order terms, the following measure of variation in the problem data is introduced in [54].

**Definition 2.4.1** (Data Oscillation). Let  $f_K$  and  $g_E$  denote piecewise constant approximations to  $f$  and  $g$  on element  $K$  and edge  $E$  respectively. The data oscillation on the subset of elements  $\omega \subset \Omega$  according to the triangulation  $\mathcal{T}_h$  is defined as,

$$\text{osc}_h(\omega)^2 := \sum_{K \subset \omega} \left\{ h_K^2 \|f - f_K\|_{L^2(K)}^2 + \sum_{E \in \mathcal{E}(K) \cap E(\Gamma_N)} h_E \|g - g_E\|_{L^2(E)}^2 \right\}.$$

The concept of data oscillation is vital in the proof of convergence of an adaptive scheme in [54].

**Corollary 2.4.8.** *The error estimator of theorem 2.4.5 satisfies,*

$$\frac{1}{MC_{\text{eff}}^2} \sum_{K \in \mathcal{T}_h} \eta_K^2 - \text{osc}_h(\bar{\Omega})^2 \leq \|e\|^2 \leq C_{\text{rel}}^2 \sum_{K \in \mathcal{T}_h} \eta_K^2. \quad (2.4.33)$$

*Proof.* Summing (2.4.23) over the elements leads to,

$$\sum_{K \in \mathcal{T}_h} \eta_K^2 \leq MC_{\text{eff}}^2 \text{osc}_h(\bar{\Omega})^2 + MC_{\text{eff}}^2 \|e\|^2, \quad (2.4.34)$$

and so the result follows by combining the upper bound (2.4.5).  $\square$

From the above corollary it is now clearer the role that data oscillation plays in the development of an adaptive algorithm. To tighten the lower bound the data oscillation must be reduced. This is one of the main features of the algorithm of Morin, Nochetto and Siebert presented next.

## The adaptive algorithm of MNS

Armed with an *a posteriori* error estimate to drive an adaptive procedure we turn to the issues regarding the design of such a procedure. With a local error indicator a marking strategy is used to determine which elements are to be refined.

**Marking strategy**

Let  $\mathcal{M}_R \subset \mathcal{T}_h$  denote the elements that are to be refined and define,

$$\mathcal{M}_R := \bigcup_{K \in \mathcal{M}_R} K.$$

For a collection of elements  $S$  define  $\eta(S) := \sum_{K \in S} \eta_K^2$ , so  $\eta(S)$  is the contribution to the error estimate of the collection. The proof of convergence provided in [54] requires that the data oscillation is tackled at each step of the algorithm. This leads to a modification of the traditional marking schemes such as that in [28], where a percentage of elements are chosen to be refined only from the information gleaned from the estimator.

**Marking Strategy MS:**


---

Given  $\theta_0, \theta_1$ ,  $0 < \theta_0, \theta_1 < 1$ :

1. Construct the minimal subset of elements  $\mathcal{M}_R \subset \mathcal{T}$  such that,

$$\eta(\mathcal{M}_R) \geq \theta_0 \eta(\mathcal{T}).$$

2. Enlarge  $\mathcal{M}_R$  so that,

$$\text{osc}_h(\mathcal{M}_R) \geq \theta_1 \text{osc}_h(\Omega).$$


---

In practice, the set  $\mathcal{M}_R$  is constructed by sequentially taking the elements with the largest error indicators, and the enlargement in step 2 proceeds analogously. The idea of the marking strategy is to first identify those elements that contribute a fraction  $\theta_0$  to the total error, and then enlarge this set by including those elements that make significant contributions to the data oscillation. The full adaptive algorithm of MNS is given below [54].

**MNS Algorithm:**


---

Choose  $\theta_0, \theta_1$ ,  $0 < \theta_0, \theta_1 < 1$ :

1. Construct  $\mathcal{T}_0$  such that coefficients are resolved as constants over the domain. Set  $k = 0$ .
2. Solve on  $\mathcal{T}_k$  for  $u_k$ .
3. Compute the estimator  $\eta$ .
4. Mark elements for refinement using marking procedure MS.
6. Refine  $\mathcal{T}_k$  to get  $\mathcal{T}_{k+1}$  using longest edge bisection.
7. Set  $k:=k+1$ . Go to step 2.

---

The proof of convergence of an AFEM using algorithm MNS is based on proving an error reduction property for the algorithm.

**Theorem 2.4.9** (Error Reduction). *Let  $\mathcal{T}_H$  be a triangulation of  $\Omega$  and let  $\mathcal{T}_h$  be a mesh achieved by interior node bisection, then there exist constants  $0 < \vartheta < 1$ ,  $\varrho > 0$  depending on  $c_a$ ,  $C_a$ ,  $\theta_0$  and the minimum angle such that for any  $\epsilon > 0$ , if,*

$$\text{osc}_H(\Omega) \leq \varrho\epsilon, \quad (2.4.35)$$

then either  $\|u - u_H\| \leq \epsilon$  or the solution  $u_h \in V^h$  satisfies,

$$\|u - u_h\| \leq \vartheta \|u - u_H\|.$$

*Proof.* See [54]. □

The above theorem states that if the data are sufficiently resolved then a refinement step results in an error reduction. The proof depends on the refinement scheme and requires the introduction of an interior node to all marked elements (see appendix A for more details). Based on the marking strategy and interior node longest edge bisection together with a reliable and efficient estimator, the following appears in [54].

**Theorem 2.4.10.** *Let  $\{u_{h_k}\}_{k \in \mathbb{N}}$  be a sequence of piecewise linear finite element approximations in nested finite element spaces  $\{V^{h_k}\}$  produced by algorithm (MNS), then there exists  $\beta \in (0, 1)$  such that,*

$$\|u - u_{h_k}\| \leq \beta^k. \quad (2.4.36)$$

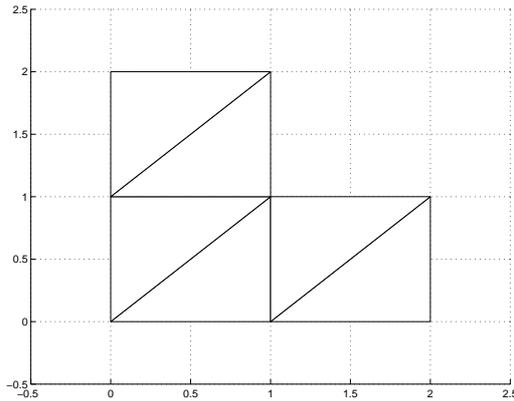


Figure 2.3: Initial triangulation  $\mathcal{T}_0$  of the domain  $\Omega$ .

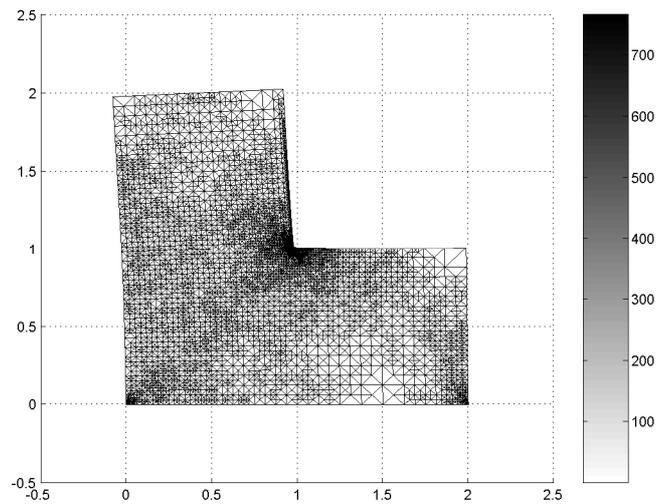
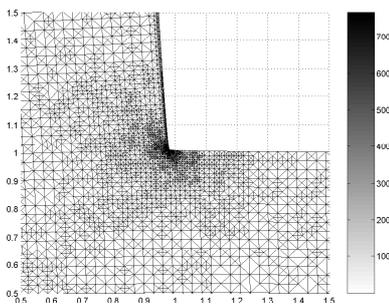
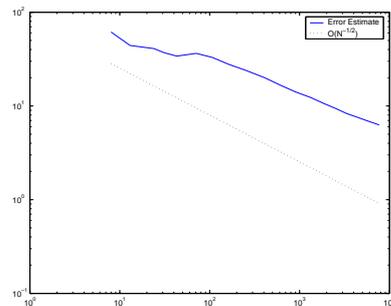
The results of [54] are presented for scalar elliptic Dirichlet problems, however, the results are valid for linear elasticity with mixed boundary conditions with little modification.

As an example, we consider a linear elastic body occupying an L-shape domain as shown in figure 2.3. The displacement is fixed along the boundary  $\Gamma_D := \{(x, y) \mid 0 < x < 2, y = 0\}$ , and the body is subjected to constant surface tractions along the boundary region  $\{(x, y) \mid x = 1, 1 < y < 2\}$ . The material parameters are,

$$E = 50000, \quad \nu = 0.3.$$

In our example we specify the number of iterations, then setting  $\theta = 0.3$  we arrive after 20 iterations at the mesh given in figure 2.4.

In this example the data oscillation term is zero. However we can still see a period where the algorithm is struggling to resolve some feature. However once it passes through it almost achieves the optimal asymptotic convergence rate of  $O(N^{-1/d})$ . Results for cases with oscillatory data are similar. Initial transitory periods persist until the data are resolved and the optimal asymptotic rate is achieved once the data are resolved.

Figure 2.4: The refined mesh,  $\mathcal{T}_{20}$ .Figure 2.5: Close up of  $\mathcal{T}_{20}$ .Figure 2.6:  $\eta_{\mathcal{T}}$  v Dofs

## 2.5 AFEM for linear systems of ordinary differential equations

In this section we consider the continuous Galerkin (cG) method as applied to linear systems of ordinary differential equations. Estep and French [34] provide an analysis of the method applied to the general system of ODEs,

$$z_t + f(z(t), t) = 0, \quad 0 < t \leq T, \quad (2.5.1)$$

$$z(0) = z_0 \in \mathbb{R}^d, \quad d \geq 1. \quad (2.5.2)$$

Using duality techniques (section 2.3) they derive a *a priori* and a *a posteriori* estimates for the approximation error  $e = z - z_h$ . For sufficiently small time step parameter  $k = \max_{1 \leq i \leq N} k_i$ , where  $k_i = t_i - t_{i-1}$ , their *a priori* estimate for the piecewise linear continuous Galerkin method takes the form,

$$|z - z_h|_{W^{m,\infty}(I)} \leq C(1 + LTe^{CLT})^{1/2} \max_{1 \leq i \leq N} k_i^{2-m} |z|_{W^{2,\infty}(I_i)}, \quad m = 0, 1, \quad (2.5.3)$$

where  $L$  is the Lipschitz constant for the function  $f(\cdot, t)$ . Their *a posteriori* estimate takes the form,

$$\max_{1 \leq i \leq N} |e|_{L^\infty(I_i)} \leq S(T) \max_{1 \leq i \leq N} k_i^2 \|Df(z_h(t), t)\|_{L^\infty(I_i)}. \quad (2.5.4)$$

The term  $S(T)$  is the stability factor of the solution to the dual problem, which in this instance is an ODE running backwards in time and the formal adjoint of (2.5.1). We return to the general non-linear situation in a later chapter, and for the time being focus on linear systems. We choose to focus on such systems since they arise in the discretisation of space and time problems, their relevance to this study increased by the fact that the internal variable formulation of linear viscoelasticity introduced in the next chapter involves such a system. It also gives us an opportunity to contrast the differing aspects of adaptivity for spatial and temporal problems and see how they can be combined for space and time adaptivity.

The problem we wish to solve is: Given the symmetric positive definite matrix  $A$ , and vector valued function  $f : \mathbb{R} \rightarrow \mathbb{R}^d$ , Find  $z : \mathbb{R} \rightarrow \mathbb{R}^d$ , such that,

$$z_t(t) + Az(t) = f(t), \quad (2.5.5)$$

$$z(0) = 0. \quad (2.5.6)$$

We assume that the problem is zero at the initial condition since any linear system of ODEs can be transformed into such a system with a suitable modification of  $f$ . In the following sections we will apply the cG(1) finite element method, derive optimal order *a posteriori* error estimates, and discuss adaptive time step selection mechanisms.

### 2.5.1 Finite element approximation

Let  $I = [0, T]$ , and let  $(\cdot, \cdot)$  denote the Euclidean inner product. Define the space  $W_0^{1,p}(I) = W^{1,p}(I) \cap \{v : v(0) = 0\}$ , then the variational formulation of (2.5.5) is: Find  $z \in W_0^{1,p}(I)$

such that,

$$\int_I (z_t + Az, w) dt = \int_I (f, w) dt \quad \forall w \in L^q(I). \quad (2.5.7)$$

Partition the time interval  $I = [0, T]$  into  $N$  subintervals  $I_i := (t_{i-1}, t_i]$  of length  $k_i :=$

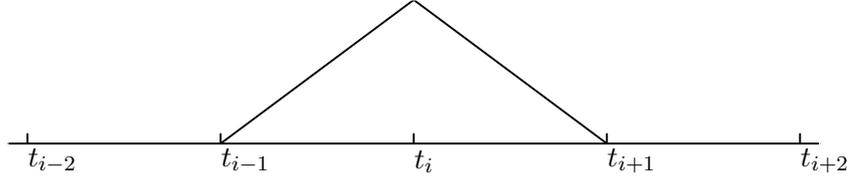


Figure 2.7: Illustration of the temporal basis function  $\phi_i(t)$ .

$t_i - t_{i-1}$ , with  $t_0 = 0$ ,  $t_N = T$  and define  $k := \max_{1 \leq i \leq N} k_i$ . Associated with this partition is the space of nodal basis functions, an example of which is shown in figure 2.5.1. The nodal basis functions are defined by,

$$\phi_0(t) = \begin{cases} \frac{t_1 - t}{k_1}, & t \in I_1, \\ 0, & \text{otherwise,} \end{cases} \quad (2.5.8)$$

$$1 \leq i \leq N - 1, \quad \phi_i(t) = \begin{cases} \frac{t - t_{i-1}}{k_i}, & t \in I_i, \\ \frac{t_i - t}{k_i}, & t \in I_{i+1}, \\ 0, & \text{otherwise,} \end{cases} \quad (2.5.9)$$

$$\phi_N(t) = \begin{cases} \frac{t - t_N}{k_N}, & t \in I_N, \\ 0, & \text{otherwise.} \end{cases} \quad (2.5.10)$$

Let  $\mathcal{T}_k$  denote the mesh parameterised by the meshwidth  $k$ . Then the trial space for the cG(1) method is  $\mathcal{S}_1(\mathcal{T}_k)$ . The test space is defined piecewise, such that on an interval  $I_i$ , it is given by  $\mathbb{P}_0(I_i)$ . Define the test space to be  $\mathcal{D}_0(\mathcal{T}_k) := \{v \mid v|_{I_i} \in \mathbb{P}_0(I_i)\}$ . The finite element problem is then: Find  $z_h \in \mathcal{S}_1(\mathcal{T}_k) \cap \{v \mid v(0) = 0\}$ , such that,

$$\int_I (z_{h,t} + Az_h, w) dt = \int_I (f, w) dt, \quad \forall w \in \mathcal{D}_0(\mathcal{T}_k). \quad (2.5.11)$$

Let  $\pi_i^0 : L^2(I_i) \rightarrow \mathbb{P}_0(I_i)$  be the  $L^2(I_i)$  projection of theorem 2.2.5, and let  $z_h^i = z_h(t_i)$ , then (2.5.11) becomes,

$$\frac{z_h^i - z_h^{i-1}}{k_i} + \frac{1}{2}A(z_h^i + z_h^{i-1}) = \pi_i^0 f, \quad i = 1, \dots, N, \quad (2.5.12)$$

$$z_h^0 = 0. \quad (2.5.13)$$

A sequence of nodal values of the approximate solution can now be generated from stepping scheme (2.5.12). In fact, let,

$$C_i = \left( I + \frac{k_i}{2}A \right)^{-1} \left( I - \frac{k_i}{2}A \right), \quad b_i = \left( I + \frac{k_i}{2}A \right)^{-1} \pi_i^0 f, \quad (2.5.14)$$

then scheme (2.5.12) can be written as,

$$z_h^i = C_i z_h^{i-1} + b_i, \quad 1 \leq i \leq N, \\ z_h^0 = 0.$$

We remark that this linear recurrence can be solved for  $z_h^n$ ,

$$z_h^n = \sum_{i=1}^{n-1} \left( \prod_{j=i+1}^n C_j \right) b_i + b_n, \quad n \geq 1, \quad (2.5.15)$$

in theory bypassing the requirement of a stepping scheme. It is more suitable for us however, to continue with the stepping scheme (2.5.12), and in the following section derive an adaptive time stepping algorithm that enables us to control the error.

### 2.5.2 A posteriori error analysis

To derive an *a posteriori* error estimate, we will use the duality technique based on the introduction of a backward dual problem. Define the dual problem to be: Find  $\chi$  such that,

$$-\chi_t(t) + A\chi(t) = g(t), \quad 0 \leq t < T, \quad (2.5.16)$$

$$\chi(T) = \psi. \quad (2.5.17)$$

We now use the dual problem to generate an error representation formula. First though, we define the residual.

**Definition 2.5.1.** Define the residual of the approximation (2.5.11) to problem (2.5.7) by,

$$\langle R(z_h), w \rangle := \int_I (e_t + Ae, w) dt = \int_I (f - z_{h,t} - Az_h, w) dt. \quad (2.5.18)$$

**Lemma 2.5.1.** Let  $\chi$  be the solution to problem (2.5.16) with (2.5.17), then the error in the approximation (2.5.11) to problem (2.5.7), defined by  $e = z - z_h$  satisfies the following relationship,

$$(e(T), \psi) + \int_I (e, g) dt = \langle R(z_h), \chi \rangle. \quad (2.5.19)$$

The error representation essentially allows us to chose a linear functional of the error on the left hand side. For given choices of  $\psi$  and  $g$ , the effect is transmitted through the dual solution appearing in the right hand side, which is a function of both  $\psi$  and  $g$ . Localising the representation (2.5.19) and using Galerkin orthogonality we get,

$$(e(T), \psi) + \int_I (e, g) dt = \sum_{i=1}^N \int_{I_i} (f - Az_h - z_{h,t}, \chi - \pi_i^0 \chi) dt. \quad (2.5.20)$$

The dual weighted residual method uses representation (2.5.20) to compute estimates of the error. By first choosing specific values of  $\psi$  and  $g$ , the method approximates the value of the error representation by computing an approximation to the dual problem (2.5.16) and evaluating (2.5.20).

In terms of deriving a computable upper bound, rather than approximating the exact value of (2.5.20), we can use the error estimate of theorem 2.2.5 for  $\pi^0$  to get,

$$\left| (e(T), \psi) + \int_I (e, g) dt \right| \leq \sum_{i=1}^N k_i \|f - Az_h - z_{h,t}\|_{L^p(I_i)} |\chi|_{W^{1,q}(I_i)}. \quad (2.5.21)$$

We can now use this estimate to drive an adaptive stepping scheme. For time dependent problems, the situation with adaptivity is very different of that for spatial discretisations. For spatial adaptivity, we can think about equidistribution of the error and various other schemes for dealing with the local error terms, all within a Solve-Estimate-Refine (SER) loop. For a time dependent problem, we perform an SER step at each time level, solving for the solution at the next time value, while reducing the step size until a criterion on the error over the step interval is met. To make this strategy easier to implement, it is common to take the maximum of the local error over all of the elements as the term we wish to

control,

$$\left| (e(T), \psi) + \int_I (e, g) dt \right| \leq \sum_{i=1}^N k_i \|f - Az_h - z_{h,t}\|_{L^p(I_i)} |\chi|_{W^{1,q}(I_i)}, \quad (2.5.22)$$

$$\leq \max_{1 \leq i \leq N} k_i \|f - Az_h - z_{h,t}\|_{L^\infty(I_i)} \sum_{i=1}^N |\chi|_{W^{1,1}(I_i)}, \quad (2.5.23)$$

$$\leq \max_{1 \leq i \leq N} k_i \|f - Az_h - z_{h,t}\|_{L^\infty(I_i)} |\chi|_{W^{1,1}(I)}. \quad (2.5.24)$$

Define the local error indicator by,

$$\eta_i = k_i \|f - Az_h - z_{h,t}\|_{L^\infty(I_i)}. \quad (2.5.25)$$

The two terms of the error estimate (2.5.24) measure two different types of error that are apparent in the solution of initial value problems by finite elements. The local error indicator  $\eta_i$  measures the error due to the local approximation properties of the discrete scheme, the semi-norm of the dual solution measures the cumulative effect of integrating the initial value problem over discrete intervals for the given choice of  $\psi$  and  $g$ .

Therefore, depending on the choice of  $\psi$  and  $g$ , we can choose to control an arbitrary linear functional of the error, so long as we can bound or compute the  $W^{1,1}(I)$  seminorm of the dual problem. Therefore  $S(T) := |\chi|_{1,1}$  is referred to as the stability factor. Suppose that we are interested in the error at final time  $T$ . We can choose  $g = 0$  and  $\psi = e(T)|e(T)|^{-1}$ , and we can determine a bound for  $S(T)$  as follows. Taking the inner product of (2.5.16) with  $-\chi_t(t)$  and integrating over  $I$  results in,

$$\begin{aligned} \|\chi_t\|_{L^2(I)}^2 - \int_I (A\chi, \chi_t) dt &= 0 \\ \|\chi_t\|_{L^2(I)}^2 - \frac{1}{2}(A\chi, \chi)|_0^T &= 0 \\ \|\chi_t\|_{L^2(I)}^2 - \frac{1}{2}(A\chi(T), \chi(T)) + \frac{1}{2}(A\chi(0), \chi(0)) &= 0. \end{aligned}$$

Since  $A$  is positive definite we discard the term involving  $\chi(0)$  to make inequality, then we can use the condition  $|\psi| = 1$  and rearrange to get,

$$\|\chi_t\|_{L^2(I)} \leq \frac{1}{\sqrt{2}} \|A\|_2. \quad (2.5.26)$$

Using  $\|\chi_t\|_{L^1(I)} \leq T^{1/2}\|\chi_t\|_{L^2(I)}$ , we can combine (2.5.24) with the explicit value for  $S(T)$  to get the fully computable *a posteriori* upper bound,

$$|e(T)| \leq S(T) \max_{1 \leq i \leq N} \eta_i. \quad (2.5.27)$$

where the stability factor  $S(T)$  is given by,

$$S(T) = \left(\frac{T}{2}\right)^{1/2} \|A\|_2. \quad (2.5.28)$$

Alternatively, suppose that we wish to control the global error  $\|e\|_{L^\infty(I)}$ . In this instance, we take  $\psi = 0$ , and we leave  $g$  arbitrary. Then assuming that we can get a bound like  $|\chi|_{W^{1,1}(I)} < S(T)\|g\|_{L^1(I)}$ , then we have,

$$\|e\|_{L^\infty(I)} = \sup_{g \in L^1(I)} \frac{(e, g)}{\|g\|_{L^1(I)}} \leq S(T) \max_{1 \leq i \leq N} k_i \|f - Az_h - z_{h,t}\|_{L^\infty(I_i)}. \quad (2.5.29)$$

To get the desired bound on the dual solution, we can use the original ODE (2.5.16) and the triangle inequality to get,

$$\|\chi_t\|_{L^1(I)} \leq \|A\chi\|_{L^1(I)} + \|g\|_{L^1(I)}. \quad (2.5.30)$$

Then using the explicit representation of the solution, we can multiply by  $A$  and get bound on  $\|\chi\|_{L^1(I)}$ ,

$$\chi(t) = \int_t^T e^{A(t-s)}g(s) ds \Rightarrow \|A\chi\|_{L^1(I)} \leq (e^{AT} - 1)\|g\|_{L^1(I)}, \quad (2.5.31)$$

which gives the upper bound for the stability factor corresponding to  $L^\infty(I)$  error control of  $S(T) = e^{AT}$ .

### 2.5.3 Adaptive time-stepping

We will now apply an adaptive solution algorithm to the discrete problem (2.5.12) using the *a posteriori* error estimator (2.5.27). Suppose that we have a global tolerance **GTOL**, and the aim is to compute  $z_h$  so that,

$$\|e\|_{L^\infty(I)} \leq \mathbf{GTOL}. \quad (2.5.32)$$

We will use adaptive error control to ensure that the computation satisfies (2.5.32). A first attempt at time stepping algorithm is given below.

Adaptive time-stepping algorithm:

---

1. Set  $t = 0$ .  $\text{LTOL} = S(T)^{-1}\text{GTOL}$ .
  2. Do:
    - i) Set  $k = \Upsilon_1(k_{\text{old}}, t, T)$ .
    - ii) Calculate  $z_h^{i+1}$  and  $\eta$ .
    - iii) While  $k\eta > \text{LTOL}$ :
      - a) Set  $k = \Upsilon_2(\text{LTOL}, \eta)$ .
      - b) Calculate  $z_h^{i+1}$  and  $\eta$ .
    - iv) Set  $z_i = z_{i+1}$ ,  $t = t + k$ ,  $i = i + 1$ .
- while  $t \leq T$ .
- 

While there is little possible variation on how such an algorithm can proceed, there are two choices to be made in respect of the functions  $\Upsilon_i$ ,  $i = 1, 2$ . The function  $\Upsilon_1$  determines what step size should be chosen for the initial solve at each time level. In space and time problems, this becomes more critical due to the computational expense involved in a solve at each time level. Therefore, we want to determine the minimum number of steps such that we meet the criteria (2.5.32). To begin with, we considered the following two approaches:

1.  $\Upsilon_1(k_{\text{old}}, t, T) = k_{\text{old}}$ . In this instance, the previous step size is taken for the initial choice at the new step.
2.  $\Upsilon_1(k_{\text{old}}, t, T) = T - t$ . In this scenario, we shoot for the final time and take as an initial guess for the time step the remainder of the interval from the current point.

In the first instance, there is a problem in that the step sizes form a non-increasing sequence. Therefore, if at some time point early in the computation we require a small step size, any future step size is at least as small which is certainly undesirable. In the second case, we find that we are being too ambitious, almost always require a refinement of the proposed

step, and then propose a new step that is much smaller than it needs to be. Based on these observations, a combination of the above two approaches seems the most appropriate way forward. We should use local information based on the previous step size, but allow for the step size to grow if there is some indication that it might be smaller than it needs to be. This subject is tackled in the book [62], where several algorithms are presented. We consider a modification of algorithm 1.24 from [62]. The parameters  $\delta_i$ ,  $i = 1, 2$  are the reduction and growth factors respectively of the step size and  $\theta$  plays the role of threshold parameter for determining when to increase the step size. Typical values are  $\delta_1 \approx 1/\sqrt{2}$ ,  $\delta_2 \approx \sqrt{2}$ ,  $\theta = 0.5$ . The time stepping algorithm is then:

Improved adaptive time-stepping algorithm:

---

1. Start with parameters  $\delta_1 \in (0, 1)$ ,  $\delta_2 > 1$  and  $\theta \in (0, 1)$ . Set  $t = 0$ .  
 $\text{LTOL} = S(T)^{-1}\text{GTOL}$ ,  $k_{\text{old}} = T$ .
  2. Do:
    - i) Set  $k = k_{\text{old}}$ .
    - ii) Calculate  $z_h^{i+1}$  and  $\eta$ .
    - iii) While  $k\eta > \text{LTOL}$ :
      - a) Set  $k = \delta_1 k$ .
      - b) Calculate  $z_h^{i+1}$  and  $\eta$ .
    - iv) If  $\eta < \theta \text{LTOL}$ , then  $k = \delta_2 k$ .
    - v) Set  $z_i = z_{i+1}$ ,  $t = t + k$ ,  $i = i + 1$ .
 while  $t \leq T$ .
- 

To demonstrate the adaptive algorithms and some of the properties discussed, we consider the one dimensional problem,

$$z_t + az = \sin(a\pi t), \quad 0 < t \leq T,$$

$$z(0) = 0,$$

with exact solution given by,

$$z(t) = \frac{1}{a^2 + 4\pi^2} \left\{ a \sin(\pi t) - 2\pi \cos(2\pi t) - 2\pi e^{-at} \right\}. \quad (2.5.33)$$

We have recorded various measures of the performance of the cG(1) method and of the *a posteriori* error estimator applied to the test problem. To examine the performance of the method, and confirm that we have the correct convergence rate in the various quantities, we consider the empirical order of convergence (EOC). Let  $k$  denote the refinement level, and let  $u_k$ ,  $k = 1, 2, \dots$ , be a sequence of approximations. The empirical order of convergence of the sequence  $\{u_i\}_{i \geq 1}$  is,

$$\text{EOC} := \frac{\ln\left(\frac{u_{k-1}}{u_k}\right)}{\ln\left(\frac{\text{dofs}_{k-1}}{\text{dofs}_k}\right)}. \quad (2.5.34)$$

Furthermore, we consider the effectivity index of the estimator  $\eta$ ,

$$\text{Eff}(\eta) := \frac{\eta}{\|z - z_h\|_{L^\infty(I)}}. \quad (2.5.35)$$

As we can see in table 2.1, we achieve the expected convergence rate in both the  $L^\infty$  norm and in the max norm at the nodes under uniform mesh refinement. Encouragingly the *a posteriori* error estimate also converges at the same rate, however as can be seen from the effectivity index, the overestimation is quite significant. To evaluate the step size selection criteria, we consider a more challenging problem where the exact solution is a function representing an impulse at time  $t_0$ ,

$$z(t) = e^{-\frac{(t-t_0)^2}{\epsilon}}, \quad 0 < \epsilon \ll 1. \quad (2.5.36)$$

We apply the adaptive algorithm presented earlier with various values of GTOL. To assess the performance of the step size selection we look at the work efficiency of the stepping scheme. This is defined as the total number time points used over the total number of solve steps required. A value of 1 implies that the proposed step size was always accepted. Thinking about the step size selector that shoots for the final time, we can see that it will have efficiency tending to 0.5, since it will always perform at least one refinement unless it is the final step. The results are presented in 2.2.

DOFS	$\ e\ _{L^\infty(I)}$	EOC	$\max  e(t_i) $	EOC	$S(T) \max k_i \eta_i$	EOC	Eff
2	2.2949e-01		2.4275e-02		7.6862e+00		23.682
4	2.3994e-01	-6.4187	9.1962e-02	-1.9216	3.9695e+00	9.5331	11.698
8	9.6944e-02	1.3074	1.5344e-02	2.5833	1.6214e+00	1.2917	11.826
16	2.5979e-02	1.8998	3.5219e-03	2.1233	4.4366e-01	1.8697	12.075
32	6.6311e-03	1.9700	9.1675e-04	1.9417	1.1340e-01	1.9681	12.092
64	1.6911e-03	1.9713	2.3049e-04	1.9918	2.8395e-02	1.9977	11.873
128	4.2434e-04	1.9947	5.7894e-05	1.9932	7.1200e-03	1.9957	11.864
256	1.0652e-04	1.9941	1.4470e-05	2.0003	1.7864e-03	1.9948	11.858
512	2.6802e-05	1.9907	3.6188e-06	1.9995	4.4752e-04	1.9970	11.807

Table 2.1: Convergence of the discrete scheme for problem with solution (2.5.33),  $a = 1.2$ .

## 2.6 Summary

In this chapter we have reviewed the theory of AFEM and considered the contrasting applications of a stationary elliptic problem and a linear system of ODEs. In the next two chapters we depart from *a posteriori* error analysis to present the reformulation of quasistatic linear viscoelasticity using internal variables, a finite element approximation and the related *a priori* error analysis. We will return to the ideas of this chapter in chapter 5, where we provide an *a posteriori* error analysis and AFEM for quasistatic linear viscoelasticity.

GTOL	Steps	$\ e\ _{L^\infty(I)}$	$\max  e(t_i) $	$\max S(T)\eta_I$	Eff	WEff
1e-00	20	0.0469602	0.00362714	0.894524	19.0486	0.606061
1e-01	44	0.00459914	0.000414067	0.0971823	21.1305	0.709677
1e-02	105	0.000514615	4.91301e-05	0.00987558	19.1902	0.826772
1e-03	282	5.73821e-05	6.29862e-06	0.000977697	17.0384	0.915584
1e-04	898	6.16935e-06	6.62939e-07	9.9942e-05	16.1998	0.967672
1e-05	2845	6.40876e-07	6.93937e-08	9.99742e-06	15.5996	0.98819
1e-06	8377	6.8948e-08	7.94936e-09	9.9958e-07	14.4976	0.995603
1e-07	27975	6.55056e-09	8.28688e-10	9.99966e-08	15.2654	0.998572

Table 2.2: Error data for a range of tolerances for problem with solution given by (2.5.36),  $a = 1.2$ ,  $\epsilon = 0.01$ .

## Chapter 3

# Finite element approximation of quasistatic linear viscoelasticity

The purpose of this chapter is to present a space and time Galerkin finite element approximation to a reformulation of the quasistatic hereditary linear viscoelasticity problem given in chapter 1, comprising of equations (1.5.1) and (1.5.4) together with the boundary conditions (1.5.2) and (1.5.3). We first present the problem in the hereditary integral formulation, which can be viewed as an abstract Volterra problem. For given  $f \in L^p(I)$ , find  $u \in L^p(I)$  such that,

$$Au(t) = f(t) + \int_0^t B(u(s); t - s)ds. \quad (3.0.1)$$

For the viscoelasticity problem,  $A$  and  $B(\cdot, t - s)$  are second order partial differential operators. Analytical and numerical solution methods for Volterra problems are described in the book by Linz [50]. For the finite element discretisation of Volterra equations, Bedivan and Fix [12] describe a continuous Galerkin approximation to the scalar problem ( $A = 1$ ,  $B(u(s), t - s) = k(t, s)u(s)$  in (3.0.1), and focus attention on the implications of quadrature errors. With specific application to viscoelasticity problems a parallel solver is formulated by Buch et al. in [17]. These works, in contrast to the time stepping approach that will be considered here, present global spacetime, one-shot solvers.

More pertinent is the work of Shaw and Whiteman ([65] [67], [68], [69]) on the quasistatic hereditary integral formulation of linear viscoelasticity and the related abstract Volterra

problem. The numerical solution using finite elements for the spatial discretisation and the trapezoidal rule applied to the Volterra integral term for equation (1.5.8) is considered in the papers [66] and [65]. A drawback of the FEM+Trapezoidal rule approach, as remarked in [65] is that the *a priori* error bounds contain the entire history of time steps and there was no obvious path to *a posteriori* error control. In [68] a discontinuous finite element approximation of (3.0.1) in the case  $A = 1$  and  $B(\cdot; t - s) = \phi(t - s) \cdot$  is presented with an *a posteriori* estimate for negative norms of the approximation error.

It is shown in [67] that the duality method for deriving *a posteriori* error estimates as outlined in section 1.2 is limited in application to Volterra problems. The limitation is that the analysis requires strong stability of the dual solution so that optimal order interpolation error estimates can be used, and hence explicit dependence on the discretisation parameter of the error estimate. That this is a limitation for Volterra problems stems from the fact that there is no way of bounding derivatives of  $u$  in terms of  $f$  alone. In general, the best one can hope for is a stability bound with the same order of time derivative appearing on both  $u$  and  $f$ . Motivated by the work of Süli and Houston [75], a negative norm is used in [67] to introduce a power of the temporal discretisation parameter, and hence controlability of the *a posteriori* error estimate.

A full extension of the results of [68] to the quasistatic linear viscoelasticity problem are presented in [69] and [70]. The results are discussed in [67] where a number of difficulties are reported. First the temporal error component of the *a posteriori* error estimator is unstable as  $h \rightarrow 0$  or, is prohibitively expensive to compute. Secondly, unless only nested refinements of the spatial mesh are permitted, jumps in the approximate stress over edges that are not in the current mesh persist. The purpose of this and following chapters is to follow up on the remarks in [67], that a representation of the solution algorithm in terms of internal variables could offer an improvement on this scenario.

### 3.1 Existence and uniqueness

In this section we show existence and uniqueness of a weak solution to the problem given by (1.5.8). Rather than use standard methods for Volterra problems, we show that under mild

restrictions, theorem 2.1.1 can be applied to a fully weak formulation of (1.5.8). Let  $v$  in equation (1.5.8) also vary in time. Then integration over  $I$  results in the abstract problem,

$$A(u, v) = L(v), \quad (3.1.1)$$

where,

$$A(u, v) := \int_I a(u(t), v(t)) - \int_0^t \partial_s \varphi(t-s) a(u(s), v(t)) ds dt, \quad (3.1.2)$$

$$L(v) := \int_I l(t; v(t)) dt. \quad (3.1.3)$$

The problem of determining the displacement can now be posed as: Find  $u \in L^p(I; V)$  such that,

$$A(u, v) = L(v), \quad \forall v \in L^q(I; V). \quad (3.1.4)$$

From theorem 2.1.1 we recognise that we require  $L^q(I; V)$  to be a reflexive Banach space. Therefore we have the immediate restriction that  $1 < q < \infty$ . We will see that a condition for existence and uniqueness is that the order of the Lebesgue spaces must satisfy the usual relationship  $p^{-1} + q^{-1} = 1$ .

We aim to apply theorem 2.1.1 to prove existence and uniqueness, therefore we must verify conditions (2.1.2) and (2.1.3) for the bilinear form (3.1.2) and show that (3.1.3) is a continuous linear functional on the  $L^q(I; V)$ . This is a non-standard method for proving existence and uniqueness, more traditional would be to use the contraction mapping theorem on the integral operator as in standard proofs of existence and uniqueness for ODEs [36]. From section 2.4 we know that the bilinear form  $a : V \times V \rightarrow \mathbb{R}$  is an inner product on the space  $V$  with associated norm  $\|\cdot\|$ . Furthermore note that  $V^* = V$  since  $V$  is a Hilbert space. Therefore,

$$\|v\|_{L^p(I; V)} := \sup_{w \in L^q(I; V)} \frac{\int_I a(v(t), w(t)) dt}{\|w\|_{L^q(I; V)}}, \quad (3.1.5)$$

where,

$$\|w\|_{L^q(I; V)} = \left( \int_0^T \|\|w(t)\|\|^q dt \right)^{1/q}. \quad (3.1.6)$$

**Lemma 3.1.1.** *The bilinear form (3.1.2) satisfies conditions (2.1.2) and (2.1.3).*

*Proof.* Starting with (2.1.2), the Cauchy-Schwarz inequality implies that,

$$\begin{aligned} A(w, v) &= \int_I a(w, v) - \int_0^t \partial_s \varphi(t-s) a(w(s), v(t)) ds dt, \\ &\geq \int_I a(w, v) dt - \int_I \int_0^t \partial_s \varphi(t-s) \|w(s)\| \cdot \|v(t)\| ds dt, \\ &\geq \int_I a(w, v) dt - \int_I \|v(t)\| \int_0^t \partial_s \varphi(t-s) \|w(s)\| ds dt, \end{aligned}$$

Applying Hölders inequality on the time integral of the second term and using Young's inequality for convolutions (1.6.6), we have,

$$A(w, v) \geq \int_I a(w, v) dt - \|v\|_{L^q(I;V)} \cdot \|\partial_t \varphi\|_{L^1(I)} \cdot \|w\|_{L^p(I;V)}. \quad (3.1.7)$$

Dividing by  $\|v\|_{L^q(I;V)} \neq 0$  and taking the supremum over  $v \in L^q(I;V)$  results in,

$$\sup_{v \in L^q(I;V)} \frac{A(w, v)}{\|v\|_{L^q(I;V)}} \geq \|w\|_{L^p(I;V)} - \|\partial_t \varphi\|_{L^1(I)} \cdot \|w\|_{L^p(I;V)}. \quad (3.1.8)$$

Division by  $\|w\|_{L^p(I;V)}$  with  $w \neq 0$  and taking the infimum over all  $w$  results in,

$$\inf_{w \in L^p(I;V)} \sup_{v \in L^q(I;V)} \frac{A(w, v)}{\|w\|_{L^p(I;V)} \|v\|_{L^q(I;V)}} \geq 1 - \|\partial_t \varphi\|_{L^1(I)}. \quad (3.1.9)$$

Then from the second part of the fading memory hypothesis (1.4.24), that  $\varphi'(t) < 0, \forall t \in I$ , it follows that,

$$\|\partial_t \varphi\|_{L^1(I)} = - \int_0^T \partial_t \varphi(t) dt = -\varphi(T) + \varphi(0) = 1 - \varphi(T). \quad (3.1.10)$$

Therefore condition (2.1.2) is satisfied with,

$$\inf_{w \in L^p(I;V)} \sup_{v \in L^q(I;V)} \frac{A(w, v)}{\|w\|_{L^p(I;V)} \|v\|_{L^q(I;V)}} \geq \varphi(T). \quad (3.1.11)$$

To show (2.1.3), assume that  $\forall w \in L^p(I;V), A(w, v) = 0$ . Then by Cauchy-Schwarz,

Hölder's and Young's inequality again, we have,

$$\begin{aligned}
& \int_I a(w(t), v(t)) - \int_0^t \partial_s \varphi(t-s) a(w(s), v(t)) ds dt = 0, \\
& \int_I a(w(t), v(t)) - \|v(t)\| \cdot \int_0^t \partial_s \varphi(t-s) \|w(s)\| ds dt \leq 0, \\
& \int_I a(w(t), v(t)) dt - \|v\|_{L^q(I;V)} \cdot \|\partial_t \varphi(t)\| \|w(t)\| \|_{L^p(I)} \leq 0, \\
& \int_I a(w(t), v(t)) dt - \|v\|_{L^q(I;V)} \cdot \|w\|_{L^p(I;V)} \cdot \|\partial_t \varphi\| \leq 0, \\
& \sup_{w \in L^p(I;V)} \frac{\int_I a(w(t), v(t)) dt}{\|w\|_{L^p(I;V)}} - \varphi(T) \|v\|_{L^q(I;V)} \leq 0, \\
& (1 - \varphi(T)) \|v\|_{L^q(I;V)} \leq 0.
\end{aligned}$$

Since  $0 < \varphi(t) < 1$ ,  $\forall t > 0$  (section 1.4), the inequality  $(1 - \varphi(T)) \|v\|_{L^q(I;V)} \leq 0$  implies that  $v = 0$  in  $L^q(I;V)$ .  $\square$

**Lemma 3.1.2.** *Let  $f \in L^p(I; L^2(\Omega))$  and  $g \in L^p(I; L^2(\Gamma_N))$ , then the linear functional  $L : L^q(I;V) \rightarrow \mathbb{R}$  is continuous on  $L^q(I;V)$ .*

*Proof.* The proof follows immediately from lemma 2.4.3 in that,

$$\begin{aligned}
|L(v)| & \leq \int_I |l(t; v(t))| dt, \\
& \leq \int_I \left( c_a^{-1} \|f(t)\|_{L^2(\Omega)} + c_a^{-1} C_{\Gamma_N} \|g(t)\|_{L^2(\Gamma_N)} \right) \|v(t)\| dt, \\
& \leq \left( c_a^{-1} \|f\|_{L^p(I;L^2(\Omega))} + c_a^{-1} C_{\Gamma_N} \|g\|_{L^p(I;L^2(\Gamma_N))} \right) \|v\|_{L^q(I;V)}, \\
\|L\|_{L^p(I;V)} & \leq c_a^{-1} \|f\|_{L^p(I;L^2(\Omega))} + c_a^{-1} C_{\Gamma_N} \|g\|_{L^p(I;L^2(\Gamma_N))}.
\end{aligned}$$

$\square$

We can now show an existence and uniqueness result for (3.1.4) by applying theorem 2.1.1.

