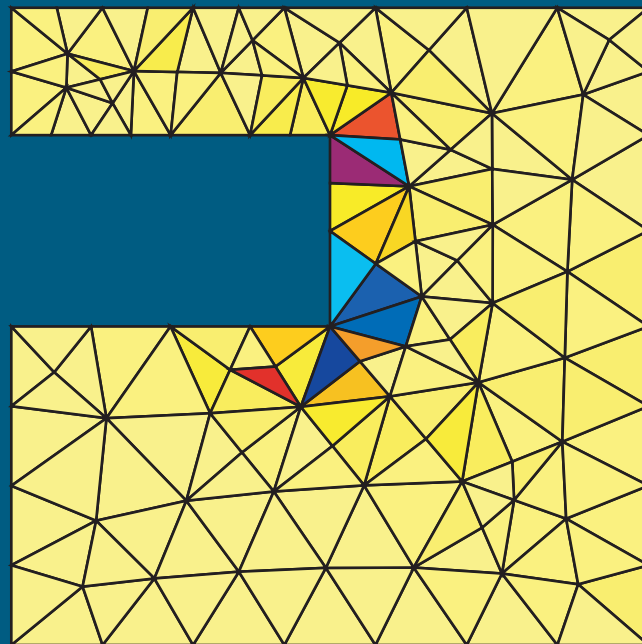


# ERROR ASSESSMENT FOR FUNCTIONAL OUTPUTS OF PDE'S:

bounds and goal-oriented adaptivity

Núria Parés Mariné

---



Doctoral Thesis  
Barcelona, January 2005



# ERROR ASSESSMENT FOR FUNCTIONAL OUTPUTS OF PDE'S:

bounds and goal-oriented adaptivity

Núria Parés Mariné

---



Doctoral Thesis

Advisors: Pedro Díez and Antonio Huerta

Barcelona, January 2005

Departament de Matemàtica Aplicada III

Programa de Doctorat de Matemàtica Aplicada



*Al Josep, la Dolors i l'Oriol,*  
*al Francesc,*  
*i molt especialment, a la padrina.*



# Acknowledgements

Voldria donar el meu més sincer agraïment a totes aquelles persones que amb el seu suport i confiança han fet possible la realització d'aquesta tesi.

En primer lloc voldria donar les gràcies a les persones sense les quals res d'això hagués estat possible, al Pedro Díez i a l'Antonio Huerta. Per la seva direcció, per les seves idees, pel seu incondicional recolzament i sobretot per haver dedicat tantes estones a formar-me no sols com a investigadora sinó també com a persona durant aquest llarg procés. *Moltíssimes gràcies!*

En segon lloc no em voldria oblidar de les persones que van fer possible que comencés aquest llarg camí: al Ramón Nolla per haver estat un gran professor durant la meua època d'estudiant d'institut i per haver-me recolzat en els moments en què havia de decidir quin havia de ser el meu futur. Als professors de la Facultat de Matemàtiques i Estadística per fer-me gaudir durant els quatre anys que vaig estar a la Facultat del món de les matemàtiques. En especial voldria agrair als professors Joan Solà-Morales, Pere Pascual i Marta València el seu recolzament. També voldria agrair-li a la Sílvia la seva inesgotable paciència durant aquelles llargues hores d'estudi i la seva amistat.

Tampoc no voldria oblidar-me dels companys d'Urgell amb qui he compartit les primeres classes i molt bons moments i dels companys del grup de recerca LaCàN per les llargues estones compartides, en especial a l'Antonio Rodríguez, al Pep, a la Sònia, a l'Agustí, al Xevi i a l'Imma.

Al Jaume Peraire per haver-me donat l'oportunitat de fer dues estades al MIT. Per les inacabables idees i les llargues converses. Gràcies per haver-me acollit tan bé junt a l'Anna, l'Anton, l'Oleg i el Jaume. Em va fer sentir com a casa tot i estar tan lluny.

A la Yolanda. Des d'aquell primer dia de residència quan ens vam conèixer a les fosques, ara ja fa més de nou anys, has estat *la* meva gran amiga. Hem compartit amics, estudis, pis, doctorat... èpoques bones, èpoques no tan bones i èpoques millors. Espero poder seguir compartint tot el que ens espera d'ara en endavant. Gràcies mils. ♪ *In these moments, moments of our lives... you and I sharing this time together, sharing the same dream, as the time goes by...* ♪

A la meva família: als meus avis Joan i Lolita i Pere i Antonieta, al meu tiet Toni i a les meves cosines Anna i Laia i evidentment a la meva padrina a qui tantes coses he d'agrair. ♪ *Son aquellas pequeñas cosas, que nos dejó un tiempo de rosas en un rincón, en un papel o en un cajón...* ♪