**Theorem 3.1.3.** *If  $f \in L^p(I; L^2(\Omega))$ ,  $g \in L^p(I; L^2(\Gamma_N))$ , then for  $1 < p < \infty$  there exists a unique solution  $u \in L^p(I;V)$  to problem (3.1.4) satisfying,*

$$\|u\|_{L^p(I;V)} \leq \varphi(T)^{-1} c_a^{-1} \|f\|_{L^p(I;L^2(\Omega))} + \varphi(T)^{-1} c_a^{-1} C_{\Gamma_N} \|g\|_{L^p(I;L^2(\Gamma_N))}. \quad (3.1.12)$$

*Proof.* From lemma 3.1.1 the bilinear form (3.1.2) satisfies conditions (2.1.2) and (2.1.2). By lemma 3.1.2 the linear functional (3.1.3) is continuous on  $L^q(I; V)$ , hence by theorem 2.1.1 there exists a unique  $u \in L^p(I; V)$  satisfying equation (3.1.4).  $\square$

For the reasons mentioned in the introduction to this chapter, we now consider a reformulation using internal variables. Internal variables for this type of problem are not new, and are presented for example in [71]. We follow the approach set out in [48] of utilising *internal variables* to rewrite the system (1.5.8), but build on it by exploiting the quasistatic assumption to formally separate the governing equations into a stationary elasticity type problem and a system of ordinary differential equations governing “internal” stresses. We believe that this representation, exploiting the quasistatic assumption is new, and while of limited applicability, covers a sufficient range of problems that it warrants further study.

## 3.2 Internal variable formulation

Internal variable formulations of linear viscoelasticity are discussed in [3], [38] and [71]. The idea is to rewrite the hereditary integral form of the stress presented above as a sum of the instantaneous strain and a sum of strains of internal variables. The internal variables can then be shown to be governed by a linear differential equation. An internal variable formulation of (1.4.26) can be achieved by defining the functions,

$$z^i(x, t) = \int_0^t \beta_i e^{-\alpha_i(t-s)} u(x, s) ds, \quad \beta = (\alpha_i \varphi_i)^{\frac{1}{2}}. \quad (3.2.1)$$

The  $\beta$  term is introduced to make the matrix in the resulting system of ordinary differential equations symmetric (see below, proposition 3.2.1). With the introduction of the internal variables, the constitutive law (1.5.4) can be written as,

$$\boldsymbol{\sigma}(x, t) = \mathbf{C}\boldsymbol{\epsilon}(u(x, t)) - \sum_{i=1}^{n_v} \beta_i \mathbf{C}\boldsymbol{\epsilon}(z^i(x, t)), \quad (3.2.2)$$

where the functions  $\{z^i\}_{i=1}^{n_v}$  satisfy the differential equations,

$$\partial_t z^i(x, t) + \alpha_i z^i(x, t) = \beta_i u(x, s), \quad i = 1, \dots, n_v. \quad (3.2.3)$$

Using internal variables, equation (1.5.8) can be written as,

$$a(u(t), v) = l(t; v) + \sum_{i=1}^{n_v} \beta_i a(z^i(t), v) \quad \forall v \in V. \quad (3.2.4)$$

Applying the operator  $\mathbf{C}\epsilon(\cdot)$  to equation (3.2.3) and taking the tensor inner product with  $\epsilon(v^i)$  for an arbitrary  $v^i \in V$ , integration over  $\Omega$  gives,

$$a(\partial_t z^i(t), v^i) + \alpha_i a(z^i(t), v^i) = \beta_i a(u(t), v^i), \quad i = 1, \dots, n_v. \quad (3.2.5)$$

Substituting equation (3.2.4) with  $v = v^i$  into the right hand side of (3.2.5) results in,

$$a(\partial_t z^i(t), v^i) + \alpha_i a(z^i(t), v^i) = \beta_i l(t; v^i) + \sum_{j=1}^{n_v} \beta_i \beta_j a(z^j(t), v^i), \quad i = 1, \dots, n_v. \quad (3.2.6)$$

Summing over  $i$  from 1 to  $n_v$  and rearranging leads to,

$$\sum_{i=1}^{n_v} a(\partial_t z^i + \sum_{j=1}^{n_v} m_{ij} z^j, v^i) dt = \sum_{i=1}^{n_v} \beta_i l(t; v^i), \quad \forall v^i \in V, \quad (3.2.7)$$

where  $m_{ij} = \delta_{ij}\alpha_i - \beta_j\beta_j$ . For notational simplification and clarity we introduce the following functional setting for the internal variables. Define an inner product with associated norm on the  $n_v$  product space  $V^{n_v} = V \times \dots \times V$  by,

$$(z, w)_{n_v} := \sum_{i=1}^{n_v} a(z^i, w^i), \quad \|\cdot\|_{n_v} := \sqrt{(\cdot, \cdot)_{n_v}}. \quad (3.2.8)$$

Also, to simplify the right hand side define the linear functional on  $V^{n_v}$ ,

$$r(t; v) = \sum_{i=1}^{n_v} \beta_i l(t; v^i). \quad (3.2.9)$$

With this new notation, the internal variable problem (3.2.7) can be written as,

$$(\partial_t z(t) + Mz(t), v)_{n_v} = r(t; v), \quad \forall v \in V^{n_v}, \quad (3.2.10)$$

$$z(0) = 0. \quad (3.2.11)$$

where  $M = (m_{ij}) \in \mathbb{R}^{n_v \times n_v}$  is the matrix with entries  $m_{ij} = \delta_{ij}\alpha_i - \beta_j\beta_j$ .

**Proposition 3.2.1** (Properties of M). *The matrix M is symmetric, furthermore if  $\varphi_0 > 0$ , then it is positive definite.*

*Proof.* Since M can be written as  $M = D - \beta\beta^T$  where  $D = (d_{ij})$  is a diagonal matrix with  $d_{ii} = \alpha_i$  and  $\beta^T = (\beta_1, \dots, \beta_{n_v})$ , symmetry follows the symmetry of diagonal matrices

and outer products and linearity of transposition,  $(D - \beta\beta^T)^T = D - \beta\beta^T$ . To see positive definiteness we show that  $x^T Mx > 0$ ,  $\forall x \in \mathbb{R}^{n_v} \setminus \{0\}$ . Then,

$$\begin{aligned}
x^T Mx &= x^T D x - x^T \beta \beta^T x, \\
&= x^T D x - (x^T \beta)^2, \\
&= \sum_{i=1}^{n_v} \alpha_i x_i^2 - \left( \sum_{i=1}^{n_v} \beta_i x_i \right)^2, \\
&\geq \sum_{i=1}^{n_v} \alpha_i x_i^2 - \left( \sum_{i=1}^{n_v} \alpha_i x_i^2 \right) \left( \sum_{i=1}^{n_v} \varphi_i \right) \\
&\geq \sum_{i=1}^{n_v} \alpha_i x_i^2 \left( 1 - \sum_{i=1}^{n_v} \varphi_i \right) \\
&= \varphi_0 \sum_{i=1}^{n_v} \alpha_i x_i^2 > 0.
\end{aligned}$$

where we have used the properties of  $\varphi$  of normalisation,  $\varphi(0) = 1$  and that  $\varphi_0 > 0$ .  $\square$

As a result of the positive definiteness of  $M$ , we have for all  $x \in \mathbb{R}^{n_v}$ ,

$$m_s \|x\|^2 \leq \|M^{1/2} x\|^2 \leq m_l \|x\|, \quad (3.2.12)$$

where  $m_s$  and  $m_l$  are the smallest and largest eigenvalues of  $M$ .

To apply the finite element method to (3.2.10) a weak formulation is achieved by integrating over the time interval. Define the space,

$$W_0^{k,p}(I; V^{n_v}) := \{v | v \in W^{k,p}(I; V^{n_v}), v(0) = 0\}. \quad (3.2.13)$$

Then define the bilinear form  $B : W_0^{1,p}(I; V^{n_v}) \times L^q(I; V^{n_v}) \rightarrow \mathbb{R}$  by,

$$B(z, v) := \int_I (\partial_t z(t) + Mz(t), v)_{n_v} dt, \quad (3.2.14)$$

and the linear functional,  $F : L^q(I; V^{n_v}) \rightarrow \mathbb{R}$ , by,

$$F(v) := \int_I r(t; v) dt, \quad (3.2.15)$$

then the weak problem is as follows. Find  $z \in W_0^{1,p}(I; V^{n_v})$  such that,

$$B(z, v) = F(v), \quad \forall v \in L^q(I; V^{n_v}). \quad (3.2.16)$$

Assuming that the internal variables are known, the problem of determining the displacement is as follows. Find  $u \in L^\infty(I; V)$  such that,

$$a(u(t), v) = l(t; v) + \sum_{i=1}^{n_v} \beta_i a(z^i(t), v), \quad \forall v \in V, \quad a.e. t \in I. \quad (3.2.17)$$

The next lemma is used to show that  $z$  is one order smoother than  $u$  in time, which is a consequence of the definition of  $z$ . The result is required in the following chapter so that *a priori* error estimates for  $u$  can be expressed completely in terms of  $u$ . Define the Volterra operator  $K$ , by,

$$K(u) = \int_0^t \kappa(t-s)u(s)ds. \quad (3.2.18)$$

**Lemma 3.2.2.** *Let  $r \geq 0$  and suppose that  $u \in W^{r,p}(I; Y)$  for some Banach space  $Y$ ,  $\kappa \in W^{r+1,1}(I)$ , then  $K : W^{r,p}(I; Y) \rightarrow W^{r+1,p}(I; Y)$  is continuous and so there exists a  $C > 0$  such that,*

$$\|K(u)\|_{W^{r+1,p}(I; Y)} \leq C \|u\|_{W^{r,p}(I; Y)}. \quad (3.2.19)$$

*Proof.* From Young's inequality for convolutions (1.6.6),  $K$  is continuous on  $L^p(I; Y)$ ,

$$\|K(u)\|_{L^p(I; Y)} \leq \|\kappa\|_{L^1(I)} \|u\|_{L^p(I; Y)}. \quad (3.2.20)$$

Differentiating (3.2.18)  $k+1$  times, results in,

$$\partial_t^{(k+1)} K(u) = \sum_{i=0}^k \partial_t^{(i)} \kappa(0) \partial_t^{(k-i)} u(t) + \int_0^t \partial_t^{(k+1)} \kappa(t-s) u(s) ds. \quad (3.2.21)$$

Taking the  $Y$  norm, raising to the power  $p$  and using the inequality  $(a+p)^p \leq C(a^p + b^p)$ , results in,

$$\begin{aligned} \|\partial_t^{(k+1)} K(u)\|_Y^p &\leq C \left( \sum_{i=0}^k |\partial_t^i \kappa(0)|^p \|\partial_t^{(k-i)} u(t)\|^p \right. \\ &\quad \left. + \left( \int_0^t \partial_t^{(k+1)} \kappa(t-s) \|u(s)\| ds \right)^p \right). \end{aligned}$$

Absorbing all constants into the first, Hölders inequality implies that,

$$\|\partial_t^{(k+1)} K(u)\|_Y^p \leq C \left( \sum_{i=0}^k \|\partial_t^{(k-i)} u(t)\|^p + \|u\|_{L^p(I; Y)}^p \right). \quad (3.2.22)$$

Integrating over  $I$ , gives,

$$\begin{aligned} \|\partial_t^{(r+1)} K(u)\|_{L^p(I;Y)}^p &\leq C \left( \sum_{i=0}^k \|\partial_t^{(k-i)} u\|_{L^p(I;Y)}^p + T \|u\|_{L^p(I;Y)}^p \right), \\ \Rightarrow |K(u)|_{W^{k+1,p}(I;Y)}^p &\leq C \|u\|_{W^{k,p}(I;Y)}^p, \quad k \geq 0. \end{aligned}$$

The result follows by summing over  $k$  from 0 to  $r$ , together with (3.2.20) and taking  $p$ -th roots.  $\square$

We can apply lemma 3.2.2 to bound  $\|z\|_{W^{r+1,p}(I;Y^{n_v})}$ .

**Corollary 3.2.3.** *Under the assumptions of lemma 3.2.2, there holds,*

$$\|z\|_{W^{r+1,p}(I;Y^{n_v})} \leq C \|u\|_{W^{r,p}(I;Y)}. \quad (3.2.23)$$

*Proof.* Setting  $z^i = K_i(u)$ , with,

$$K_i(u) = \int_0^t \beta_i e^{-\alpha_i(t-s)} u(s) ds. \quad (3.2.24)$$

Then it follows that,

$$\begin{aligned} \|z\|_{W^{r+1,p}(I;Y^{n_v})} &= \left( \sum_{i=1}^{n_v} \|z^i\|_{W^{r+1,p}(I;Y)}^p \right)^{1/p}, \\ &= \left( \sum_{i=1}^{n_v} \|K_i(u)\|_{W^{r+1,p}(I;Y)}^p \right)^{1/p}, \\ &\leq C \left( \sum_{i=1}^{n_v} \|u\|_{W^{r,p}(I;Y)}^p \right)^{1/p}, \\ &\leq C \|u\|_{W^{r,p}(I;Y)}. \end{aligned}$$

Furthermore, it follows from the definition of the internal variables that if  $u \in L^\infty(I;Y)$  for some Banach space  $Y$ , then,

$$\|z(t)\|_{Y^{n_v}} \leq C \|u\|_{L^\infty(I;Y)} \quad \forall t \in I. \quad (3.2.25)$$

$\square$

### 3.3 Finite element approximation

For the approximation in the spatial variables, a conforming Galerkin simplicial Lagrange finite element method as described in section 2.2 is used. For the approximation of the internal variables both a space and time approximation is required. For the temporal approximation we will use the continuous Galerkin method and the discretisation presented in section 2.5.

#### 3.3.1 Displacement

For the approximation of the displacement we will construct an approximate solution in the space  $V$ , defined in (1.5.5), at each time  $t_i$ , from functions belonging to  $S_1(\mathcal{T}_i; \mathbb{R}^d)$  (2.2.11). Denote the finite element space at time  $t_i$  by  $V_h^i$  and take it to be,

$$V_h^i = S_1(\mathcal{T}_i; \mathbb{R}^d) \cap V. \quad (3.3.1)$$

Denote the set of basis functions of  $V_h^i$  by  $\{\eta_j^i\}_{j=1}^{\dim V_h^i}$  and write functions from  $V_h^i$  as,

$$v_h^i = \sum_{j=1}^{\dim V_h^i} v_j^i \eta_j^i, \quad (3.3.2)$$

where  $v_j^i$  is the component of the vector  $\mathbf{v}^i \in \mathbb{R}^{\dim V_h^i}$  corresponding to the  $j$ -th basis function of  $V_h^i$ .

#### 3.3.2 Internal Variables

The approximation of the internal variables is in both space and time. For the approximation by the cG(1) method, the trial space at a given time point  $t_i$  is  $V_h^i$  as given in 3.3.1. Then on each time level, we consider the space of linear polynomials of the form,

$$v_h|_{I_i} = v_h^i \phi_i(t) + v_h^{i-1} \phi_{i-1}(t), \quad (3.3.3)$$

where  $\phi_i(t)$  are the functions given by (2.5.8), (2.5.9) and (2.5.10), and denote this space by  $\mathbb{P}_1(I_i; V_h^{i-1} + V_h^i)$ . To cope with the fact that there are  $n_v$  internal variables, we form the  $n_v$  product space  $\mathbb{P}_1(I_i; V_h^{i-1} + V_h^i)^{n_v} = \mathbb{P}_1(I_i; V_h^{i-1} + V_h^i) \times \dots \times \mathbb{P}_1(I_i; V_h^{i-1} + V_h^i)$ , and denote the piecewise composition of this space by  $Z_h$ ,

$$Z_h := \{z \in C^0(I; V^{n_v}), z(0) = 0 : z|_{I_i} \in \mathbb{P}_1(I_i; V_h^{i-1} + V_h^i)^{n_v}\}. \quad (3.3.4)$$

The test space is taken to be the space of piecewise constant mappings of the time intervals into  $V_h^i$  defined by,

$$Y_h := \{v \in L^\infty(I; V^{n_v}) : v|_{I_i} \in \mathbb{P}_0(I_i; V_h^i)^{n_v}\}, \quad (3.3.5)$$

where  $\mathbb{P}_0(I_i; V_h^i)^{n_v}$  is the  $n_v$  product space of piecewise constant polynomials such that  $v_h|_{I_i} \in V_h^{i, n_v}$ . The finite element problem corresponding to the internal variable equation is: Find  $z_h \in Z_h$  such that,

$$B(z_h, v) = F(v), \quad \forall v \in Y_h. \quad (3.3.6)$$

Assuming that the internal variables are known, the finite element approximation for the displacement problem then takes the familiar form of that for an elliptic problem. The finite element problem for the displacement is: Find  $u_h(t_i) \in V_h^i$  such that,

$$a(u_h(t_i), v) = l(t_i; v) + \sum_{j=1}^{n_v} \beta_j a(z_h^{j,i}, v), \quad \forall v \in V_h^i, \quad i = 0, 1, \dots, N. \quad (3.3.7)$$

Given that we only solve for the displacement at the time nodes, we consider  $u_t(t)$  to be a piecewise linear function of the form,

$$u_h|_{I_i} = u_h^i \phi_i(t) + u_h^{i-1} \phi_{i-1}(t). \quad (3.3.8)$$

### 3.3.3 The discrete scheme

Since the internal variable problem is a linear system, when the substitutions to derive the discrete scheme are made, there will be two types of vectors and matrices, those associated with the original system, and those associated with the discretisation. Therefore we have a system of finite element approximations and we briefly explain some notation relating to this situation. The finite element approximation to the  $k$ -th component of  $z$  on time interval  $I_i$  is given by,

$$z_h^k(t) = \phi_{i-1}(t) \sum_{j=1}^{\dim V_h^{i-1}} z_j^{k, i-1} \eta_j^{i-1} + \phi_i(t) \sum_{j=1}^{\dim V_h^i} z_j^{k, i} \eta_j^i. \quad (3.3.9)$$

Denote the  $n_v$  dimensional vector of internal variable approximations by

$$z_h(t) = (z_h^1(t), \dots, z_h^{n_v}(t)). \quad (3.3.10)$$

Since there are inner products of functions from different meshes, define the matrices,

$$A = (A_{mn}) = a(\eta_m^i, \eta_n^i), \quad \tilde{A} = (\tilde{A}_{mn}) = a(\eta_m^i, \eta_n^{i-1}). \quad (3.3.11)$$

Choosing  $v(t) = \chi_{I_i}(t)v^i$  in (3.3.6), where  $\chi_{I_i}$  is the indicator function for the interval  $I_i$ , and  $v^i \in V_h^{i, n_v}$ , (3.3.6) becomes,

$$(z_h^i - z_h^{i-1}, v^i)_{n_v} + \frac{k_i}{2}(M(z_h^i + z_h^{i-1}), v^i)_{n_v} = \int_{I_i} r(t; v^i) dt, \quad \forall v^i \in V_h^{i, n_v}. \quad (3.3.12)$$

Let  $\mathbf{z}^i = (\mathbf{z}^{1,i}, \dots, \mathbf{z}^{n_v,i})^T$  be the vector representations of the internal variables with respect to the basis of  $V_h^i$ , and let  $A_{n_v}$  and  $\tilde{A}_{n_v}$  denote the block matrices,

$$A_{n_v} = \begin{pmatrix} A & 0 & \dots & 0 \\ 0 & A & \dots & \dots \\ \dots & \dots & \dots & 0 \\ 0 & \dots & 0 & A \end{pmatrix}, \quad \tilde{A}_{n_v} = \begin{pmatrix} \tilde{A} & 0 & \dots & 0 \\ 0 & \tilde{A} & \dots & \dots \\ \dots & \dots & \dots & 0 \\ 0 & \dots & 0 & \tilde{A} \end{pmatrix}, \quad (3.3.13)$$

where  $A$  and  $\tilde{A}$  are matrices from (3.3.11). Furthermore, let  $M_{n_v}$  denote the block matrix created by expanding each entry of  $M$  into a square block of size  $\dim V_h^i$ . Sampling (3.3.12) at each basis function of  $V_h^i$  implies that,

$$A_{n_v} \mathbf{z}^i - \tilde{A}_{n_v} \mathbf{z}^{i-1} + \frac{k_i}{2}(A_{n_v} M_{n_v} \mathbf{z}^i + \tilde{A}_{n_v} M_{n_v} \mathbf{z}^{i-1}) = \mathbf{r}^i, \quad (3.3.14)$$

where,

$$\mathbf{r}^i = \int_{I_i} \mathbf{r}(t) dt, \quad \text{where } \mathbf{r}(t) = (r(t; \eta_1^i), \dots, r(t; \eta_{\dim V_h^i}^i))^T. \quad (3.3.15)$$

Rearranging (3.3.14) we get,

$$A_{n_v} \left( I_{n_v} + \frac{k_i}{2} M_{n_v} \right) \mathbf{z}^i - \tilde{A}_{n_v} \left( I_{n_v} - \frac{k_i}{2} M_{n_v} \right) \mathbf{z}^{i-1} = \mathbf{r}^i, \quad (3.3.16)$$

where  $I_{n_v}$  is the identity matrix of dimension  $n_v$ . Let  $\mathbf{u}^i \in \mathbb{R}^{\dim V_h^i}$  denote the vector representation of  $u_h(t_i)$ . Then (3.3.7) becomes,

$$A \mathbf{u}^i = \mathbf{l}^i + \sum_{j=1}^{n_v} A \mathbf{z}^{j,i}, \quad (3.3.17)$$

where,

$$\mathbf{l}^i = (l(t_i; \eta_1^i), \dots, l(t_i; \eta_{\dim V_h^i}^i))^T. \quad (3.3.18)$$

The solution algorithm is then as follows:

Solution Algorithm:

---

Determine elasticity solution,  $\mathbf{u}^0 = A^{-1}\mathbf{b}^0$ . For  $i = 0, 1, \dots$ :

1. Determine internal variables. Find  $\mathbf{z}^i$  such that,

$$\left(I_{n_v} + \frac{k_i M_{n_v}}{2}\right) A_{n_v} \mathbf{z}^i - \left(I_{n_v} - \frac{k_i M_{n_v}}{2}\right) \tilde{A}_{n_v} \mathbf{z}^{i-1} = \mathbf{r}^i, \\ \mathbf{z}^0 = \mathbf{0}.$$

2. Determine displacement: Find  $\mathbf{u}^i$  such that,

$$A\mathbf{u}^i = \mathbf{I}^i + \sum_{j=1}^{n_v} \beta_j A\mathbf{z}^{j,i}. \quad (3.3.19)$$


---

While the algorithm has been presented so that the internal variables and displacement are solved sequentially, this is not a requirement.

We remark that the notation introduced in this section for the internal variables allows us to write the stepping scheme in a familiar form, however, the matrices  $A_{n_v}$  and  $\tilde{A}_{n_v}$  would not be formed in practice.

**Proposition 3.3.1.** *Assume that  $V_h^{i-1} \subset V_h^i$ ,  $\forall i = 1, 2, \dots$ , then there exists a constant  $C$  independent of the discretisation such that the solution to the discrete problem (3.3.6) satisfies the stability estimate,*

$$\|z_h(t_n)\|^2 + \frac{1}{m_s} \int_0^{t_n} \|\partial_t z_h(t)\|^2 dt \leq \frac{C}{m_s} \int_0^{t_n} \|l(t)\|_{\tilde{V}^*}^2 dt. \quad (3.3.20)$$

*Proof.* We have,

$$\int_{I_i} (\partial_t z_h(t) + M z_h(t), v)_{n_v} dt = \int_{I_i} r(t; v) dt, \quad \forall v \in \mathbb{P}_0(I_i; V_h^i)^{n_v}.$$

Since  $V_h^{i-1} \subset V_h^i$ , we can take  $v = \partial_t z_h$  to get,

$$\int_{I_i} \|\partial_t z_h(t)\|^2 + (M z_h(t), \partial_t z_h(t))_{n_v} dt = \int_{I_i} r(t; \partial_t z_h(t)) dt. \quad (3.3.21)$$

Recalling the definition of  $r(t; \cdot)$  (3.2.9), we can bound the right hand side of (3.3.21) using lemma 3.1.2 and Youngs inequality,

$$\left| \int_{I_i} r(t; \partial_t z_h(t)) dt \right| \leq \frac{C}{2} \int_0^{t_n} \|l(t)\|_{\tilde{V}^*}^2 dt + \frac{1}{2} \int_{I_i} \|\partial_t z\|_{n_v}^2 dt. \quad (3.3.22)$$

Substituting into the right hand side of (3.3.21) an rearranging, we get,

$$\int_{I_i} \frac{1}{2} \|\partial_t z_h(t)\|^2 + (M z_h(t), \partial_t z_h(t))_{n_v} dt \leq \frac{C}{2} \int_0^{t_n} \|l(t)\|_{V^*}^2 dt. \quad (3.3.23)$$

Now, because of the symmetry of  $M$  (lemma 3.2.1) it follows that,

$$(M z_h(t), \partial_t z_h(t))_{n_v} = \frac{1}{2} \partial_t \|M^{1/2} z_h(t)\|^2,$$

so (3.3.23), after performing the integration and clearing the factor of a half becomes,

$$\int_{I_i} \|\partial_t z_h(t)\|^2 dt + \|M^{1/2} z_h(t_i)\|^2 - \|M^{1/2} z_h(t_{i-1})\|^2 \leq C \int_0^{t_n} \|l(t)\|_{V^*}^2 dt. \quad (3.3.24)$$

Summation over  $i$  from 0 to  $n$  and using the zero initial condition for  $z$ , results in,

$$\int_0^{t_n} \frac{1}{2} \|\partial_t z_h(t)\|^2 dt + \|M^{1/2} z_h(t_n)\|^2 \leq \frac{C}{2} \int_0^{t_n} \|l(t)\|_{V^*}^2 dt. \quad (3.3.25)$$

Using (3.2.12) implies the result.  $\square$

### 3.4 Summary

In this chapter we have considered the question of existence and uniqueness for the quasi-static hereditary integral formulation of linear viscoelasticity. Under standard assumptions for elliptic problems and a mild but physically reasonable condition on the relaxation function, this problem fits into the abstract framework of section 2.1. Nevertheless, we present an alternative formulation to deal with certain theoretical and computational difficulties outlined for example by Shaw and Whiteman in [67]. The reformulation results in a modified elliptic problem together with a system of ordinary differential equations in the energy space of the problem. In the following two chapters we derive *a priori* and *a posteriori* error estimates for the finite element approximation of the reformulated problem. Given the reformulation, most of the techniques surveyed in chapter 2 are available, in the sense that the elliptic component is similar to the linear elasticity problem of section 2.4, and the system of internal variables similar to the ODE system of section 2.5.

## Chapter 4

# *A priori* error analysis

In this chapter we present an *a priori* analysis for the approximation error in the displacement and internal variables. Given the formulation presented in the previous chapter the analysis of the error for the internal variables can proceed independently of that for the displacement. However we begin with displacement and show an *a priori* upper bound for the energy norm of the error  $e_u(t) = u(t) - u_h(t)$  at time  $t_N$  which depends on the errors in the internal variables  $\{e_{z_j}(t)\}_{j=1}^{n_v}$  and a term typical of *a priori* error estimates for Galerkin approximations to elliptic problems. Our attention then turns to *a priori* upper bounds on the  $n_v$ -energy norm of the error in the internal variable approximation,  $e_z(t) := z(t) - z_h(t)$ . We apply the method described in section 1.2 of using a discrete dual problem and show that the proposed finite element method is of optimal order.

Since we are primarily concerned with determining the order of convergence of the method, in contrast to the treatment of *a posteriori* error estimates, for the remainder of this chapter we will absorb all constants that do not depend on the discretisation parameters  $h$  and  $k$  into a generic constant  $C$ .

### 4.1 Displacement

In this section an *a priori* estimate for the error in the displacement at time  $t_i$  is derived. Since we are in effect dealing with a Galerkin approximation to an elliptic problem, we pursue the direction set in section (2.1) of using Galerkin orthogonality to show a proof of

a variation on Céa's lemma. The result and proof are based on Céa's lemma but there is an additional term due to the internal variables.

**Lemma 4.1.1** (Galerkin Orthogonality). *The approximation error in the displacement  $e_u(t)$  and internal variables  $\{e_{z^j}(t)\}_{j=1}^{n_v}$  satisfies the relationship,*

$$a(e_u(t_i), v) - \sum_{j=1}^{n_v} \beta_j a(e_{z^j}(t_i), v) = 0, \quad \forall v \in V^i, \forall i = 0, 1, 2, \dots \quad (4.1.1)$$

*Proof.* Choose  $v \in V_h^i$  in equation (3.2.17) and subtract from (3.3.7).  $\square$

With Galerkin orthogonality, a modified version of Céa's lemma (lemma 2.1.4) can be proven.

**Lemma 4.1.2.** *Let  $u$  be the solution of (3.2.17),  $u_h$  the solution to (3.3.7), and let  $z$  be the solution of (3.2.16),  $z_h$  the solution to (3.3.6), then there exists a  $C > 0$  such that the errors  $e_u(t) = u(t) - u_h(t)$ ,  $e_z(t) = z(t) - z_h(t)$  satisfy,*

$$\|e_u(t_i)\| \leq C \left\{ \inf_{w_h^i \in V_h^i} \|u(t_i) - w_h^i\| + |\varphi'(0)|^{1/2} \cdot \|e_z(t_i)\|_{n_v} \right\}. \quad (4.1.2)$$

*Proof.* It is only necessary to bound the portion of error contained within the discrete space  $V_h^i$ , since,

$$\|u(t_i) - u_h^i\| \leq \|u(t_i) - w_h^i\| + \|w_h^i - u_h^i\|, \quad \forall w_h^i \in V_h^i. \quad (4.1.3)$$

To bound  $w_h^i - u_h^i$ , Galerkin Orthogonality (lemma 4.1.1) followed by Cauchy-Schwarz implies that,

$$\begin{aligned} \|w_h^i - u_h^i\|^2 &= a(w_h^i - u_h^i, w_h^i - u_h^i), \\ &= a(w_h^i - u(t_i) + u(t_i) - u_h^i, w_h^i - u_h^i), \\ &= a(w_h^i - u(t_i), w_h^i - u_h^i) + a(e(t_i), w_h^i - u_h^i), \\ &= a(w_h^i - u(t_i), w_h^i - u_h^i) + \sum_{j=1}^{n_v} \beta_j a(e_{z^j}(t_i), w_h^i - u_h^i), \\ &\leq \|w_h^i - u(t_i)\| \cdot \|w_h^i - u_h^i\| + \|w_h^i - u_h^i\| \sum_{j=1}^{n_v} \beta_j \|e_{z^j}(t_i)\|, \\ &\leq \|w_h^i - u(t_i)\| \cdot \|w_h^i - u_h^i\| + \|w_h^i - u_h^i\| \cdot |\varphi'(0)|^{1/2} \|e_z(t_i)\|_{n_v}. \end{aligned}$$

Young's inequality then gives,

$$\begin{aligned} \|w_h^i - u_h^i\|^2 &\leq \frac{1}{2\epsilon} \|w_h^i - u(t_i)\|^2 + \epsilon \|w_h^i - u_h^i\| + \frac{1}{2\epsilon} |\varphi'(0)| \|e_z(t_i)\|_{n_v}^2, \\ \|w_h^i - u_h^i\|^2 &\leq \frac{1}{2\epsilon(1-\epsilon)} \left\{ \|w_h^i - u(t_i)\|^2 + |\varphi'(0)| \cdot \|e_z(t_i)\|_{n_v}^2 \right\}. \end{aligned}$$

From which it follows that,

$$\|w_h^i - u_h^i\| \leq \frac{1}{\sqrt{2}} \|w_h^i - u(t_i)\| + \frac{1}{\sqrt{2}} |\varphi'(0)|^{1/2} \cdot \|e_z(t_i)\|_{n_v}, \quad (4.1.4)$$

where we have taken  $\epsilon = 1/2$ . Combining (4.1.3) with (4.1.3) implies the result. Applying the Cauchy-Schwarz inequality to the sum and using the fact,

$$\sum_{j=1}^{n_v} \beta_j^2 = \sum_{j=1}^{n_v} \alpha_j \varphi_j = -\varphi'(0), \quad (4.1.5)$$

then the result follows.  $\square$

Given the result of lemma 4.1.2 the problem is to now derive an *a priori* upper bound on the error in the internal variables.

## 4.2 Internal variables

To determine an *a priori* error estimate for the internal variables, we carry out the scheme from [30], described in section 2.3. Estep and French [34] provide an analysis of the continuous Galerkin method applied to a general system of ODEs as described in section 2.5. We anticipate achieving a similar result, at least in the convergence with respect to the time discretisation, however we must cope with the added complication due to the spatial approximation. We will see that this added complication is not too severe and that under reasonable assumptions for an *a priori* analysis we are able to derive suitable estimates.

Once again we need to derive an upper bound on the discrete portion of the error. Let  $\Pi : Z \mapsto Z_h$  be an arbitrary mapping into the discrete space  $Z_h$ . Given that we can write,

$$e_z = z - z_h = (I - \Pi)z + \Pi z - z_h, \quad (4.2.1)$$

then so long as we can derive an error estimate for the mapping  $\Pi$ , we need only bound  $\Pi z - z_h$  to achieve an *a priori* error estimate for the error  $e_z$ . We define  $e_z^h := \Pi z - z_h$  and

refer to the term  $e_z^h$  as the discrete portion of the error. Once again we will make use of Galerkin orthogonality. Recall the bilinear form of the internal variable problem,

$$B(z, v) := \int_I (\partial_t z(t) + Mz(t), v)_{n_v} dt. \quad (4.2.2)$$

**Lemma 4.2.1** (Galerkin Orthogonality). *The error  $e_z = z - z_h$  satisfies the orthogonality relationship,*

$$B(e_z, w_h) = 0, \quad \forall w_h \in Y_h. \quad (4.2.3)$$

*Proof.* Choose  $v = w_h \in Y_h$  in the full problem (3.2.16) and subtract the finite element equations (3.3.6).  $\square$

To achieve a bound on  $e_z^h$ , a discrete dual problem with solution  $\chi_h$  is introduced. The discrete dual problem is designed to satisfy,

$$\| \| e_z^h(t_N) \| \|_{n_v}^2 = B(e_z^h, \chi_h). \quad (4.2.4)$$

Suppose that the equation (4.2.4) holds, then with Galerkin Orthogonality(lemma 4.2.1), we can write,

$$\begin{aligned} \| \| e_z^h(t_N) \| \|_{n_v}^2 &= B(e_z^h, \chi_h), \\ &= B(e_z^h, \chi_h) - B(e, \chi_h), \\ &= B(\Pi z - z, \chi_h). \end{aligned}$$

which characterises the discrete portion of the error in terms of the discrete dual solution and an arbitrary interpolant. In the next section we present the discrete dual problem and the properties of it that we require.

#### 4.2.1 The dual problem

Let  $Y_h$  be the space given in (3.3.5), then consider the discrete dual problem: Find  $\chi_h \in Y_h$  such that,

$$-\sum_{i=1}^N (v(t_i), \llbracket \chi_h \rrbracket_i)_{n_v} + \sum_{i=1}^N \int_{t_{i-1}}^{t_i} (v, M\chi_h)_{n_v} dt = 0, \quad \forall v \in Z_h, \quad (4.2.5)$$

$$(v(t_N), \chi_h(t_N^+))_{n_v} = (v(t_N), e_z^h(t_N))_{n_v} \quad \forall v \in Z_h. \quad (4.2.6)$$

**Proposition 4.2.2.** *There exists a unique solution to problem (4.2.5) with initial condition (4.2.6).*

*Proof.* To see uniqueness consider two solutions  $\chi_{h,1}, \chi_{h,2}$  of (4.2.5) supplemented with the condition (4.2.6). Then their difference  $\theta_h = \chi_{h,1} - \chi_{h,2}$  satisfies the problem

$$-\sum_{i=1}^N (v(t_i), \llbracket \theta_h \rrbracket_i)_{n_v} + \sum_{i=1}^N \int_{t_{i-1}}^{t_i} (v^h, M\theta_h)_{n_v} dt = 0, \quad \forall v \in Z_h \quad (4.2.7)$$

$$(v(t_N^+), \theta_h(t_N^+))_{n_v} = 0, \quad \forall v \in Z_h. \quad (4.2.8)$$

Let  $\theta^i$  denote the vector representation of  $\theta$  at time  $i$  with respect to the basis of  $V_h^{i,n_v}$ . The stepping scheme resulting from (4.2.7) is given (using notation developed in subsection 3.3.3 by,

$$A_{n_v} \left( I_{n_v} - \frac{k_i}{2} M_{n_v} \right) \theta^i = \tilde{A}_{n_v} \left( I_{n_v} + \frac{k_i}{2} M_{n_v} \right) \theta^{i+1}. \quad (4.2.9)$$

Then it follows from (4.2.8), that if  $\theta^N = \mathbf{0}$ , then  $\theta^i = \mathbf{0}$ ,  $\forall i = 0, \dots, N$ . Therefore  $\chi_{h,1} = \chi_{h,2}$  and the solution to problem (4.2.5) is unique. Since problem (4.2.5) is finite dimensional, existence follows from uniqueness.  $\square$

**Lemma 4.2.3.** *There exist constants  $C > 0$  and  $C' > 0$  such that the solution  $\chi_h$  of the discrete dual problem (4.2.5) and (4.2.6) satisfies,*

$$\|\|\chi_h^i\|\|_{n_v} \leq C \|\|e_z^h(t_N)\|\|_{n_v}, \quad \forall i \in \{1, \dots, N\}. \quad (4.2.10)$$

*Furthermore,*

$$\|\|\chi_h\|\|_{L^q(I; V^{n_v})} \leq C' \|\|e_z^h(t_N)\|\|_{n_v}. \quad (4.2.11)$$

*Proof.* Choose  $v = \phi_i \chi_h^i$  in (4.2.5), where  $\phi_i$  is the basis function for the  $i$ -th node in the temporal discretisation. Then since  $\llbracket \chi_h \rrbracket_i = \chi_h^{i+1} - \chi_h^i$ , equation (4.2.5) becomes,

$$\begin{aligned} & -(\chi_h^i, \llbracket \chi_h \rrbracket_i)_{n_v} + \int_{t_{i-1}}^{t_i} (\chi_h^i, M\chi_h)_{n_v} \phi_i dt + \int_{t_i}^{t_{i+1}} (\chi_h^i, M\chi_h)_{n_v} \phi_i dt = 0, \\ \Rightarrow & \|\|\chi_h^i\|\|^2 - (\chi_h^i, \chi_h^{i+1})_{n_v} + \frac{k_i}{2} (\chi_h^i, M\chi_h^i)_{n_v} + \frac{k_{i+1}}{2} (\chi_h^i, M\chi_h^{i+1})_{n_v} = 0. \end{aligned}$$

Using the fact that  $M$  is positive definite (proposition 3.2.1), we can remove the non-negative term  $(\chi_h^i, M\chi_h^i)_{n_v}$  to get,

$$\|\|\chi_h^i\|\|^2 - (\chi_h^i, \chi_h^{i+1})_{n_v} + \frac{k_{i+1}}{2} (\chi_h^i, M\chi_h^{i+1})_{n_v} \leq 0. \quad (4.2.12)$$

Dividing by  $\|\chi_h^i\|_{n_v}$  and using the Cauchy-Schwarz inequality results in,

$$\begin{aligned} \|\chi_h^i\|_{n_v} &\leq \left(1 + \frac{k_{i+1}m_l}{2}\right) \|\chi_h^{i+1}\|_{n_v}, \\ \Rightarrow \|\chi_h^i\|_{n_v} &\leq \prod_{j=i}^{N-1} \left(1 + \frac{k_j m_l}{2}\right) \|e_z^h(t_N)\|_{n_v}. \end{aligned}$$

From lemma 1.6.5 we have,

$$\begin{aligned} \|\chi_h^i\|_{n_v} &\leq \prod_{j=i}^{N-1} \left(1 + \frac{k_j m_l}{2}\right) \|e_z^h(t_N)\|_{n_v}, \\ &\leq \exp\left(\frac{m_l}{2} \sum_{j=i}^{N-1} k_j\right) \|e_z^h(t_N)\|_{n_v}, \\ &\leq \exp\left(\frac{m_l}{2} t_N\right) \|e_z^h(t_N)\|_{n_v}. \end{aligned}$$

To show (4.2.11), we can use the above result in the following way,

$$\begin{aligned} \|\chi_h\|_{L^q(I; V^{n_v})} &= \left(\sum_{i=1}^N \|\chi_h\|_{L^q(I_i; V^{n_v})}^q\right)^{1/q} = \left(\sum_{i=1}^N k_i \|\chi_h^i\|_{n_v}^q\right)^{1/q}, \\ &= \left(\sum_{i=1}^N k_i\right)^{1/q} \exp\left(\frac{m_l}{2} t_N\right) \|e_z^h(t_N)\|_{n_v}, \\ &= T^{1/q} \exp\left(\frac{m_l}{2} t_N\right) \|e_z^h(t_N)\|_{n_v}. \end{aligned}$$

□

In the next lemma we show that the discrete portion of the error can be represented in terms of the solution to the discrete dual problem and the interpolation error.