Als meus pares: Dolors i Josep que han estat al meu darrera sempre encoratjant-me a seguir endavant en els moments de desànim i donant-me l'oportunitat d'arribar fins aquí. Evidentment sense el seu gran esforç tampoc no hagués pogut estar possible res d'això. Us estaré sempre infinitament agraïda. ♪ *I jo que m'adormia entre els teus braços amb la boca enganxada en el teu pit. L'amor d'un home ja ens havia unit abans d'aquell matí d'hivern en què vaig néixer. El record d'aquell temps, el vent no l'arrossega: quan estalviaves pa per a donar-me mantega...* ♪

Al meu germà, l'Ori, que ha sigut l'alegria més gran de la meva vida. Espero que em doni moltes més alegries. ♪ *Vull que em molestis, sóc aquí, parla'm o plora que no tinc cap altra feina aquesta nit...* ♪

I finalment, al Francesc. Per estar al meu costat en tot moment. Per fer possible aquesta i tantes altres coses... amb la seva paciència... ♪ *So many 25th's of December, just as many 4th of July's, and we're still holding it together It only comes down to you and I. I know you can still remember Things we said right from the start When we said that this could be special, I'm keeping those words deep down in my heart.* ♪

*Perquè tampoc és el mateix...*

Barcelona  
December, 2004

Núria Parés Mariné



# Contents

|  |           |
|--|-----------|
| <b>Acknowledgements</b>  | <b>v</b>  |
| <b>1 Introduction</b>  | <b>1</b>  |
| 1.1 Motivation . . . . .   | 1         |
| 1.2 Objectives of the thesis . . . . .   | 2         |
| 1.3 Overview . . . . .   | 6         |
| <b>2 Bounds for linear outputs of interest for a general variational problem</b> | <b>7</b>  |
| 2.1 Model problem . . . . .  | 8         |
| 2.1.1 Operator decomposition . . . . .   | 9         |
| 2.2 Energy reformulation . . . . .   | 10        |
| 2.2.1 Bounds for self-adjoint model problems . . . . .                           | 13        |
| 2.2.2 Bounds for nonself-adjoint model problems . . . . .                        | 14        |
| 2.3 Summary . . . . .  | 17        |
| <b>3 Bounds for the energy norm of solutions of self-adjoint model problems</b>  | <b>19</b> |
| 3.1 Model problem . . . . .  | 20        |
| 3.2 Upper bounds for the energy . . . . .  | 21        |
| 3.2.1 Sufficient condition for the upper bound property . . . . .                | 22        |
| 3.2.2 The equilibrated residual method . . . . .                                 | 23        |
| 3.3 Lower bounds for the energy . . . . .  | 27        |
| 3.3.1 Sufficient condition for the lower bound property . . . . .                | 28        |

|                       |  |            |
|-----------------------|--|------------|
| 3.3.2                 | Lower bounds by post-processing . . . . .  | 29         |
| 3.3.3                 | Optimization of the lower bounds . . . . .   | 31         |
| 3.4                   | Implementation issues . . . . .  | 35         |
| 3.4.1                 | Remarks on the derivation of upper bounds for the energy .                             | 35         |
| 3.4.2                 | Remarks on the derivation of lower bounds for the energy .                             | 37         |
| <b>4</b>              | <b>Subdomain-based <i>flux-free</i> a posteriori error estimator</b>                   | <b>41</b>  |
| 4.1                   | Model problem . . . . .  | 42         |
| 4.2                   | Upper bound for the energy . . . . .   | 43         |
| 4.3                   | Lower bound for the energy . . . . .   | 46         |
| 4.4                   | Computational aspects . . . . .  | 46         |
| 4.5                   | Numerical examples . . . . .   | 48         |
| <b>5</b>              | <b>Strict bounds for the energy norm of weak solutions to the elasticity equations</b> | <b>53</b>  |
| 5.1                   | Model problem . . . . .  | 54         |
| 5.2                   | Strict upper bound for the energy . . . . .  | 55         |
| 5.2.1                 | Sufficient condition for the upper bound property . . . . .                            | 58         |
| 5.3                   | Certification . . . . .  | 59         |
| 5.4                   | Numerical examples . . . . .   | 59         |
| 5.4.1                 | Square plate . . . . .   | 61         |
| 5.4.2                 | $J$ -integral . . . . .  | 64         |
| <b>6</b>              | <b>Conclusions</b>   | <b>71</b>  |
| 6.1                   | Future developments . . . . .  | 73         |
| <br><b>APPENDICES</b> |  |            |
| <b>A</b>              | <b>Bounds for linear outputs of interest for a general variational problem</b>         | <b>A–1</b> |
| A.1                   | Energy reformulation . . . . .   | A–1        |
| A.1.1                 | Minimization reformulation . . . . .   | A–2        |