**Lemma 4.2.4** (Error Representation). *Let  $\chi_h$  be the solution to the discrete dual problem (4.2.5) and (4.2.6), then the following holds*

$$\|e_z^h(t_N)\|_{n_v}^2 = B(\Pi z - z, \chi_h). \quad (4.2.13)$$

*Proof.* To begin, we show that,

$$\|e_z^h(t_N)\|_{n_v}^2 = B(e_z^h, \chi_h), \quad (4.2.14)$$

since if this holds, the conclusion follows from Galerkin orthogonality,

$$\begin{aligned} \|e_z^h(t_N)\|_{n_v}^2 &= B(e_z^h, \chi_h), \\ &= B(e_z^h, \chi_h) - B(e, \chi_h), \\ &= B(\Pi z - z, \chi_h). \end{aligned}$$

Integration by parts and rearranging the endpoint evaluations to form jumps implies that,

$$\begin{aligned} B(e_z^h, \chi_h) &= \sum_{i=1}^N \int_{I_i} (\partial_t e_z^h + M e_z^h, \chi_h)_{n_v} dt, \\ &= \sum_{i=1}^N (e_z^h, \chi_h)_{n_v} \Big|_{t_{i-1}^+}^{t_i^-} + \sum_{i=1}^N \int_{t_{i-1}}^{t_i} (e_z^h, M \chi_h)_{n_v} dt, \\ &= (e_z^h(t_N), \chi_h(t_N^-))_{n_v} - \sum_{i=1}^{N-1} (e_z^h(t_i), \llbracket \chi_h \rrbracket_i)_{n_v} \\ &\quad + \sum_{i=1}^N \int_{t_{i-1}}^{t_i} (e_z^h, M \chi_h)_{n_v} dt. \end{aligned}$$

Adding and subtracting the term  $(e_z^h(t_N), \chi_h(t_N^+))$  we get,

$$\begin{aligned} B(e_z^h, \chi_h) &= (e_z^h(t_N), \chi_h(t_N^+))_{n_v} - \sum_{i=1}^N (e_z^h(t_i), \llbracket \chi_h \rrbracket_i)_{n_v} \\ &\quad + \sum_{i=1}^N \int_{t_{i-1}}^{t_i} (e_z^h, M \chi_h)_{n_v} dt. \end{aligned}$$

Since  $\chi_h$  is the solution to (4.2.5) with (4.2.6), the first term becomes the norm of the error at time  $t_N$  and the second two vanish. From the argument at the beginning of the proof, the result follows.  $\square$

Now that we have the error representation, we use approximation properties of local interpolation operators to derive an upper bound on the discrete portion of the error. However, first we must specify the operator  $\Pi$  and derive suitable error estimates. Let  $R_i : V^{n_v} \mapsto V_h^{i, n_v}$  be the orthogonal projection with respect to the inner product  $(\cdot, \cdot)_{n_v}$ . Then from section 2.2, due to the best approximation property, we have the estimate,

$$\|(I - R_i)w\|_{n_v} \leq Ch_i |w|_{H^2(\Omega)^{n_v}}. \quad (4.2.15)$$

**Definition 4.2.1.** Let  $\Pi : Z \mapsto Z_h$  be defined by,

$$(\Pi w)(t)|_{I_i} := R_{i-1}w(t_{i-1})\phi_{i-1}(t) + R_i w(t_i)\phi_i(t), \quad (4.2.16)$$

where  $\{\phi_{i-1}, \phi_i\}$  are the Lagrange basis functions for the time discretisation (2.5.9), and  $R_i$  is the orthogonal projection with respect to the inner product  $(\cdot, \cdot)_{n_v}$ .

We can now show an upper bound for the discrete portion of the error in terms of the interpolation error. In doing so we require the assumption that the spaces are nested moving forwards in time, that is  $V_h^{i-1, n_v} \subset V_h^{i, n_v}$ . In general, this is undesirable for a space and time adaptive algorithm since the likelihood is that the meshes will be different, due to refinement and coarsening from one step to the next. However in the context of an *a priori* analysis there is little sense in pursuing a more general result since what we are looking for is a result as the discretisation is uniformly refined, i.e., as  $h \rightarrow 0$ .

**Lemma 4.2.5.** *Assume that  $V_h^{i-1, n_v} \subset V_h^{i, n_v}$ , then there is a constant  $C > 0$  such that the discrete portion  $e_z^h$  of the error  $e_z$  satisfies,*

$$\|e_z^h(t_N)\|_{n_v} \leq C \|(I - \Pi)z\|_{L^p(I; V^{n_v})}. \quad (4.2.17)$$

*Proof.* From the result of the lemma (4.2.13) integration by parts over the time domain gives,

$$\begin{aligned} \|e_z^h(t_N)\|_{n_v}^2 &= \int_I (\partial_t(\Pi z - z) + M(\Pi z - z), \chi_h)_{n_v} dt, \\ &= \sum_{i=1}^N \left[ (\Pi z - z, \chi_h)_{n_v} \Big|_{t_{i-1}^+}^{t_i^-} + \int_{t_{i-1}^+}^{t_i^-} (\Pi z - z, M\chi_h)_{n_v} dt \right]. \end{aligned}$$

Note that  $\chi_h(t_i^-) = \chi_h(t_{i-1}^+) = \chi_h^{i-1} \in V_h^{i-1, n_v}$ . Evaluating the first term on the right hand side at the endpoints of the interval results in,

$$\begin{aligned} (\Pi z - z, \chi_h)_{n_v} \Big|_{t_{i-1}^+}^{t_i^-} &= (\Pi z(t_i) - z(t_i), \chi_h(t_i^-))_{n_v} - (\Pi z(t_{i-1}) - z(t_{i-1}), \chi_h(t_{i-1}^+))_{n_v}, \\ &= ((R_i - I)z(t_i), \chi_h^{i-1})_{n_v}, \\ &= 0, \end{aligned}$$

due to  $R_{i-1}$  being the orthogonal projection and the assumption that  $V_h^{i-1, n_v} \subset V_h^{i, n_v}$ . This leaves,

$$\begin{aligned}
\|e_z^h(t_N)\|_{n_v}^2 &= \sum_{i=1}^N \int_{t_{i-1}}^{t_i} (M(\Pi z - z), \chi_h^i)_{n_v} dt, \\
&\leq \sum_{i=1}^N \int_{t_{i-1}}^{t_i} \|M(I - \Pi)z\|_{n_v} \|\chi_h^i\|_{n_v} dt, \\
&\leq m_l \sum_{i=1}^N \|\chi_h^i\|_{n_v} \int_{t_{i-1}}^{t_i} \|(I - \Pi)z\|_{n_v} dt, \\
&\leq m_l \sum_{i=1}^N k_i^{1/q} \|\chi_h^i\|_{n_v} \|(I - \Pi)z\|_{L^p(I_i; V^{n_v})}.
\end{aligned}$$

From the stability of the discrete dual problem, lemma 4.2.10, we arrive at the bound for the discrete error in terms of the interpolation error as follows,

$$\begin{aligned}
\|e_z^h(t_N)\|_{n_v}^2 &\leq m_l \sum_{i=1}^N k_i^{1/q} \|\chi_h^i\|_{n_v} \|\Pi z - z\|_{L^p(I_i; V^{n_v})}, \\
&\leq m_l \|e_z^h(t_N)\|_{n_v} \exp\left(\frac{m_l t_N}{2}\right) \sum_{i=1}^N k_i^{1/q} \|(I - \Pi)z\|_{L^p(I_i; V^{n_v})}, \\
&\leq m_l \|e_z^h(t_N)\|_{n_v} \exp\left(\frac{m_l t_N}{2}\right) t_N^{1/q} \|(I - \Pi)z\|_{L^p(I; V^{n_v})}.
\end{aligned}$$

□

The purpose of the next lemma is to derive an error estimate in the  $L^p(I; V^{n_v})$  norm for the space time interpolation operator  $\Pi$ . However, we present the main proposition in a partial state since it gives us some indications of what to expect in the later *a posteriori* analysis. We then show that under particular assumptions on the sequence of spaces and the underlying discretisation, different estimates are available.

**Proposition 4.2.6.** *There is a constant  $C > 0$ , such that the interpolation operator defined by (4.2.16) satisfies the error estimate,*

$$\begin{aligned}
\|(I - \Pi)z\|_{L^p(I_i; V^{n_v})} &\leq C \left( \|(I - R_i)z\|_{L^p(I; V^{n_v})} + \|(I - R_{i-1})z\|_{L^p(I_i; V^{n_v})} \right. \\
&\quad \left. + k_i \|(R_i - R_{i-1})\partial_t z\|_{L^p(I_i; V^{n_v})} + k^2 \|\partial_t^2 z\|_{L^p(I_i; V^{n_v})} \right).
\end{aligned}$$

*Proof.* Restricting  $t$  to the interval  $I_i := (t_{i-1}, t_i)$ , Taylor's Theorem (1.6.7) gives,

$$R_i z(t_i) = R_i z(t) + (t_i - t)R_i \partial_t z(t) + \int_t^{t_i} R_i \partial_t^2 z(s)(t_i - s) ds, \quad (4.2.18)$$

$$R_{i-1} z(t_{i-1}) = R_{i-1} z(t) - (t - t_{i-1})R_{i-1} \partial_t z(t) + \int_{t_{i-1}}^t R_{i-1} \partial_t^2 z(s)(s - t_{i-1}) ds. \quad (4.2.19)$$

Using the identity  $\phi_i(t) + \phi_{i-1}(t) = 1$ ,  $\forall t \in I_i$ , then using (4.2.18) and (4.2.19) and the definitions of  $\phi_i$  and  $\phi_{i-1}$  (2.5.9), it follows that,

$$(I - \Pi)z(t) = (z(t) - R_i z(t_i))\phi_i(t) + (z(t) - R_{i-1} z(t_{i-1}))\phi_{i-1}(t) \quad (4.2.20)$$

$$= (I - R_i)z(t)\phi_i(t) + (I - R_{i-1})z(t)\phi_{i-1}(t), \quad (4.2.21)$$

$$+ \phi_i(t) \int_{t_{i-1}}^t R_i \partial_t^2 z(s)(t_i - s) ds \quad (4.2.22)$$

$$+ \phi_{i-1}(t) \int_t^{t_i} R_{i-1} \partial_t^2 z(s)(s - t_{i-1}) ds \quad (4.2.23)$$

$$+ \psi_i(t)(R_i - R_{i-1})\partial_t z(t), \quad (4.2.24)$$

where  $\psi_i(t) = k_i^{-1}(t - t_{i-1})(t_i - t)$ . Now consider each line (4.2.21), (4.2.22), (4.2.23) and (4.2.24) in turn. The  $L^p(I_i; V^{nv})$  norm of the right hand side of (4.2.21), and the triangle inequality results in,

$$\|(I - R_i)z\phi_i + (I - R_{i-1})z\phi_{i-1}\|_{L^p(I_i; V^{nv})} \quad (4.2.25)$$

$$\leq \|\phi_i\|_{L^\infty(I_i)}\|(I - R_i)z\|_{L^p(I_i; V^{nv})} + \|\phi_{i-1}\|_{L^\infty(I_i)}\|(I - R_{i-1})z\|_{L^p(I_i; V^{nv})},$$

$$\leq \|(I - R_i)z\|_{L^p(I_i; V^{nv})} + \|(I - R_{i-1})z\|_{L^p(I_i; V^{nv})}. \quad (4.2.26)$$

For (4.2.22), from Hölder's inequality and that  $\|\phi_i\|_{L^p(I_i)} = (p+1)^{-1/p}k_i^{1/p}$  we have,

$$\left\| \phi_i(t) \int_{t_{i-1}}^t R_i \partial_t^2 z(s)(t_i - s) ds \right\|_{L^p(I_i; V^{nv})} \leq \|\phi_i(t)\|_{L^p(I_i)} \|R_i \partial_t^2 z(t)(t_i - t)\|_{L^1(I_i; V^{nv})}, \quad (4.2.27)$$

$$\leq (p+1)^{-1/p}k_i^{1/p} \|R_i \partial_t^2 z\|_{L^p(I_i; V^{nv})} \|t_i - t\|_{L^q(I_i)}. \quad (4.2.28)$$

Then since  $\|t_i - t\|_{L^q(I_i)} = (q-1)^{-1/q}k_i^{1+1/q}$ , and  $\|R_i\| \leq 1$  it follows that,

$$\begin{aligned} \left\| \phi_i(t) \int_{t_{i-1}}^t R_i \partial_t^2 z(s)(t_i - s) ds \right\|_{L^p(I_i; V^{nv})} &\leq \frac{k_i^{1/p} k_i^{1+1/q}}{(p+1)^{1/p} (q+1)^{1/q}} \|\partial_t^2 z\|_{L^p(I_i; V^{nv})}, \\ &\leq C k_i^2 \|\partial_t^2 z\|_{L^p(I_i; V^{nv})}. \end{aligned} \quad (4.2.29)$$

The exact same bound can be found for (4.2.23). For (4.2.24) since,

$$\|\psi_i\|_{L^\infty(I_i)} = \frac{k_i}{4},$$

it follows that,

$$\|\psi_i(R_i - R_{i-1})\partial_t z\|_{L^p(I_i; V^{n_v})} \leq \frac{k_i}{4} \|(R_i - R_{i-1})\partial_t z\|_{L^p(I_i; V^{n_v})}. \quad (4.2.30)$$

Combining (4.2.26), (4.2.28) and (4.2.30) implies the result.  $\square$

**Corollary 4.2.7.** *Suppose that  $V_h^i = V_h^j$ ,  $\forall i, j$ , then there exists a constant  $C > 0$  independent of the discretisation parameters such that the approximation operator  $\Pi$  defined in 4.2.1 satisfies the error estimate,*

$$\|(I - \Pi)w\|_{L^p(I; V^{n_v})} \leq C \left\{ h\|w\|_{L^p(I; H^2(\Omega)^{n_v})} + k^2 \|\partial_t^2 w\|_{L^p(I; V^{n_v})} \right\}, \quad l \geq 0. \quad (4.2.31)$$

*Proof.* From the result of proposition 4.2.6, it follows that if  $V_h^i = V_h^j = V_h$ , then  $R_i = R$  for all  $i$ . Therefore the term involving  $R_i - R_{i-1}$  vanishes, and since  $R$  is the orthogonal projection onto  $V_h$ , it satisfies the best approximation estimate (4.2.15). Summing over the intervals proves the result.  $\square$

**Corollary 4.2.8.** *Suppose that there exist constants  $c, c' > 0$  such that,*

$$h_{i-1} \leq ch_i, \quad (4.2.32)$$

$$h_i \leq c'k_i, \quad (4.2.33)$$

*then there exists a constant  $C > 0$  independent of the discretisation parameters such that the approximation operator  $\Pi$  defined in 4.2.1 satisfies the error estimate,*

$$\|(I - \Pi)w\|_{L^p(I; V^{n_v})} \leq C \left\{ h\|w\|_{L^p(I; H^2(\Omega))} + k^2 \|\partial_t w\|_{L^p(I; H^2(\Omega)^{n_v})} + k^2 \|\partial_t^2 w\|_{L^p(I; V^{n_v})} \right\}.$$

*Proof.* From the result of proposition 4.2.6, we need only to bound the terms,

$$\|(I - R_i)w\|_{L^p(I_i; V^{n_v})}, \quad \|(I - R_{i-1})w\|_{L^p(I_i; V^{n_v})}, \quad \text{and} \quad \|(R_i - R_{i-1})\partial_t w\|_{L^p(I_i; V^{n_v})}. \quad (4.2.34)$$

First use the triangle inequality on the last term,

$$\|(R_i - R_{i-1})\partial_t z\|_{L^p(I_i; V^{nv})} \leq \|(I - R_i)\partial_t w\|_{L^p(I_i; V^{nv})} + \|(I - R_{i-1})\partial_t w\|_{L^p(I_i; V^{nv})}, \quad (4.2.35)$$

Then due to assumption (4.2.32), we have,

$$\begin{aligned} \|(R_i - R_{i-1})\partial_t z\|_{L^p(I_i; V^{nv})} &\leq Ch_i \|\partial_t w\|_{L^p(I; H^2(\Omega))} + Ch_{i-1} \|\partial_t w\|_{L^p(I; H^2(\Omega))}, \\ &\leq C(1 + c)h_i \|\partial_t w\|_{L^p(I; H^2(\Omega))}. \end{aligned}$$

Then similarly using (4.2.32) to collect the two orthogonal projection error terms together results in,

$$\begin{aligned} \|(I - \Pi)w\|_{L^p(I_i; V^{nv})} &\leq C \left\{ h_i \|w\|_{L^p(I_i; H^2(\Omega)^{nv})} + k_i h_i \|\partial_t w\|_{L^p(I_i; H^2(\Omega)^{nv})} \right. \\ &\quad \left. + k_i^2 \|\partial_t^2 w\|_{L^p(I_i; V^{nv})} \right\}. \end{aligned}$$

We can use assumption (4.2.33) to take a power of  $h_i$  away from the middle term to get an extra power of  $k_i$ , so that,

$$\begin{aligned} \|(I - \Pi)w\|_{L^p(I_i; V^{nv})} &\leq C \left\{ h_i \|w\|_{L^p(I_i; H^2(\Omega)^{nv})} + k_i^2 \|\partial_t w\|_{L^p(I_i; H^2(\Omega)^{nv})} \right. \\ &\quad \left. + k_i^2 \|\partial_t^2 w\|_{L^p(I_i; V^{nv})} \right\}. \end{aligned}$$

Raising to the power  $p$  summing and taking  $p$ -th roots gives the result.  $\square$

The previous two results, corollaries 4.2.7 and 4.2.8 provide specialisations of theorem 4.2.6 that we will require in the sequel. In particular, they show that under the specified assumptions, the error for a  $\mathbb{P}_1$  in space  $\mathbb{P}_1$  in time interpolation is  $O(h + k^2)$ . We can now combine these results with lemmas 4.1.2 and 4.2.5 to provide *a priori* upper bounds on the discretisation error.

### 4.3 *A priori* estimates

We now combine the work of the previous sections to show *a priori* error estimates for the approximation to the displacement and the internal variables.

**Theorem 4.3.1.** *Assume that  $V_h^{i-1, n_v} \subset V_h^{i, n_v}$ ,  $\forall i \in 1, \dots, N$  and suppose that there exist constants  $c, c' > 0$  such that there holds,*

$$h_{i-1} \leq ch_i, \quad h_i \leq c'k_i, \quad (4.3.1)$$

then the error in the internal variables  $e_z(t) = z(t) - z_h(t)$  satisfies the a priori estimate,

$$\begin{aligned} \| \|e_z(t_N)\| \|_{n_v} \leq C \left\{ h|z(t_N)|_{H^2(\Omega)} + h\|z\|_{L^p(I; H^2(\Omega)^{n_v})} \right. \\ \left. + k^2\|\partial_t^2 z\|_{L^p(I; V^{n_v})} + k^2\|\partial_t z\|_{L^p(I; H^2(\Omega)^{n_v})} \right\}. \end{aligned}$$

*Proof.* The triangle inequality combined with (4.2.15) and lemma 4.2.5 implies that,

$$\begin{aligned} \| \|e_z(t_N)\| \|_{n_v} &\leq \| \|z(t_N) - \Pi z(t_N)\| \|_{n_v} + \| \|e_z^h(t_N)\| \|_{n_v}, \\ &\leq C \left\{ h|z(t_N)|_{H^2(\Omega)^{n_v}} + \|(I - \Pi)z\|_{L^p(I; V^{n_v})} \right\}. \end{aligned}$$

Then using the result of 4.2.8 completes the proof.  $\square$

Given the a priori estimate for the internal variables, we can use the result of corollary 3.2.3 which allows us to bound Sobolev norms of  $z$  in terms of a weaker Sobolev norm of  $u$ , together with lemma (4.1.2) to derive an estimate completely in terms of  $u$ .

**Theorem 4.3.2.** *There exist constants  $C > 0$ ,  $C' > 0$  independent of the discretisation parameters such that the error  $e_u(t_N)$  between the solution of (3.3.7) and the solution to (3.2.17) satisfies,*

$$\| \|e_u(t_N)\| \| \leq Ch \left\{ \|u(t_N)\|_{H^2(\Omega)} + \|u\|_{L^p(I; H^2(\Omega))} + \|u\|_{L^p(I; V)} \right\} + C'k^2\|u\|_{W^{1,p}(I; V)}. \quad (4.3.2)$$

*Proof.* Lemmas 4.1.2 and 4.2.15 taken together with theorem 4.3.1 imply that,

$$\begin{aligned} \| \|e_u(t_N)\| \| &\leq \inf_{w_h^N \in V_h^N} \| \|u(t_N) - w_h^N\| \| + |\varphi'(0)| \| \|e_z(t_N)\| \|_{n_v}, \\ &\leq C \left\{ h\|u(t_N)\|_{H^2(\Omega)} + h|z(t_N)|_{H^2(\Omega)^{n_v}} + h\|z\|_{L^p(I; H^2(\Omega)^{n_v})} \right. \\ &\quad \left. + k^2\|\partial_t^2 z\|_{L^p(I; V^{n_v})} + k^2\|\partial_t z\|_{L^p(I; H^2(\Omega)^{n_v})} \right\}. \end{aligned}$$

From corollary 3.2.3 we have  $|z(t_N)|_{H^{t+1}(\Omega)^{n_v}} \leq C\|u\|_{L^\infty(I;H^{t+1}(\Omega))}$  and  $\|z\|_{W^{r+1,p}(I;X^{n_v})} \leq C\|u\|_{W^{r,p}(I;X)}$ , then,

$$\|e_u(t_N)\| \leq C \left\{ h\|u\|_{L^\infty(I;H^2(\Omega))} + h\|u\|_{L^p(I;H^2(\Omega))} + k^2\|\partial_t u\|_{L^p(I;V)} + k^2\|u\|_{L^p(I;H^2(\Omega))} \right\}.$$

□

## 4.4 Numerical results

The purpose of this section is to computationally demonstrate the convergence properties of the approximation scheme presented in section 3.3 and analysed earlier in this chapter.

### 4.4.1 Implementation

The implementation of the algorithms has been carried out using MATLAB. The code we use is derived from that presented in [5] which contains a short implementation of the  $\mathbb{P}_1$  Galerkin finite element method for linear elasticity in two and three dimensions. As we have already mentioned one of the advantages of the internal variable formulation of linear viscoelasticity is the reformulation results in a modified linear elasticity problem, together with a system of ODEs. Therefore the only modifications to the linear elasticity code are to enable the addition of a vector to the right hand side of the linear elasticity system which contains the internal variable information. In addition to the internal variable solver code, a mesh refinement/coarsening suite has also been implemented (see appendix A for details of the refinement process). Given the requirement for adaptive spatial mesh, we have limited ourselves to the two dimensional case. In addition to the above software, we note some alterations made to the original code of [5] (for example higher order quadrature) in assembling the right hand side, that are required for interpreting the results. To be more specific:

- Integration over elements: sixth order quadrature for triangles [51].
- Integration over edges and time intervals: sixth order one dimensional Gaussian quadrature [51].

#### 4.4.2 Convergence diagnostics

To test the convergence properties of the method, we consider problems for which we know the exact solution, achieved by first choosing the solution and applying the required operators to get the problem data such as the volume force  $f$  and boundary traction  $g$ . As a first test we confirm that the problem generates the exact solution when possible.

To examine separately the convergence rates as either  $h \rightarrow 0$  or  $k \rightarrow 0$ , we choose test problems such that they are exact in either space or time. Given a problem that has a solution that is linear in time, we can examine the spatial convergence of the method as  $h \rightarrow 0$ , since any error is due to the spatial discretisation. Similarly for a given problem that has a solution that is linear in the spatial variables, we can examine the temporal convergence as  $k \rightarrow 0$ . For both of these scenarios we will have exact solutions available so that we can measure the true value of the error. To establish orders of convergence we use the empirical order of convergence (EOC) given in (2.5.34).

According to the results of the previous section, in particular theorems 4.3.1 and 4.3.2, we expect the energy norm of the errors at the time nodes to converge linearly as  $h \rightarrow 0$  and quadratically as  $k \rightarrow 0$ .

#### Convergence in space

In this section we will examine the convergence of the approximation as  $h \rightarrow 0$ . Given that the problem is designed so that the approximation is exact in time, we can take a fixed number of timesteps and examine the error as the spatial mesh is refined. Let  $\Omega = [0, 1] \times [0, 1]$  and  $I = [0, 1]$  be the problem domain. Let  $\Gamma_D = \{(x, 0), x \in (0, 1)\} \cup \{(x, 1), x \in (0, 1)\}$  and let  $\Gamma_N = \{(0, y), y \in (0, 1)\} \cup \{(1, y), y \in (0, 1)\}$ .

For this test problem let the displacement and internal variable be given by,

$$u = \begin{pmatrix} \kappa(1 + \alpha_1 t) \sin(2\pi y) \\ 0 \end{pmatrix}, \quad (4.4.1)$$

$$z = \begin{pmatrix} \kappa\beta_1 t \sin(2\pi y) \\ 0 \end{pmatrix}. \quad (4.4.2)$$

	Degrees of freedom			
t	25	81	289	1089
0	6.246518e-002	3.283594e-002	1.676324e-002	8.419381e-003
0.2	7.245961e-002	3.808969e-002	1.944536e-002	9.766482e-003
0.4	8.245404e-002	4.334344e-002	2.212748e-002	1.111358e-002
0.6	9.244847e-002	4.859719e-002	2.480960e-002	1.246068e-002
0.8	1.024429e-001	5.385094e-002	2.749172e-002	1.380778e-002
1.0	1.124373e-001	5.910469e-002	3.017384e-002	1.515489e-002

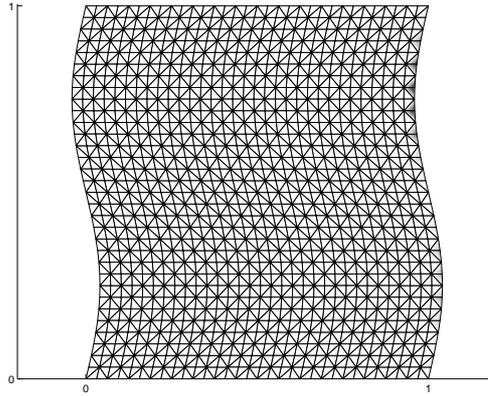
Table 4.1: Energy norm errors of the displacement at time points.

	Degrees of freedom			
t	25	81	289	1089
0.0	0	0	0	0
0.2	7.901290e-003	4.153454e-003	2.120401e-003	1.064977e-003
0.4	1.580258e-002	8.306908e-003	4.240803e-003	2.129954e-003
0.6	2.370387e-002	1.246036e-002	6.361204e-003	3.194930e-003
0.8	3.160516e-002	1.661382e-002	8.481605e-003	4.259907e-003
1.0	3.950645e-002	2.076727e-002	1.060201e-002	5.324884e-003

Table 4.2: Energy norm errors of the internal variable at selected time points.

Dofs	$\max_i \ e_u(t_i)\ $	EOC	$\max_i \ e_z(t_i)\ _{n_v}$	EOC
25	2.3148e+00	0	5.1439e-01	0
81	1.1055e+00	1.2573e+00	2.4566e-01	1.2573e+00
289	5.4677e-01	1.1070e+00	1.2150e-01	1.1070e+00
1089	2.7218e-01	1.0517e+00	6.0484e-02	1.0517e+00
4225	1.3591e-01	1.0244e+00	3.0203e-02	1.0244e+00

Table 4.3: Maximum energy norms and EOC values for the displacement and internal variable at different refinement levels.

Figure 4.1: Solution at  $t=1.0$ .

The stress is then given by,

$$\boldsymbol{\sigma} = \begin{pmatrix} 0 & 2\mu\pi\kappa(1 + \alpha_i t - \beta_1^2 t) \cos(2\pi y) \\ 2\mu\pi\kappa(1 + \alpha_i t - \beta_1^2 t) \cos(2\pi y) & 0 \end{pmatrix}, \quad (4.4.3)$$

and the body force and surface tractions can be calculated accordingly. Tables 4.1 and 4.3 show the energy norm of the errors in the approximation for various mesh widths at different points in time. For calculating the EOC value, we use the modified form of (2.5.34),

$$EOC = \frac{\ln(e_{k-1}/e_k)}{\ln(h_{k-1}/h_k)}, \quad (4.4.4)$$

where  $h_k$  is the mesh width at refinement level  $k$ . In the case of adaptive mesh refinements, it is more useful to consider the EOC in terms of degrees of freedom. Generally a refinement of a two dimensional triangle reduces the meshwidth by a factor of  $2^{-\frac{1}{2}}$

Calculating the experimental order of convergence we find that the method does indeed display linear convergence as  $h \rightarrow 0$ , and this is visible in the log-log plot of figure 4.2.

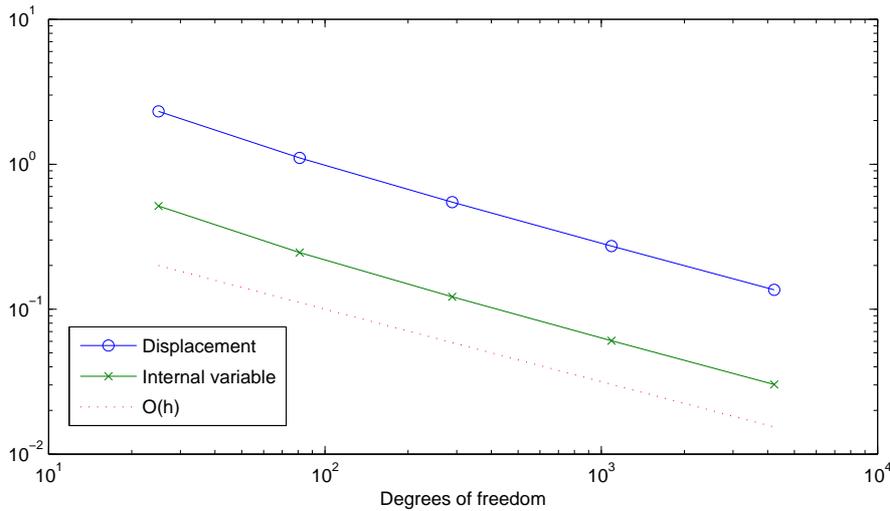


Figure 4.2: Log-log plot of the maximum energy norm of the errors against the degrees of freedom.

### Convergence as $k \rightarrow 0$

In this section we will examine the convergence of the approximation as  $k \rightarrow 0$ . The problem given below is designed so that the approximation is exact in the spatial variable, and therefore we can examine the behaviour of the error due to the time discretisation as the number of time steps is increased. Let  $\Omega = [0, 1] \times [0, 1]$  and  $I = [0, 1]$  be the problem domain. Let  $\Gamma_D = \{(x, 0), x \in (0, 1)\}$  and let  $\Gamma_N = \{(0, y), y \in (0, 1)\} \cup \{(1, y), y \in (0, 1)\} \cup \{(x, 1), y \in (0, 1)\}$ . For this test problem the displacement and internal variable are given by,

$$u = \begin{pmatrix} \kappa y \sin(2\pi t) \\ 0 \end{pmatrix}, \quad (4.4.5)$$

$$z = \begin{pmatrix} \frac{\kappa\beta_1}{\alpha_1^2 + 4\pi^2} \left\{ \alpha_1 \sin(2\pi t) - 2\pi \cos(2\pi t) + 2\pi e^{-\alpha_1 t} \right\} y \\ 0 \end{pmatrix}. \quad (4.4.6)$$

The solution is a linear shearing motion applied sinusoidally in time. The initial mesh is shown in figure 4.3.

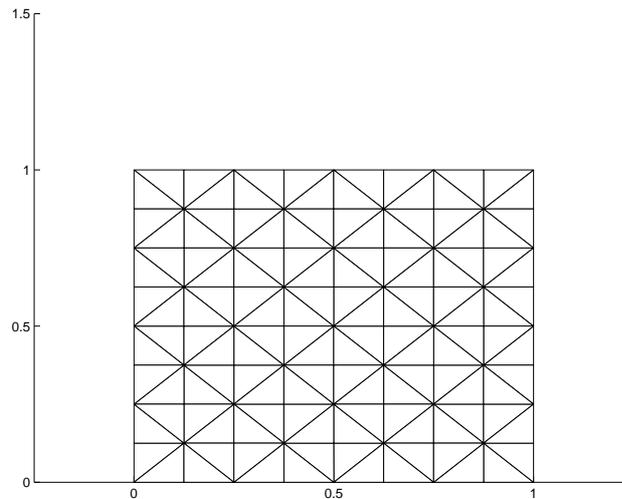


Figure 4.3: Initial mesh.

As in the previous example the stress and surface tractions can be calculated according to the strong form of the PDE and the boundary conditions, furthermore it follows that  $f = 0$ . Table 4.4 shows the energy norm of the approximation errors at the final time  $t = T$ . Calculating the experimental order of convergence we find that the method does display the expected second order convergence rate as  $k \rightarrow 0$ . This is further illustrated by the log-log plot in figure 4.4.

## 4.5 Summary

In this chapter we have provided an *a priori* error analysis of finite element approximation to quasistatic linear viscoelasticity. We proved and displayed by numerical examples that the spatial approximation convergences linearly in the energy norm, and that we have quadratic convergence in the temporal discretisation. In the next chapter we consider an *a posteriori* error analysis and derive an adaptive in space and time algorithm for generating solutions.

Timesteps	$\max_i \ e_u(t_i)\ $	EOC	$\max_i \ e_z(t_i)\ _{n_v}$	EOC
4	3.3289e-04	0	8.3221e-04	0
8	8.0831e-05	2.0420e+00	2.0208e-04	2.0420e+00
16	2.0429e-05	1.9843e+00	5.1073e-05	1.9843e+00
32	5.1435e-06	1.9898e+00	1.2859e-05	1.9898e+00
64	1.2856e-06	2.0003e+00	3.2140e-06	2.0003e+00
128	3.2169e-07	1.9987e+00	8.0422e-07	1.9987e+00
256	8.0420e-08	2.0000e+00	2.0105e-07	2.0000e+00
512	2.0105e-08	2.0000e+00	5.0262e-08	2.0000e+00

Table 4.4: Maximum energy norms and EOC values for the displacement and internal variable at different time refinement levels.

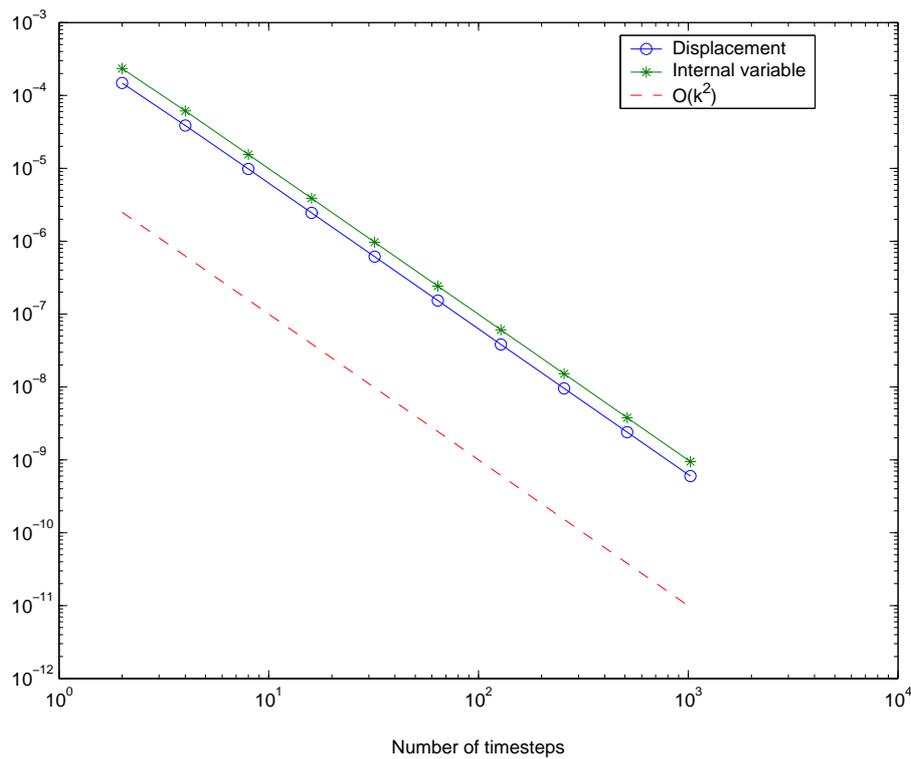
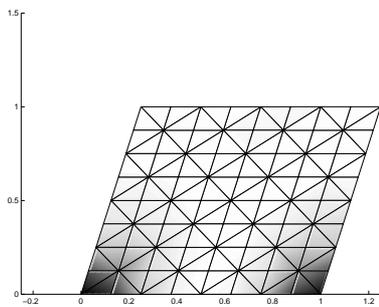
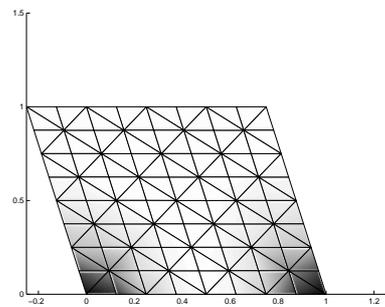


Figure 4.4: Log-log plot of the errors against the degrees of freedom.

Figure 4.5: Solution at  $t=0.25$ .Figure 4.6: Solution at  $t=0.75$ .

## Chapter 5

# *A posteriori* error analysis

The chapter presents an *a posteriori* error analysis of the finite element approximation described in section 3.3. In contrast to the previous chapter where the error in the approximation is bounded in terms of the discretisation parameters and norms of the true solution, this chapter focuses on deriving computable upper bounds of the approximation error in terms of the known approximate solution and the problem data. The methods and results of chapter 2 are applied to derive reliable and efficient residual based *a posteriori* error estimates for both the displacement and internal variable problem. A slightly more general approach is taken, in that the results regarding reliability are, conditional on obtaining stability of the dual solution, applicable for errors in arbitrary linear functionals, however as remarked earlier we are required to specialise to energy norm estimates at the time nodes for proofs of efficiency, since such results are not available for arbitrary linear functionals.

In both the displacement and internal variable problems, we begin by deriving representations of linear functionals of the error in terms of the residual acting upon the solution of a suitably defined dual problem. After confirming stability properties of the dual solution, we proceed to derive computable upper and lower bounds on the error using optimal order interpolation and quasi-interpolation error estimates from section 2.2 and the bubble functions of section 2.4.

We first show an *a posteriori* estimate for the error in the displacement  $e_u(t) = u(t) - u_h(t)$ . The estimate comprises of two parts, one is a residual type term for the discretisation,

the other is the error in the approximation of the internal variables. The results of section 2.4 are used to analyse the residual type term since it has a form similar to that of the linear elasticity residual, with minor modifications to due to the internal variables.

To characterise the error in the internal variables  $e_z(t) = z(t) - z_h(t)$ , we look to the results presented in [34] and [49] and the review article [30] for error estimation of time and space-time problems. Taking a lead from the sketch for ODEs given in section 2.5, we attempt to derive results similar to those given in [34], however there is an added complication in that there is also a spatial discretisation at play.

For both error estimators we consider the properties of reliability and efficiency. Recall that an error indicator is reliable and efficient if it is both an upper and lower bound for the true error modulo higher order data terms. For time dependent problems efficiency appears to be more difficult to prove than in stationary problems, and it is common to call the estimator efficient if the *a posteriori* error estimate is bounded above by an optimal order *a priori* estimate. This is different from genuine efficiency as our numerical results will show. For the spatial discretisation, the theory presented in 2.4 is used to show that local error indicators are bounded above by a projection of the true error together with a term representing the time variation of the error and higher order data terms. For the temporal discretisation, we must settle for showing that the temporal error indicator satisfies an optimal order *a priori* upper bound.

We close the chapter with some numerical experiments confirming our theoretical results, however they also expose a weakness with the temporal indicator. While it is reliable, and satisfies an optimal order *a priori* upper bound, it is deficient in that it is non-zero when the approximation is exact, or equivalently when the true solution is piecewise linear in time. Therefore any stepping scheme based on this error indicator will require time steps that are smaller than actually required. We consider this problem in greater detail in chapter 6.

## 5.1 Displacement

To derive *a posteriori* error estimates we will follow the approach set out in section 2.3 of using a dual problem to facilitate an error representation formula. Let  $V$  be the space

defined in (1.5.5) and assume that  $J : V \rightarrow \mathbb{R}$  is a bounded linear functional. Define the dual problem: Find  $\chi(t) \in V$  such that,

$$a(v, \chi(t)) = J(v), \quad \forall v \in V, \quad a.e. t \in I, \quad (5.1.1)$$

where  $a(\cdot, \cdot)$  is the positive definite symmetric bilinear form of section 2.4. By theorem 2.1.1 there exists a unique solution to equation (5.1.1) satisfying the *a priori* bound  $\|\chi(t)\| \leq \|J\|_{V^*}$ . Define the stability constant  $C_{\text{stab}} = \|J\|_{V^*}$ .

**Definition 5.1.1.** The residual of the approximation (3.3.7) to problem (3.2.17) is a linear functional on  $V$ , defined by,

$$\langle R(u_h(t), z_h(t)), v \rangle = l(t; v) + \sum_{j=1}^{n_v} \beta_j a(z_h^j(t), v) - a(u_h(t), v), \quad \forall v \in V, \quad (5.1.2)$$

where  $l(t; \cdot)$  is the linear functional defined by (1.5.7).

We also have the Galerkin orthogonality relationship for each  $t_n$ ,  $n = 1, \dots, N$ ,

$$\langle R(u_h(t_n), z_h(t_n)), v \rangle = 0, \quad \forall v \in V_n^h, \quad (5.1.3)$$

which follows directly from the finite element equations (3.3.7) and the definition of the residual (5.1.2). Since  $J$  is linear, the error in the functional is equal to the functional acting on the error,  $J(u) - J(u_h) = J(e_u)$ . We can now derive a representation of the error in the functional  $J$  in terms of the solution to the dual problem  $\chi$  (5.1.1), and the residual (5.1.2).

**Lemma 5.1.1.** *If  $\chi$  is the solution to the problem (5.1.1) then the functional  $J$  acting on the error in the displacement  $e_u(t)$  satisfies the representation,*

$$J(e_u) = \langle R(u_h(t), z_h(t)), \chi \rangle + \sum_{j=1}^{n_v} \beta_j a(e_{z^j}, \chi). \quad (5.1.4)$$

*Proof.* Starting with the dual problem (5.1.1) we have,

$$\begin{aligned} J(e_u(t)) &= a(e_u(t), \chi(t)), \\ &= a(u(t), \chi(t)) - a(u_h(t), \chi(t)), \\ &= l(t; \chi(t)) + \sum_{j=1}^{n_v} \beta_j a(z_h^j(t), \chi(t)) - a(u_h(t), \chi(t)), \end{aligned}$$

where we use (3.2.17) in the last step. Adding and subtracting  $\sum_{j=1}^{n_v} \beta_j a(z_h^j(t), \chi(t))$  from the right hand side, we can organise terms so that,

$$J(e_u(t)) = l(t; \chi(t)) + \sum_{j=1}^{n_v} \beta_j a(z_h^j(t), \chi(t)) - a(u_h(t), \chi(t)) + \sum_{j=1}^{n_v} \beta_j a(e_{z^j}(t), \chi(t)).$$

The first three terms on the right hand side can be recognised as those forming the residual (5.1.2) □

The derivation of this error representation may seem odd at first since the typical approach would be to write the full dual problem, collecting together both the displacement equation and the internal variable equation. However, it is our aim to attempt to separate the problems as much as possible, therefore we stick with treating the displacement problem as a linear elasticity problem, with the internal variables as part of the problem data. Inclusion of some dual internal variables would lead back to the fully coupled situation for the error estimates, and the separation would not be possible.

### 5.1.1 Reliability

We now focus on deriving an upper bound for the residual term (5.1.2). Based on the observation that it has similarities to the residual for the linear elasticity problem we derive a localised representation, followed by an upper bound using similar arguments as those of section 2.4. Recall that the residual based *a posteriori* error estimate for the linear elasticity problem involved the discrete divergence of the stress, culminating in jumps over edges in the mesh. For our problem, there is a complication since the stress now depends on both the displacement and the internal variables. This will have implications later on when developing an algorithm, but for now we proceed by observing that the discrete stress has the form,

$$\boldsymbol{\sigma}(u_h(t), z_h(t)) = \mathbf{C}\boldsymbol{\epsilon}(u_h(t)) - \sum_{j=1}^{n_v} \beta_j \mathbf{C}\boldsymbol{\epsilon}(z_h^j(t)). \quad (5.1.5)$$

The following lemma uses integration by parts to arrive at the localised representation of the residual.

**Lemma 5.1.2.** *The residual of the approximation of problem (3.3.7) to problem (3.2.17) at time  $t_n$  has the localised representation,*

$$\langle R(u_h(t_n), z_h(t_n)), v \rangle = \sum_{K \in \mathcal{T}_n} \left\{ (f(t_n), v)_K + \sum_{E \in \mathcal{E}(K)} (R_E(t_n), v)_E \right\}, \quad (5.1.6)$$

where,

$$R_E(t_n) = \begin{cases} \frac{1}{2} \llbracket \boldsymbol{\sigma}(u_h(t_n), z_h(t_n)) \rrbracket, & \text{on } E \subset \Omega, \\ g(t_n) - \boldsymbol{\sigma}(u_h(t_n), z_h(t_n)) n_E, & \text{on } E \subset \Gamma_N. \end{cases} \quad (5.1.7)$$

*Proof.* From the definition of the residual (5.1.2), integration by parts over  $\Omega$  gives,

$$\begin{aligned} \langle R(u_h(t_n), z_h(t_n)), v \rangle &= l(t_n; v) + \sum_{j=1}^{n_v} \beta_j a(z_h^j(t_n), v) - a(u_h(t_n), v), \\ &= \sum_{K \in \mathcal{T}_n} (f(t_n), v)_{L^2(K)} + \sum_{E \in \Gamma_N} (g(t_n), v)_{L^2(E)} \\ &\quad + \sum_{K \in \mathcal{T}_n} \left\{ (\operatorname{div} \boldsymbol{\sigma}(u_h(t_n), z_h(t_n)), v)_{L^2(K)} \right. \\ &\quad \left. - (\boldsymbol{\sigma}(u_h(t_n), z_h(t_n)) n_{\partial K}, v)_{L^2(\partial K)} \right\}. \end{aligned}$$

Since  $u_h(t_n)$  and each  $z_h(t_n)$  are piecewise linear with respect to the triangulation  $\mathcal{T}_n$  of affine finite elements, and since  $\mathbf{C}$  is constant over  $\Omega$ , the divergence of the stress on each element is zero. As in lemma 2.4.12 we can collect the internal edge terms together and assign half to each element on either side of the edge to arrive at the result.  $\square$

Using lemma 5.1.2, we can now bound the error in the displacement as we did in the linear elasticity problem by using the optimal order quasi-interpolation error estimates of section 2.3.

**Lemma 5.1.3.** *There exists a constant  $C_{u,\text{rel}}$  depending on the domain  $\Omega$ , the coercivity constant  $c_a$ , the minimum angle in the domain, and the maximum number of overlapping element neighbourhoods  $c_M$  such that the residual (5.1.6) satisfies the bound,*

$$\| \|R(u_h(t_n), z_h(t_n))\| \|_* \leq C_{u,\text{rel}} \left( \sum_{K \in \mathcal{T}_n} \eta_K^2 \right)^{1/2}, \quad (5.1.8)$$

where,

$$\eta_K^2 := h_K^2 \|f(t_n)\|_{L^2(K)}^2 + \sum_{E \in \mathcal{E}(K)} h_E \|R_E(t_n)\|_{L^2(E)}^2. \quad (5.1.9)$$

*Proof.* Let  $\mathcal{I}_n^1 : V \rightarrow V_h^n$  be the quasi-interpolation operator of section 2.3, then by Galerkin orthogonality (5.1.3) and (5.1.6) there holds,

$$\begin{aligned} \langle R(u_h(t_n), z_h(t_n)), v \rangle &= \langle R(u_h(t_n), z_h(t_n)), v - \mathcal{I}_n^1 v \rangle, \\ &= \sum_{K \in \mathcal{T}} \left\{ (f(t_n), v - \mathcal{I}_n^1 v)_{L^2(K)} + \sum_{E \in \mathcal{E}(K)} (R_E(t_n), v - \mathcal{I}_n^1 v)_{L^2(E)} \right\}. \end{aligned}$$

Applying the Cauchy-Schwarz inequality, then using the interpolation estimates (2.2.18), (2.2.19) gives,

$$\begin{aligned} |\langle R(u_h(t_n), z_h(t_n)), v - \mathcal{I}_i^1 v \rangle| &\leq \sum_{K \in \mathcal{T}_n} \left\{ c_K h_K \|f(t_n)\|_{L^2(K)} |v|_{H^1(\tilde{\omega}_K)} \right. \\ &\quad \left. + \sum_{E \in \mathcal{E}(K)} c_E h_E^{1/2} \|R_E(t_n)\|_{L^2(K)} |v|_{H^1(\tilde{\omega}_E)} \right\}. \end{aligned}$$

The Cauchy-Schwarz inequality once more and pulling through the constants results in,

$$\begin{aligned} &|\langle R(u_h(t_n), z_h(t_n)), v - \mathcal{I}_i^1 v \rangle| \\ &\leq \max\{c_K, c_E\} \left\{ \sum_{K \in \mathcal{T}_n} h_K^2 \|f(t_n)\|_{L^2(K)}^2 + \sum_{E \in \mathcal{E}} h_E \|R_E(t_n)\|_{L^2(E)}^2 \right\}^{1/2} \\ &\quad \times \left\{ \sum_{K \in \mathcal{T}_n} |v|_{H^1(\tilde{\omega}_K)}^2 + \sum_{E \in \mathcal{E}} |v|_{H^1(\tilde{\omega}_E)}^2 \right\}^{1/2}, \\ &\leq c_M \max\{c_K, c_E\} \left( \sum_{K \in \mathcal{T}_n} \eta_K^2 \right)^{1/2} |v|_{H^1(\Omega)}. \end{aligned}$$

Using the equivalence of the  $H^1(\Omega)$  and energy norms,  $|v|_{H^1(\Omega)} \leq \|v\|_{H^1(\Omega)} \leq \frac{1}{c_a} \|v\|$ , the result follows with reliability constant,

$$C_{u,\text{rel}} = \frac{c_M \max\{c_K, c_E\}}{c_a}. \quad (5.1.10)$$

□

We can now bound the linear functional of the error in the displacement  $e_u$  in the following way.

**Theorem 5.1.4.** *The linear functional  $J : V \rightarrow \mathbb{R}$  evaluated at the error in the displacement satisfies the upper bound,*

$$|J(e_u(t_n))| \leq C_{\text{stab}} \left\{ C_{u,\text{rel}} \left( \sum_{K \in \mathcal{T}_i} \eta_K^2 \right)^{1/2} + |\varphi(0)|^{1/2} \cdot \|e_z(t_n)\| \right\}. \quad (5.1.11)$$

*Proof.* From the definition 5.1.1, the Cauchy-Schwarz inequality and lemma 5.1.3 imply that,

$$\begin{aligned} |J(e_u(t_n))| &\leq |\langle R(u_h(t_n), z_h(t_n)), \chi \rangle| + \left| \sum_{j=1}^{n_v} \beta_j a(e_{z^j}, \chi) \right|, \\ &\leq \| \|R(u_h(t_n), z_h(t_n))\|_* \| \chi(t_n) \| + \sum_{j=1}^{n_v} \beta_j \| e_{z^j}(t_n) \| \cdot \| \chi(t_n) \|, \\ &\leq C_{u,\text{rel}} \left( \sum_{K \in \mathcal{T}_n} \eta_K^2 \right)^{1/2} \| \chi(t_n) \| + |\varphi(0)|^{1/2} \cdot \| e_z(t_n) \|_{n_v} \cdot \| \chi(t_n) \|. \end{aligned}$$

Then under the assumption that  $\| \chi(t) \| \leq C_{\text{stab}}$ , the result follows.  $\square$

A case of primary interest for the remainder of this chapter is when the linear functional  $J : V \mapsto \mathbb{R}$  is given by,

$$J(v) = \frac{a(v, e)}{\| e \|}. \quad (5.1.12)$$

In this instance, since  $\chi$  solves problem (5.1.1), and satisfies the *a priori* bound  $\| \chi \| \leq \| J \|_*$ , then,

$$\| \chi \| \leq \| J \|_* \leq \sup_{v \neq 0} \frac{J(v)}{\| v \|} = \sup_{v \neq 0} \frac{a(v, e)}{\| e \| \| v \|} = 1. \quad (5.1.13)$$

Therefore in this instance,  $C_{\text{stab}} = 1$ .

### 5.1.2 Efficiency

We will extend the result of Theorem 2.4.6 to show that the error indicator for the displacement is efficient in the sense that it describes the true error up to the error in the internal variables and the data oscillation. The reason that the internal variables appear in the inequality is that they are contained in the discrete stress. The extension shows that the internal variable terms pose no extraordinary difficulty, and that we have another instance of where results for the linear elasticity problem are applicable with slight modification.

Unfortunately there does not appear to be a route through to showing the efficiency of the indicator (5.1.9) for the error in linear functionals. The technique of the proof relies on using the localised form of the residual (5.1.6) acting on particular bubble functions described in section 2.4.

Given the previous remarks, we focus on showing that the error estimator of theorem 5.1.4 is efficient with respect to the energy norm error at a fixed time  $t_n$ . The arguments

contained in the proof of the following theorem are those of theorem 2.4.6, however we include it to confirm the previous statement that the presence of the internal variable causes no major complication. Let  $f_K(t_n)$  and  $g_E(t_n)$  denote piecewise constant approximations at time  $t_n$  to  $f(t_n)$  and  $g(t_n)$  on the element  $K$  and on the edge  $E$  respectively.

**Theorem 5.1.5.** *Suppose that  $K \in \mathcal{T}_n$ , then there is a constant  $C_{u,\text{eff}}$ , independent of the mesh parameter  $h_K$  such that the local error indicator of theorem 5.1.4 for the element  $K$  satisfies the following upper bound,*

$$\begin{aligned} \eta_K^2 \leq C_{u,\text{eff}}^2 & \left\{ \|e_u(t_n)\|_{\omega_K}^2 + |\varphi'(0)| \cdot \|e_z(t_n)\|_{n_v, \omega_K}^2 \right. \\ & \left. + h_K^2 \|f(t_n) - f_K(t_n)\|_{L^2(K)}^2 + h_E \|g(t_n) - g_E(t_n)\|_{L^2(E)}^2 \right\}. \end{aligned}$$

*Proof.* Step 1. Volume terms. Inequality (2.4.18) together with the localised form of the residual (5.1.6) implies that,

$$\begin{aligned} \|f_K(t_n)\|_{L^2(K)}^2 & \leq \varepsilon_1^2 (f_K(t_n), f_K(t_n) b_K)_K, \\ & \leq \varepsilon_1^2 \left\{ (f(t_n), f_K(t_n) b_K)_K + (f_K(t_n) - f(t_n), f_K(t_n) b_K)_K \right\}, \\ & \leq \varepsilon_1^2 \left\{ a(e(t_n), f_K(t_n) b_K) - \sum_{j=1}^{n_v} \beta_j a(e_{z^j}(t_n), f_K(t_n) b_K) \right. \\ & \quad \left. + (f_K(t_n) - f(t_n), f_K(t_n) b_K)_K \right\}. \end{aligned}$$

The Cauchy-Schwarz inequality implies,

$$\begin{aligned} \|f_K(t_n)\|_{L^2(K)}^2 & \leq \varepsilon_1^2 \left\{ \|e(t_n)\|_K + |\phi'(0)| \cdot \|e_z(t_n)\|_{n_v, K} \right\} \|f_K(t_n) b_K\|_K \\ & \quad + \varepsilon_1^2 \|f_K(t_n) - f(t_n)\|_{L^2(K)} \|f_K(t_n) b_K\|_{L^2(K)}. \end{aligned}$$

We now use the inverse estimate (2.4.19) in conjunction with the continuity of the bilinear form  $a(\cdot, \cdot)$ , to get  $\|f_K(t_n) b_K\|_K \leq C_a^{1/2} \varepsilon_2 h_K^{-1} \|f_K(t_n)\|_{L^2(K)}$ , where  $C_a$  is the continuity constant of  $a(\cdot, \cdot)$ . Then,

$$\begin{aligned} \|f_K(t_n)\|_{L^2(K)}^2 & \leq C_a^{1/2} \varepsilon_1^2 \varepsilon_2 h_K^{-1} \left\{ \|e(t_n)\|_K + |\phi'(0)|^{1/2} \cdot \|e_z(t_n)\|_{n_v, K} \right\} \|f_K(t_n)\|_{L^2(K)} \\ & \quad + \varepsilon_1^2 \|f_K(t_n) - f(t_n)\|_{L^2(K)} \|f_K(t_n)\|_{L^2(K)}, \end{aligned}$$

where we have also used  $\|f_K(t_n)b_K\|_{L^2(K)} \leq \|f_K(t_n)\|_{L^2(K)}$ . Then using Young's inequality (1.6.1) and multiplying by  $h_K^2$  we can rearrange the above so that we have,

$$h_K^2 \|f_K(t_n)\|_{L^2(K)}^2 \leq C_1 \left\{ \|e(t_n)\|_K^2 + |\phi'(0)| \cdot \|e_z(t_n)\|_{n_v, K}^2 + h_K^2 \|f_K(t_n) - f(t_n)\|_{L^2(K)}^2 \right\},$$

where we have absorbed the messy constant expression into  $C_1 > 0$ .

Step 2. Edge terms. Edges on the Neumann boundary are treated as follows. Define  $\tilde{R}_E := g_E(t_n) - \sigma(u_h(t_n), z_h(t_n))$ . Then since  $\tilde{R}_E$  is piecewise constant, the estimate (2.4.20) combined with the definition of the local form of the residual (5.1.6) leads to,

$$\begin{aligned} \|\tilde{R}_E\|_{L^2(E)}^2 &\leq \varepsilon_3^2 (\tilde{R}_E, \tilde{R}_E b_E)_E, \\ &\leq \varepsilon_3^2 \left\{ (R_E, \tilde{R}_E b_E)_E + (g_E(t_n) - g(t_n), \tilde{R}_E b_E)_E \right\}, \\ &\leq \varepsilon_3^2 \left\{ a(e(t_n), \tilde{R}_E b_E) - \sum_{K \in \omega_E} (R_K, \tilde{R}_E b_E)_K - \sum_{j=1}^{n_v} \beta_j a(e_m(t_n), \tilde{R}_E b_E) \right. \\ &\quad \left. + (g_E(t_n) - g(t_n), \tilde{R}_E b_E)_E \right\}. \end{aligned}$$

Once again, the Cauchy-Schwarz inequality implies that,

$$\begin{aligned} \|\tilde{R}_E\|_{L^2(E)}^2 &\leq \varepsilon_3^2 \left\{ \|e(t_n)\|_{\omega_E} \|b_E \tilde{R}_E\|_{\omega_E} + \sum_{K \in \omega_E} \|R_K\|_{L^2(K)} \|b_E \tilde{R}_E\|_{L^2(K)} \right. \\ &\quad \left. + |\phi'(0)| \cdot \|e_z(t_n)\|_{n_v, \omega_E} \|b_E \tilde{R}_E\|_{\omega_E} \right. \\ &\quad \left. + \|g_E(t_n) - g(t_n)\|_{L^2(E)} \|b_E \tilde{R}_E\|_{L^2(E)} \right\}. \end{aligned}$$

The bubble function inequalities (2.4.20) and (2.4.22) lead to,

$$\begin{aligned} \|\tilde{R}_E\|_{L^2(E)}^2 &\leq \varepsilon_3^2 \left\{ \varepsilon_4 h_E^{-1/2} \|e_u(t_n)\|_{\omega_E} + \sum_{K \in \omega_E} \varepsilon_5 h_E^{1/2} \|R_K\|_{L^2(K)} \right. \\ &\quad \left. + \varepsilon_4 h_E^{-1/2} |\phi'(0)| \cdot \|e_z(t_n)\|_{n_v, \omega_E} + \|g_E(t_n) - g(t_n)\|_{L^2(E)} \right\} \|\tilde{R}_E\|_{L^2(E)}. \end{aligned}$$

Applying Young's inequality as before, though this time multiplying by  $h_E$  and absorbing the constants into  $C_2 > 0$ , we can rearrange to get, for  $E \subset \Gamma_N$ ,

$$\begin{aligned} h_E \|\tilde{R}_E\|_{L^2(E)}^2 &\leq C_2 \left\{ \|e_u(t_n)\|_{\omega_E}^2 + |\phi'(0)| \cdot \|e_z(t_n)\|_{n_v, \omega_E}^2 \right. \\ &\quad \left. + h_E \|g_E(t_n) - g(t_n)\|_{L^2(E)}^2 + \sum_{K \in \omega_E} h_K^2 \|f(t_n)\|_{L^2(K)}^2 \right\}. \end{aligned}$$

Using the previous estimate above for the element based term  $h_K^2 \|f(t_n)\|_{L^2(K)}^2$  and absorbing the extra constant factors into  $C_2$ , we have,

$$h_E \|\tilde{R}_E\|_{L^2(E)}^2 \leq C_2 \left\{ \|e_u(t_n)\|_{\omega_E}^2 + |\varphi'(0)| \cdot \|e_z(t_n)\|_{n_v, \omega_E}^2 + h_E \|g_E(t_n) - g(t_n)\|_{L^2(E)}^2 + \sum_{K \in \omega_E} h_K^2 \|f_K(t_n) - f(t_n)\|_{L^2(K)}^2 \right\}.$$

For internal edges the proof is similar to that of the Neumann edges above, and that given in the proof of theorem (2.4.6). Following the same pattern and absorbing the constants into  $C_3 > 0$  results in,

$$h_E \|R_E\|_{L^2(E)}^2 \leq C_3 \left\{ \|e_u(t_n)\|_{\omega_E}^2 + |\varphi'(0)| \cdot \|e_z(t_n)\|_{n_v, \omega_E}^2 + \sum_{K \in \omega_E} h_K^2 \|f_K(t_n) - f(t_n)\|_{L^2(K)}^2 \right\}.$$

From the definition of  $\eta_K^2$ , we can now use the three inequalities above for the element, boundary edge and internal edges, and the result follows with  $C_{u, \text{eff}}^2 = \max\{C_1, C_2, C_3\}$ .  $\square$

## 5.2 Internal variables

In this section we consider the *a posteriori* error analysis of the approximation to the internal variables. The aim is to derive computable upper and lower bounds on the error  $e_z(t) = z(t) - z_h(t)$ . The internal variable problem involves both space and time discretisations so ideally we would like to determine error indicators that represent errors originating from the different parts of the discretisation. Once we derive upper bounds, we must also show that the estimators are efficient. The concept of efficiency is clear for spatial discretisations, typically it means to show, as we have done in sections 2.4 and 5.1 that the local error indicator for a given element is bounded above by the error in a covering neighbourhood of the element. This implies that the error indicator provides local information regarding the error. For discretisations involving both space and time, it is more common in the literature to show that the *a posteriori* error estimators are bounded above by optimal order *a priori* error estimates.

Considering what we want in a space and time adaptive algorithm - that at each time step, calculate the error due to the time discretisation, adjusting the time step until the time

error is smaller than a given tolerance, followed by an adaptive procedure for the spatial discretisation until a similar tolerance is reached - leads us to the following requirements of space and time *a posteriori* error indicators. Firstly, we would like to be able to separate the error effects of the two discretisations, and furthermore, that altering the spatial discretisation at a fixed time point has no undesirable effect on the time error.

In the remainder of this section we derive an *a posteriori* error estimate that is made up of two terms, one measuring the spatial error, the other the temporal error. We then carry out some numerical tests and evaluate the estimators in light of the above remarks. From the result of the previous section, we require only the energy norm of the error in the internal variables at a given time  $t_n$ . Since we have a general interest in methods for space time problems, we take a more general approach than is required showing that the upper bounds are suitable for errors in linear functionals of the solution, however, once again we revert to energy norms when we look to show efficiency.

The approach of this section is similar to that of the last, we start with the residual and show that the error in a linear functional satisfies a representation in terms of the residual acting upon the solution to a suitably defined dual problem. After confirming the strong stability of the dual solution we proceed to derive upper bounds. Let  $I = [0, t_n]$ , with  $t_n$  not necessarily equal to the final time  $T$ .

**Definition 5.2.1.** The residual of the approximation of problem (3.2.16) by (3.3.6) is,

$$\langle R(z_h), v \rangle := \int_I r(t; v) - (\partial_t z_h + M z_h, v)_{n_v} dt. \quad (5.2.1)$$

We introduce the dual problem: Find  $\chi \in W^{1,q}(I; V^{n_v})$  such that,

$$\int_I (w, -\partial_t \chi + M \chi)_{n_v} dt = \int_I \phi(t; w) dt, \quad \forall w \in L^p(I; V^{n_v}), \quad (5.2.2)$$

$$\chi(t_n) = \psi. \quad (5.2.3)$$

Note that the dual problem is a backward problem, with  $t$  running from  $t_n$  to zero.

We now derive a representation of the error in terms the dual problem and the residual defined in (5.2.1).

**Lemma 5.2.1.** *Let  $\chi$  be the solution of problem (5.2.2) with (5.2.3), then the error in the finite element approximation to the internal variables satisfies,*

$$(e_z(t_n), \psi)_{n_v} + \int_I \phi(t; e_z) dt = \langle R(z_h), \chi \rangle. \quad (5.2.4)$$

*Proof.* Setting  $w = e_z$  in the dual problem, integration by parts over  $I$  implies,

$$\begin{aligned} \int_I \phi(t; e_z) &= \int_I (e_z, -\partial_t \chi + M\chi)_{n_v} dt, \\ &= -(e_z(t_n), \chi(t_n))_{n_v} + (e_z(0), \chi(0))_{n_v} + \int_I (\partial_t e_z + M e_z, \chi)_{n_v} dt, \\ &= -(e_z(t_n), \chi(t_n))_{n_v} + \langle R(z_h), \chi \rangle, \end{aligned}$$

where we have used the property that  $e_z(0) = 0$ . Rearranging the above and using condition (5.2.3) implies the result.  $\square$

In terms of deriving an error estimate for  $e_z(t_n)$  we can consider the dual problem with  $\phi = 0$  and  $\psi = e_z(t_n) \| \| e_z(t_n) \| \|^{-1}$ . In this case, (5.2.4) becomes,

$$\| \| e_z(t_n) \| \|_{n_v} = \langle R(z_h), \chi \rangle, \quad (5.2.5)$$

in which case the following lemma provides the required strong stability estimate.

**Lemma 5.2.2.** *The solution of the dual problem (5.2.2) with  $\phi = 0$  and (5.2.3) with  $\psi = e_z(t_n) \| \| e_z(t_n) \| \|^{-1}$  satisfies,*

$$\| \chi \|_{L^1(I; V^{n_v})} \leq S_0(t_n), \quad (5.2.6)$$

$$\| \partial_t \chi \|_{L^1(I; V^{n_v})} \leq S_1(t_n). \quad (5.2.7)$$

where,

$$S_0(t_n) := \left( \frac{1 - e^{-m_s t_n}}{m_s} \right) \quad (5.2.8)$$

$$S_1(t_n) := m_l S_0(t_n), \quad (5.2.9)$$

and  $m_s$  and  $m_l$  are the smallest and largest eigenvalues of  $M$ .

*Proof.* Define  $\hat{\chi}(t) := \chi(t_n - t)$ , then the equation under the time integral in the dual problem (5.2.2) can be written as,

$$(v, \partial_t \hat{\chi} + M \hat{\chi})_{n_v} = 0, \quad \forall v \in V^{n_v}. \quad (5.2.10)$$

Choosing  $v = \hat{\chi}$  it follows that,

$$(\hat{\chi}, \partial_t \hat{\chi} + M \hat{\chi})_{n_v} = \frac{1}{2} \frac{d}{dt} (\hat{\chi}, \hat{\chi})_{n_v} + (\hat{\chi}, M \hat{\chi})_{n_v}.$$

From (3.2.12) it follows that,

$$\frac{1}{2} \frac{d}{dt} (\hat{\chi}, \hat{\chi})_{n_v} + m_s (\hat{\chi}, \hat{\chi})_{n_v} \leq 0.$$

Set  $p(t) := (\hat{\chi}(t), \hat{\chi}(t))_{n_v}$ , then  $p(t)$  satisfies the inequality,

$$\frac{dp}{dt} \leq -2m_s p(t), \quad (5.2.11)$$

and Gronwall's inequality implies,

$$p(t) \leq e^{-2m_s t} p(0). \quad (5.2.12)$$

Replacing  $p(t)$  with  $\|\hat{\chi}(t)\|_{V^{n_v}}^2$ , since  $\|\hat{\chi}(0)\|_{V^{n_v}} = 1$ , taking square roots implies that,

$$\|\hat{\chi}(t)\|_{V^{n_v}} \leq e^{-m_s t}.$$

Integrating over  $[0, t_n]$  and the fact that  $\|\hat{\chi}\|_{L^1(I; V^{n_v})} = \|\chi\|_{L^1(I; V^{n_v})}$  implies (5.2.8). To show (5.2.9), we use the fact that  $\|\partial_t \chi(t)\| \leq \|M^{1/2} \chi(t)\| \leq m_l \|\chi(t)\|$  together with (5.2.8).  $\square$

The stability result assures us that the accumulating error as we march through time does not grow in any undesirable way. In fact, the stability factors are  $O(1)$  as  $t \rightarrow \infty$ .

The aim now is to localise the residual to give an indication of the source and size of the error, be it the spatial or temporal discretisation. Furthermore we would like to localise these indicators to time intervals for the temporal indicator, and space mesh elements for the spatial indicator. The question is how do we derive these separate quantities from the residual? For the spatial indicator we proceed as we have before in sections 2.4 and 5.1 by integrating by parts over  $\Omega$  so that we can use the optimal order quasi-interpolation error estimates. For the temporal indicator it is less clear about whether to integrate by parts over  $\Omega$  or not. The main point is that we want the right form of the residuals, so that we can then use interpolation error estimate (2.2.24) to build powers of  $k_i$  into the indicator, and the use of that estimate does not require that we integrate by parts over  $\Omega$ .

Under the assumption that the mesh only gets finer or stays the same at each step, there are no problems integrating by parts over  $\Omega$ . However, if we wish to relax this assumption and allow coarsening when moving forward in time, there is a slight difficulty due to the fact that we will integrate the function  $z_h \in \mathbb{P}_1(I_i; V_h^{i-1} + V_h^i)^{n_v}$  over  $\Omega$ . In this instance we must cope with the fact that this function is in general defined over different meshes corresponding to the different discretisations at the different time points. Things are made much easier because, as discussed in section 2.2 (with more details in appendix A) we only consider hierarchic families of meshes. Let  $\mathcal{T}_{i-1,i}$  denote the union mesh, that can be visualised as the mesh formed by overlaying both  $\mathcal{T}_{i-1}$  and  $\mathcal{T}_i$ . With the idea of the union mesh, we can localise the residual in the following way.

**Lemma 5.2.3.** *The residual (5.2.4) has the following localised form,*

$$\langle R(z_h), v \rangle = \sum_{i=1}^n \int_{I_i} \sum_{K \in \mathcal{T}_{i-1,i}} \left\{ (\beta f, v)_{L^2(K)^{n_v}} + \sum_{E \in \mathcal{E}(K)} (R_E, v)_{L^2(E)^{n_v}} \right\} dt, \quad (5.2.13)$$

where,

$$R_E = \begin{cases} \frac{1}{2} \llbracket \mathbf{C}\epsilon(\partial_t z_h + M z_h) \rrbracket, & \text{on } E \subset \Omega, \\ \beta g(t) - \mathbf{C}\epsilon(\partial_t z_h + M z_h) n_E, & \text{on } E \subset \Gamma_N. \end{cases} \quad (5.2.14)$$

*Proof.* The result follows immediately by integration by parts of the right hand side of (5.2.1) and using the fact that  $z_h$  is piecewise linear and  $\mathbf{C}$  is constant in space. Elemental divergence terms are zero and jumps over internal edges are grouped together to form jumps.  $\square$

While the introduction of the union mesh allows the correct representation, it does not solve the problem of how to introduce the quasi-interpolation error estimates. The volume terms can be grouped together to form integrals over the elements in the current mesh  $\mathcal{T}_i$ , however the difficulty remains with jumps over edges that are not edges of the current mesh, i.e., those edges that were in  $\mathcal{T}_{i-1}$  but have been removed due to coarsening. To solve this problem we require the following trace result (see [18] and references therein for details).

**Proposition 5.2.4.** *If  $K$  is a bounded Lipschitz domain, then there exists a constant  $c_K$  which depends only on the shape of the domain  $K$  but not on its size  $h_K = \text{diam}(K)$ , such*

that, for all  $f \in W^{1,2}(K)$ ,

$$\|f\|_{L^2(\partial K)}^2 \leq c_K \left\{ h_K^{-1} \|f\|_{L^2(K)}^2 + h_K |f|_{W^{1,2}(K)}^2 \right\}. \quad (5.2.15)$$

We can use the result of proposition 5.2.4 in conjunction with the error estimate (2.2.18) to bound the  $L^2(E)$  norm of the error in approximation by the quasi-interpolant  $\mathcal{I}_i^1 : H^1(\Omega) \rightarrow \mathbb{P}_1(\mathcal{T}_i)$ , where  $E$  is not an edge of the triangulation  $\mathcal{T}_i$  but bisects an element  $K \in \mathcal{T}_i$ . Consider an element  $K \in \mathcal{T}_i$  made up of two smaller elements  $K_1$  and  $K_2$  from  $\mathcal{T}_{i-1}$ , with  $K = K_1 \cup K_2$  and observe that  $u - \mathcal{I}_i^1 u \in W^{1,2}(K)$  implies  $u - \mathcal{I}_i^1 u \in W^{1,2}(K_i)$ ,  $i = 1, 2$ . Proposition (5.2.4) can be applied on  $K_1$  and  $K_2$ , then averaging the inequalities, it follows from the mesh regularity assumption,  $h_{K_1} \leq ch_{K_2}$  and the fact that  $h_K \leq h_{K_1} + h_{K_2}$  that there is a constant depending on the shapes  $K_1$  and  $K_2$ , the mesh regularity and the error estimate (2.2.18) such that,

$$\|u - \mathcal{I}_i^1 u\|_{L^2(E)} \leq c_{K_1, K_2} h_K^{1/2} |u|_{W^{1,2}(K)}. \quad (5.2.16)$$

We now press on with developing the *a posteriori* error estimates. Let  $\pi_i^0$  denote the  $L^2(I_i)$  projection onto piecewise constant functions on the interval  $I_i$  introduced in section 2.2, and let  $\pi^0$  denote the global form such that  $\pi^0|_{I_i} = \pi_i^0$ . Similarly let  $\mathcal{I}^1$  denote the global quasi-interpolant  $\mathcal{I}^1|_{I_i} := \mathcal{I}_i^1$  where  $\mathcal{I}_i^1$  is the quasi-interpolant introduced in section 2.2. Since  $\pi_i^0 \mathcal{I}_i^1 \chi \in V_h^{i, n_v}$ , by Galerkin orthogonality, we can write,

$$\begin{aligned} \langle R(z_h), \chi \rangle &= \langle R(z_h), \chi - \pi^0 \mathcal{I}^1 \chi \rangle, \\ &= \langle R(z_h), \chi - \pi^0 \chi \rangle + \langle R(z_h), \pi^0 (I - \mathcal{I}^1) \chi \rangle. \end{aligned}$$

Based on the interpolation operators appearing in the right hand sides of the duality brackets, we will make the first term the time error and the second the spatial error. For the time error we choose the original form of the residual (5.2.1). For the spatial term we choose the local representation provided by lemma 5.2.3. In the following subsections we consider the spatial and temporal error indicators in turn. Since we have different terms for the space and time errors, we introduce an additional subscript on the reliability constant  $C_{\text{rel}}$ , so that  $C_{\text{s,rel}}$  and  $C_{\text{t,rel}}$  and are the reliability constants for the spatial and temporal error indicators respectively.

### 5.2.1 Spatial Residual

In this section we derive an upper and lower bound for the spatial indicator.

**Lemma 5.2.5.** *There exists a constant  $C_{s,\text{rel}}$  depending on the domain  $\Omega$ , the coercivity constant  $c_a$ , the minimum angle in the triangulation, the maximum number of overlapping element neighbourhoods, such that the residual (5.1.6) satisfies the bound,*

$$|\langle R(z_h), \pi_0(1 - \mathcal{I})\chi \rangle| \leq C_{s,\text{rel}} \max_{1 \leq i \leq n} \left( \sum_{K \in \mathcal{T}_i} \eta_K^2 \right)^{1/2} \|\chi\|_{L^1(I; V^{nv})},$$

where,

$$\eta_K^2 = h_K^2 \|\pi_i^0 \beta f\|_{L^2(K)}^2 + \sum_{E \in \mathcal{E}(K)} h_E \|\pi_i^0 R_E\|_{L^2(E)}^2 + \sum_{E \subset K, E \in \mathcal{E}(\mathcal{T}_{i-1})} h_K \|\pi_i^0 R_E\|_{L^2(E)}^2, \quad (5.2.17)$$

where,

$$R_E(t) = \begin{cases} \frac{1}{2} [\mathbf{C}\epsilon(\partial_t z_h(t) + M z_h(t))], & \text{on } E \in \Omega, \\ \beta g(t) - \mathbf{C}\epsilon(\partial_t z_h(t) + M z_h(t)) n_E, & \text{on } E \in \Gamma_N. \end{cases} \quad (5.2.18)$$

*Proof.* Starting with the representation (5.2.13), we can use the  $L^2$  projection  $\pi_i^0$  that  $(\pi_i^0 w, v)_{L^2(I_i)} = (\pi_i^0 w, \pi_i^0 v)_{L^2(I_i)}$  to make the expression constant over the interval  $I_i$ . Then since the expression is constant over  $I_i$  carry out the integration to introduce a factor of  $k_i$ .

This gives,

$$\begin{aligned} \langle R(z_h), \pi^0(I - \mathcal{I}^1)\chi \rangle &= \sum_{i=1}^n k_i \sum_{K \in \mathcal{T}_{i-1,i}} \left\{ (\pi_i^0 \beta f, \pi_i^0(I - \mathcal{I}_i^1)\chi)_{L^2(K)^{nv}} \right. \\ &\quad \left. + \sum_{E \in \mathcal{E}(K)} (\pi_i^0 R_E, \pi_i^0(I - \mathcal{I}_i^1)\chi)_{L^2(E)^{nv}} \right\}. \end{aligned}$$

Split the jump terms into those on the current mesh and those on the mesh at the previous time step,

$$\begin{aligned} \langle R(z_h), \pi^0(I - \mathcal{I}^1)\chi \rangle &= \sum_{i=1}^n k_i \sum_{K \in \mathcal{T}_i} \left\{ (\pi_i^0 \beta f, \pi_i^0(I - \mathcal{I}_i^1)\chi)_{L^2(K)^{nv}} \right. \\ &\quad \left. + \sum_{E \in \mathcal{E}(K)} (\pi_i^0 R_E, \pi_i^0(I - \mathcal{I}_i^1)\chi)_{L^2(E)^{nv}} \right. \\ &\quad \left. + \sum_{E \subset K, E \in \mathcal{E}(\mathcal{T}_{i-1})} (\pi_i^0 R_E, \pi_i^0(I - \mathcal{I}_i^1)\chi)_{L^2(E)^{nv}} \right\}. \end{aligned}$$

We can now use the interpolation estimates (2.2.18), (2.2.19), together with the estimate (5.2.16) for the edge terms from the previous mesh. Then iterated application of the Cauchy-Schwarz inequality results in,

$$|\langle R(z_h), \pi^0(I - \mathcal{I}^1)\chi \rangle| \leq C_{s,\text{rel}} \sum_{i=1}^n k_i \left( \sum_{K \in \mathcal{T}_i} \eta_K^2 \right)^{1/2} \|\pi_i^0 \chi\|_{n_v}. \quad (5.2.19)$$

Then using  $\|\pi_i^0 \chi\|_{n_v} \leq k_i^{-1} \|\chi\|_{L^1(I_i; V^{n_v})}$  we get,

$$|\langle R(z_h), \pi^0(I - \mathcal{I}^1)\chi \rangle| \leq C_{s,\text{rel}} \max_{1 \leq i \leq n} \left( \sum_{K \in \mathcal{T}_i} \eta_K^2 \right)^{1/2} \|\chi\|_{L^1(I; V^{n_v})}. \quad (5.2.20)$$

□

We now consider the efficiency of the spatial error estimator. The statement of the theorem and pattern of proof is similar to that in theorem 2.4.6 and theorem 5.1.5, where we utilise the bubble functions presented in section 2.4. However, there is a notable difference in the term that provides the upper bound.