---

|          |  |            |
|----------|--|------------|
| A.1.2    | Lagrange multiplier . . . . .  | A-3        |
| A.1.3    | Strong duality . . . . .   | A-4        |
| A.2      | Bounds for the error in the quantity of interest . . . . .                     | A-5        |
| A.2.1    | Bounds for self-adjoint model problems . . . . .                               | A-6        |
| A.2.2    | Bounds for nonself-adjoint model problems . . . . .                            | A-10       |
| <b>B</b> | <b>Comparison between subdomain-based <i>flux-free</i> residual methods</b>    | <b>B-1</b> |
| <b>C</b> | <b>Strict bounds for outputs of interest: bounds for the <i>J</i>-integral</b> | <b>C-1</b> |

## APPENDED PAPERS

### Paper A

Díez, Parés and Huerta (2003), ‘Recovering lower bounds of the error by postprocessing implicit residual a posteriori error estimates’, *International Journal for Numerical Methods in Engineering* 56 (10), 1465-1488.

### Paper B

Parés, Díez and Huerta (2004), ‘Subdomain-based flux-free a posteriori error estimators’, *Computer Methods in Applied Mechanics and Engineering* 194. In press.

### Paper C

Parés, Bonet, Huerta and Peraire (2004), ‘The computation of bounds for linear-functional outputs of weak solutions to the two-dimensional elasticity equations’, *Computer Methods in Applied Mechanics and Engineering* 194. In press.



# Chapter 1

## Introduction

### 1.1 Motivation

Computational methods are today an inherent process in engineering design. Numerical simulation of physical phenomena allows to reduce the cost of product development at the same time that assessing the quality of the final product, providing faster, flexible and inexpensive alternative to experiments.

The systems under consideration are usually modelled using a set of partial differential equations (mathematical model) which are then discretized and solved using numerical methods. However, no matter how sophisticated or how appropriate the mathematical model can be for characterizing the physical phenomena and how accurate the numerical methods involved in the discretization and approximation process might be, all computational results will be in error. Therefore, before using any numerical simulation in a decision making process, one has to decide whether the computational results are reliable or not: are we solving the right model? and are we solving the model right?

Validation and verification are thus crucial in order to use computer simulations as reliable tools for engineering design. Validation accounts for the modelling error, that is, the assessment of the error introduced by approximating the physical problem by a mathematical model, while verification accounts for the error introduced by solving the continuum models using numerical approximations.

Engineering applications typically consist in studying a physical phenomena in order to predict certain quantities relevant to the analysis such as averages of the

solution, flow rates, velocities or shear stress at a given critical point in the domain. It is frequent to study if a design meets the security requirements or to study how to modify a design in order to improve its performance requirements, which are the quantities of analysis. To approximate these quantities, numerical approximations of the physical phenomena are used. Therefore, the accuracy of the numerical results is given by its capacity to provide reliable quantitative information about the quantities of interest also called outputs. Obtaining an approximated solution with a global prescribed accuracy is not the main goal but rather the control of the error in the output, which represents the relevant engineering quantity.

In this context, the validation and verification techniques must be capable of ensuring that the quantities predicted by the mathematical model agree with the quantities obtained in experimental results and also that the discretized mathematical model is solved with sufficient precision so that the error in the predicted quantities meet the accuracy requirements. Moreover, in the case where the error in the output is too large, it has to provide information on how to modify the mathematical model or the discretization procedure in order to be able to achieve the desired accuracy.

## 1.2 Objectives of the thesis

This thesis is focused on the verification of numerical results, that is, in the evaluation of the errors introduced in the discretization process of transforming the mathematical model into a numerical problem. The goal is to control and assess the discretization error, which is the difference between the exact solution of the mathematical model,  $u$ , and the solution of its discretized approximation,  $u_H$ . Validation, that is assessing the difference between the exact solution of the mathematical model and the physical phenomenon under consideration is out of the scope of this thesis. In particular, special interest is paid in the assessment of the discretization error not only in a global norm but in a particular quantity of interest.

The quantities of interest are typically functionals of the field variables  $u$ ,  $\ell^{\mathcal{O}}(u)$ . These field variables are approximated using numerical schemes, such as the finite element method, yielding an approximated solution  $u_H$ , which in turn yields the approximation of the quantity of interest  $\ell^{\mathcal{O}}(u_H)$ . The goal is then to be able to assess

the error in this approximation  $\ell^\mathcal{O}(u) - \ell^\mathcal{O}(u_H)$  to decide if the current approximation is accurate enough, and, if not, to be able to provide quantitative information on how to improve the approximation  $u_H$  to achieve the desired accuracy.

Usually, the assessment of the error is achieved using a hierarchy of numerical approximations. The first discretization or “working” coarse mesh provides the approximated solution  $u_H$  and the approximated output  $\ell^\mathcal{O}(u_H)$  with a relatively low effort. The second discretization or “reference” fine mesh produces an approximation  $u_h$  and an output  $\ell^\mathcal{O}(u_h)$  for which it is assumed that  $|\ell^\mathcal{O}(u) - \ell^\mathcal{O}(u_h)|$  is negligible with respect to  $|\ell^\mathcal{O}(u) - \ell^\mathcal{O}(u_H)|$ , at least in the asymptotic range of convergence. The reference discretization serves to assess the error in the output associated with the coarse discretization since it is assumed that  $\ell^\mathcal{O}(u) - \ell^\mathcal{O}(u_H) \approx \ell^\mathcal{O}(u_h) - \ell^\mathcal{O}(u_H)$ . Obviously, the “reference” mesh calculations are too expensive, or impossible, to be performed. However, it is possible to obtain bounds for  $\ell^\mathcal{O}(u_h)$ , with an extra cost similar to the output calculation on the “working” mesh, using the approximation  $u_H$ . One would expect that if the “reference” mesh is *fine enough*, the bounds for  $\ell^\mathcal{O}(u_h)$  would also hold for  $\ell^\mathcal{O}(u)$ . Although a priori error estimates techniques might be used to have information on the asymptotic range of convergence of  $u_h$ , one can not answer the question: are the bounds for  $\ell^\mathcal{O}(u_h)$  still valid when assessing the exact output  $\ell^\mathcal{O}(u)$ ?

The final goal of the thesis is to obtain bounds for  $\ell^\mathcal{O}(u)$ , using only the coarse mesh approximation  $u_H$ , which are uniformly valid, that is, they are valid regardless of the size of the underlying coarse discretization. The proposed techniques provide a certificate of precision for a predicted output with a cost that does not overwhelm the cost of computation of  $u_H$ . With each approximation  $\ell^\mathcal{O}(u_H)$  bounds  $s_l$  and  $s_u$  are provided guaranteeing that  $\ell^\mathcal{O}(u) \in [\ell^\mathcal{O}(u_H) + s_l, \ell^\mathcal{O}(u_H) + s_u]$ . Furthermore, the procedure also provides local information on the contributions to the error in the output  $\ell^\mathcal{O}(u) - \ell^\mathcal{O}(u_H)$  which might be used in adaptive methods to drive a solution to an arbitrary precision.

Note that, although attention is not paid to the error introduced from the choice of a mathematical model, the bounds for the output can also be used to invalidate models in the cases where the experimental data does not fit the predicted intervals of approximation of the output  $\ell^\mathcal{O}(u)$ .

The aforementioned goal of obtaining strict bounds for outputs of interest depending on the exact weak solution of a set of partial differential equations is probably one of the most important open problems in error verification. Since the 1990s significant advances toward obtaining bounds for quantities of interest have been made but it still remains a great deal of work to be done.

Having in mind this “so-ambitious” goal, attention has been focused on several partial goals, the solution of which provides concrete advances toward the resolution of the main problem.

First of all the problem has been simplified to obtaining bounds for *linear-functional outputs*, thus leaving aside the treatment of non-linear outputs. Only a particular quadratic-functional output has been considered, the evaluation of the  $J$ -integral in linear fracture mechanics see (Chapter 5). In this scenario the following partial goals are considered:

- ▷ **Obtaining bounds for linear-functional outputs depending on the solution of self-adjoint coercive boundary value problems (thermal model problem and elasticity):** bounds for quantities of interest for self-adjoint coercive problems may be obtained using techniques providing upper bounds for the error measured in the energy norm. Moreover, these bounds may be enhanced if also techniques providing lower bounds for the error measured in the energy norm are available. Many implicit a-posteriori residual type error estimators provide upper bounds for the energy norm of the error, thus, the first goal is to develop simple and inexpensive techniques to *obtain sharp lower bounds from the available information (finite element approximation and upper bound estimate)*. This work is detailed in (Díez, Parés and Huerta 2003) and also in Section 3.3.
- ▷ **Obtaining a methodology providing sharp upper and lower bounds for linear-functional outputs and being competent for 3D applications:** as mentioned before, obtaining bounds for a quantity of interest rely on obtaining bounds for the energy. Although there are many available techniques to obtain upper bounds for the energy, either they provide bounds which are not really accurate or they rely on equilibration techniques which are complex to implement, specially in 3D applications. Thus, in order to be able to consider



complex practical applications in three dimensions one can either obtain too pessimistic bounds or must spend an immense amount of effort in the implementation of equilibrating techniques. Thus, the second goal is to *develop an estimation technique providing accurate estimates and being simple to implement*. This work is detailed in (Parés, Díez and Huerta 2005) and in Chapter 4.