**Theorem 5.2.6.** *There exists a constant  $C_{s,\text{eff}}$  depending on the domain  $\Omega$ , the coercivity constant  $c_a$ , the minimum angle in the triangulation, the maximum number of overlapping element neighbourhoods and the degree of mesh coarsening, such that the local error indicator  $\eta_K \in \mathcal{T}_i$  of (5.2.17) satisfies the bound,*

$$\begin{aligned} \eta_K^2 \leq C_{s,\text{eff}} \left\{ \left\| \left\| \frac{e_z(t_i) - e_z(t_{i-1})}{k_i} \right\| \right\|_{n_v}^2 + \|\pi_i^0 e\|_{n_v} \right. \\ \left. + h_K^2 \|\pi_i^0(\beta f_K - \beta f)\|_{L^2(K)} + h_E \|\pi_i^0(\beta g_E - \beta g)\|_{L^2(E)} \right\}. \end{aligned}$$

*Proof.* We will break the proof into three steps which follows closely the proofs provided in theorems 2.4.6 and 5.1.5.

Step 1. Volume terms. Let  $f_K$  denote the  $L^2(K)$  projection onto piecewise constants on  $K$ , then by (2.4.18),

$$\begin{aligned} \|\beta \pi_i^0 f_K\|_{L^2(K)^{n_v}}^2 &\leq \varepsilon_1^2 (\beta \pi_i^0 f_K, b_K \beta \pi_i^0 f_K)_{L^2(K)^{n_v}}, \\ &\leq \varepsilon_1^2 (\beta \pi_i^0 f, b_K \pi_i^0 f_K)_{L^2(K)^{n_v}} + \varepsilon_1^2 (\beta \pi_i^0 (f_K - f), b_K \beta \pi_i^0 f_K)_{L^2(K)^{n_v}}. \end{aligned}$$

Using the definition of the residual its local form (5.2.13) together with the definition of  $\pi_i^0$  we have the relationship,

$$\left( \frac{e_z(t_i) - e_z(t_{i-1})}{k_i} + M\pi_i^0 e_z, b_K \right)_{n_v} = (\pi_i^0 \beta f, b_K)_{n_v}, \quad b_K \in \mathbb{P}_0(I_i). \quad (5.2.21)$$

Using (5.2.21) to replace the first term on the right hand side in the previous upper bound, the Cauchy-Schwarz inequality leads to ,

$$\begin{aligned} & \|\beta \pi_i^0 f_K\|_{L^2(K)^{n_v}}^2 \\ & \leq \varepsilon_1^2 \left\{ \left( \frac{e_z(t_i) - e_z(t_{i-1})}{k_i} + M\pi_i^0 e_z, b_K \beta \pi_i^0 f_K \right)_{n_v} + (\pi_i^0 \beta (f_K - f), b_K \beta \pi_i^0 f_K)_{L^2(K)^{n_v}} \right\}, \\ & \leq \varepsilon_1^2 \left\{ \left( \left\| \frac{e_z(t_i) - e_z(t_{i-1})}{k_i} \right\|_{n_v}^2 + \|M\pi_i^0 e_z\|_{n_v} \right) \|b_K \pi_i^0 \beta f_K\|_{n_v} \right. \\ & \quad \left. + \|\pi_i^0 \beta (f_K - f)\|_{L^2(K)} \|b_K \beta \pi_i^0 f_K\|_{L^2(K)^{n_v}} \right\}. \end{aligned}$$

Then using the inverse estimate  $\|v\|_{n_v} \leq C_a^{1/2} \varepsilon_2 h_K^{-1} \|v\|_{L^2(K)^{n_v}}$ , and (2.4.19)

$$\begin{aligned} & \|\beta \pi_i^0 f_K\|_{L^2(K)^{n_v}}^2 \\ & \leq \varepsilon_1^2 \left\{ C_a^{1/2} \varepsilon_2 h_K^{-1} \left( \left\| \frac{e_z(t_i) - e_z(t_{i-1})}{k_i} \right\|_{n_v}^2 + \|M\pi_i^0 e_z\|_{n_v} \right) \|\pi_i^0 \beta f_K\|_{L^2(K)^{n_v}} \right. \\ & \quad \left. + \|\pi_i^0 \beta (f_K - f)\|_{L^2(K)} \|\beta \pi_i^0 f_K\|_{L^2(K)^{n_v}} \right\}. \end{aligned}$$

Multiplying through by  $h_K$  and using Young's inequality to get squares on the right hand side and rearranging implies that there is a constant  $C_1$  independent of  $h_K$  such that,

$$\begin{aligned} & h_K \|\beta \pi_i^0 f_K\|_{L^2(K)^{n_v}} \\ & \leq C_1 \left\{ \left\| \frac{e_z(t_i) - e_z(t_{i-1})}{k_i} \right\|_{n_v, K} + \|\pi_i^0 e_z\|_{n_v, K} + \|\pi_i^0 \beta (f_K - f)\|_{L^2(K)^{n_v}} \right\}. \end{aligned}$$

Step 2. Internal edges. We first consider the case of internal edges  $E \in \mathcal{T}_i$ , and then for  $E \in \mathcal{T}_{i-1}$ . Recall that on internal edges,  $R_E$  is piecewise constant in space, therefore, (2.4.20) gives,

$$\begin{aligned} \|\pi_i^0 R_E\|_{L^2(E)^{n_v}}^2 & \leq \varepsilon_3^2 (\pi_i^0 R_E, \pi_i^0 R_E)_{L^2(E)^{n_v}} \\ & \leq \varepsilon_3^2 \left\{ \left( \frac{e_z(t_i) - e_z(t_{i-1})}{k_i} + M\pi_i^0 e_z, b_E \pi_i^0 R_E \right)_{n_v, \omega_E} \right. \\ & \quad \left. - (\pi_i^0 \beta R_E, \pi_i^0 R_E)_{L^2(\omega_E)^{n_v}} \right\}, \end{aligned}$$

where we have used the residual equation in the second step. We can now proceed as in the previous similar proofs, utilising (2.4.21) and (2.4.22) to get to the stage where there is a constant  $C_2$  independent of  $h_K$ , such that,

$$h_E^{1/2} \|\pi_i^0 R_E\|_{L^2(E)^{n_v}} \leq C_2 \left\{ \left\| \left\| \frac{e_z(t_i) - e_z(t_{i-1})}{k_i} \right\| \right\|_{n_v} + \|\pi_i^0 e_z\|_{n_v, \omega_E} + \sum_{K \in \omega_E} h_K \|\pi_i^0 \beta(f_K - f)\|_{L^2(K)^{n_v}} \right\}.$$

For the case of an edge from the previous mesh, the argument is no different, since for the hierarchic mesh structure,  $E \in \mathcal{T}_{i-1}$  implies that  $\omega_E = K \in \mathcal{T}_i$ . Therefore the same inequality holds for those edges.

Step 3. Neumann edges. As in previous proofs, we introduce the piecewise constant approximation  $g_E$  to  $g$  on a given edge  $E$ , and as in step 2 of the proof of theorem 5.1.5 introduce the approximation to  $R_E$ , denoted by  $\hat{R}_E := \beta g_E - \mathbf{C}\epsilon(\partial_t z_h + M z_h)$ . Using a similar argument to the above step, it follows that there is a constant  $C_3$  independent of  $h_E$  such that for  $E \in \Gamma_N$ ,

$$h_E^{1/2} \|\pi_i^0 R_E\|_{L^2(E)^{n_v}} \leq C_3 \left\{ \left\| \left\| \frac{e_z(t_i) - e_z(t_{i-1})}{k_i} \right\| \right\|_{n_v} + \|\pi_i^0 e_z\|_{n_v} + h_E \|\pi_i^0 \beta(g_E - g)\|_{L^2(E)^{n_v}} + h_K \|\pi_i^0 \beta(f_K - f)\|_{L^2(K)^{n_v}} \right\}.$$

To show the result, we combine the estimates from the three steps.  $\square$

### 5.2.2 Temporal residual

In this section we show upper and lower bounds for the temporal indicator. However, while these results ensure that the indicator converges to zero as the error does, there is a problem in that the indicator is non-zero when the approximation is exact. We deal with this issue in the next chapter.

**Lemma 5.2.7.** *There exists a constant  $C_{t,\text{rel}}$  depending on the domain  $\Omega$  and the coercivity constant  $c_a$  such that the temporal residual satisfies the upper bound,*

$$|\langle R(z_h), \chi - \pi^0 \chi \rangle| \leq C_{t,\text{rel}} \max_{1 \leq i \leq n} k_i \xi_i \|\partial_t \chi\|_{L^1(I; V^{n_v})}. \quad (5.2.22)$$

where

$$\xi_i := \|(1 - \pi_0^i)\boldsymbol{\beta}f\|_{L^\infty(I_i; L^2(\Omega)^{nv})} + \|(1 - \pi_0^i)\boldsymbol{\beta}g\|_{L^\infty(I_i; L^2(\Gamma_N)^{nv})} + \|(1 - \pi_0^i)Mz_h\|_{L^\infty(I_i; V^{nv})}. \quad (5.2.23)$$

*Proof.* Recalling the definition of the residual (5.2.1), we have,

$$\begin{aligned} \langle R(z_h), (1 - \pi_0)\chi \rangle &= \int_I r(t; (1 - \pi_0)\chi) - (\partial_t z_h + Mz_h, (1 - \pi_0)\chi)_{n_v} dt, \\ &= \sum_{i=1}^n \int_{t_{i-1}}^{t_i} (\boldsymbol{\beta}f, (1 - \pi_0^i)\chi)_{L^2(\Omega)^{nv}} \\ &\quad + (\boldsymbol{\beta}g, (1 - \pi_0^i)\chi)_{L^2(\Gamma_N)^{nv}} \\ &\quad - (\partial_t z_h + Mz_h, (1 - \pi_0^i)\chi)_{n_v} dt. \end{aligned}$$

Note that the term involving the time derivative of  $z_h$  is zero by the orthogonality of  $\pi_i^0$ . Exploiting this property once again by subtracting off terms involving  $\pi_i^0\boldsymbol{\beta}f$ ,  $\pi_i^0\boldsymbol{\beta}g$  and  $\pi_i^0Mz_h$  gives,

$$\begin{aligned} |\langle R(z_h), (1 - \pi_0)\chi \rangle| &= \sum_{i=1}^n \int_{t_{i-1}}^{t_i} ((1 - \pi_0^i)\boldsymbol{\beta}f, (1 - \pi_0^i)\chi)_{L^2(\Omega)^{nv}} \\ &\quad + ((1 - \pi_0^i)\boldsymbol{\beta}g, (1 - \pi_0^i)\chi)_{L^2(\Gamma_N)^{nv}} \\ &\quad - ((1 - \pi_0^i)Mz_h, (1 - \pi_0^i)\chi)_{n_v} dt. \end{aligned}$$

Applying the Cauchy-Schwarz inequality, the trace theorem (theorem 2.4.2), and using the fact that  $\|v\|_{H^1(\Omega)} \leq c_a^{-1}\|v\|$ , together with the error estimates for  $\pi_0^i$  (2.2.24) results in,

$$\begin{aligned} \langle R(z_h), (1 - \pi_0)\chi \rangle &\leq C_{t,\text{rel}} \sum_{i=1}^n k_i \left\{ \|(1 - \pi_0^i)\boldsymbol{\beta}f\|_{L^\infty(I_i; L^2(\Omega)^{nv})} \right. \\ &\quad + \|(1 - \pi_0^i)\boldsymbol{\beta}g\|_{L^\infty(I_i; L^2(\Gamma_n)^{nv})} \\ &\quad \left. + \|(1 - \pi_0^i)Mz_h\|_{L^\infty(I_i; V^{nv})} \right\} \|\partial_t \chi\|_{L^1(I_i; V^{nv})} \end{aligned}$$

where  $C_{t,\text{rel}} = c_1^{-1} \max\{1, C_{\Gamma_N}\}$  where  $C_{\Gamma_N}$  is the constant from the trace theorem. A final application of Hölders inequality proves the result.  $\square$

To show that the estimator given in lemma 5.2.7 is efficient, we show that it satisfies an optimal order *a priori* upper bound.

**Theorem 5.2.8.** *There exists a constant  $C_{t,\text{eff}}$  independent of the discretisation parameters such that the temporal error estimator of theorem 5.2.7 satisfies,*

$$k_i \xi_i \leq C_{t,\text{eff}} \left\{ k_i \|e_z(t_i)\|_{n_v} + k_i \|e_z(t_{i-1})\|_{n_v} + k_i^2 \|\partial_t z\|_{L^\infty(I_i; V^{n_v})} + k_i^2 \|\partial_t f\|_{L^\infty(I_i; L^2(\Omega))} + k_i^2 \|\partial_t g\|_{L^\infty(I_i; L^2(\Gamma_N))} \right\}$$

*Proof.* For the data terms, we use estimate the error estimate for  $\pi_i^0$  (2.2.22). For the term involving  $z_h$ , we again use error estimate (2.2.22),

$$\begin{aligned} \|(1 - \pi_0^i) M z_h\|_{L^\infty(I_i; V^{n_v})} &\leq k_i \|M \partial_t z_h\|_{L^\infty(I_i; V^{n_v})}, \\ &\leq \|M\| \|z_h(t_i) - z_h(t_{i-1})\|_{n_v}, \\ &\leq \|M\| \|z_h(t_i) - z(t_i) + z(t_i) - z(t_{i-1}) + z(t_{i-1}) - z_h(t_{i-1})\|_{n_v}, \\ &\leq \|M\| \left\{ \|e_z(t_i)\|_{n_v} + \|z(t_i) - z(t_{i-1})\|_{n_v} + \|e_z(t_{i-1})\|_{n_v} \right\}. \end{aligned}$$

Then since,

$$\begin{aligned} \|z(t_i) - z(t_{i-1})\|_{n_v} &= \left\| \int_{t_{i-1}}^{t_i} \partial_t z(t) dt \right\|_{n_v}, \\ &\leq \int_{t_{i-1}}^{t_i} \|\partial_t z(t)\|_{n_v} dt, \\ &\leq k_i \|\partial_t z\|_{L^\infty(I_i; V^{n_v})}, \end{aligned}$$

the result follows.  $\square$

The previous result shows that the temporal error indicator converges like  $O(k^2)$  subject to the regularity of the solution and the data. We bring the results of this section together to give the following *a posteriori* error estimate for the internal variable problem.

**Theorem 5.2.9.** *There exist constants  $C_{s,\text{rel}}$  and  $C_{t,\text{rel}}$ , depending on the domain  $\Omega$ , the Neumann boundary  $\Gamma_N$ , the coercivity constant  $c_a$ , the minimum angle in the triangulation, the maximum number of overlapping element neighbourhoods, such that the error in the finite element approximation to the internal variables satisfies,*

$$|(e_z(T), \psi)_{n_v} + \int_I \phi(t; e_z) dt| \leq \Theta(z_h) \|\chi\|_{L^1(I; V^{n_v})} + \Phi(z_h) \|\partial_t \chi\|_{L^1(I; V^{n_v})}, \quad (5.2.24)$$

where,

$$\Theta(z_h) = C_{s,\text{rel}} \max_{1 \leq i \leq n} \left( \sum_{K \in \mathcal{T}_i} \eta_K^2 \right),$$

$$\Phi(z_h) = C_{t,\text{rel}} \max_{1 \leq i \leq n} k_i \xi_i,$$

where  $\eta_K$  and  $\xi_i$  are given by (5.2.17) and (5.2.23) respectively.

*Proof.* From lemma 5.2.1, Galerkin orthogonality implies that,

$$\begin{aligned} (e_z(T), \psi)_{n_v} + \int_I \phi(t; e_z) dt &= \langle R(z_h), \chi - \pi^0 \mathcal{I}^1 \chi \rangle, \\ &= \langle R(z_h), \chi - \pi^0 \chi \rangle + \langle R(z_h), \pi^0 (I - \mathcal{I}) \chi \rangle. \end{aligned}$$

Taking absolute values, the triangle inequality and lemmas 5.2.5 and 5.2.7 imply the result.  $\square$

### 5.3 *A posteriori* estimates

We now bring together the results of this chapter. Recall from section 2.4 definition (2.4.1) of the data oscillation. Then extending the definition to deal with the mesh at time  $t_n$ , we have,

$$\text{osc}_h(\mathcal{T}_n)^2 := \sum_{K \subset \omega} \left\{ h_K^2 \|f(t_n) - f_K(t_n)\|_{L^2(K)}^2 + \sum_{E \in \mathcal{E}(K) \cap E(\Gamma_N)} h_E \|g(t_n) - g_E(t_n)\|_{L^2(E)}^2 \right\}. \quad (5.3.1)$$

Combining the upper and lower bounds of this section, and incorporating the data oscillation, we have the following result for the displacement.

**Corollary 5.3.1.** *The energy norm of the error in the displacement  $e_u(t_n) = u(t_n) - u_h(t_n)$  satisfies the relationship,*

$$\begin{aligned} C_{u,\text{eff}}^{-1} \left\{ \left( \sum_{K \in \mathcal{T}_i} \eta_K^2 \right)^{1/2} - |\varphi(0)| \cdot \|e_z(t_n)\|_{n_v} - \text{osc}_h(\mathcal{T}_n) \right\} \\ \leq \|e_u(t_n)\| \\ \leq C_{u,\text{rel}} \left\{ \left( \sum_{K \in \mathcal{T}_i} \eta_K^2 \right)^{1/2} + |\varphi(0)| \cdot \|e_z(t_n)\|_{n_v} \right\}. \end{aligned}$$

*Proof.* Given the definition of the data oscillation then the result follows by combining the results of theorems 5.1.4 and 5.1.5.  $\square$

This result shows that if the oscillation in the data and the error in the internal variables are sufficiently small, then we have computable upper and lower bounds that reflect the true value of the error in the displacement. We have already discussed how data oscillation can be controlled in section 2.4. We have the following result which deals with the problem of controlling the error in the internal variables.

**Corollary 5.3.2.** *There exist constants  $C_{s,\text{rel}}$  and  $C_{t,\text{rel}}$ , depending on the domain  $\Omega$ , the Neumann boundary  $\Gamma_N$ , the coercivity constant  $c_a$ , the minimum angle in the triangulation, the maximum number of overlapping element neighbourhoods, such that the error in the finite element approximation to the internal variables satisfies,*

$$\|e_z(t_n)\|_{n_v} \leq S_0(t_n)\Theta(z_h) + S_1(t_n)\Phi(z_h), \quad (5.3.2)$$

where  $\Theta(z_h)$  and  $\Phi(z_h)$  are given in theorem 5.2.9.

*Proof.* Combining theorem 5.2.9 with the stability result of lemma 5.2.2 gives the result.  $\square$

## 5.4 Adaptive algorithms

In chapter 2 we presented the standard Solve-Estimate-Refine (SER) process typical of AFEM for stationary problems. We also presented an adaptive time stepping algorithm for an ODE in time. Adaptive algorithms for problems in space and time are discussed in [62], and a classification of space time adaptive algorithms is presented (see references in [62]). We will follow the approach they refer to as *Implicit strategy A*. This strategy uses the mesh from the previous time step as the starting mesh on the new time step, The SER algorithm is then applied until the error criteria on the current mesh are met.

In this section we look at how these algorithms together with the *a posteriori* error estimate of the previous sections can be combined for space and time problems. We will assume that the computation is to proceed under the condition that,

$$\|u(t_n) - u_h(t_n)\| \leq \text{GTOL}, \quad (5.4.1)$$

for some final time  $t_n$ . There are a number of issues to tackle, mainly involving the partition of the tolerance. First we consider the algorithm we would use if we were only solving the internal variable problem, which in itself is a space time problem. Recalling the *a posteriori* error estimate of theorem 5.2.9 and the result of corollary 5.3.2, we would typically want to satisfy the condition,

$$S_0(t_n)\Theta(z_h) + S_1(t_n)\Phi(z_h) \leq \text{GTOL}. \quad (5.4.2)$$

Setting  $\text{LTOL} = \max\{S_0(t_n)^{-1}, S_1(t_n)^{-1}\}\text{GTOL}$ , then at each time point we can ensure that the condition will be satisfied by making the sum of the error indicators less than  $\text{LTOL}$ . Recall that  $\Theta(z_h)$  is the spatial error indicator and  $\Phi(z_h)$  is the temporal error indicator. We propose the following space and time adaptive algorithm for generating solutions for the internal variable problems.

Space and time adaptive algorithm:

---

1. Start with parameters  $\delta_1 \in (0, 1)$ ,  $\delta_2 > 1$ ,  $\theta \in (0, 1)$  and  $\gamma \in (0, 1)$ . Set  $t = 0$ ,  $\text{LTOL} = \max\{S_0(t_n)^{-1}, S_1(t_n)^{-1}\}\text{GTOL}$ ,  $k_{\text{old}} = T$ , and set  $i = 0$ .
2. Do:
  - i) Set  $k = k_{\text{old}}$ ,  $\mathcal{T}_{i+1} = \mathcal{T}_i$
  - ii) Calculate  $z_h^{i+1}$  and  $\Theta(z_h)$  and  $\Phi(z_h)$ .
  - iii) While  $\Theta(z_h) + \Phi(z_h) > \text{LTOL}$ :
    - a) [Refine timestep] While  $\Phi(z_h) > \gamma\text{LTOL}$  and  $\Theta(z_h) + \Phi(z_h) > \text{LTOL}$ :
      - Set  $k = \delta_1 k$ .
      - Calculate  $z_h^{i+1}$ ,  $\Theta(z_h)$  and  $\Phi(z_h)$ .
    - b) [Refine spatial mesh] While  $\Theta(z_h) > (1 - \gamma)\text{LTOL}$ :
      - Adapt  $\mathcal{T}_{i+1}$ .
      - Calculate  $z_h^{i+1}$  and  $\Theta$ .
  - iv) If  $\Phi(z_h) < \theta\text{LTOL}$ , then  $k = \delta_2 k$ .
  - v) Set  $z_h^i = z_h^{i+1}$ ,  $t = t + k$ ,  $i = i + 1$ .

while  $t \leq T$ .

---

The algorithm comprises of two main parts, corresponding to the temporal and spatial refinement steps. For each new proposal solution, we calculate the full error estimate. If the local error condition is not satisfied, then we refine the time step until the temporal error term is less than a fraction (determined by  $\gamma$ ) of the permitted tolerance or the total error satisfies the criterion. Once satisfied, we move to the spatial error and apply an SER algorithm. However, when adapting the spatial discretisation we must consider both refinement and coarsening of the mesh. In fact it has been shown in [13] and [72] that we should consider coarsening even for stationary problems.

We described the marking strategy for mesh refinements in section 2.4. When coarsening the mesh, additional terms will be added to the error estimate at the next time step (see lemma 5.2.5). Care must be taken not to enter into a coarsen-refine oscillation where the error introduced by coarsening leads to a refinement at the next step and so on. To mark elements that are to be coarsened, we form as large a set as possible of elements that make up a fraction  $\theta_c$  of the error. The full marking algorithm is given below.

Marking Strategy:

---

Given  $\theta_r$ ,  $\theta_o$  and  $\theta_c$ ,  $0 < \theta_o, \theta_r, \theta_c < 1$ :

1. Construct a subset of elements  $\mathcal{M}_R \subset \mathcal{T}_i$  such that,

$$\eta(\mathcal{M}_R) \geq \theta_r \eta(\mathcal{T}_i).$$

2. Enlarge  $\mathcal{M}_R$  so that,

$$\text{osc}_h(\mathcal{M}_R) \geq \theta_o \text{osc}_h(\mathcal{T}_i).$$

3. Construct a subset of elements  $\mathcal{M}_C \subset (\mathcal{T}_i \setminus \mathcal{M}_R)$ , such that,

$$\eta(\mathcal{M}_C) \leq \theta_c \eta(\mathcal{T}_i).$$

---

As discussed in chapter 2, there is a growth condition on the time steps so that if the condition is “over ”satisfied, the time steps are increased. From this perspective we can see that it is essential for the adaptive algorithm to allow mesh coarsening, otherwise the mesh would simply become finer after each time step that required a refinement. Furthermore, it is a central requirement that  $\Phi(z_h)$  is stable with respect to spatial refinements and  $\Theta(z_h)$  is stable with respect to step size adjustments. There are two choices for generating the displacement, either step by step alongside the internal variables, or we could in principle solve the internal variables over the whole domain, and then follow up by solving the displacement at the desired time points. In the first instance we perform an SER loop for the displacement for the displacement and use the resulting mesh as the initial guess for the mesh at the next time step for the internal variables.

## 5.5 Numerical results

We now computationally examine the properties of the error estimators. We consider the same group of test problems as those presented in section 4.4, where we had model problems that were linear in either space or time. The utility of that choice here is that we can examine the error indicators when the error present in the approximation is exactly that which it should measure, for example for the spatial indicator, the approximation will be exact in time, so the only error in the approximation will be due to the spatial approximation.

The constants that appear in the estimates are all set to 1. This is because we want to focus on the convergence properties of the estimators rather than definite values of the upper bounds, however the determination of the constants is an important issue in practice. Johnson and Hansbo [44] consider calculating the constants in their *a posteriori* bounds for linear elasticity, and there are a few publications dedicated to determining the constants for quasi-interpolation error estimates. A further constant that we would require is that from Korn’s inequality (lemma 2.4.1). In summary, the constants would have to calculate and references to papers dealing with the calculation are:

1. Coercivity constant  $c_a$ , [46] and [47].

2. Quasi-interpolation constants  $c_K$ ,  $c_E$  and  $c_{K_1, K_2}$ , [43], [19] and [79].

### 5.5.1 Spatial convergence

To demonstrate the spatial convergence of the estimators we consider the problem described in section 4.4 with displacement given by (4.4.1) and internal variable given by (4.4.2). With a fixed number of timesteps we examine the behaviour of the estimators as the spatial discretisation is refined. From figure 5.1 we can see that the indicator of the displacement error converges at the same rate as the true error.

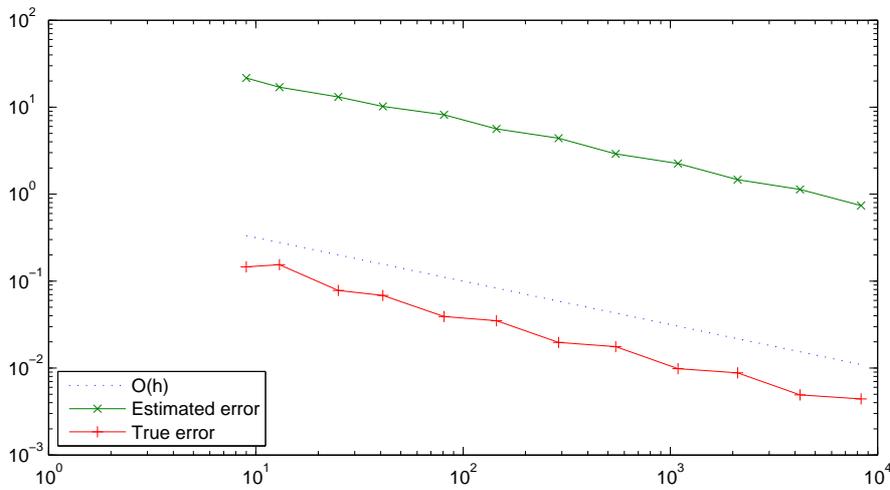


Figure 5.1: Log-log plot of the error estimate and the maximum energy norm of the error in the displacement against the degrees of freedom.

Similarly, from figure 5.2 we can see that the spatial indicator of the internal variables  $\Theta(z_h)$  behaves like the true error. However, it is interesting to note that the temporal indicator displays convergence as the spatial discretisation is refined. This is a deficiency of the temporal indicator that we return to in the discussion at the end of this chapter and deal with in the next.

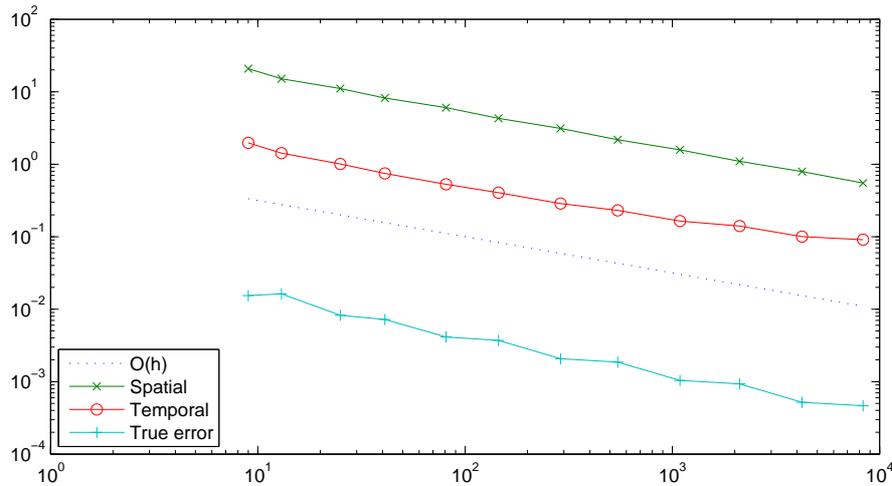


Figure 5.2: Log-log plot of the spatial and temporal error indicators and the maximum energy norm of the error in the internal variable against the degrees of freedom.

### 5.5.2 Temporal convergence

To demonstrate the temporal convergence of the estimators we consider the problem described in section 4.4 with displacement given by (4.4.5) and internal variable given by (4.4.6). With a fixed spatial mesh we examine the behaviour of the estimators as the number of timesteps increases over the fixed time interval  $I$ . From figure 5.3 we can see that the temporal indicator of the internal variables error converges at the same rate as the true error. Also note that the indicator for the spatial error in the displacement is zero. Our attention however, is drawn to the increasing quantity, the spatial indicator for the internal variables. Since the solution is piecewise linear in space, the jump terms should be zero, as are the volume contributions, the only possible contribution should be from the Neumann edges. However, a numerical check confirms that this growth appears in the edge terms.

So is the spatial indicator for the internal variables unstable? It certainly appears to grow as  $k \rightarrow 0$ . However this is in the case where there is no spatial error and its value is very small, and in theory the jumps are zero. If we consider a problem that is not exact in space, then the growth as  $k \rightarrow 0$  does not occur, therefore we conclude that the spatial indicator is not unstable and the growth is simply due to round-off error.

Dofs	$\max_i \ e_u(t_i)\ $	$\max_i \eta_{\mathcal{T}_i}$	$\max_i \ e_z(t_i)\ _{n_v}$	$\Theta(z_h)$
9	1.4570e-01	2.1632e+01	1.5360e-02	2.0814e+01
25	7.8142e-02	1.3130e+01	8.2375e-03	1.1042e+01
81	3.9190e-02	8.1723e+00	4.1312e-03	6.0365e+00
289	1.9668e-02	4.3877e+00	2.0733e-03	3.1191e+00
1089	9.8377e-03	2.2458e+00	1.0371e-03	1.5774e+00
4225	4.9190e-03	1.1318e+00	5.1854e-04	7.9182e-01

Table 5.1: Convergence in space of the true error and the spatial error indicators for the displacement and internal variables.

Dofs	$\max_i \ e_u(t_i)\ $	$\max_i \eta_{\mathcal{T}_i}$	$\max_i \ e_z(t_i)\ _{n_v}$	$\Theta(z_h)$	$\Phi(z_h)$
5	2.3539e-05	4.1244e-14	5.8846e-05	8.6804e-15	5.1531e-02
9	5.7156e-06	3.4350e-14	1.4289e-05	1.6460e-14	1.4742e-02
17	1.4446e-06	3.9690e-14	3.6114e-06	3.1340e-14	3.8368e-03
33	3.6370e-07	7.8891e-14	9.0924e-07	7.0600e-14	9.6916e-04
65	9.0905e-08	9.3981e-14	2.2726e-07	1.1846e-13	2.4295e-04
129	2.2747e-08	9.4151e-14	5.6867e-08	8.5692e-13	6.0785e-05
257	5.6866e-09	9.5532e-14	1.4216e-08	1.3219e-12	1.5200e-05
513	1.4216e-09	9.3990e-14	3.5541e-09	3.9165e-12	3.8003e-06

Table 5.2: Convergence in time of the true error and the spatial and temporal error indicators for the displacement and internal variables.

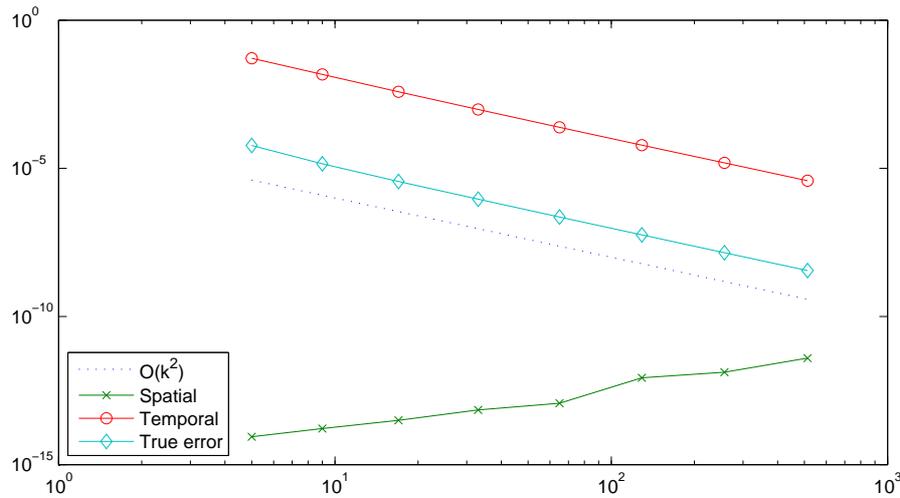


Figure 5.3: Log-log plot spatial and temporal error estimators and the maximum energy norm of the error in the internal variable against the degrees of freedom.

## 5.6 Summary

In this chapter we carried out an *a posteriori* error analysis of the finite element approximation presented in chapter 3. After showing that the estimators are reliable and efficient we described an adaptive algorithm and presented some numerical confirmation of the theoretical results.

However, while theorem 5.2.9 provides a reliable and efficient computable upper bound of the approximation error, the numerical results show that it is questionable how informative it is regarding the actual error, and therefore how useful is it for driving an adaptive routine. The flaw lies in the fact that if the approximation was exact, the temporal error indicator,  $\Phi(z_h)$  is still non-zero.

Where does the analysis go wrong? Essentially we didn't give the estimator a chance to be zero because there is no measure of difference, unlike in the ODE (2.5) case where we use the difference of the approximate ODE with the true problem data. We were unable to do this in our case since we did not integrate by parts over  $\Omega$ . However, integrating by parts means introducing jump terms over edges, which are far from desirable in the temporal

indicator, since these are non-zero when there is error due to the spatial discretisation, irrespective of whether there is error due to the time approximation.

Perhaps if we had tried an alternative splitting to,

$$\chi - \pi^0 \mathcal{I} \chi = \chi - \pi^0 \chi + \pi^0 (I - \mathcal{I}) \chi, \quad (5.6.1)$$

we might be able to manipulate the residual to get something more like what we want. Since  $\pi^0 \mathcal{I}^1 = \mathcal{I}^1 \pi^0$ , why not choose,

$$\chi - \pi^0 \mathcal{I} \chi = \chi - \mathcal{I}^1 \chi + \mathcal{I}^1 (I - \pi^0) \chi, \quad (5.6.2)$$

instead? Indeed we would get a slightly different estimator if we were to pursue the alternative splitting (5.6.2). However, it turns out that the resulting spatial error indicator depends on the temporal error. In fact splitting (5.6.2) was the first splitting used in this study, though on discovery of the dependence on the temporal error, the alternative (5.6.1) was used, though as we mention, still flawed. We presented the *a posteriori* estimator above because in fact the spatial indicator is the desired one. However, we remark that the alternative splitting (5.6.2) does lead to the correct temporal indicator but we postpone presenting it until the next chapter, with the bias in presentation due to chronological factors. In the next section we consider a more refined analysis of how space and time errors interact to determine what the error indicators should look like, and in turn present an appropriate estimator. We find that the splitting,

$$\chi - \pi^0 \mathcal{I} \chi = \pi^0 (I - \mathcal{I}) \chi + (I - \pi^0) \mathcal{I} + (I - \pi^0) (I - \mathcal{I}) \chi, \quad (5.6.3)$$

is a more natural splitting, and this allows us to derive error indicators that are exact in the sense that they are zero for the exact solution.

## Chapter 6

# Exact *a posteriori* error estimators

In the previous chapter we presented an *a posteriori* error estimate which satisfied most of the criteria typically required of a useful *a posteriori* error estimator. However, there was an issue with the temporal residual in that when the approximation was exact, i.e., when the true solution belonged to the approximating space, the error estimator was non-zero. In this chapter we tackle this weakness by designing exact error indicators, where by exact, we mean that they are zero when the approximation is exact, or equivalently when the true solution belongs to the approximation space. This is not a new idea, since the definition of efficiency of an estimator (2.3.16) implies that an efficient estimator is exact.

In this chapter, we first consider a standard  $L^2$  projection in space and time as a model scenario. We show that the error in such an approximation can and should be decomposed into three terms, one corresponding to error in space, one to error in time, and a third which contains the mixed effects. We then cast *a posteriori* error estimators in the same light, and derive a new *a posteriori* error estimate for the internal variable problem. However while the new error estimate has certain desirable theoretical properties, there are practical issues involved. Those issues are discussed fully in the numerical experiments section. We close the chapter with a summary of the *a posteriori* error estimates that we have considered and make recommendations for practical implementation.

## 6.1 Space-time projections

The problems outlined in the previous chapter are due to the fact that we have not adequately separated the errors due to the different parts, spatial and temporal, of the discretisation. To see this, we consider a simple space time  $L^2$  projection. Let  $\Omega \subset \mathbb{R}^d$  be the spatial domain and  $I \subset \mathbb{R}$  the time domain. Suppose that given a function  $f \in L^2(\Omega \times I; \mathbb{R})$ , we want to compute an approximation  $f_h$  by projecting  $f$  onto a finite dimensional subspace of  $S_h(\Omega \times I; \mathbb{R}) \subset L^2(\Omega \times I; \mathbb{R})$ . The problem can be posed: Find  $f_h \in S_h$  such that,

$$(f_h, s)_{\Omega \times I} = (f, s)_{\Omega \times I}, \quad \forall s \in S_h(\Omega \times I; \mathbb{R}).$$

This is an orthogonal projection in space and time. The approximation  $f_h$  satisfies the best approximation property,

$$\|f - f_h\|_{L^2(\Omega \times I)} = \inf_{s \in S_h} \|f - s\|_{L^2(\Omega \times I)}.$$

Now suppose that we required the error to satisfy an explicit error bound  $\epsilon$ , i.e.,

$$\|f - f_h\|_{L^2(\Omega \times I)} \leq \epsilon.$$

Since we know  $f$  and  $f_h$  then we can compute this quantity and keep adapting the space  $S_h$  until the bound is satisfied. However there is no information telling us whether we should increase the dimension of the finite dimensional space with regard to the spatial or temporal discretisation.

For the discretisation, partition the time interval  $I = [0, T]$  into  $N$  subintervals  $\{I_i\}_{i=1}^N$  with  $I_i := (t_{i-1}, t_i]$  and let  $\mathcal{T}_i$  denote a partition of the domain  $\Omega$  into simplices. Let  $P_i^1 : L^2(\Omega) \mapsto \mathbb{P}_1(\mathcal{T}_i)$  be the  $L^2$ -projection onto the space of continuous piecewise linear functions defined by,

$$(P_i^1 f, v)_{\Omega} = (f, v)_{\Omega}, \quad \forall v \in \mathbb{P}_1(\mathcal{T}_i).$$

Let  $\pi_i^0 : L^2(I_i) \mapsto \mathbb{P}_0(I_i)$  be the  $L^2$ -projection onto piecewise constants defined by,

$$(\pi_i^0 f, v)_{I_i} = (f, v)_{I_i}, \quad \forall v \in \mathbb{P}_0(I_i).$$

The discrete problem is then to find  $f_h \in \oplus_{i=1}^N \mathbb{P}_0(I_i; \mathbb{P}_1(\mathcal{T}_i))$  such that,

$$(f_h, v)_{\Omega \times I} = (f, v)_{\Omega \times I}, \quad \forall v \in \oplus_{i=1}^N \mathbb{P}_0(I_i; \mathbb{P}_1(\mathcal{T}_i)).$$

The restriction of the problem to a particular time step gives: Find  $f_h^i \in \mathbb{P}_0(I_i; \mathbb{P}_1(\mathcal{T}_i))$  such that,

$$(f_h^i, v)_{\Omega \times I_i} = (f, v)_{\Omega \times I_i}, \quad \forall v \in \mathbb{P}_0(I_i; \mathbb{P}_1(\mathcal{T}_i)).$$

The approximation can be written as the composition of the two  $L^2$  projections introduced earlier, i.e.,  $f_h^i = \pi_i^0 P_i^1 f$ . If we want to measure the error incurred due to the different parts of the approximation, then consider the difference,

$$f - f_h.$$

If the functions agreed spatially, then this term does not have to be zero, since  $f$  is variable in time, while  $f_h$  is piecewise constant in time. This leads to the idea that to capture the spatial error we should consider the difference,

$$\pi_i^0 f - \pi_i^0 P_i^1 f, \tag{6.1.1}$$

since now if  $f = P_i^1 f$  then 6.1.1 is zero. Similarly for the time error we should consider the difference,

$$P_i^1 f - \pi_i^0 P_i^1 f, \tag{6.1.2}$$

which will be non-zero only if the functions disagree temporally. We summarise these observations in the following proposition.