- ▷ **Obtaining strict bounds for linear-functional outputs in elasticity:** most residual type error estimators introduce a reference or fine mesh in order to derive upper bounds for linear-functional outputs. This yields to the loss of the upper bound property with respect to the exact output,  $\ell^{\mathcal{O}}(u)$ , but still returns an upper bound with respect to the reference output,  $\ell^{\mathcal{O}}(u_h)$ . In the asymptotic range of convergence, the reference solution will approach the exact solution and the upper bounds for  $\ell^{\mathcal{O}}(u_h)$  are hoped to provide also an upper bound for  $\ell^{\mathcal{O}}(u)$ . However, given a working mesh, it is crucial to be able to assess the error with respect to the exact solution and not with respect to a reference solution. Sauer-Budge, Bonet, Huerta and Peraire (2004) and Sauer-Budge and Peraire (2004) derive strict bounds for linear-functional outputs for scalar model problems. An extension to the elasticity equations is presented in (Parés, Bonet, Huerta and Peraire 2005) and in Chapter 5.
  
- ▷ **Obtaining a general framework to derive sharp bounds for linear-functional outputs in nonself-adjoint model problems:** finally, the idea of enhancing the bounds for the output by means of obtaining lower bounds for the error measured in the energy norm, appearing when dealing with self-adjoint problems, has been extended for nonself-adjoint model problems. The goal is to be able to obtain accurate bounds for linear-functional outputs in convection-dominated convection-diffusion-reaction model problems. Here only the general framework to obtain the bounds is presented (Chapter 2). The application of these techniques to numerical examples has only been tested for one dimensional model problems and is now being tested in a two dimensional setting.

### 1.3 Overview

The thesis is divided in three main parts (three-layer presentation): the exposition, the appendices referring to the exposition and finally the three main contributions of the thesis enclosed in form of published or accepted papers.

The exposition part aims at presenting the contributions of this thesis in a clear and concise manner providing the main ideas and core concepts. The full details of the deduction and implementation of the methodologies are appended either in the appendices A, B and C, which support the exposition part, or in the final papers.

The exposition part is divided in 5 chapters: Chapters 2 present a general framework to obtain sharp bounds for outputs of interest depending on solutions of mathematical problems dealing with both self-adjoint and nonself-adjoint operators. The most important concept in this chapter is the relation between bounding outputs of interest and bounding energy norms. It is sufficient to obtain upper and lower bounds for the errors measured in the energy norm to be able to recover sharp bounds for the output. Thus, Chapter 3 is concerned with obtaining upper and lower bounds for the energy norm. It presents the main ingredients present in any implicit residual type a posteriori error estimation technique which allows to compute upper and lower bounds for the energy norm of the error.

Chapters 4 and 5 (and its correspondent appendices B and C) present a brief description of the methods detailed in the papers (Parés, Díez and Huerta 2005) and (Parés, Bonet, Huerta and Peraire 2005) respectively. Finally Chapter 6 presents the conclusions and future developments.

The three appended papers at the end of the thesis correspond to the references Díez et al. (2003), Parés, Díez and Huerta (2005) and Parés, Bonet, Huerta and Peraire (2005) respectively. Throughout the thesis these papers are cited using the corresponding reference.

## Chapter 2

# Bounds for linear outputs of interest for a general variational problem

The goal of many finite element computations is to determine a specific quantity  $\ell^{\mathcal{O}}(u)$ , such as the prediction of system characteristics in engineering design. Since the quantity of interest depends on the unknown solution  $u$  it is not possible in general to know the value of the output, and the finite element solution is used to approximate it, that is,  $\ell^{\mathcal{O}}(u)$  is approximated by  $\ell^{\mathcal{O}}(u_H)$ . In this context, the classical assessment of the energy norm of the error in the field solution  $u$  does not provide any information on the accuracy of the approximation of the output which is the aim of goal oriented error estimation techniques.

This chapter presents a general framework aimed at obtaining both upper and lower bounds for the error in the quantity of interest. Thus, if  $s := \ell^{\mathcal{O}}(u) - \ell^{\mathcal{O}}(u_H)$  is the error in the quantity of interest, scalar quantities  $s_l$  and  $s_u$  will be computed satisfying

$$s_l \leq s \leq s_u.$$

An immediate consequence of these error bounds for the output is that upper and lower bounds may be obtained for the quantity of interest itself, namely

$$\ell^{\mathcal{O}}(u_H) + s_l \leq \ell^{\mathcal{O}}(u) \leq \ell^{\mathcal{O}}(u_H) + s_u.$$

In order to develop a general framework to obtain bounds for  $s$ , it is convenient to regard the quantity of interest as a bounded, linear functional  $\ell^{\mathcal{O}}(\cdot)$  acting on the space  $\mathcal{V}$  of admissible functions for the problem at hand. In this case, the error in the quantity of interest may be written in the form  $s = \ell^{\mathcal{O}}(u) - \ell^{\mathcal{O}}(u_H) = \ell^{\mathcal{O}}(e)$

where  $e = u - u_H$  is the error in the finite element approximation of  $u$ . The treatment of non-linear outputs is considered by Larsson, Hansbo and Runesson (2002) and Sarrate, Peraire and Patera (1999) and also by Xuan, Parés and Peraire (2005) where bounds for a particular output, the  $J$ -integral, are found without introducing a linearization of the output.

A notable feature of the presented theory is that it is possible to obtain upper and lower bounds for  $s$  using techniques developed for estimating the error in the global energy norm for self-adjoint model problems. These techniques are not discussed in this chapter which is primarily concerned with establishing the general framework for assessing the error in quantities of interest. A survey of techniques to estimate the error in the global energy norm for self-adjoint model problems can be found in the book by Ainsworth and Oden (2000). Moreover, Chapters 4 and 5 describe and discuss two new error estimation procedures yielding bounds for the error measured in the energy norm.

This chapter is structured as follows: first, the model problem is presented. Then upper and lower bounds for the error in the output  $s$  are derived for both self-adjoint and nonself-adjoint coercive model problems. The theory presented here is a survey of the most popular techniques to obtain bounds for output quantities from energy norm estimates. However, for nonself-adjoint model problems a novelty is introduced. Although the presented approach only introduces a slightly modification on the existing theory, it allows to enhance the bounds for the output sidestepping the degradation of the bounds appearing in some problems.

## 2.1 Model problem

Consider a general variational problem: find  $u \in \mathcal{V}$  such that

$$a(u, v) = \ell(v) \quad \forall v \in \mathcal{V}, \quad (2.1)$$

where  $\mathcal{V}$  is a Hilbert space, the functional  $\ell \in \mathcal{V}'$  is a continuous linear functional over  $\mathcal{V}$ , and  $a : \mathcal{V} \times \mathcal{V} \rightarrow \mathbb{R}$  is a continuous, coercive bilinear form not necessarily symmetric. For ease of exposition the Dirichlet boundary conditions are taken to be homogeneous, in which case the solution and test space coincide. Non-

homogeneous Dirichlet conditions may be dealt with in an analogous fashion (see Parés, Bonet, Huerta and Peraire 2005).

The Lax-Milgram Theorem guarantees both existence and uniqueness of the solution to (2.1) (see Brenner and Scott 1994). The exact solution is approximated by the finite element solution  $u_H$  lying in a finite dimensional subspace  $\mathcal{V}^H \subset \mathcal{V}$ . The finite element solution also known as Ritz-Galerkin approximation is the solution of equation (2.1) with  $\mathcal{V}$  replaced by  $\mathcal{V}^H$ , namely

$$a(u_H, v) = \ell(v) \quad \forall v \in \mathcal{V}^H. \quad (2.2)$$

The error  $e = u - u_H$  belongs to the space  $\mathcal{V}$  and satisfies the residual equation

$$a(e, v) = \ell(v) - a(u_H, v) =: R^P(v) \quad \forall v \in \mathcal{V}, \quad (2.3)$$

where  $R^P(\cdot)$  is the residue associated to the finite element solution  $u_H$ . Moreover, the standard orthogonality condition for the error in the Galerkin projection holds

$$a(e, v) = 0 \quad \forall v \in \mathcal{V}^H. \quad (2.4)$$

### 2.1.1 Operator decomposition

The bilinear form  $a(\cdot, \cdot)$  can be split into its symmetric,  $a^s(v, w) = a^s(w, v)$ , and antisymmetric (or skew-symmetric),  $a^{ss}(v, w) = -a^{ss}(w, v)$ , contributions

$$a^s(v, w) = \frac{1}{2} (a(v, w) + a(w, v)), \quad a^{ss}(v, w) = \frac{1}{2} (a(v, w) - a(w, v)), \quad (2.5)$$

namely

$$a(v, w) = a^s(v, w) + a^{ss}(v, w). \quad (2.6)$$

Since  $a^s(\cdot, \cdot)$  is a continuous, coercive, symmetric bilinear form, it defines an inner product and the associated norm  $\|v\|_s = \sqrt{a^s(v, v)}$ . This immediately implies that both the Cauchy-Schwarz inequality

$$|a^s(v, w)| \leq \|v\|_s \|w\|_s, \quad (2.7)$$

and the parallelogram identity

$$a^s(v, w) = \frac{1}{4} \|v + w\|_s^2 - \frac{1}{4} \|v - w\|_s^2, \quad (2.8)$$

hold for the scalar product  $a^s(\cdot, \cdot)$ . Additionally, denoting by  $\|v\| = \sqrt{a(v, v)}$  the energy norm of a function in  $\mathcal{V}$ , from the definition of the symmetric operator,  $\|v\|_s = \|v\|$ . Thus, equations (2.7) and (2.8) are still valid when the symmetric norms,  $\|\cdot\|_s$ , are replaced by the energy norm,  $\|\cdot\|$ . In fact, from now on, instead of the symmetric norm  $\|\cdot\|_s$  the energy norm  $\|\cdot\|$  is used, that is, the norm associated with the symmetric bilinear form  $a^s(\cdot, \cdot)$  is taken to be directly the energy norm,  $\|\cdot\|$ .