**Proposition 6.1.1.** *The error in the  $L^2(I; L^2(\Omega))$  projection onto  $\mathbb{P}_0(I_i; \mathbb{P}_1(\mathcal{T}_i))$ , where  $f_h$  is the approximation to  $f$  satisfies,*

$$\|f - f_h\|_{L^2(\Omega \times I_i)}^2 = \mathcal{E}_i^s + \mathcal{E}_i^t + \mathcal{E}_i^m, \tag{6.1.3}$$

where,

$$\mathcal{E}_i^s := \|\pi_i^0(1 - P_i^1)f\|_{L^2(\Omega \times I_i)}^2, \tag{6.1.4}$$

$$\mathcal{E}_i^t := \|(1 - \pi_i^0)P_i^1 f\|_{L^2(\Omega \times I_i)}^2, \tag{6.1.5}$$

$$\mathcal{E}_i^m := \|(1 - \pi_i^0)(1 - P_i^1)f\|_{L^2(\Omega \times I_i)}^2. \tag{6.1.6}$$

*Proof.* The proof follows by decomposing the space  $L^2(I_i; L^2(\Omega))$  using the projections. Since  $L^2(\Omega) = P_i^1 L^2(\Omega) \oplus (I - P_i^1)L^2(\Omega)$  and  $L^2(I_i) = \pi_i^0 L^2(I_i) \oplus (I - \pi_i^0)L^2(I_i)$ , we can

decompose  $L^2(I_i; L^2(\Omega))$  in the following way,

$$\begin{aligned} L^2(I_i; L^2(\Omega)) &= L^2(I_i; P_i^1 L^2(\Omega)) \oplus L^2(I_i; (1 - P_i^1) L^2(\Omega)) \\ &= \pi_i^0 L^2(I_i; P_i^1 L^2(\Omega)) \oplus (1 - \pi_i^0) L^2(I_i; P_i^1 L^2(\Omega)) \\ &\quad \oplus \pi_i^0 L^2(I_i; (1 - P_i^1) L^2(\Omega)) \oplus (1 - \pi_i^0) L^2(I_i; (1 - P_i^1) L^2(\Omega)) \end{aligned}$$

From this, we recognise that the function  $f$  can be expressed as,

$$f = \pi_i^0 P_i^1 f + (1 - \pi_i^0) P_i^1 f + \pi_i^0 (1 - P_i^1) f + (1 - \pi_i^0) (1 - P_i^1) f.$$

Therefore we have the representation,

$$f - \pi_i^0 P_i^1 f = (1 - \pi_i^0) P_i^1 f + \pi_i^0 (1 - P_i^1) f + (1 - \pi_i^0) (1 - P_i^1) f.$$

Taking the  $L^2(\Omega \times I_i)$  norm gives,

$$\begin{aligned} \|f - \pi_i^0 P_i^1 f\|_{L^2(\Omega \times I)}^2 &= \|(1 - \pi_i^0) P_i^1 f\|_{L^2(\Omega \times I_i)}^2 \\ &\quad + \|\pi_i^0 (1 - P_i^1) f\|_{L^2(\Omega \times I_i)}^2 \\ &\quad + \|(1 - \pi_i^0) (1 - P_i^1) f\|_{L^2(\Omega \times I_i)}^2 \\ &\quad + 2((1 - \pi_i^0) P_i^1 f, (1 - \pi_i^0) (1 - P_i^1) f)_{L^2(\Omega \times I_i)} \\ &\quad + (\pi_i^0 (1 - P_i^1) f, (1 - \pi_i^0) (1 - P_i^1) f)_{L^2(\Omega \times I_i)}. \end{aligned}$$

The last two terms of the above expression are zero by orthogonality.  $\square$

We can see from (6.1.3) that if the approximation is exact in the spatial variables, in the sense that the function being approximated is contained inside the approximating space, then  $\mathcal{E}_i^s = \mathcal{E}_i^m = 0$ . Similarly then for an approximation which is exact in time, we have  $\mathcal{E}_i^t = \mathcal{E}_i^m = 0$ .

## 6.2 Error indicators

We now consider the derivation of an *a posteriori* error estimator based on the idea of exactness introduced in the previous section. A key observation of the previous section is that to measure the spatial error, we had to look in the temporally discrete space, i.e.,

the spatial error is the result of a projection of the full error onto the temporally discrete space. The roles of the spaces and projections are reversed for the temporal error. Recall the internal variable problem,

$$\int_I (\partial_t z(t) + Mz(t), v)_{n_v} dt = \int_I r(t; v) dt, \quad \forall v \in L^2(I; V^{n_v}), \quad (6.2.1)$$

$$z(0) = 0. \quad (6.2.2)$$

Motivated by the previous section, we note that when we have the representation,

$$(e_z(t_n), \psi)_{n_v} + \int_I \phi(t; e_z) dt = \langle R(z_h), \chi \rangle, \quad (6.2.3)$$

rather than subtracting off the residual evaluated on the discrete space (Galerkin orthogonality), we can decompose the dual solution using the splitting,

$$\chi - \pi^0 \mathcal{I} \chi = \pi^0 (I - \mathcal{I}) \chi + (I - \pi^0) \mathcal{I} + (I - \pi^0) (I - \mathcal{I}) \chi. \quad (6.2.4)$$

Galerkin orthogonality then implies,

$$\begin{aligned} \langle R(z_h), \chi \rangle &= \langle R(z_h), \pi^0 (I - \mathcal{I}) \chi \rangle \\ &\quad + \langle R(z_h), (I - \pi^0) \mathcal{I} \chi \rangle \\ &\quad + \langle R(z_h), (I - \pi^0) (I - \mathcal{I}) \chi \rangle. \end{aligned}$$

We recognise that we have already bounded the first term on the right hand side as it is the spatial residual of chapter 5. We call the second and third terms on the right hand side the temporal and mixed residuals respectively.

Before progressing to derive upper bounds on the temporal and mixed residuals, we first consider what we want them to look like by examining the exact problem on the discrete spaces. This analysis confirms that we have the correct form for the error indicators.

In the following lemma we consider equations that the true solution satisfies when the approximation is exact in space, projected onto the temporally discrete space.

**Lemma 6.2.1.** *If the solution  $z$  of (3.2.16) is such that  $z(t_i) \in V_h^i$ , then it satisfies the*

relationships,

$$-\operatorname{div}\mathbf{C}\epsilon(\pi_i^0(\partial_t z + Mz)) = \pi_i^0\beta f, \quad \forall K \in \mathcal{T}_i, \quad (x, t) \in K \times I_i, \quad (6.2.5)$$

$$\frac{1}{2}[\![\mathbf{C}\epsilon(\pi_i^0(\partial_t z + Mz))]\!] = 0, \quad \forall E \subset \Omega, \quad (x, t) \in E \times I_i, \quad (6.2.6)$$

$$\pi_i^0\beta g - \mathbf{C}\epsilon(\pi_i^0(\partial_t z + Mz))n_E = 0, \quad \forall E \subset \Gamma_N, \quad (x, t) \in E \times I_i. \quad (6.2.7)$$

*Proof.* The proof follows by inspecting the weak form of the internal variable problem (3.2.16) on the discrete space  $\mathbb{P}_0(I_i; V^{n_v})$ . Integrating by parts over  $\Omega$  the equation,

$$\int_{I_i} (\partial_t z + Mz, v)_{n_v} = \int_{I_i} r(t; v) dt, \quad \forall v \in \mathbb{P}_0(I_i; V^{n_v}), \quad (6.2.8)$$

we get the spatially strong form,

$$\begin{aligned} \sum_{i=1}^N \int_{I_i} \sum_{K \in \mathcal{T}_i} \left\{ (\pi_i^0\beta f - \operatorname{div}\mathbf{C}\epsilon(\pi_i^0(\partial_t z + Mz)), v)_{L^2(K)^{n_v}} \right. \\ \left. + \sum_{E \in \mathcal{E}(K)} (\pi_0^i R_E, v)_{L^2(E)^{n_v}} \right\} dt = 0, \quad \forall v \in \mathbb{P}_0(I_i; V^{n_v}), \end{aligned}$$

where we have used  $(\pi_i^0 w, v)_{L^2(I_i)} = (w, v)_{L^2(I_i)}$ ,  $\forall v \in \mathbb{P}_0(I_i)$ . Now since  $v$  is arbitrary in  $\mathbb{P}_0(I_i; V^{n_v})$ , we can choose  $v$  to be a bubble function (2.4.16) so that it is non-zero only on the interior of the element  $K$ . This implies (6.2.5). Then given (6.2.5), (6.2.6) and (6.2.7) follow by choosing  $v$  first to be zero on the boundary of  $\Omega$ , and then an arbitrary  $V^{n_v}$  function.  $\square$

The previous lemma states local projected forms of the strong form of the original equations and helps us characterise the spatial error in the approximation. We note that the individual quantities (6.2.5), (6.2.6) and (6.2.7) are exactly those that appear in the spatial error indicator presented in lemma 5.2.5.

### 6.2.1 Temporal residual

Before embarking on the same route for the temporal error as we have above for the spatial, we first introduce a special operator which is analogous to the discrete Laplacian in other situations. However, we include the Neumann boundary condition in the mapping. Define,  $\operatorname{div}\mathbf{C}\epsilon : H^1(\Omega) \rightarrow V_h^i$  by,

$$((\operatorname{div}\mathbf{C}\epsilon)_h(w), v)_{L^2(\Omega)^{n_v}} := (\beta g, v)_{L^2(\Gamma_N)^{n_v}} - (\mathbf{C}\epsilon(w), \epsilon(v))_{L^2(\Omega)^{n_v}}, \quad \forall v \in V_h^{i, n_v}. \quad (6.2.9)$$

To determine the discrete divergence of a function  $w \in H^1(\Omega)$ , we solve the linear system,

$$L\mathbf{w}_h = \mathbf{g}, \quad (6.2.10)$$

where  $L$  is the mass matrix of the space  $V_h^{i,n_v}$ ,  $\mathbf{g}$  is the vector formed by sampling the right hand side of (6.2.9) at the basis functions  $\{\eta_i\}_{i=1}^{\dim V_h^{i,n_v}}$  of  $V_h^{i,n_v}$ , and  $\mathbf{w}_h$  is the representation of  $w$  with respect to the basis of  $V_h^{i,n_v}$ . We note however, that for  $w \in V_h^{i,n_v}$ , the system (6.2.10) becomes,

$$L\mathbf{w}_h = \mathbf{b} - A\mathbf{w}_h, \quad \mathbf{b} = (b_1, \dots, b_{\dim V_h^{i,n_v}})^T, \quad b_i = \int_{\Gamma_N} \beta g \cdot \eta_i \, dx. \quad (6.2.11)$$

where  $A$  is the stiffness matrix of the basis  $\{\eta_i\}_{i=1}^{\dim V_h^{i,n_v}}$ .

**Lemma 6.2.2.** *The solution  $z$  of (3.2.16) satisfies the relationships,*

$$-(\operatorname{div} \mathbf{C}\boldsymbol{\epsilon})_h(\partial_t z + Mz) = P_i^1 \boldsymbol{\beta} f, \quad (x, t) \in \Omega \times I_i. \quad (6.2.12)$$

*Proof.* Once again we inspect the weak form of the equations on the partially discrete space. Consider the equation,

$$\int_{I_i} (\partial_t z + Mz, v)_{n_v} = \int_{I_i} r(t; v) dt, \quad \forall v \in L^2(I_i; V_h^i). \quad (6.2.13)$$

Since  $v \in \mathbb{P}_1(\mathcal{T}_i)$ , we can introduce the  $L^2(\Omega)$  projection onto piecewise linear functions  $P_i^1$ ,

$$\int_{I_i} (\mathbf{C}\boldsymbol{\epsilon}(\partial_t z + Mz), \boldsymbol{\epsilon}(v))_{L^2(\Omega)^{n_v}} dt = \int_{I_i} (P_i^1 \boldsymbol{\beta} f, v)_{L^2(\Omega)^{n_v}} + (\boldsymbol{\beta} g, v)_{L^2(\Gamma_N)^{n_v}} dt. \quad (6.2.14)$$

Then, it follows from the definition of (6.2.9) that,

$$\begin{aligned} \int_{I_i} (-(\operatorname{div} \mathbf{C}\boldsymbol{\epsilon})_h(\partial_t z + Mz), v)_{L^2(\Omega)^{n_v}} + (\boldsymbol{\beta} g, v)_{L^2(\Gamma_N)^{n_v}} \\ = \int_{I_i} (P_i^1 \boldsymbol{\beta} f, v)_{L^2(\Omega)^{n_v}} + (\boldsymbol{\beta} g, v)_{L^2(\Gamma_N)^{n_v}}. \end{aligned}$$

Rearranging, implies,

$$\int_{I_i} ((\operatorname{div} \mathbf{C}\boldsymbol{\epsilon})_h(\partial_t z + Mz) + P_i^1 \boldsymbol{\beta} f, v)_{L^2(\Omega)^{n_v}} = 0. \quad (6.2.15)$$

Then since  $v$  is arbitrary in  $L^2(I_i; V_h^i)$ , the result follows.  $\square$

Once again, all we have done here is show that the exact solution satisfies the projected strong form of the equations. We now look to deriving an upper bound on a temporal residual, the form of which is influenced by lemma 6.2.2.

**Lemma 6.2.3.** *There exists a constant  $C_{t,\text{rel}}$  depending on the domain  $\Omega$  and the coercivity constant  $c_a$  such that the temporal residual satisfies the upper bound,*

$$|\langle R(z_h), (I - \pi^0)\mathcal{I}\chi \rangle| \leq C_{t,\text{rel}} \max_{1 \leq i \leq N} k_i \xi_i \|\chi\|_{L^1(I; V^{nv})}. \quad (6.2.16)$$

where

$$\xi_i := \|P_i^1 \beta f + (\text{div} \mathbf{C} \epsilon)_h(\partial_t z_h + M z_h)\|_{L^\infty(I_i; L^2(\Omega)^{nv})} \quad (6.2.17)$$

*Proof.* From the representation of the residual, lemma 6.2.2 implies,

$$\langle R(z_h), (I - \pi^0)\mathcal{I}\chi \rangle = \sum_{i=0}^N \int_{I_i} (P_i^1 \beta f + (\text{div} \mathbf{C} \epsilon)_h(\partial_t z_h + M z_h), (I - \pi_i^0)\mathcal{I}_i^1 \chi)_{L^2(\Omega)} dt.$$

Using Cauchy-Schwarz and the stability of  $\mathcal{I}_i^1$  in the  $L^2(\Omega)$  norm (2.2.17) implies that,

$$\begin{aligned} & |\langle R(z_h), (I - \pi^0)\mathcal{I}\chi \rangle| \\ & \leq C_{\mathcal{I}} \sum_{i=0}^N \int_{I_i} \|P_i^1 \beta f + (\text{div} \mathbf{C} \epsilon)_h(\partial_t z_h + M z_h)\|_{L^2(\Omega)^{nv}} \| (I - \pi_i^0)\chi \|_{n_v} dt. \end{aligned}$$

By Hölder's inequality, the error estimates for  $\pi_i^0$  (2.2.24) then give,

$$\begin{aligned} & |\langle R(z_h), (I - \pi^0)\mathcal{I}\chi \rangle| \\ & \leq C_{\mathcal{I}} \sum_{i=0}^N k_i \|P_i^1 \beta f + (\text{div} \mathbf{C} \epsilon)_h(\partial_t z_h + M z_h)\|_{L^p(I_i; L^2(\Omega)^{nv})} \cdot \|\chi\|_{L^1(I_i; V^{nv})} \end{aligned}$$

Then using Hölders inequality once more,

$$|\langle R(z_h), (I - \pi^0)\mathcal{I}\chi \rangle| \max_{1 \leq i \leq N} \leq k_i \xi_i \|\chi\|_{L^1(I; V^{nv})}. \quad (6.2.18)$$

□

The upper bound for the temporal residual as expressed above involves assembling the mass matrix and solving a linear system for each evaluation. This is obviously too expensive an estimator to consider for practical purposes.

### Alternative temporal residual

Out of curiosity we can consider what would happen if we were to use the localised form of the residual (5.2.13) to try to estimate the temporal error. We get the following.

**Proposition 6.2.4.** *There exists a constant  $C_{alt,rel}$  depending on the domain  $\Omega$ , the coercivity constant  $c_a$ , the minimum angle in the triangulation and the maximum number of overlapping element neighbourhoods, such that the error in the finite element approximation satisfies,*

$$|\langle R(z_h), (I - \pi^0)\mathcal{I}\chi \rangle| \leq C_{alt,rel} \max_{1 \leq i \leq n} k_i \left( \sum_{K \in \mathcal{T}_i} \eta_K^2 \right)^{1/2} \|\partial_t \chi\|_{L^1(I; V^{nv})} \quad (6.2.19)$$

$$\begin{aligned} \eta_K^2 &= \|(I - \pi_i^0)\beta f\|_{L^\infty(I_i; L^2(K)^{nv})}^2 \\ &\quad + \sum_{E \in \mathcal{E}(K)} h_K^{-1/2} \|(I - \pi_i^0)R_E\|_{L^\infty(I_i; L^2(E)^{nv})}^2 \\ &\quad + \sum_{E \subset K} h_K^{-1/2} \|(I - \pi_i^0)R_E\|_{L^\infty(I_i; L^2(E)^{nv})}^2 \end{aligned}$$

*Proof.* Starting with the representation (5.2.13), we can use the properties of  $\pi_i^0$ , to subtract off projected terms,

$$\begin{aligned} \langle R(z_h), (I - \pi^0)\mathcal{I}\chi \rangle &= \sum_{i=1}^N \int_{I_i} \sum_{K \in \mathcal{T}_{i-1,i}} \left\{ (\beta f, (I - \pi_i^0)\mathcal{I}_i^1 \chi)_{L^2(K)^{nv}} \right. \\ &\quad \left. + \sum_{E \in \mathcal{E}(K)} (R_E, (I - \pi_i^0)\mathcal{I}_i^1 \chi)_{L^2(E)^{nv}} \right\} dt, \\ &= \sum_{i=1}^N \int_{I_i} \sum_{K \in \mathcal{T}_{i-1,i}} \left\{ ((I - \pi_i^0)\beta f, (I - \pi_i^0)\mathcal{I}_i^1 \chi)_{L^2(K)^{nv}} \right. \\ &\quad \left. + \sum_{E \in \mathcal{E}(K)} ((I - \pi_i^0)R_E, (I - \pi_i^0)\mathcal{I}_i^1 \chi)_{L^2(E)^{nv}} \right\} dt. \end{aligned}$$

Combining the trace result (proposition 5.2.4) with the inverse estimate  $|v|_{W^{1,2}(K)} \leq ch_K^{-1} \|v\|_{L^2(K)}$ ,  $\forall v \in \mathbb{P}_1(K)$  we have  $\|v\|_{L^2(E)} \leq ch_K^{-1/2} \|v\|_{L^2(K)}$ . This can be used to get an upper bound

in the following way. The Cauchy-Schwarz inequality implies,

$$\begin{aligned}
& |\langle R(z_h), (I - \pi^0)\mathcal{I}\chi \rangle| \\
& \leq \sum_{i=1}^n \int_{I_i} \sum_{K \in \mathcal{T}_i} \left\{ \|(I - \pi_i^0)\beta f\|_{L^2(K)^{n_v}} \|(I - \pi_i^0)\mathcal{I}_i^1 \chi\|_{L^2(K)^{n_v}} \right. \\
& \quad + \sum_{E \in \mathcal{E}(K)} \|(I - \pi_i^0)R_E\|_{L^2(E)^{n_v}} \|(I - \pi_i^0)\mathcal{I}_i^1 \chi\|_{L^2(E)^{n_v}} \\
& \quad \left. + \sum_{E \subset K} \|(I - \pi_i^0)R_E\|_{L^2(E)^{n_v}} \|(I - \pi_i^0)\mathcal{I}_i^1 \chi\|_{L^2(E)^{n_v}} \right\} dt,
\end{aligned}$$

Now using the inverse estimate  $\|v\|_{L^2(E)} \leq ch_K^{-1/2}\|v\|_{L^2(K)}$  on the edge terms and pulling the constant through, we get,

$$\begin{aligned}
& |\langle R(z_h), (I - \pi^0)\mathcal{I}\chi \rangle| \\
& \leq c \sum_{i=1}^N \int_{I_i} \sum_{K \in \mathcal{T}_i} \left\{ \|(I - \pi_i^0)\beta f\|_{L^2(K)^{n_v}} \right. \\
& \quad + \sum_{E \in \mathcal{E}(K)} h_K^{-1/2} \|(I - \pi_i^0)R_E\|_{L^2(E)^{n_v}} \\
& \quad \left. + \sum_{E \subset K} h_K^{-1/2} \|(I - \pi_i^0)R_E\|_{L^2(E)^{n_v}} \right\} \|(I - \pi_i^0)\mathcal{I}_i^1 \chi\|_{L^2(K)^{n_v}} dt.
\end{aligned}$$

Using the stability of  $\mathcal{I}_i$  (2.2.17), the Cauchy-Schwarz inequality for the sum over the elements followed by the Hölder inequality for the time integral results in,

$$|\langle R(z_h), (I - \pi^0)\mathcal{I}\chi \rangle| \leq c \sum_{i=1}^N \left( \sum_{K \in \mathcal{T}_i} \eta_K^2 \right)^{1/2} \|(I - \pi_i^0)\chi\|_{L^1(I_i; V^{n_v})}.$$

Using the error estimate for  $\pi_i^0$  (2.2.22) and Hölders inequality, we get,

$$|\langle R(z_h), (I - \pi^0)\mathcal{I}\chi \rangle| \leq c \max_{1 \leq i \leq n} k_i \left( \sum_{K \in \mathcal{T}_i} \eta_K^2 \right)^{1/2} \|\partial_t \chi\|_{L^1(I; L^2(\Omega)^{n_v})}. \quad (6.2.20)$$

□

### 6.2.2 Mixed residual

We have bounded the spatial and temporal errors, we now bound the mixed term.

**Lemma 6.2.5.** *There exists a constant  $C_{m,rel}$  depending on the domain  $\Omega$ , the coercivity constant  $c_a$ , the minimum angle in the triangulation and the maximum number of overlapping element neighbourhoods, such that the error in the finite element approximation satisfies,*

$$|\langle R(z_h), (I - \pi^0)(1 - \mathcal{I})\chi \rangle| \leq C_{m,rel} \max_{1 \leq i \leq N} k_i \left( \sum_{K \in \mathcal{T}_i} \zeta_K^2 \right) \|\partial_t \chi\|_{L^1(I; V^{nv})} \quad (6.2.21)$$

where,

$$\begin{aligned} \zeta_K^2 &= h_K^2 \|(I - \pi_i^0)\beta f\|_{L^\infty(I_i; L^2(K)^{nv})}^2 \\ &+ \sum_{E \in \mathcal{E}(K)} h_E \|(1 - \pi_i^0)R_E\|_{L^\infty(I_i; L^2(E)^{nv})}^2 \\ &+ \sum_{E \subset K, E \in \mathcal{T}_{i-1}} h_K \|(1 - \pi_i^0)R_E\|_{L^\infty(I_i; L^2(E)^{nv})}^2. \end{aligned}$$

*Proof.* Starting with the representation (5.2.13), we can use the properties of  $\pi_i^0$ , to subtract off projected terms,

$$\begin{aligned} \langle R(z_h), (I - \pi^0)(1 - \mathcal{I})\chi \rangle &= \sum_{i=1}^N \int_{I_i} \sum_{K \in \mathcal{T}_{i-1,i}} \left\{ (\beta f, (I - \pi_i^0)(I - \mathcal{I}_i^1)\chi)_{L^2(K)^{nv}} \right. \\ &\quad \left. + \sum_{E \in \mathcal{E}(K)} (R_E, (I - \pi_i^0)(I - \mathcal{I}_i^1)\chi)_{L^2(E)^{nv}} \right\} dt, \\ &= \sum_{i=1}^N \int_{I_i} \sum_{K \in \mathcal{T}_{i-1,i}} \left\{ ((I - \pi_i^0)\beta f, (I - \pi_i^0)(I - \mathcal{I}_i^1)\chi)_{L^2(K)^{nv}} \right. \\ &\quad \left. + \sum_{E \in \mathcal{E}(K)} ((I - \pi_i^0)R_E, (I - \pi_i^0)(I - \mathcal{I}_i^1)\chi)_{L^2(E)^{nv}} \right\} dt. \end{aligned}$$

As in the earlier proofs, we use the Cauchy-Schwarz and Hölder inequalities together with the estimates (2.2.18), (2.2.19) and (5.2.16), together with the error estimate for  $\pi_i^0$  (2.2.24).  $\square$

We now present an alternative *a posteriori* error estimator to that of theorem 5.2.9 based on the splitting (6.2.4). The estimate contains three terms, one corresponding to the spatial error, which is the same as that given in theorem 5.2.9, another corresponding to the temporal error, and a third term that is a mixed spatial-temporal indicator.

**Theorem 6.2.6.** *The error in the internal variables satisfies the bound,*

$$\begin{aligned} |(e_z(t_n), \psi)_{n_v} + \int_I \phi(t; e_z) dt| \\ \leq \Theta(z_h) \|\chi\|_{L^1(I; V^{n_v})} + \Phi(z_h) \|\partial_t \chi\|_{L^1(I; V^{n_v})} + \Psi(z_h) \|\partial_t \chi\|_{L^1(I; V^{n_v})}, \end{aligned}$$

where,

$$\begin{aligned} \Theta(z_h) &= C_{s,\text{rel}} \max_{1 \leq i \leq N} \left( \sum_{K \in \mathcal{T}_i} \eta_K^2 \right) \\ \Phi(z_h) &= C_{t,\text{rel}} \max_{1 \leq i \leq N} k_i \xi_i \\ \Psi(z_h) &= C_{m,\text{rel}} \max_{1 \leq i \leq N} k_i \left( \sum_{K \in \mathcal{T}_i} \zeta_K^2 \right) \end{aligned}$$

where  $\Theta$  is given in lemma 5.2.5,  $\Phi$  is given in lemma 6.2.3 and  $\Psi$  is given in lemma 6.2.5.

*Proof.* Using (6.2.4) in the residual equation results in,

$$\begin{aligned} \langle R(z_h), \chi - \pi^0 \mathcal{I} \chi \rangle &= \langle R(z_h), \pi^0 (I - \mathcal{I}) \chi \rangle \\ &\quad + \langle R(z_h), (I - \pi^0) \mathcal{I} \chi \rangle \\ &\quad + \langle R(z_h), (I - \pi^0) (I - \mathcal{I}) \chi \rangle. \end{aligned}$$

Taking absolute values and combining the results of lemmas 5.2.5, 6.2.3 and 6.2.5 gives the result.  $\square$

We can use the results of lemma 5.2.2 together with the result of the above theorem to determine a computable *a posteriori* upper bound.

### 6.3 Adaptive algorithms

We have already discussed adaptive algorithms for space time problems in section 5.4, however we need to make a modification to deal with the mixed term, though we note this change need only be made to the internal variable section. Consider the algorithm we would

use if we were only solving the internal variable problem, Recalling the *a posteriori* error estimate of theorem 6.2.6, we would typically want to satisfy the condition,

$$S_0(t_n)\Theta(z_h) + S_1(t_n)\left\{\Phi(z_h) + \Psi(z_h)\right\} \leq \text{GTOL}. \quad (6.3.1)$$

Setting  $\text{LTOL} = \max\{S_0(t_n)^{-1}, S_1(t_n)^{-1}\}\text{GTOL}$ , then at each time point we can ensure that the condition will be satisfied by making the sum of the error indicators less than  $\text{LTOL}$ . Recall that  $\Theta(z_h)$  is the spatial error indicator,  $\Phi(z_h)$  the temporal error indicator and  $\Psi(z_h)$  is the mixed term. We propose the following space and time adaptive algorithm for generating solutions.

Space and time adaptive algorithm:

---

1. Start with parameters  $\delta_1 \in (0, 1)$ ,  $\delta_2 > 1$ ,  $\theta_1 \in (0, 1)$ , and  $\theta_2 \in (0, \theta_1)$ ,  $\gamma_1, \gamma_2 \in (0, 1)$ . Set  $t = 0$ ,  $\text{LTOL} = \max\{S_0(t_n)^{-1}, S_1(t_n)^{-1}\}\text{GTOL}$ ,  $k_{\text{old}} = T$ .
2. Do:
  - i) Set  $k = k_{\text{old}}$ ,  $\mathcal{T}_{i+1} = \mathcal{T}_i$
  - ii) Calculate  $z_h^{i+1}$  and  $\Theta$  and  $\Phi$ .
  - iii) While  $\Theta + \Psi + \Phi > \text{LTOL}$ :
    - a) While  $\Phi + (1 - \gamma_1)\Psi > \gamma_2\text{LTOL}$ :
      - Set  $k = \delta_1 k$ .
      - Calculate  $z_h^{i+1}$ ,  $\Phi$  and  $\Psi$ .
    - b) While  $\Theta + \gamma_1\Psi > (1 - \gamma_2)\text{LTOL}$ :
      - Adapt  $\mathcal{T}_{i+1}$  according to the localisation of  $\Theta$  and  $\Psi$ .
      - Calculate  $z_h^{i+1}$ ,  $\Theta$  and  $\Psi$ .
  - iv) If  $\Theta + \Phi + \Psi < \theta_2\text{LTOL}$ , then  $k = \delta_2 k$ .
  - v) Set  $z_i = z_{i+1}$ ,  $t = t + k$ ,  $i = i + 1$ .

while  $t \leq T$ .

---

The algorithm comprises of two main parts, corresponding to the temporal and spatial refinement steps. For each new proposal solution, we calculate the full error estimate. If the local error condition is not satisfied, then we refine the time step until the temporal error term together with a fraction (determined by  $\gamma_1$ ) of the mixed error term is less than a fraction (determined by  $\gamma_2$ ) of the permitted tolerance. Once satisfied, we move to the spatial error. Once again the term we seek to reduce contains the mixed error, though this time with the spatial error. After these refinement steps we can be sure that we satisfy the local tolerance condition.

## 6.4 Numerical experiments

We consider some numerical experiments to demonstrate the theoretical work of this chapter. Since it is the temporal indicator that motivated the extra effort, we start by looking at the temporal convergence. We use the problem described in section 4.4 with displacement given by (4.4.5) and internal variable given by (4.4.6). With a fixed spatial mesh we consider the behaviour of the estimators as the number of timesteps increases over the fixed time interval  $I$ . As can be seen in figure (6.1) we get the right order of convergence as the time discretisation is refined.

Now we examine the behaviour of the temporal indicators in the situation where there is no time error using the problem with internal variable given by (4.4.2). We can see from figure (6.2) that each indicator behaves in a different way. The temporal indicator of chapter 5 converges like  $O(h)$  and the alternative indicator of proposition 6.2.4 appears to remain bounded. However, the exact temporal indicator of lemma 6.2.3 is essentially zero, but as in the case for the exact spatial indicator appears to grow linearly in  $h$ .

Now that we are confident that non-zero values of the space and time error indicators imply the presence of errors in either space or time, we can put them to use in a fully adaptive algorithm. To display the adaptive properties we will consider a problem on a unit square domain with internal variable given by,

$$z(x, y, t) = \begin{pmatrix} q(x, y)r(t) \\ 0 \end{pmatrix}. \quad (6.4.1)$$

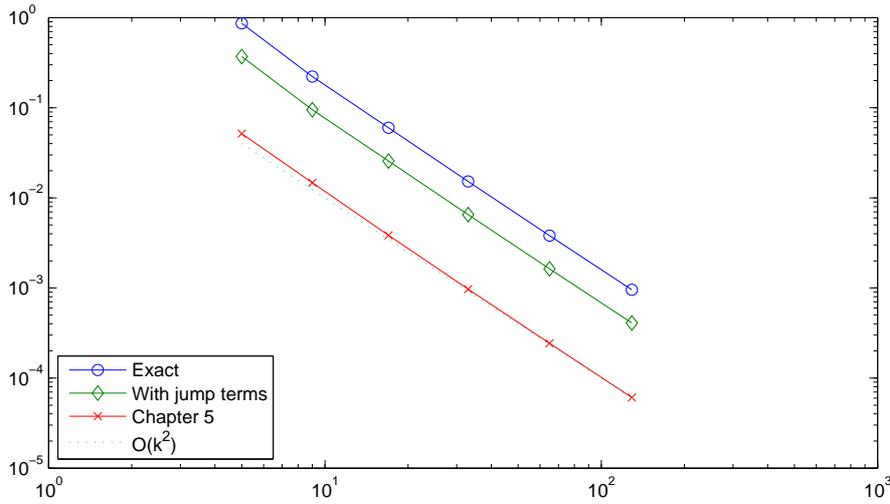


Figure 6.1: Log-log plot of the various temporal indicators against the degrees of freedom in the time discretisation.

where,

$$q(x, y) = xy(1-x)(1-y)e^{-\frac{(x-1/2)^2}{2\epsilon_1}}, \quad r(t) = te^{-\frac{(t-1/2)^2}{2\epsilon_1}}. \quad (6.4.2)$$

The choice of exact solution is motivated by the desire to provide challenging features for the adaptive algorithm to resolve — Gaussian spikes in both space and time. Given the internal variable, all other functions can be determined and so we have a Dirichlet boundary value problem for  $u$ , where  $u$  is given by,

$$u(x, y, t) = \begin{pmatrix} (\beta^{-1}r_t(t) + (\alpha\beta^{-1} - \beta)r(t))q(x, y) \\ 0 \end{pmatrix}. \quad (6.4.3)$$

Before presenting the results, we remark on the two parameters involved in tuning the adaptive algorithm,  $\gamma_1$  and  $\gamma_2$ . The parameter  $\gamma_1$  controls the split of the mixed term between spatial and temporal error, and  $\gamma_2$  controls what proportion of the local tolerance the spatial and temporal error indicators are allowed to take up, where  $\gamma_2$  is taken by the temporal indicator, and  $1 - \gamma_2$  is taken up by the spatial indicator. For example  $(\gamma_1, \gamma_2) = (0.5, 0.5)$  would mean that the mixed term has been split equally and both the spatial and temporal discretisation are allowed to consume half of the local tolerance each.

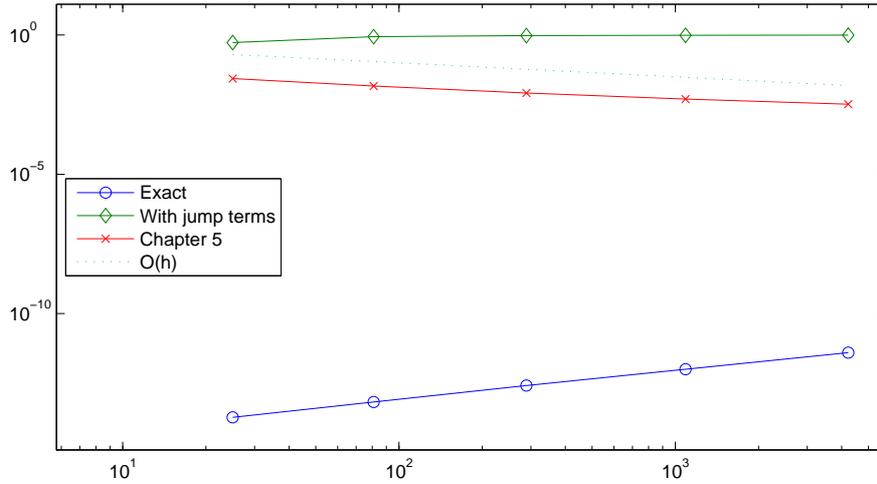


Figure 6.2: Log-log plot of the various temporal indicators against the degrees of freedom in the time discretisation.

It is our experience that  $\gamma_1 = 0.5$  is sensible unless it is the case that either the spatial or temporal indicators are zero, in which instance the mixed term is simply added to whatever indicator is non-zero. Furthermore, from now on we talk only of the spatial and temporal indicators, and assume that the mixed term has been divided up and split between them.

For  $\gamma_2$ , it is not clear what value should be chosen. To illustrate the adaptive time stepping algorithm, we present results with  $\gamma_2 = 0.01$  and for the spatial adaptivity  $\gamma_2 = 0.2$ . However, in illustrating the performance of adaptive algorithms over uniform algorithms, we split the tolerance according to the proportion of the total error that each indicator takes up when first calculated. Since  $\gamma_2$  controls the amount of tolerance that the temporal indicator can take up, we set  $\gamma_2$  equal to the percentage of the estimated temporal error (recall 6.3.1),

$$\gamma_2 = \frac{\Phi + (1 - \gamma_1)\Psi}{\Theta + \Phi + \Psi}. \quad (6.4.4)$$

This means that if the spatial indicator is large relative to the temporal indicator, it is allowed to consume more tolerance. We opt for this choice because it allows for a more balanced refinement strategy. For example, let  $\gamma_2 = 0.5$  be fixed, then at some stage in the

run, there could be scenarios where the temporal error is only 1% of the total error. With  $\gamma_2 = 0.5$  there would be heavy refinement of the spatial mesh because while the spatial error is 99% of the total error, the adaptive routine will enforce refinement until it is only half of the total error. Therefore it is more efficient to let the slack from the temporal proportion of the local tolerance be used by the spatial indicator, and that is what the choice (6.4.4) allows.

All viscoelasticity and elasticity parameters are the same as those in section 4.4, and we take Gaussian spike control parameters  $\epsilon_1 = \epsilon_2 = 0.01$ .

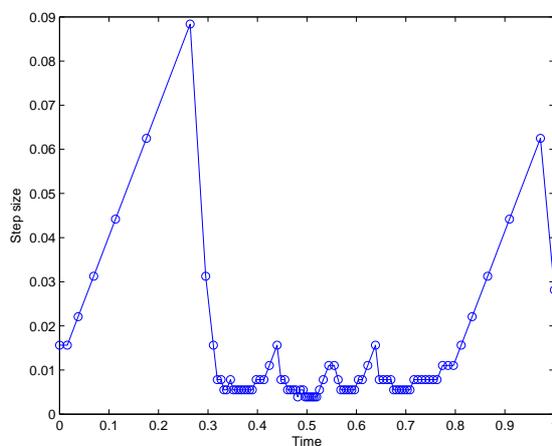


Figure 6.3: Time step sizes over the interval,  $\gamma_2 = 0.01$ .

To illustrate the temporal adaptivity, figure 6.3 shows the variation of the time steps over the time interval. As we can see the steps grow until the feature is met, they then shrink, resolving the effects of the spike and grow afterwards where the solution is once again flat. The drop in step size at the end of the interval is due to the condition ensuring that the stepper hits the final time  $T = 1$ .

To illustrate the spatial adaptivity we show a sequence of the solution and the meshes at various points throughout the time interval in figure 6.12. To balance the spatial and temporal adaptivity we take  $\gamma_2 = 0.2$ . The sequence starts with a coarse mesh, and begins refining during the first quarter of the time interval, derefines towards the center of the time

interval and then refines again before derefining towards the end of the interval (recall that  $u$  depends on  $r$  and  $r_t$  from 6.4.3).

To demonstrate the benefit of using adaptive schemes over non-adaptive ones, we look at the true error as a function of the total number of degrees of freedom used in the solver run, and compare the error for uniform and adaptive refinements in space and time. To calculate the uniform approximation, the number of spatial elements is doubled after each run, and the number of timesteps is double every other run. For the adaptive scheme we control the local tolerance, reducing by a fixed amount with each run. The results are depicted in figure 6.13. We can see that the error with the adaptive approximation, after early transients, is uniformly lower than the error with the uniform approximation with comparable degrees of freedom.

## 6.5 Summary

In this chapter we extended the work of chapter 5 to deal with the shortcomings of the temporal estimator of lemma 5.2.7 and presented numerical results demonstrating the utility of the *a posteriori* error estimates and the adaptive algorithm. To motivate the analysis we considered the error estimate of a simple space time projection and showed that the error can be decomposed in a particular way. However, the alternative temporal indicators that we presented also suffer from certain shortcomings. The indicator of lemma 6.2.3 is expensive to compute, and the indicator of 6.2.4 contains jumps over the edges. However, the indicator of lemma 6.2.3 does have the attractive property of indicating whether there is any error due to the temporal discretisation or not. We have also demonstrated the utility of adaptive methods in delivering approximation with smaller error with fewer degrees of freedom than standard uniform approximations.