Obviously, if the bilinear form  $a(\cdot, \cdot)$  is symmetric it coincides with its symmetric contribution,  $a(v, w) = a^s(v, w)$ , and  $a^{ss}(v, w) = 0 \forall v, w \in \mathcal{V}$ .

Symmetric bilinear forms  $a(\cdot, \cdot)$  derive from self-adjoint model problems which are named after symmetric model problems (though slightly abuse of language) referring to the associated bilinear form being symmetric. Similarly, nonsymmetric bilinear forms derive from nonself-adjoint model problems named after nonsymmetric model problems.

## 2.2 Energy reformulation

Attention is usually centered in bounding the error in the finite element approximation of a quantity of interest  $s = \ell^{\mathcal{O}}(e)$  where  $\ell^{\mathcal{O}} \in \mathcal{V}'$  is a bounded continuous linear functional over  $\mathcal{V}$  (see for instance Paraschivoiu, Peraire and Patera 1997, Ma-day, Patera and Peraire 1999, Prudhomme and Oden 1999, Oden and Prudhomme 2001, Patera and Peraire 2003). These strategies introduce a dual (or adjoint) problem with respect to the selected output. The weak form of the dual problem reads: find  $\psi \in \mathcal{V}$  verifying

$$a(v, \psi) = \ell^{\mathcal{O}}(v) \quad \forall v \in \mathcal{V}. \quad (2.9)$$

The finite element approximation of the dual problem is  $\psi_H \in \mathcal{V}^H$  such that

$$a(v, \psi_H) = \ell^{\mathcal{O}}(v) \quad \forall v \in \mathcal{V}^H, \quad (2.10)$$

and the error in the finite element approximation or dual error is  $\varepsilon = \psi - \psi_H \in \mathcal{V}$  solution of the dual residual problem

$$a(v, \varepsilon) = \ell^{\mathcal{O}}(v) - a(v, \psi_H) =: R^{\mathcal{D}}(v) \quad \forall v \in \mathcal{V}, \quad (2.11)$$

$R^D(\cdot)$  being the dual weak residue associated to the finite element approximation  $\psi_H$ .

The definition of the dual problem (2.9) along with the aid of the Galerkin orthogonality property (2.4) allows to rewrite the error in the output of interest as

$$s = \ell^O(e) = a(e, \psi) = a(e, \psi - \psi_H) = a(e, \varepsilon). \quad (2.12)$$

If  $a(\cdot, \cdot)$  is a symmetric bilinear form, then the parallelogram identity (2.8) yields to the well known alternative representation for the output

$$s = a(e, \varepsilon) = a^s(e, \varepsilon) = \frac{1}{4} \|\kappa e + \frac{1}{\kappa} \varepsilon\|^2 - \frac{1}{4} \|\kappa e - \frac{1}{\kappa} \varepsilon\|^2, \quad (2.13)$$

where  $\kappa \in \mathbb{R}$  is a nonzero arbitrary scalar parameter (see Babuška and Miller 1984, Prudhomme and Oden 1999, Ainsworth and Oden 2000, and Theorem A.2.1 of the present thesis).

Equation (2.13) allows to compute the error in the output  $s$  simply computing the energy norm of linear combinations of the primal and dual errors  $e$  and  $\varepsilon$ . Moreover, it reduces the problem of bounding the error in the output to the derivation of upper and lower bounds for the energy norm of the linear combinations  $\kappa e \pm \frac{1}{\kappa} \varepsilon$ . Indeed, if  $\|\kappa e \pm \frac{1}{\kappa} \varepsilon\|_{\text{UB}}$  and  $\|\kappa e \pm \frac{1}{\kappa} \varepsilon\|_{\text{LB}}$  are upper and lower bounds of  $\|\kappa e \pm \frac{1}{\kappa} \varepsilon\|$  respectively, namely

$$\|\kappa e \pm \frac{1}{\kappa} \varepsilon\|_{\text{LB}} \leq \|\kappa e \pm \frac{1}{\kappa} \varepsilon\| \leq \|\kappa e \pm \frac{1}{\kappa} \varepsilon\|_{\text{UB}},$$

the error in the output is readily bounded as

$$\frac{1}{4} \|\kappa e + \frac{1}{\kappa} \varepsilon\|_{\text{LB}}^2 - \frac{1}{4} \|\kappa e - \frac{1}{\kappa} \varepsilon\|_{\text{UB}}^2 \leq s \leq \frac{1}{4} \|\kappa e + \frac{1}{\kappa} \varepsilon\|_{\text{UB}}^2 - \frac{1}{4} \|\kappa e - \frac{1}{\kappa} \varepsilon\|_{\text{LB}}^2. \quad (2.14)$$