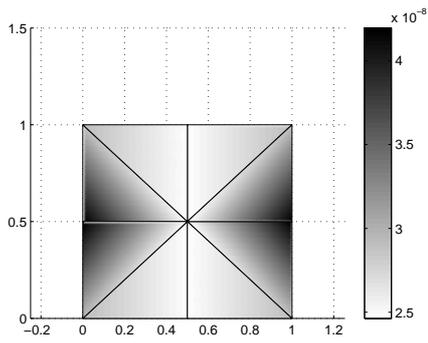


Figure 6.4:  $t = 0$ .

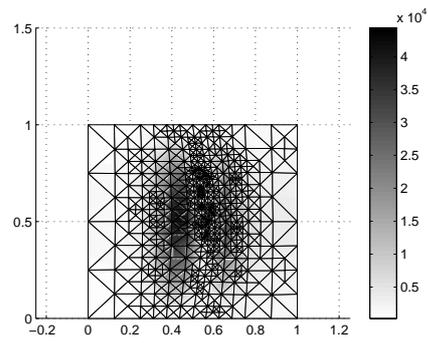


Figure 6.5:  $t = 0.3352$ .

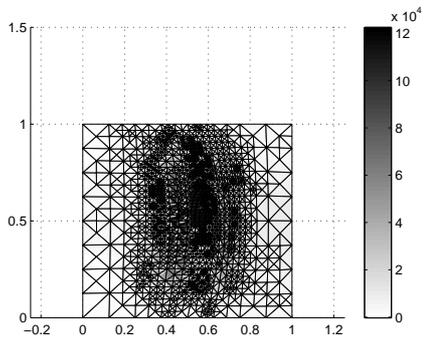


Figure 6.6:  $t = 0.3793$ .

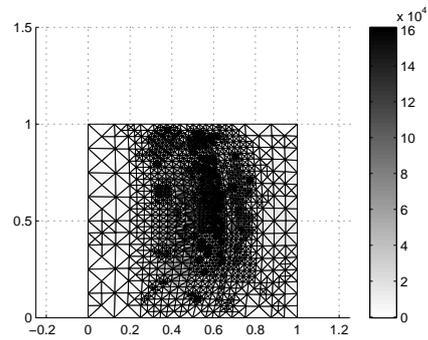


Figure 6.7:  $t = 0.4014$ .

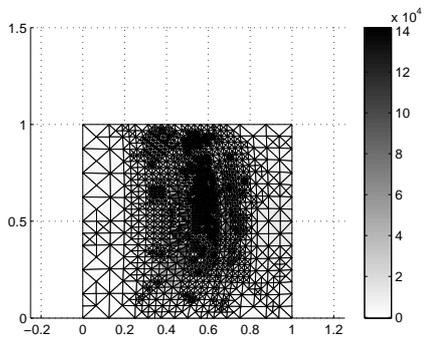


Figure 6.8:  $t = 0.4580$ .

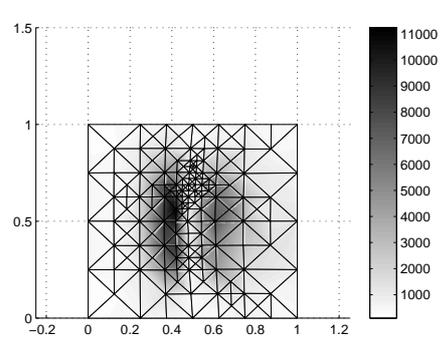


Figure 6.9:  $t = 0.5888$ .

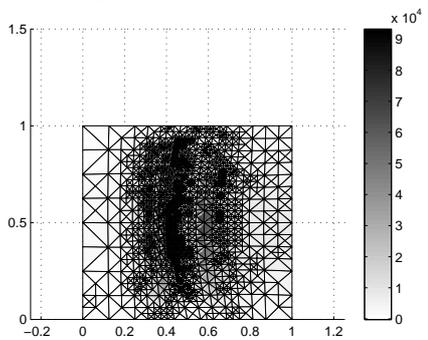


Figure 6.10:  $t = 0.7361$ .

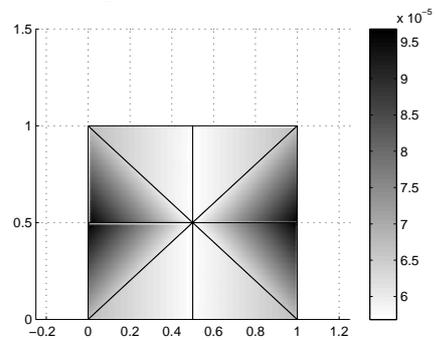


Figure 6.11:  $t = 1.0$ .

Figure 6.12: Sequence of adaptive meshes.

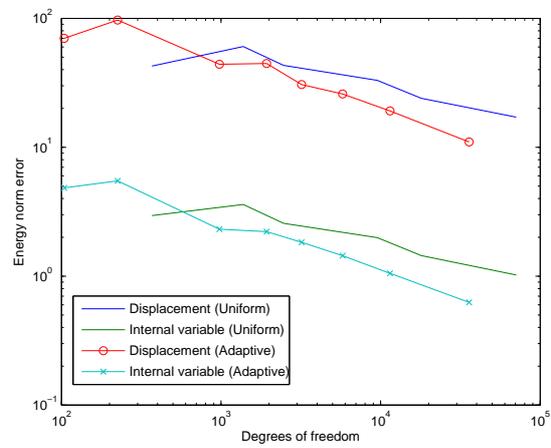


Figure 6.13: Comparison of adaptive refinement against uniform refinements.

## Chapter 7

# Nonlinear viscoelasticity

In this chapter we consider the extension of the work in the previous chapters to a problem from nonlinear viscoelasticity. The set up of the boundary value problem is the same as in section 1.5, however the constitutive equation (1.5.4) is replaced with a nonlinear law of the form,

$$\boldsymbol{\sigma}(x, t) = \mathbf{C}\boldsymbol{\epsilon}(u(x, t)) - \int_0^t \partial_s \varphi(u; t, s) \mathbf{C}\boldsymbol{\epsilon}(u(x, s)) ds. \quad (7.0.1)$$

The resulting weak form is an elliptic nonlinear Volterra problem, however in keeping with the work of previous sections, we look towards using internal variables rather than dealing with the Volterra problem directly. This is not without difficulty since firstly, we cannot remove the displacement from the internal variable equations as we did in the linear case, so not only must a simultaneous system must be solved at each time point, but the shape of the analysis presented before must change. On the plus side, we are once again left with a modified elliptic problem for the displacement, so the results of the previous chapters for the displacement equation are valid. However, since it is the nonlinearity of the internal variable equations that provide the difficulty, the focus of this chapter is on how to make progress in terms of *a posteriori* error estimation for the internal variable equation.

We begin by introducing the non-linear constitutive law, followed by the reformulation in terms of internal variables. While we first motivate the nonlinearity in the context of a displacement problem for a viscoelastic body, the main feature of the nonlinearity is the reduced time component, which features in other problems involving viscoelasticity such as

diffusion models. Therefore to isolate the difficulties provided by the nonlinearity, we leave the displacement (vector) setting and pose a simplified scalar problem that contains the essential features, allowing us to concentrate on dealing with the non-linearity.

## 7.1 Schapery-Knauss-Emri constitutive model

The Schapery-Knauss-Emri (SKE) constitutive model invokes a reduced time based upon a nonlinear functional of the strain. The constitutive law is then,

$$\boldsymbol{\sigma}(u) = \mathbf{C}\boldsymbol{\epsilon}(u) - \int_0^t \varphi_s(\rho(u; t) - \rho(u; s))\mathbf{C}\boldsymbol{\epsilon}(u(s)) ds. \quad (7.1.1)$$

where,

$$\rho(u; t) = \int_0^t \exp\left\{\frac{b\text{tr}\boldsymbol{\epsilon}(u(\tau))}{f_0 + \gamma\text{tr}\boldsymbol{\epsilon}(u(\tau))}\right\} d\tau, \quad (7.1.2)$$

and  $\varphi$  is the Prony series relaxation (1.4.25) function from chapter 1.5. The quantities  $f_0$ ,  $\gamma$  and  $b$  are positive fixed scalars. We must also ensure that the denominator in the integrand of (7.1.2) is defined, so require the constraint,

$$\text{tr}\boldsymbol{\epsilon}(u) > -\frac{f_0}{\gamma}, \quad (7.1.3)$$

which in words, means that the volumetric strain, which measures the ratio of change of a body's volume is bounded below. A negative ratio implies compression and a positive ratio implies expansion. The constraint then implies that the maximum level of compression achievable at a given point is finite. However there are no guarantees that a computed solution would maintain this constraint so it must be dealt with explicitly in any numerical treatment.

Since we are more interested in the error analysis for a reformulation in terms of internal variables we move swiftly onto the internal variables, pausing only to say the displacement problem resulting from the constitutive law (7.1.1) is a nonlinear Volterra problem which can be analysed using the methods of for example chapter 5 from [7].

### 7.1.1 Internal variable formulation

As in chapter 3, we consider an internal variable reformulation to remove the history integral. With the aim of keeping things as simple as possible we consider the case of only one internal

variable. For the time being we suppress the  $u$  argument from the function  $\rho$ , and denote its time derivative with a prime, as in  $\partial_t \rho(t) = \rho'(t)$ . Define an internal strain  $\mathbf{q}$ , by,

$$\mathbf{q} = \int_0^t \alpha_1 \varphi_1 \rho'(s) e^{-\alpha_1(\rho(t)-\rho(s))} \boldsymbol{\epsilon}(u(s)) ds, \quad (7.1.4)$$

and observe that (7.1.1) can now be written as,

$$\boldsymbol{\sigma}(u) = \mathbf{C}(\boldsymbol{\epsilon}(u) - \mathbf{q}). \quad (7.1.5)$$

It follows that the equation for the displacement in this nonlinear problem will be of similar form to that of the linear problem considered in chapter 3. However, the nonlinearity results in a more complex ODE for the internal stresses than that of previous sections.

Differentiating (7.1.4) with respect to  $t$  results in,

$$\begin{aligned} \partial_t \mathbf{q} &= \alpha_1 \varphi_1 \rho'(t) \boldsymbol{\epsilon}(u(t)) + \int_0^t \alpha_1 \varphi_1 \rho'(s) e^{-\rho(s)} \partial_t e^{-\alpha_1 \rho(t)} \boldsymbol{\epsilon}(u(s)) ds, \\ \partial_t \mathbf{q} &= \alpha_1 \varphi_1 \rho'(t) \boldsymbol{\epsilon}(u(t)) - \alpha_1 \rho'(t) \int_0^t \alpha_1 \varphi_1 \rho'(s) e^{-\alpha_1(\rho(t)-\rho(s))} \boldsymbol{\epsilon}(u(s)) ds, \\ \partial_t \mathbf{q} &= \alpha_1 \varphi_1 \rho'(t) \boldsymbol{\epsilon}(u(t)) - \alpha_1 \rho'(t) \mathbf{q} \\ \partial_t \mathbf{q} + \alpha_1 \rho'(t) \mathbf{q} &= \alpha_1 \varphi_1 \rho'(t) \boldsymbol{\epsilon}(u(t)). \end{aligned}$$

Therefore the internal strain  $\mathbf{q}$  satisfies the initial value problem,

$$\partial_t \mathbf{q} + \alpha_1 \rho'(t) \mathbf{q} = \alpha_1 \varphi_1 \rho'(t) \boldsymbol{\epsilon}(u(t)), \quad (7.1.6)$$

$$\mathbf{q}(0) = 0. \quad (7.1.7)$$

To derive a weak form for the internal variable problem, we note that the weak form for the displacement takes the form,

$$a(u(t), v) = b(t; v) + (\mathbf{q}, \mathbf{C}\boldsymbol{\epsilon}(v))_{L^2(\Omega)}, \quad \forall v \in V, \quad (7.1.8)$$

where  $a(\cdot, \cdot)$  is the bilinear form of linear elasticity introduced in section 2.4. Starting with the bilinear form  $a(\cdot, \cdot)$  and using the right hand side of equation (7.1.6) to introduce the

internal strain, we have,

$$\begin{aligned}
a(u(t), v) &= (\mathbf{C}\boldsymbol{\epsilon}(u), \boldsymbol{\epsilon}(v))_{L^2(\Omega)}, \\
&= (\boldsymbol{\epsilon}(u), \mathbf{C}\boldsymbol{\epsilon}(v))_{L^2(\Omega)}, \\
&= ((\alpha_1\varphi_1\rho'(t))^{-1}(\alpha_1\varphi_1\rho'(t))\boldsymbol{\epsilon}(u), \mathbf{C}\boldsymbol{\epsilon}(v))_{L^2(\Omega)}, \\
&= ((\alpha_1\varphi_1\rho'(t))^{-1}(\partial_t\mathbf{q} + \alpha_1\rho'(t)\mathbf{q}), \mathbf{C}\boldsymbol{\epsilon}(v))_{L^2(\Omega)}, \\
&= ((\alpha_1\varphi_1\rho'(t))^{-1}\partial_t\mathbf{q} + \varphi_1^{-1}\mathbf{q}, \mathbf{C}\boldsymbol{\epsilon}(v))_{L^2(\Omega)}.
\end{aligned}$$

Then by (7.1.8),

$$\begin{aligned}
((\alpha_1\varphi_1\rho'(t))^{-1}\partial_t\mathbf{q} + \varphi_1^{-1}\mathbf{q}, \mathbf{C}\boldsymbol{\epsilon}(v))_{L^2(\Omega)} &= b(t; v) + (\mathbf{q}, \mathbf{C}\boldsymbol{\epsilon}(v))_{L^2(\Omega)}, \\
((\alpha_1\varphi_1\rho'(t))^{-1}\partial_t\mathbf{q} + (\varphi_1^{-1} - 1)\mathbf{q}, \mathbf{C}\boldsymbol{\epsilon}(v))_{L^2(\Omega)} &= b(t; v).
\end{aligned}$$

To tidy up the previous expression, define,

$$\gamma(u; t) = \alpha_1\varphi_1\rho'(u; t), \quad \psi = \frac{1}{\varphi_1} - 1, \quad (7.1.9)$$

and observe that the constraint  $\varphi(0) = 1$  implies that,

$$\psi = \frac{1 - \varphi_1}{\varphi_1} = \frac{\varphi_0}{\varphi_1} > 0. \quad (7.1.10)$$

The spatially weak form of the internal strain ODE is then,

$$(\gamma(u; t)^{-1}\partial_t\mathbf{q} + \psi\mathbf{q}, \mathbf{C}\boldsymbol{\epsilon}(v))_{L^2(\Omega)} = b(t; v). \quad (7.1.11)$$

To try to get the problem in as similar form as that which we have dealt with earlier, we assume that there is a  $z : [0, t] \rightarrow V$  such that,

$$(\gamma(u; t)^{-1}\partial_t\boldsymbol{\epsilon}(z) + \psi\boldsymbol{\epsilon}(z), \mathbf{C}\boldsymbol{\epsilon}(v))_{L^2(\Omega)} = b(t; v). \quad (7.1.12)$$

By reversing the derivation above, the assumption of such a  $z$  is equivalent to assuming that,

$$\boldsymbol{\epsilon}(z) = \int_0^t \alpha_1\varphi_1\rho'(s)e^{-\alpha_1(\rho(t)-\rho(s))}\boldsymbol{\epsilon}(u(s)) ds. \quad (7.1.13)$$

As we can see, the new problem poses a stern challenge, though there are similarities with our earlier work. Since our focus is on deriving *a posteriori* error estimates, it is important

to isolate where the difficulties arise in deriving estimates, rather than the technicalities associated with the actual problem. To achieve this, we consider the following scalar problem. Let  $\Omega$  be a bounded Lipschitz domain in  $\mathbb{R}^d$  with Dirichlet boundary  $\Gamma_D$  and Neumann boundary  $\Gamma_N$ , and let  $I = [0, T]$ . Define the space  $H_D^1(\Omega)$  to be those functions from  $H^1(\Omega)$  that are zero on the Dirichlet boundary. The model problem in its weak form is then: Find  $u : I \rightarrow H_D^1(\Omega)$  such that,

$$\int_0^T (\gamma(u)^{-1} \nabla \partial_t u + \psi \nabla u, \nabla v)_{L^2(\Omega)} dt = \int_0^T l(t; v) dt, \quad (7.1.14)$$

$$\nabla u(x, 0) = 0, \quad (7.1.15)$$

where

$$l(t; v) = (f, v)_{L^2(\Omega)} + (g, v)_{L^2(\Gamma_N)}. \quad (7.1.16)$$

The differences between this proposed model problem and the actual problem are:

1. Reduction to a scalar problem. The generalisation of the problem from the scalar  $H^1(\Omega)$  problem back up to the original problem is not difficult since abstractly they are both elliptic operators acting on a temporal differential equation. However, a consequence of this is that there is no strain equivalent that can be put into the nonlinearity. We must however replace the nonlinearity with something that bears a similar structure to the original. To at least match the order of the derivatives, we replace the trace of the strain with  $|\nabla u|$  to get,

$$\gamma(u) = \exp \left\{ \frac{b|\nabla u|}{f_0 + \gamma|\nabla u|} \right\}. \quad (7.1.17)$$

This choice also inadvertently circumvents the difficulty that might arise when the volumetric strain becomes negative.

2. Elimination of the displacement variable. This simplification actually renders a completely different problem from that originally posed. However, once we can establish results in this scenario, we should be more aware of the difficulties that might present themselves in the full problem. We should stress though, that the problem given by (7.1.14) is not completely artificial without its own merits. In fact it is a special case of a nonlinear Sobolev equation as considered in [37].

Integrating by parts the equation (7.1.14) we get the following strong form: Find  $u : [0, T] \rightarrow H^1(\Omega)$  such that,

$$\begin{aligned} -\nabla \cdot (\gamma(u)^{-1} \nabla \partial_t u(x, t)) - \Delta u(x, t) &= f(x, t), & (x, t) \in \Omega \times (0, T], \\ u(x, t) &= 0, & (x, t) \in \Gamma_D \times [0, T], \\ (\gamma(u)^{-1} \nabla \partial_t u(x, t) + \nabla u(x, t)) \cdot n(x) &= g(x, t), & (x, t) \in \Gamma_N \times [0, T], \\ \nabla u(x, 0) &= 0, & x \in \Omega. \end{aligned}$$

Assuming that a solution to this problem exists (see [37] for references), we can derive the following stability result.

**Proposition 7.1.1.** *Let  $u : I \rightarrow H_D^1(\Omega)$  be a solution of the problem defined by (7.1.14) and (7.1.15) and suppose that there is a constant  $\bar{\gamma}$  such that  $0 < \gamma(u) \leq \bar{\gamma}$ , then there holds,*

$$\psi \|\nabla u(t)\|_{L^2(\Omega)}^2 + \int_0^t \|\gamma(u)^{-1/2} \nabla \partial_s u(s)\|_{L^2(\Omega)}^2 ds \leq C\bar{\gamma} \int_0^t \|l(s)\|_{V^*}^2 ds. \quad (7.1.18)$$

*Proof.* Choose  $v = \partial_t u$  in (7.1.14). □

## 7.2 Finite element approximation

We consider a continuous piecewise linear Galerkin finite element approximation in both space and time to the function  $u$ . Let  $1 < p < \infty$ , and define the space,

$$W^p(I; H_D^1(\Omega)) = \{u : I \rightarrow H_D^1(\Omega), u \in L^p(I; H_D^1(\Omega)), \partial_t u \in L^p(I; H_D^1(\Omega))\}, \quad (7.2.1)$$

It is well known [33] that functions from  $W^p(I; H_D^1(\Omega))$  are continuous in time, so we define  $W_0^p(I; H_D^1(\Omega)) = W^p(I; H_D^1(\Omega)) \cap \{u : u(x, 0) = 0\}$ . We write the weak problem (7.1.14) in the abstract form: Find  $u \in W_0^1(I; H_D^1(\Omega))$  such that,

$$A(u; v) = L(v), \quad \forall v \in L^p(I; H_D^1(\Omega)), \quad (7.2.2)$$

where,

$$\begin{aligned} A(u; v) &= \int_0^T (\gamma(u)^{-1} \nabla \partial_t u + \psi \nabla u, \nabla v)_{L^2(\Omega)} dt, \\ L(v) &= \int_0^T l(t; v) dt. \end{aligned}$$

Note that when writing  $A(u; v)$ , we mean that the form is nonlinear in arguments to the left of the semi-colon, and linear in those to the right. For the finite element approximation we use a continuous Galerkin  $\mathbb{P}_1$  Lagrange finite element in space as described in section 2.2, together with a continuous Galerkin  $\mathbb{P}_1$  Lagrange finite element in time, described in 2.5. Denote the trial space by  $V^h$  and the test space by  $W^h$ . Then the finite element approximation is: Find  $u_h \in V^h$ , such that,

$$A(u_h; v_h) = L(v_h), \quad \forall v_h \in W^h. \quad (7.2.3)$$

Equation (7.2.3) is a finite dimensional system of nonlinear equations for which a wide variety of algorithms based on Newton's method can be applied [7].

### 7.3 Towards *a posteriori* error analysis

In this section we consider the derivation of an *a posteriori* error estimate for a functional of the error  $e = u - u_h$ . Rannacher proposed a general scheme for doing so [61], however we prefer to proceed directly by designing the dual problem to suit our needs. The argument is a generalisation of that of earlier chapters, in that a dual problem is used to derive an error representation formula. However, the linear dual problem is based on a linearisation of the original nonlinear problem. An additional difference for this nonlinear context, is that Galerkin orthogonality takes the form,

$$A(u; v_h) - A(u_h; v_h) = 0, \quad \forall v_h \in W^h, \quad (7.3.1)$$

where the nonlinearity prevents us forming an error term in the left-hand argument, however, we note that

$$A(u; v_h) - A(u_h; v_h) = \int_0^1 \frac{d}{ds} A(su + (1-s)u_h; v_h) ds = \int_0^1 \frac{d}{ds} A(u_h + se; v_h) ds. \quad (7.3.2)$$

We now form the dual problem, effectively reverse engineering it from the result we want

to achieve. Integration by parts over the time interval results in,

$$\begin{aligned}
A(u; \chi) - A(u_h; \chi) &= \int_0^T (\gamma(u)^{-1} \nabla \partial_t u, \nabla \chi)_{L^2(\Omega)} dt - \int_0^T (\gamma(u_h)^{-1} \nabla \partial_t u_h, \nabla \chi)_{L^2(\Omega)} dt \\
&\quad + (\psi \nabla(u - u_h), \nabla \chi)_{L^2(\Omega)} dt \\
&= (\gamma(u)^{-1} \nabla u, \nabla \chi)_{L^2(\Omega)} \Big|_0^T - (\gamma(u_h)^{-1} \nabla u_h, \nabla \chi)_{L^2(\Omega)} \Big|_0^T \\
&\quad - \int_0^T (\nabla u, \partial_t (\gamma(u)^{-1} \nabla \chi))_{L^2(\Omega)} dt + \int_0^T (\nabla u_h, \partial_t (\gamma(u_h)^{-1} \nabla \chi))_{L^2(\Omega)} dt \\
&\quad + \int_0^T (\psi \nabla(u - u_h), \nabla \chi)_{L^2(\Omega)} dt
\end{aligned}$$

Using the initial condition  $\nabla u(x, 0) = 0$  and setting  $\nabla \chi(x, T) = 0$ , we get,

$$\begin{aligned}
A(u; \chi) - A(u_h; \chi) &= - \int_0^T \left\{ (\nabla u, \partial_t (\gamma(u)^{-1} \nabla \chi))_{L^2(\Omega)} - (\nabla u_h, \partial_t (\gamma(u_h)^{-1} \nabla \chi))_{L^2(\Omega)} \right\} dt \\
&\quad + \int_0^T (\psi \nabla(u - u_h), \nabla \chi)_{L^2(\Omega)} dt
\end{aligned}$$

If we now define,

$$P(u; v) := - \int_0^T (\nabla u, \partial_t (\gamma(u)^{-1} \nabla v))_{L^2(\Omega)} dt, \quad (7.3.3)$$

then,

$$\begin{aligned}
A(u; \chi) - A(u_h; \chi) &= P(u; \chi) - P(u_h; \chi) + \int_0^T (\psi \nabla(u - u_h), \nabla \chi)_{L^2(\Omega)} dt, \\
&= \int_0^1 \frac{d}{ds} P(su + (1-s)u_h; \chi) ds + \int_0^T (\psi \nabla(u - u_h), \nabla \chi)_{L^2(\Omega)} dt \\
&= \int_0^1 P'(u_h + se; e, \chi) ds + \int_0^T (\psi \nabla(u - u_h), \nabla \chi)_{L^2(\Omega)} dt.
\end{aligned}$$

We therefore define the dual bilinear form  $A^*(u_h; v, \chi)$  to be,

$$A^*(u_h; v, \chi) = \int_0^1 P'(u_h + sv; v, \chi) ds + \int_0^T (\psi \nabla v, \nabla \chi)_{L^2(\Omega)} dt.$$

The dual problem is then: Find  $\chi \in W$  such that,

$$A^*(u_h; v, \chi) = \int_I J(v) dt, \quad \forall v \in V. \quad (7.3.4)$$

Leaving aside the technical issues for the moment, we want to outline the extension of our previous results for linear problems to the current setting. Define  $J$  to be,

$$J(v) := (\nabla g, \nabla v)_{L^2(\Omega)}, \quad (7.3.5)$$

where  $g$  is the right hand side of the strong form of the dual problem. The dual problem (7.3.4) gives us,

$$\begin{aligned} J(e) &= \int_I (\nabla g, \nabla e)_{L^2(\Omega)} dt = A^*(u_h; e, \chi), \\ &= \int_0^1 P'(u_h + se; e, \chi) ds + \int_0^T (\psi \nabla e, \nabla \chi)_{L^2(\Omega)} dt, \\ &= A(u; \chi) - A(u_h; \chi). \end{aligned}$$

By Galerkin orthogonality (7.3.1) we have for all  $\phi_h \in W^h$ ,

$$\begin{aligned} J(e) &= A(u; \chi - \phi_h) - A(u_h; \chi - \phi_h), \\ &= L(\chi - \phi_h) - A(u_h; \chi - \phi_h), \\ &= \langle R(u_h), \chi - \phi_h \rangle. \end{aligned}$$

To get an error estimate, we follow the approach of chapter 5. From the representation, we can localise to the time slab and element level as follows. Integrating by parts over  $\Omega$  results in,

$$J(e) = \sum_{i=1}^n \int_{I_i} \sum_{K \in \mathcal{T}_i} (f, \chi - \phi_h)_{L^2(K)} + \sum_{E \in K} (R_E, \chi - \phi_h)_{L^2(E)} + \sum_{E \subset K} (R_E, \chi - \phi_h)_{L^2(E)}, \quad (7.3.6)$$

where,

$$R_E := \begin{cases} \llbracket \gamma(u_h)^{-1} \nabla (\partial_t u_h + \psi u_h) \rrbracket, & E \in \Omega, \\ g - \gamma(u_h)^{-1} \nabla (\partial_t u_h + \psi u_h) \cdot n, & E \in \Gamma_N. \end{cases} \quad (7.3.7)$$

We can now proceed as in chapter 5 to derive an upper bound. Since  $\phi_h \in W^h$  is arbitrary, take  $(\chi - \phi_h)|_{I_i} = \chi - \pi_i^0 \chi + \pi_i^0 (I - \mathcal{I}_i) \chi$  where  $\pi_i^0$  is the  $L^2$  projection onto piecewise constants from theorem 2.2.5 and  $\mathcal{I}_i$  is the piecewise linear quasi-interpolant from theorem 2.2.3. Then,

$$\langle R(u_h), \chi - \phi_h \rangle = \langle R(u_h), \chi - \pi^0 \rangle + \langle R(u_h), \pi^0 (I - \mathcal{I}) \chi \rangle. \quad (7.3.8)$$

Take the first term to be the time error, then we have,

$$\begin{aligned} \langle R(u_h), \chi - \pi^0 \chi \rangle &= \sum_{i=1}^n \int_{I_i} (f, \chi - \pi^0 \chi)_{L^2(\Omega)} + (g, \chi - \pi^0 \chi)_{L^2(\Gamma_N)} \\ &\quad - (\gamma(u_h)^{-1} \nabla (\partial_t u_h + \psi u_h), (I - \pi^0) \nabla \chi)_{L^2(\Omega)} dt. \end{aligned}$$

Arguing as we have in the previous chapters, we can subtract off projected forms of  $f$ ,  $g$  and the discrete equation to get the following upper bound,

$$|\langle R(u_h), \chi - \pi^0 \chi \rangle| \leq C_{t,\text{rel}} \max_{1 \leq i \leq n} k_i \xi_i \|\nabla \partial_t \chi\|_{L^1(I; L^2(\Omega))}, \quad (7.3.9)$$

where  $\xi_i$  is defined by,

$$\begin{aligned} \xi_i = & \|(I - \pi_i)f\|_{L^\infty(I_i; L^2(\Omega))} + \|(I - \pi_i)g\|_{L^\infty(I_i; L^2(\Omega))} \\ & + \|(I - \pi_i)\gamma(u_h)^{-1} \nabla(\partial_t u_h + \psi u_h)\|_{L^\infty(I_i; L^2(\Omega))}. \end{aligned}$$

For the spatial residual, we get the following,

$$|\langle R(u_h), \pi^0(I - \mathcal{I})\chi \rangle| \leq C_{s,\text{rel}} \left( \sum_{K \in \mathcal{T}_i} \eta_K^2 \right)^{1/2} \|\nabla \chi\|_{L^1(I; L^2(\Omega))}, \quad (7.3.10)$$

where,

$$\eta_K^2 = h_K^2 \|\pi_i^0 f\|_{L^2(K)}^2 + \sum_{E \in K} h_E \|\pi_i^0 R_E\|_{L^2(E)}^2 + \sum_{E \subset K} h_E \|\pi_i^0 R_E\|_{L^2(E)}^2. \quad (7.3.11)$$

If we define the stability factors,

$$\begin{aligned} S_0(I) &:= \frac{\|\nabla \chi\|_{L^1(I; L^2(\Omega))}}{\|\nabla g\|_{L^2(I; L^2(\Omega))}}, \\ S_1(I) &:= \frac{\|\nabla \partial_t \chi\|_{L^1(I; L^2(\Omega))}}{\|\nabla g\|_{L^2(I; L^2(\Omega))}}, \end{aligned}$$

then we have the *a posteriori* error estimate,

$$\begin{aligned} \|e\|_{L^2(I; L^2(\Omega))} &= \sup_{g \in L^2(I; H^1(\Omega))} \frac{J(e)}{\|\nabla g\|_{L^2(I; L^2(\Omega))}} \\ &\leq C_{s,\text{rel}} \max_{1 \leq i \leq n} \left( \sum_{K \in \mathcal{T}_i} \eta_K^2 \right)^{1/2} S_0(I) + C_{t,\text{rel}} \max_{1 \leq i \leq n} k_i \xi_i S_1(I). \end{aligned}$$

## 7.4 Summary

The mechanics of this extension appear straightforward with the exception of the stability analysis of the dual problem. It would appear that this is highly non-trivial, but it could form the focus of further work. Alternatives to analysis are the computational approaches of determining the stability factors  $S_0(T)$  and  $S_1(T)$  [35] or preserving the error representation and computing the dual solution as in the dual weighted residual method [61], [75].

## Chapter 8

# Summary and recommendations for further work

In this chapter we provide a summary of the work that has been presented and make some recommendations for future lines of research.

### 8.1 Summary

We have considered the adaptive finite element method (AFEM) for the solution of the displacement problem for quasistatic linear viscoelasticity and made inroads into an extension towards dealing with non-linear viscoelasticity. Our main focus has been *a posteriori* error estimation and the resulting adaptive algorithms. Due to a novel reformulation utilising internal variables, the resulting problem became two almost distinct problems. One, an augmented elliptic boundary value problem, the other, an elliptic operator acting on a system of ODEs. The reformulation led us to review the work in the areas of AFEM for elliptic boundary value problems and ODEs and relevant results and algorithms were presented in chapter 2.

After showing that the proposed finite element method for the solution of the reformulated problem converged (chapter 4), we considered the *a posteriori* error analysis for the problem and presented an *a posteriori* error estimate in chapter 5. Given the deficiencies of the proposed estimator, we proposed an alternative in chapter 6 together with numeri-

cal examples and an adaptive algorithm. In performing the numerical experiments, it was necessary for us to implement a mesh adaptation package, a summary of which appears in appendix A.

After dealing with the linear problem, we turned to a nonlinear generalisation of the viscoelastic model involving reduced time in chapter 7. Noting the obvious difficulties we reduced the problem to a simpler one where we could more accurately evaluate the methodology of deriving *a posteriori* error estimates for nonlinear problems. An *a posteriori* error estimate was presented and a comments regarding its extension were made.

## 8.2 Recommendations for future work

Based on the contents of this thesis, we now make some recommendations for further work, and outline possible extensions of what we have done already.

1. The execution of the exact indicators did not go through cleanly, however it seems that the idea of decomposing the space as we did for the  $L^2$  space time projection is the right idea, and further study into this area to circumvent the computational difficulties might yield improved results.
2. Algorithm analysis and tuning parameters. In the adaptive algorithms that we have presented, there are a number of tuning parameters that need to be selected. A lot of experimentation would be required to determine some optimal values, and in practice probably only sensible ones will be used. However, given that there is a general structure to the adaptive algorithms, irrespective of the underlying problem, there is scope to try to adapt the tuning parameters during the computation. As an example, in adaptive time step selection, we had the fixed constants  $\delta_1$  and  $\delta_2$  scaling the proposed timestep in a switch like manner. However, it would be possible to build functionality so that  $\delta_1$  and  $\delta_2$  were modified depending on whether the time step coming from their proposal was accepted or not. The theoretical question would be to determine methods of determining the parameters that offered gains in efficiency of the algorithms.

3. While the form of relaxation function taken here is frequently used in the literature, it would be interesting to consider alternative relaxation functions that are also popular.
4. Nonlinear viscoelasticity. We have considered the extension of the results for linear problems to the nonlinear case in chapter 7. Work now should move forward by analysing the dual problem, and or computing via the DWR method. With success, the next step from there would be to consider a coupled problem, where the main variable is contained in a functional that is a coefficient of the internal variable equation. Such problems exist in the field of polymer diffusion.
5. Modelling error. While the development of an *a posteriori* error estimate for the finite element error in the approximation to the nonlinear problem is further work in its own right, there is an interesting thread of analysis that uses the same methodology for estimating modelling errors. The model of nonlinear viscoelasticity that we presented in 7 contains the linear problem that was the focus of our earlier work as a special case. In this respect we have a two layer hierarchy of models. An interesting future line of research would be adaptive modelling, where a computation would start with the linear model, then during the computation determine whether the model is still appropriate, and if not switch to the nonlinear model. Ideally this would happen locally rather than across the whole domain. Modelling error estimation has been considered in [58] [14].

# Appendix A

## Adaptive mesh refinement in MATLAB

### A.1 Introduction

This paper is stimulated by the opening remarks of [5]. We aim to offer an open-box MATLAB implementation of local adaptive mesh generation that is simple, easy to understand and easy to modify. The offering is a set of MATLAB routines that are fully compatible with the output of the mesh generation software given in [59], and with the mesh storage mechanism underlying the finite element implementations given in [4], [5] and [20]. Our aim is not to provide a full package such as MClite [45], but to illustrate a basic set of elementary routines that can easily be combined with other freely available MATLAB codes to provide a flexible problem solving environment.

A generator of geometrically conforming meshes is presented and implemented in MATLAB in [59]. The software provided there allows the user to specify the structure of the mesh *a priori*, and while useful, there is no mechanism for local adaptivity. The Solve - Estimate - Refine (SER) approach of adaptive finite element methods (see [28], [54] and [77]) is based on approximations on locally adapted meshes. The local adaptation is driven by *a posteriori* error estimates, and while the MATLAB implementations of the finite element method described in [4], [5], and [20] discuss *a posteriori* error estimation, they stop short of dealing with the issue of adaptive mesh refinement based on such estimators.

In this paper we consider the generation of locally adapted meshes for geometrically conforming triangulations, suitable for example, for conforming  $\mathbb{P}_1$  finite element approximations. However it is worth pointing out that this is a more difficult consideration than the non-conforming case, and the code presented can be easily adjusted to deal with such an instance. Our method of storing the mesh is taken from [4], [5], [20], and [59], and we show in examples how the code presented here can easily be used alongside the software presented in those papers as part of an adaptive solution procedure.

The choice of refinement procedure is motivated by the results presented in [54], and uses at its starting point the discussion of hierarchic meshes in [62]. However, rather than introduce a new mesh structure to imitate implementations of the refinement algorithm in other programming languages, the algorithm and code are designed to make use of the powerful array handling facilities provided in standard MATLAB.

### A.1.1 Notation

We adopt the usual misleading nomenclature of referring to triangular subdomains as elements rather than element domains. A mesh is then a triple  $(\mathcal{T}, \mathcal{E}, \mathcal{N})$ , of elements, edges and nodes respectively. Let  $\Omega$  denote a polygonal domain with boundary  $\Gamma$ . The set  $\mathcal{T}$  forms a partition of  $\Omega$  into closed triangles  $K$  satisfying:

1. The partition covers the closure of  $\Omega$ ,  $\bar{\Omega} = \bigcup_{K \in \mathcal{T}} K$ .
2. The intersection of the interiors of distinct elements is empty, so that if  $K \neq K'$  then  $\overset{\circ}{K} \cap \overset{\circ}{K}' = \emptyset$ .
3. The intersection of any 2 distinct elements results in an edge, a node or is empty.

That is if  $K \neq K'$ , then

$$K \cap K' = \begin{cases} E \in \mathcal{E}, & \text{or,} \\ z \in \mathcal{N}, & \text{or,} \\ \emptyset. \end{cases}$$

The symbol  $\mathcal{E}(\Gamma)$  denotes those edges on the boundary, and if  $K \in \mathcal{T}$  then the edges belonging to  $K$  are denoted by  $\mathcal{E}(K)$ . While not required for the discussion on mesh

refinement, they are required later, so we introduce the gridsize function  $h$  measuring the size of elements and edges, defined by,

$$h(S) := \text{diam}(S),$$

and the shape regularity property,

$$\frac{h_K}{\rho_K} \leq c, \quad \rho_K := \sup\{r \mid B_r(x_0) \subset K, \forall x_0 \in K\}, \quad (\text{A.1.1})$$

where  $B_r(x_0)$  represents the ball of radius  $r$  centered at  $x_0$ , and  $c$  is independent of  $K$  and  $h$ . In two dimensions, if the smallest angle in the triangulation is bounded below then the shape regularity condition holds.

For each  $K \in \mathcal{T}$  denote the longest edge of  $K$  by  $\bar{E}_K := \arg \max_{E \in \mathcal{E}(K)} h(E)$ , and the element adjacent to the longest edge of  $K$  by  $\alpha_K$ ,

$$\alpha_K := K', \text{ where } K \cap K' = \bar{E}_K.$$

The set of elements adjacent to an element  $K$  is denoted by

$$\Lambda_K := \{K' \in \mathcal{T} : \mathcal{E}(K) \cap \mathcal{E}(K') \neq \emptyset\}.$$

The set of elements marked to be refined is denoted by  $\mathcal{M}$ .

## A.2 Local mesh refinement

Locally refined meshes can be generated in a number of ways (see [77] for example). All of them essentially involve the refinement of individual elements, together with a principle regarding the geometric conformity of the mesh. In the non-conforming case, where hanging nodes are admitted, local refinement proceeds by refining the desired element with no regard for the surrounding elements. However, if geometric conformity is to be maintained, the refinement algorithm must include the refinement of neighbouring elements.

### A.2.1 Element refinement

The refinement procedure used for marked elements subdivides the element into 6 smaller triangles as shown in figure A.1. This is achieved by longest edge bisection of the parent element, and the resulting child elements. The last step is the addition of the internal node which is achieved by refining the two child elements whose longest edge is on the interior. An important consideration for finite element approximations is the resulting structure of the mesh. Constants appearing in error estimates for quasi-interpolants depend on the smallest angle appearing in the mesh [77]. It is known (see [77] and references within) that the smallest angle in a mesh arrived at by longest edge bisections is bounded below by half of the smallest angle in the initial mesh.

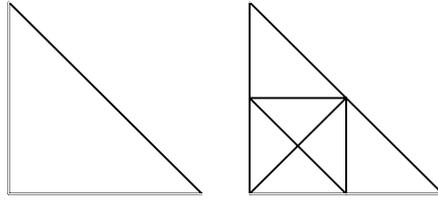


Figure A.1: Element domain before and after refinement.

The key difference between this type of refinement and most others is that it introduces a node into the interior of the domain. See [54] for details regarding the introduction of an interior node during refinement and convergence of adaptive finite element methods.

### A.2.2 Closure

Let  $K$  denote an element marked for refinement. The introduction of new nodes on the edges of  $K$  means that to ensure geometric conformity, all elements in  $\Lambda_K := \{K', K'', K'''\}$  must be refined. If  $\mathcal{E}(K) = \{\bar{E}_{K'}, \bar{E}_{K''}, \bar{E}_{K'''}\}$ , refinement of the marked element can take place with each adjacent element undergoing a longest edge bisection as shown in figure A.2.

The difficulty occurs if an element adjacent to  $K$  does not share its longest edge with  $K$ . Let  $K'$  be the element such that  $K' \cap K \neq \bar{E}_{K'}$ . Then  $K'$  must be refined until the resulting element adjacent to  $K$  shares its longest edge with  $K$ . Suppose that  $\alpha_{K'} \cap K = \bar{E}_{\alpha_{K'}}$ ,

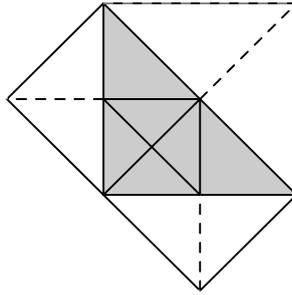


Figure A.2: The shaded element has been refined. The dotted lines represent the edges added to maintain conformity.

then since these elements share their longest edge they can be bisected simultaneously. The resulting subelement neighbouring  $K$  can be bisected and conformity is maintained.

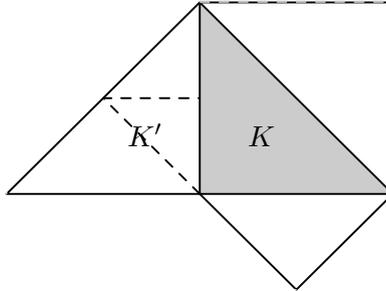


Figure A.3: Configuration of elements showing a case where refinement must take place before the marked element can be refined. The dotted lines show the required edges.

The process of finding the longest edge neighbours results in a path for each element  $K' \in \Lambda_K$ , where a sequence  $\{K_0 = K', K_1, K_2, \dots, K_n\}$  is generated from the rule where  $K_{i+1} = \alpha_{K_i}$ . The path terminates when  $\alpha_{K_n} = K_{n-1}$  or  $\bar{E}_{K_n} \in \mathcal{E}(\Gamma)$ . Once the paths are determined, a backwards recursive bisection algorithm is applied from the end of the path back to  $K$  so that it can be refined. Pathological constructions for which the path for refinement of a single element is the entire mesh are rare in practice. A discussion of the properties of these paths appears in [74].

The implementation presented here takes a different approach from the path based recursive algorithm discussed above. To avoid recursion we introduce an indicator array

which stores the required information from the propagation paths. However rather than work back up the paths individually, the indicator array allows for a vectorized treatment of the refinement process, so that elements that are refined almost simultaneously. This allows us to exploit the strength of MATLAB in handling arrays. The whole process can be broken down into three stages. The first, we call the closure process, this is the enlargement of the marked elements to include those that are to be refined to preserve the geometric conformity of the mesh. The information from this process is stored in the indicator array. Then refinement takes place. The implementation of the refining process is sloppy in the sense that duplicate nodes are created during the refinement process but removed at the end during the administration of the degrees of freedom stage. The final stage is the application of the MATLAB `unique` function which handles the administration of the degrees of freedom, ensuring that node numberings are consistent.

The MATLAB code takes as input the following sets of data.

1. A list of coordinates, `coordinates`, storing the Euclidean coordinate position of nodes. The global node number of a node corresponds to the position within this list.
2. A list of elements, `elements`, storing the three global node numbers of its constituent nodes. The nodes should be numbered in the anticlockwise direction with the longest edge as the first two numbers.
3. The Dirichlet and Neumann boundary condition regions, stored as lists of edges `dirichlet` and `neumann`.
4. A list of indices corresponding to which of the elements are to be refined, `marked`.

### A.2.3 Closure algorithm

To describe the closure algorithm we introduce the indicator array that is the key tool in the implementation presented here. The indicator array is a binary matrix of dimension (number of elements)  $\times$  3 for which the usual rules of boolean logic (and  $\wedge$ , or  $\vee$ , etc...) are applied componentwise, e.g.,

$$\begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} \vee \begin{pmatrix} 1 & 0 \\ 1 & 0 \end{pmatrix} = \begin{pmatrix} 1 & 0 \\ 1 & 1 \end{pmatrix}$$

The closure algorithm consists of two stages. First the set of elements to be refined for geometric conformity is found. Denote this set by  $\mathcal{C}$ . The elements that are on the longest edge of elements in  $\mathcal{C}$  are added to  $\mathcal{C}$ , until all elements in  $\mathcal{C}$  share their longest edge with another element in  $\mathcal{C}$ , or their longest edge is a boundary edge. While this proceeds the indicator array records what elements are to be refined and how. Recall that elements are stored with their longest edge appearing first and numbered in the anti-clockwise direction. The row of the indicator matrix will contain 1's or 0's depending on whether the first second or third edge is to be refined. Denote the set of elements marked for refinement by  $\mathcal{M}$ .

Closure Algorithm:

---

Input  $\mathcal{M}$ :

1. Define the set  $\mathcal{C} := \cup_{K \in \mathcal{M}} \Lambda_K$  and the indicator matrix  $I$  such that

$$I_{ij} = \begin{cases} 1 & E_j \in \mathcal{E}(K_i) \cap \mathcal{E}(\mathcal{M}) \\ 0 & \text{otherwise} \end{cases}$$

2. Set  $I^{new} := I^{old}$  and mark the longest edge of any element that has an edge to be refined

$$I_{i1}^{new} := I_{i1}^{old} \vee I_{i2}^{old} \vee I_{i3}^{old}.$$

3. While  $I^{new} \neq I^{old}$

$$I^{old} := I^{new}$$

$$I_{ij}^{test} := \begin{cases} 1 & E_j \in \mathcal{E}(K_i) \cap \left( \cup_{K \in \mathcal{C}} \bar{E}_K \right) \\ 0 & \text{otherwise} \end{cases}$$

$$I_{i1}^\alpha = I_{i1}^\alpha \vee I_{i2}^\alpha \vee I_{i3}^\alpha$$

$$I^{new} = I^{new} \vee I^{test}$$

$$\mathcal{C} := \left( \bigcup_{K \in \mathcal{C}} \alpha_K \right)$$

The initial indicator array is constructed in MATLAB in the following way. Firstly the arrays of edges and marked edges are created by the following lines

```
edges=[elements(:,[1,2]);elements(:,[2,3]);elements(:,[3,1])];
markededges=[elements(marked,[1,2]);
             elements(marked,[2,3]);
             elements(marked,[3,1])];
markededges=[markededges;markededges(:,[2,1])];
```

The arrays are then compared using the `ismember` function.

```
tf=ismember(edges,markededges,'rows');
```

The return value is a logical array the same length as `edges` containing a 1 or 0 depending on whether the entry of `edges` in that row is contained in `markededges`. To create the matrix this list is just reshaped using the MATLAB `reshape` function.

```
indicator=reshape(tf,size(elements));
```

**Example** Suppose that the central element 543 in figure A.4 is marked for refinement.

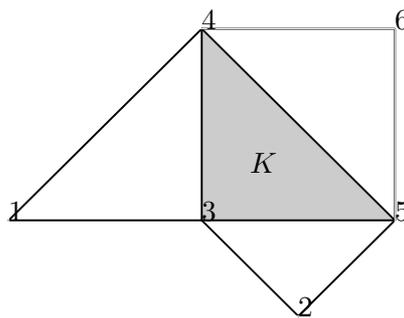


Figure A.4: Element configuration to illustrate usage of the indicator array.

The corresponding arrays are

$$\begin{array}{rcc} & 5 & 4 & 3 & & 1 & 1 & 1 \\ \text{elements} = & 4 & 5 & 6 & \text{indicator} = & 1 & 0 & 0 \\ & 2 & 5 & 3 & & 0 & 1 & 0 \\ & 4 & 1 & 3 & & 0 & 0 & 1 \end{array}$$

If an element has an edge that must be bisected, then its longest edge must also be bisected. The longest edges of elements with an edge to be refined is marked by the following lines of code:

```
newindicator=indicator;
newindicator(:,1)=any(indicator,2);
```

The MATLAB function `any` returns an array with entry 1 if that row in `indicator` contains non-zero entries and 0 otherwise. If `newindicator` and `indicator` are equal then the closure algorithm is complete. If not, all elements that are adjacent to the longest edges of elements in  $\mathcal{M}$  must be found and marked accordingly. Below is the code segment that performs the closure algorithm:

```
while ~isequal(newindicator,indicator)
    indicator=newindicator;
    longestedges=edges(indicator(:,1),[2 1]);
    tf=ismember(edges,longestedges,'rows');
    longestedgeindicator=reshape(tf,size(elements));
    longestedgeindicator(:,1)=any(longestedgeindicator,2);
    newindicator=or(newindicator,longestedgeindicator);
end
```

The output of the closure process is an array `indicator`, with a row for each element containing 1's and 0's which refer to the edges to be refined.

### A.2.4 Refinement

The refinement procedure is based on the information contained in the indicator array. Together with the `unique` function this allows for a vectorized approach to the refinement. Rather than considering each element and how it is to be refined, whole lists of elements can be refined at the same time. The whole procedure comprises of six components, five types of element refinement, and boundary refinement. To minimize memory usage, the refinement functions do not change `elements`, `dirichlet` or `neumann` directly. MATLAB passes by value only those arguments that a function modifies, so since they are only referred to in the function, not modified, they are passed by reference. The refinement functions only use values from `elements`, `dirichlet` and `neumann` to construct new members. The main routine is responsible for removing those that have been altered. Each step adds new coordinates to the coordinate list. While these coordinates may already appear on the list, what is important is that the new elements and edges know where their nodes appear on that list, since in the administration of the degrees of freedom stage the `unique` function will take care of these details.

#### Refinement of the boundary

Refinement of the boundary edges where `dirichlet` and `neumann` conditions are in place begins by finding what edges that are being refined are in `dirichlet` or `neumann`. The arrays `D` and `N` are the indices of these edges in the lists `dirichlet` and `neumann` respectively.

```
% Boundary Refinement
markededges=edges([indicator(:,1);indicator(:,2);indicator(:,3)],:);

[tfd,D]=ismember(markededges,dirichlet,'rows');
[tfn,N]=ismember(markededges,neumann,'rows');
D(D==0)=[];
N(N==0)=[];
[newdirichlet,newcoordinates]=boundaryrefine(D,dirichlet,coordinates);
[newneumann,newcoordinates]=boundaryrefine(N,neumann,newcoordinates);
```

```

dirichlet(D,:)=[];
neumann(N,:)=[];
newdirichlet=[dirichlet;newdirichlet];
newneumann=[neumann;newneumann];

```

The function `boundaryrefine` takes in turn the index vectors `D` and `N` with the respective sets of edges `dirichlet` and `neumann` and performs the bisection on the required edges. The first call reads in the list `coordinates` while the second reads in `newcoordinates`. Each refinement step after the refinement of the Dirichlet boundary takes as an input `newcoordinates` as it is important that the edges or elements created maintain connectivity to their correct nodes. This is important for the approach to administration of the degrees of freedom. The edges that have been bisected are then removed and the new lists formed.

The function `boundaryrefine` illustrates the principles behind the implementation of the element mesh refining functions so we explain it in detail here. The size of the input coordinates array measured is denoted by `S`, as it is from this number that the new nodes will be numbered from. The new coordinate array is created for all edges simultaneously and the size of it is used to create an array running from 1 to `n`. This array is to used number the new coordinates, by adding to each entry the number of original coordinates, therefore giving each new node a, a unique number. The connectivity between the new edges and the nodes is achieved by using the `index` array as a row selector.

```

function [newedges, newcoordinates]=
    ...boundaryrefine(index,edges,coordinates)

% Boundary Refinement
S=size(coordinates,1); newcoordinates=[];newedges=[];

if ~isempty(index)
    newcoordinates=0.5*(coordinates(edges(index,1),:))...
        +coordinates(edges(index,2),:);
    n=[1:size(newcoordinates,1)]';
    newedges=[edges(index,1), S+n;

```

```

        S+n, edges(index,2)];
end

newcoordinates=[coordinates; newcoordinates];

```

### Element refinement

The refinement of the elements uses five functions each carrying out the type of refinement specified in the indicator matrix. The indices of these elements are found via the MATLAB `find` function. For example

```
find(indicator(:,1)==1 & indicator(:,2)==0 & indicator(:,3)==1)
```

finds the numbers of those elements that are to have their longest and third edges bisected. Once the lists have been resolved, they are each passed into their respective refinement functions:

1. `refine_marked`

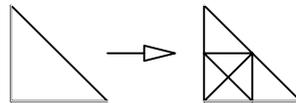


Figure A.5: Six triangle longest edge bisection.

2. `refine_111`



Figure A.6: Four triangle longest edge bisection.

3. `refine_110`

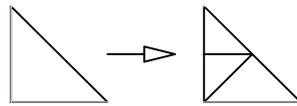


Figure A.7: First (longest) and second edge bisection.

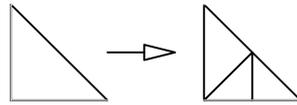


Figure A.8: First and third edge bisection.

4. refine\_101

5. refine\_100

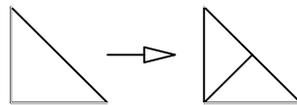


Figure A.9: First edge bisection.

Each function takes in and adds more coordinates to the list `newcoordinates`. This ensures that the final node numbering will remain consistent. The refining process is carried out in a vectorized format, for example,

```
function [nelements,newcoordinates]=...
    refine_111(marked,elements,coordinates)

%REFINE_111 Four triangle longest edge bisection.
% New coordinates
newcoords=0.5*(coordinates(elements(marked,[1 2 3]),:)+...
               coordinates(elements(marked,[2 3 1]),:));
newcoordinates=[coordinates;newcoords];

% New node numbers
```

```

t=size(newcoords,1)/3;
z=[1:t]';
N=size(coordinates,1);

% New elements
newelements=[elements(marked,1),      N+z      ,      N+2*t+z;
             N+z,                      elements(marked,2),  N+t+z;
             elements(marked,3),      N+z      ,      N+t+z;
             N+z,                      elements(marked,3),  N+2*t+z;];

```

As with the boundary refinement procedure, the new coordinates are created, the linear array that will give the new numbers is used, then the explicit construction of the new elements is carried out.

### A.2.5 Administration of the degrees of freedom

As already mentioned, the administration of the degrees of freedom is taken care of via the `unique` function. The elements that have been refined are removed, and new ones assembled. The call to the MATLAB function `unique` returns a sorted, unique list of coordinates, together with an index vector `J`. The vector `J` is the same size as the input array, and each entry contains the position that each entry from the input, occupies in the sorted unique list. Applying `J(·)` to the arrays `elements`, `dirichlet` and `neumann` results in the node numbers corresponding to the rows of the sorted unique coordinates list replacing the previous ones.

```

elements(indicator(:,1),:)=[];
newelements=[elements;newmarked;newelem111;
             newelem110;newelem101;newelem100];
[newcoordinates,I,J]=unique(newcoordinates,'rows');
newelements=J(newelements);
newdirichlet=J(newdirichlet);
newneumann=J(newneumann);

```

### A.2.6 Towards full adaptivity

A truly adaptive meshing package should contain the facility to coarsen or de-refine. We have not included such procedures here, since it appears to be necessary to implement additional data structures similar to those in [62], which we outline in the following.

It is reasonable to assume the existence of a mesh  $(\mathcal{T}_0, \mathcal{E}_0, \mathcal{N}_0)$  which can be thought of as the coarsest possible triangulation capturing the relevant geometrical features of the domain. In practice  $(\mathcal{T}_0, \mathcal{E}_0, \mathcal{N}_0)$  is probably a mesh supplied from a CAD package. If we are to coarsen, we must still ensure that geometrical conformity is maintained. It seems that the most suitable way of doing this is to treat coarsening as an inverse of refinement, in the sense that coarsening operations are undoing refinements that have already taken place. Considering the type of coarsening that can take place, it appears that there are only two scenarios where edges can be removed:

1. The configuration shown in figure A.10. When two elements with an edge on the boundary share an edge which has a node on the boundary not contained in  $\mathcal{N}_0$ , the node  $z$  can be removed without any concerns over breaking conformity of the mesh.

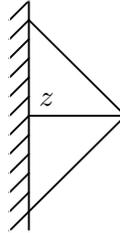


Figure A.10: Configuration of elements on the boundary for which coarsening can take place.

2. The configuration shown in figure A.11. When four elements share a node  $z$  not in  $\mathcal{N}_0$ , then  $z$  can be removed. However there are two possibilities for the resulting configuration. Either the line  $db$ , or  $ac$  is removed.

To ensure that a sequence of refinements, followed by a sequence of coarsening arrives at  $(\mathcal{T}_0, \mathcal{E}_0, \mathcal{N}_0)$ , it is necessary to in some way to store the refinements that have taken place.

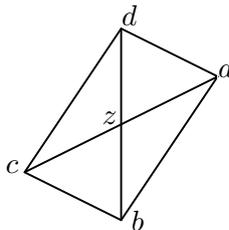


Figure A.11: Configuration of elements which coarsening can take place.

The process then of coarsening is then the reconstruction of elements from the previous mesh that have previously been refined. A solution to this problem is to introduce a hierarchic tree to describe the mesh. The root is the mesh  $(\mathcal{T}_0, \mathcal{E}_0, \mathcal{N}_0)$ , then a refinement generates branches from each element that has been refined. For further details see [62].

### A.3 Summary

In summary, we have shown how local mesh refinement can be implemented MATLAB in a clear and flexible manner. The code is fully compatible with the mesh generation and finite element codes already presented in the literature. For time dependent problems and optimal solution of stationary problems, mesh coarsening is also a requirement. This is the subject of further work. The closure algorithm is independent of the spatial dimension of the domain, and we believe that with suitable adjustments, the approach might also prove suitable for tetrahedral meshes.

### A.4 Main routine

```

1  function [newcoordinates,newelements,newdirichlet,newneumann]=...
2          refinemesh(coordinates,elements,dirichlet,neumann,marked)
3
4  %REFINEMESH Mesh Refinement algorithm
5  %   Refines a two dimensional geometrically conforming triangular mesh
6  %   using longest edge bisection with new interior node. MARKED is an
7  %   index array into ELEMENTS of those elements to undergo interior
8  %   node refinement.
9  %   Requires Inputs:

```

```
10 % COORDINATES - [x,y] is an array of coordinates, each row
11 %             corresponding to a node.
12 % ELEMENTS - [n1,n2,n3] is an array containing a row for each element.
13 %             The entries are the row numbers of each node in
14 %             coordinates listed in an anticlockwise manner, longest
15 %             edge first.
16 % DIRICHLET - [n1,n2] is an array of edges with each row representing
17 %             an edge, the entries are the node numbers.
18 % NEUMANN - [n1,n2] is an array of edges with each row representing an
19 %             edge, the entries are the node numbers.
20
21 % See report Adaptive Mesh Refinement in MATLAB for full details
22
23 % Author H.R. Hill
24 % Date 10 October 2005
25
26 %CLOSURE
27
28 % Create arrays containing the mesh edges and those that are marked
29 % and compare them to make the initial indicator matrix.
30 edges=[elements(:, [1,2]);elements(:, [2,3]);elements(:, [3,1])];
31 markededges=[elements(marked, [1,2]);
32             elements(marked, [2,3]);
33             elements(marked, [3,1])];
34 markededges=[markededges;markededges(:, [2,1])];
35 tf=ismember(edges,markededges,'rows');
36 indicator=reshape(tf,size(elements));
37
38 % Mark the longest edge of adjacent elements to those already marked
39 newindicator=indicator;
40 newindicator(:,1)=any(indicator,2);
41
42 % Closure loop
43 while ~isequal(newindicator,indicator)
44     indicator=newindicator;
45     longestedges=edges(indicator(:,1),[2 1]);
46     tf=ismember(edges,longestedges,'rows');
```

```
47     longestedgeindicator=reshape(tf,size(elements));
48     longestedgeindicator(:,1)=any(longestedgeindicator,2);
49     newindicator=or(newindicator,longestedgeindicator);
50 end
51
52 %REFINEMENT
53
54 % Boundary Refinement
55 % Look for marked edges in dirichlet and neumann
56 markededges=edges([indicator(:,1);indicator(:,2);indicator(:,3)],:);
57 [tfd,D]=ismember(markededges,dirichlet,'rows');
58 [tfn,N]=ismember(markededges,neumann,'rows');
59 D(D==0)=[];
60 N(N==0)=[];
61
62 % Refine the required edges using trhe boundary refine function
63 [newdirichlet,newcoordinates]=...
64     boundaryrefine(D,dirichlet,coordinates);
65 [newneumann,newcoordinates]=...
66     boundaryrefine(N,neumann,newcoordinates);
67
68 % Remove the old edges and add the new ones.
69 dirichlet(D,:)=[];
70 neumann(N,:)=[];
71 newdirichlet=[dirichlet;newdirichlet];
72 newneumann=[neumann;newneumann];
73
74 % Element Refinement.
75 % Sort into groups for the different types of refinement
76 m111=setdiff(find(all(indicator,2)),marked);
77 m110=find(indicator(:,1)==1 & indicator(:,2)==1 & indicator(:,3)==0);
78 m101=find(indicator(:,1)==1 & indicator(:,2)==0 & indicator(:,3)==1);
79 m100=find(indicator(:,1)==1 & indicator(:,2)==0 & indicator(:,3)==0);
80
81 % Refine the groups
82 [newmarked,newcoordinates]=...
83     refine_marked(marked,elements,newcoordinates);
```

---

```

84 [newelem111,newcoordinates]=refine_111(m111,elements,newcoordinates);
85 [newelem110,newcoordinates]=refine_110(m110,elements,newcoordinates);
86 [newelem101,newcoordinates]=refine_101(m101,elements,newcoordinates);
87 [newelem100,newcoordinates]=refine_100(m100,elements,newcoordinates);
88
89 %ADMINISTRATION OF DEGREES OF FREEDOM
90 % Remove elements that have been refined
91 elements(indicator(:,1),:)=[];
92 % Create new element list
93 newelements=[elements;newmarked;newelem111;
94             newelem110;newelem101;newelem100];
95 [newcoordinates,I,J]=unique(newcoordinates,'rows');
96 newelements=J(newelements); newdirichlet=J(newdirichlet);
97 newneumann=J(newneumann);
98
99 % BUG FIX Stops 1x2 arrays being switched to 2x1
100 if size(newdirichlet,2)==1
101     newdirichlet=newdirichlet';
102 end
103 if size(newneumann,2)==1
104     newneumann=newneumann';
105 end
106
107 %SUBFUNCTIONS
108 % These are the individual refinement functions for the different
109 % types of refinement and the boundary refinement. All essentially
110 % follow the same process of creating the new coordinates, the new
111 % node numbers and assembling the new elements.
112
113 %REFINE_MARKED New interior node refinement
114 function [newelements,newcoordinates]=...
115         refine_marked(marked,elements,coordinates)
116
117 % New coordinates
118 newedgecoords=0.5*(coordinates(elements(marked,[1 2 3]),:))...
119                 +coordinates(elements(marked,[2 3 1]),:));
120 newintcoords=0.5*(coordinates(elements(marked,3),:))...

```

```
121         +newedgecoords([1:end/3,:]);
122 newcoordinates=[coordinates;newedgecoords;newintcoords];
123
124 % New node numbers
125 t=size(newedgecoords,1)/3;
126 z=[1:t]';
127 N=size(coordinates,1);
128
129 % New elements
130 newelements=[elements(marked,1),      N+z      ,      N+2*t+z;
131             N+z,                      elements(marked,2),  N+t+z;
132             N+t+z,                    elements(marked,3),  N+3*t+z;
133             elements(marked,3),      N+2*t+z      ,      N+3*t+z;
134             N+2*t+z,                  N+z          ,      N+3*t+z;
135             N+z,                      N+t+z          ,      N+3*t+z];
136
137
138 %REFINE_111 Four triangle longest edge bisection.
139 function
140 [newelements,newcoordinates]=...
141     refine_111(marked,elements,coordinates)
142
143 % New coordinates
144 newcoords=0.5*(coordinates(elements(marked,[1 2 3]),:))...
145     +coordinates(elements(marked,[2 3 1]),:));
146 newcoordinates=[coordinates;newcoords];
147
148 % New node numbers
149 t=size(newcoords,1)/3;
150 z=[1:t]';
151 N=size(coordinates,1);
152
153 % New elements
154 newelements=[elements(marked,1),      N+z      ,      N+2*t+z;
155             N+z,                      elements(marked,2),  N+t+z;
156             elements(marked,3),      N+z          ,      N+t+z;
157             N+z,                      elements(marked,3),  N+2*t+z];
```

```
158
159
160 %REFINE_110 Bisects longest and second edges
161 function [newelements,newcoordinates]=...
162         refine_110(marked,elements,coordinates)
163
164 % New coordinates
165 newcoords=0.5*(coordinates(elements(marked,[1 2]),:))
166         +coordinates(elements(marked,[2 3]),:));
167 newcoordinates=[coordinates;newcoords];
168
169 % New node numbers
170 t=size(newcoords,1)/2;
171 z=[1:t]';
172 N=size(coordinates,1);
173
174 % New elements
175 newelements=[elements(marked,3),         elements(marked,1),   N+z;
176             elements(marked,3),         N+z             , N+t+z;
177             N+z,             elements(marked,2),   N+t+z];
178
179
180 %REFINE_101 Bisects longest and third edges
181 function [newelements,newcoordinates]=...
182         refine_101(marked,elements,coordinates)
183
184 % New coordinates
185 newcoords=0.5*(coordinates(elements(marked,[1 3]),:))...
186         +coordinates(elements(marked,[2 1]),:));
187 newcoordinates=[coordinates;newcoords];
188
189 % New node numbers
190 t=size(newcoords,1)/2;
191 z=[1:t]';
192 N=size(coordinates,1);
193
194 % New elements
```

```

195 newelements=[N+z,          elements(marked,3),  N+t+z;
196           elements(marked,1),      N+z,      N+t+z;
197           elements(marked,2), elements(marked,3),  N+z];
198
199 %REFINE_100 Bisects longest edges
200 function [newelements,newcoordinates]=...
201           refine_100(marked,elements,coordinates)
202
203 % New coordinates
204 newcoords=0.5*(coordinates(elements(marked,1),:)...
205           +coordinates(elements(marked,2),:));
206 newcoordinates=[coordinates;newcoords];
207
208 % New node numbers
209 t=size(newcoords,1);
210 z=[1:t]';
211 N=size(coordinates,1);
212
213 % New elements
214 newelements=[elements(marked,3),  elements(marked,1), N+z;
215           elements(marked,2),  elements(marked,3),  N+z];
216
217 %BOUNDARY_REFINE
218 function [newedges,newcoordinates]=...
219           boundaryrefine(index,edges,coordinates)
220
221 S=size(coordinates,1);
222 newcoordinates=[];
223 newedges=[];
224
225 if ~isempty(index)
226     newcoordinates=0.5*(coordinates(edges(index,1),:)...
227           +coordinates(edges(index,2),:));
228     n=[1:size(newcoordinates,1)]';
229     newedges=[edges(index,1), S+n;
230           S+n, edges(index,2)];
231 end

```

232

233 newcoordinates=[coordinates; newcoordinates];

# Bibliography

- [1] Adams R. A. and J.J.F. Fournier, *Sobolev Spaces*, 2nd ed., Pure and Applied Mathematics, vol. 140, Academic Press, Amsterdam, 2003.
- [2] M. Ainsworth and J. T. Oden, *A Posteriori Error Estimation in Finite Element Analysis*, Pure and Applied Mathematics, Wiley, New York, 2000.
- [3] H-D. Alber, *Materials with Memory: Initial-Boundary Value Problems for Constitutive Equations with Internal Variables*, Lecture Notes in Mathematics, vol. 1682, Springer, Berlin, 1998.
- [4] J. Albery, C. Carstensen, and S. A. Funken, *Remarks around 50 lines of Matlab: short finite element implementation*, Numerical Algorithms **20** (1999), 117–137.
- [5] J. Albery, C. Carstensen, S.A. Funken, and R. Klose, *Matlab implementation of the finite element method in elasticity*, Computing **69** (2002), no. 3, 239–263.
- [6] S.S. Antman, *Nonlinear Problems of Elasticity*, Applied mathematical sciences, vol. 107, Springer-Verlag, New York, 1995.
- [7] K. Atkinson and W. Han, *Theoretical Numerical Analysis*, 2nd ed., Springer, New York, 2005.
- [8] A. K. Aziz and P. Monk, *Continuous finite elements in space and time for the heat equation*, Mathematics of Computation **52** (1989), no. 186, 255–274.
- [9] I. Babuška, R. Durán, and R. Rodríguez, *Analysis of the efficiency of an a posteriori error estimator for linear triangular finite elements*, SIAM J. Numer. Anal **29** (1992), 947–964.

- [10] I. Babuška and M. Vogelius, *Feedback and adaptive finite element solution of one dimensional boundary value problems*, Numerische Mathematik **44** (1984), 75–102.
- [11] W. Bangerth and R. Rannacher, *Adaptive Finite Element Methods for Differential Equations*, Lectures in Mathematics, Birkhäuser Verlag, Basel, 2003.
- [12] D. M. Bedivan and G. J. Fix, *Analysis of finite element approximation and quadrature of Volterra integral equations*, Numerical methods Partial Differential Equations (1997), no. 13, 663–672.
- [13] P. Binev, W. Dahmen, and R. DeVore, *Adaptive finite element methods with convergence rates*, Numerische Mathematik **97** (2004), 219 – 268.
- [14] M. Braack and A. Ern, *A posteriori control of modeling errors and discretization errors*, SIAM Multiscale Modeling and Simulation **1** (2003), no. 2, 221–238.
- [15] S. C. Brenner and C. Carstensen, *Finite Element Methods*, ch. 4, John Wiley and Sons, New York, 2004.
- [16] S.C. Brenner and L. R. Scott, *The Mathematical Theory of Finite Element Methods*, Texts in Applied Mathematics, Springer, New York, 2002.
- [17] M. Buch, A. Idesman, R. Niekamp, and E. Stein, *Finite elements in space and time for parallel computing of viscoelastic deformation*, Comput. Mech. **24** (1999), 386–395.
- [18] C. Carstensen, *Quasi-interpolation and a posteriori estimates in finite element methods*, Mathematical Modelling and Numerical Analysis **33** (1999), no. 6, 1187–1202.
- [19] C. Carstensen and S.A. Funken, *Constants in Clément-interpolation error and residual based a posteriori estimates in finite element methods*, East-West Journal of Numerical Analysis **8** (2000), no. 3, 153–175.
- [20] C. Carstensen and R. Klose, *Elastoviscoplastic finite element analysis in 100 lines of Matlab*, J. Numer. Math. **10** (2002), no. 3, 157–192.
- [21] R. M. Christensen, *Theory of Viscoelasticity, An Introduction*, Academic Press, New York, 1971.

- [22] P.G. Ciarlet, *Mathematical Elasticity. Volume 1: Three Dimensional Elasticity*, Studies in Mathematics and its Applications, vol. 20, North-Holland, Amsterdam, 1993.
- [23] ———, *The Finite Element Method for Elliptic Problems*, Classics in Applied Mathematics, vol. 40, SIAM, Philadelphia, 2002.
- [24] P. Clément, *Approximation by finite element functions using local regularization*, RAIRO, Série rouge **Analyse Numérique** (1975), no. R-2, 77–84.
- [25] B.D. Coleman and W. Noll, *Foundations of linear viscoelasticity*, Reviews of Modern Physics **33** (1961), 239–249.
- [26] R. Courant, *Variational methods for the solution of problems of equilibrium and vibrations*, Bulletin of the American Mathematical Society **49** (1943), 1–23.
- [27] W. A. Day, *The Thermodynamics of Simple Materials*, Springer Tracts in Natural Philosophy, vol. 22, Springer - Verlag, New York, 1972.
- [28] W Dörfler, *A convergent adaptive algorithm for Poisson's equation*, SIAM J. Numer. Anal. **33** (1996), 1106–1124.
- [29] G. Duvaut and J. L. Lions, *Inequalities in Mechanics and Physics*, Springer, Berlin, 1976.
- [30] K. Eriksson, D. Estep, P. Hansbo, and C. Johnson, *Introduction to adaptive methods for differential equations*, Acta Numerica (1995), 105–158.
- [31] K. Eriksson and C. Johnson, *Adaptive finite element methods for parabolic problems. I: a linear model problem*, SIAM J. Numer. Anal **28** (1991), 43–77.
- [32] ———, *Adaptive finite element methods for parabolic problems. II: optimal error estimates in  $L_\infty L_2$  and  $L_\infty L_\infty$* , SIAM J. Numer. Anal **32** (1995), 706–740.
- [33] A. Ern and J-L. Guermond, *Theory and Practice of Finite Elements*, Applied Mathematical Sciences, vol. 40, Springer-Verlag, New York, 2002.

- [34] D. Estep and D. French, *Global error control for the continuous Galerkin method for ordinary differential equations*, Mathematical modelling and numerical analysis **28** (1994), no. 7, 815–852.
- [35] D. J. Estep, M. G. Larson, and R. D. Williams, *Estimating the error of numerical solutions of systems of reaction-diffusion equations*, Memoirs A.M.S. **146** (2000), 1–109.
- [36] L. C. Evans, *Partial Differential Equations*, Graduate Studies in Mathematics, vol. 19, AMS, Rhode Island, 1998.
- [37] R.E. Ewing, *Time-stepping Galerkin methods for nonlinear Sobolev partial differential equations*, SIAM Journal on Numerical Analysis **15** (1978), no. 6, 1125–1150.
- [38] M. Fabrizio and A. Morro, *Mathematical Problems in Linear Viscoelasticity.*, Studies in Applied Mathematics, vol. 12, SIAM, Philadelphia, 1992.
- [39] D. French and T. E. Peterson, *A continuous space-time finite element method for the wave equation*, Mathematics of Computation **65** (1996), no. 214, 491–506.
- [40] J. M. Golden and G. A. C. Graham, *Boundary Value Problems in Linear Viscoelasticity*, Springer-Verlag, Berlin, 1988.
- [41] G. H. Golub and C. F Van Loan, *Matrix Computations*, 3rd ed., The John Hopkins University Press, Baltimore, 1996.
- [42] P. Grisvard, *Elliptic Problems in Nonsmooth Domains*, 2nd ed., Pitman Advanced Publishing program, Pitman, Boston, 1985.
- [43] D. Handscomb, *Errors of linear interpolation on a triangle*, Tech. report, Oxford University, 1995.
- [44] P. Hansbo and C. Johnson, *Adaptive finite element methods in computational mechanics*, Computer Methods in Applied Mechanics and Engineering. (1992), no. 101, 143–181.

- [45] M. Holst, *Mclite: An adaptive multilevel finite element matlab package for scalar nonlinear elliptic equations in the plane*, Tech. report, UCSD, <http://cam.ucsd.edu/~mholst/pubs/dist/mclite.pdf>, 2000.
- [46] C. O. Horgan, *Korn's inequalities and their applications in continuum mechanics*, SIAM Review **37** (1995), no. 4, 491–511.
- [47] C. O. Horgan and J.K. Knowles, *Eigenvalue problems associated with Korn's inequalities*, Archives for Rational Mechanics and Analysis **40** (1971), 384–402.
- [48] A. R. Johnson and A. Tessler, *A viscoelastic high order beam finite element*, pp. 333–345, Wiley, Chichester, 1997, In J. R. Whiteman, Editor, The Mathematics of Finite Elements and its Applications. MAFELAP.
- [49] C. Johnson, *Error estimates and adaptive time-step control for a class of one-step methods for stiff ordinary differential equations*, SIAM J. Numer. Anal **25** (1988), 908–926.
- [50] P. Linz, *Analytical and Numerical Methods for Volterra Equations*, Studies in Applied Mathematics, SIAM, Philadelphia, 1985.
- [51] Křížek M. and P Neittaanmäki, *Finite Element Approximation of Variational Problems and Applications*, Pitman Monographs and Surveys in Pure and Applied Mathematics, 1990, New York, 2005.
- [52] K. Mekchay and R. H. Nochetto, *Convergence of adaptive finite element methods for general second order linear elliptic pde*, Preprint, <http://www.math.umd.edu/~rhn/publications.html>.
- [53] K-S. Moon, E. von Schwerin, A. Szepessy, and R. Tempone, *Convergence rates for an adaptive dual weighted residual finite element algorithm*, Preprint, <http://www.nada.kth.se/~szepessy/papers.html>.
- [54] P. Morin, R. H. Nochetto, and K. G. Siebert, *Data oscillation and convergence of adaptive FEM*, SIAM J. Numer. Anal. **38** (2000), no. 2, 466–488.

- [55] ———, *Convergence of adaptive finite element methods*, SIAM Review **44** (2002), no. 4, 631–658.
- [56] R.H. Nochetto, *Removing the saturation assumption in a posteriori error analysis*, Istit. Lombardo Accad. Sci. Lett. Rend. A **127** (1993), 67–82.
- [57] W. Noll, *A new mathematical theory of simple materials*, Archives for Rational Mechanics and Analysis **48** (1972), 1–50.
- [58] J. T. Oden and S. Prudhomme, *Adaptive modeling in computational mechanics*, Journal of Computational Physics **182** (2002), 496–515.
- [59] P.-O. Persson and G. Strang, *A simple mesh generator in matlab*, SIAM Review **46** (2004), no. 2, 329–345.
- [60] A. Quarteroni and A. Valli, *Numerical Approximation of Partial Differential Equations*, 2nd ed., Springer Series in Computational Mathematics, vol. 23, Springer, Berlin, 1997.
- [61] R. Rannacher, *Adaptive Galerkin finite element methods for partial differential equations*, Journal of Computational and Applied Mathematics **128** (2001), 205–233.
- [62] A. Schmidt and K. G. Kunibert, *Design of Adaptive Finite Element Software: The Finite Element Toolbox ALBERTA*, Lecture note on computational science and engineering, Springer, Berlin, 2005.
- [63] C. Schwab, *p- and hp-Finite Element Methods. Theory and Applications to Solid and Fluid Mechanics*, Oxford University Press, Oxford, 1998.
- [64] R. Scott and Zhang S., *Finite element interpolation of non smooth functions satisfying boundary conditions*, Mathematics of Computation **54** (1990), no. 190, 483–493.
- [65] S. Shaw, M. K. Warby, and J. R. Whiteman, *Error estimates with sharp constants for a fading memory Volterra problem in linear solid viscoelasticity*, SIAM Journal of Numerical Analysis (1997), no. 34, 1237–1254.

- [66] S. Shaw, M. K. Warby, J. R. Whiteman, C. Dawson, and M. F. Wheeler, *Numerical techniques for the treatment of quasistatic viscoelastic stress problems in linear isotropic solids*, *Comput. methods in Appl. Mech. Engrg.* (1994), no. 118, 211–237.
- [67] S. Shaw and J. R. Whiteman, *Adaptive space-time finite element solution for Volterra equations arising in viscoelasticity problems*, *J. Comput. Appl. Math.* (2000), no. 125, 337–345.
- [68] ———, *Negative norm error control for second-kind convolution Volterra equations*, *J. Comput. Appl. Math.* (2000), no. 125, 337–345.
- [69] ———, *Numerical solution of linear quasistatic hereditary viscoelasticity problems*, *SIAM Journal of Numerical Analysis* (2000), no. 38, 80–97.
- [70] ———, *A posteriori error estimates for space-time finite element approximation of quasistatic hereditary linear viscoelasticity problems*, *Comput. Methods Appl. Mech. Engrg.* (2004), no. 193, 5551–5572.
- [71] J. C. Simo and T.J.R. Hughes, *Computational Inelasticity*, *Interdisciplinary Applied Mathematics*, vol. 7, Springer, New York, 1998.
- [72] R. Stevenson, *An optimal adaptive finite element method*, *SIAM Journal of Numerical Analysis* **42** (2005), no. 5, 2188–2217.
- [73] G. Strang, *Variational crimes in the finite element method.*, Academic Press, New York, 1972, In: A. Aziz ed., *The Mathematical Foundations of the Finite Element Method with Applications to Partial Differential Equations*.
- [74] J. P. Suárez, A. Plaza, and G.F. Carey, *Propagation path properties in iterative longest-edge refinement*, *Proceedings, 12th International Meshing Roundtable*, Sandia National Laboratories (2003), 79–90.
- [75] E. Süli and P. Houston., *Finite element methods for hyperbolic problems: a posteriori error analysis and adaptivity.*, pp. 441–471, Oxford University Press, Oxford, 1997, In: I. Duff and G.A. Watson, eds., *State of the Art in Numerical Analysis*.

- 
- [76] C. Truesdell and W. Noll, *The non-linear field theories of mechanics*, vol. III/1, Springer, Berlin, 1961.
- [77] R. Verfürth, *A Review of a posteriori Error Estimation and Adaptive Mesh Refinement Techniques*, Teubner-Wiley, Stuttgart, 1996.
- [78] ———, *A review of a posteriori error estimation techniques for elasticity problems.*, In: *On new Advances in Adaptive Computational Methods in Mechanics* (P. Ladeve, J. T. Oden; eds.) (1998), 257 – 274.
- [79] ———, *Error estimates for some quasi-interpolation operators*, *Mathematical Modelling and Numerical Analysis* **33** (1999), no. 4, 695–713.
- [80] ———, *On the constants in some inverse inequalities for some finite element functions*, Tech. report, Universität Bochum, 1999.
- [81] M. Šilhavý, *The Mechanics and Thermodynamics of Continuous Media*, Texts and Monographs in Physics, Springer, Berlin, 2002.