Unfortunately, this result is no longer valid for nonsymmetric bilinear forms. In fact, for a general bilinear form containing nonself-adjoint terms, the idea of developing upper and lower bounds for  $s$  from energy norm estimates must be in some sense extended, since there is no natural energy norm in which to measure the error. Is there any reason to suspect that bounds for  $s$  can be obtained for a general model problem using only available techniques for estimating the error measured in the energy norm of solutions of symmetric model problems?

With the aid of the symmetric part of the bilinear form  $a(\cdot, \cdot)$ , a symmetric analogous of the dual residual problem (2.11) can be defined as: find  $\varepsilon^s \in \mathcal{V}$  verifying

$$a^s(v, \varepsilon^s) = R^D(v) \quad \forall v \in \mathcal{V}. \quad (2.15)$$

Then, replacing  $v = e$  first in equation (2.11) and then in equation (2.15) and using the Cauchy-Schwarz inequality (2.7), it follows that

$$|s| = |a(e, \varepsilon)| = |R^D(e)| = |a^s(e, \varepsilon^s)| \leq \|e\| \|\varepsilon^s\|.$$

Moreover, if the symmetric analogous of the primal residual problem (2.3) is introduced: find  $e^s \in \mathcal{V}$  such that

$$a^s(e^s, v) = R^P(v) \quad \forall v \in \mathcal{V}, \quad (2.16)$$

taking  $v = e$  first in equation (2.3) and then in equation (2.16), it follows that  $\|e\|^2 = R^P(e) = a^s(e^s, e)$  and then a routine application of the Cauchy-Schwarz inequality yields  $\|e\| \leq \|e^s\|$ . Hence

$$|s| \leq \|e^s\| \|\varepsilon^s\| \implies -\|e^s\| \|\varepsilon^s\| \leq s \leq \|e^s\| \|\varepsilon^s\|. \quad (2.17)$$

Consequently, it is possible to obtain computable bounds for  $s$  in terms of upper bounds for the energy norm of solutions of symmetric boundary value problems. If  $\|e^s\|_{\text{UB}}$  and  $\|\varepsilon^s\|_{\text{UB}}$  are upper bounds for the energy norm of  $e^s$  and  $\varepsilon^s$  respectively, namely

$$\|e^s\| \leq \|e^s\|_{\text{UB}}, \quad \|\varepsilon^s\| \leq \|\varepsilon^s\|_{\text{UB}},$$

the output of the error is readily bounded as

$$-\|e^s\|_{\text{UB}} \|\varepsilon^s\|_{\text{UB}} \leq s \leq \|e^s\|_{\text{UB}} \|\varepsilon^s\|_{\text{UB}}. \quad (2.18)$$

Note that since both  $e^s$  and  $\varepsilon^s$  are solutions of symmetric boundary value problems, the upper bounds  $\|e^s\|_{\text{UB}}$  and  $\|\varepsilon^s\|_{\text{UB}}$  may be computed using standard a posteriori error estimation techniques producing upper bounds for the energy norm.

The repeated use of the Cauchy-Schwarz inequality to deduce the bounds given in equation (2.18) generally produces rather pessimistic bounds. In the symmetric case the use of the Cauchy-Schwarz inequality can be sidestepped taking advantage of the parallelogram identity. However, the parallelogram identity is no longer valid for nonsymmetric bilinear forms  $a(\cdot, \cdot)$  and other techniques have to be used in order to refine the bounds.

The rest of this chapter investigates to what extent estimates for the energy norm may be used to find accurate bounds for outputs of interest depending on solutions



of nonself-adjoint problems. First the bounds for the symmetric model problem are revised. Then, bounds for nonsymmetric model problems found in the earlier literature are presented and finally, an improvement of the bounds is introduced which allows to obtain sharp bounds for the error in the output. Here only a brief account of how the bounds are derived is provided. A complete deduction of the proposed bounds can be found in Appendix A.

### 2.2.1 Bounds for self-adjoint model problems

If the bilinear form  $a(\cdot, \cdot)$  is symmetric,  $a^s(v, w) = a(v, w)$ , as seen in the previous section, an exact representation of the quantity of interest  $s = \ell^{\mathcal{O}}(e)$  can be deduced using the parallelogram identity, namely

$$s = \frac{1}{4} \left\| \kappa e + \frac{1}{\kappa} \varepsilon \right\|^2 - \frac{1}{4} \left\| \kappa e - \frac{1}{\kappa} \varepsilon \right\|^2. \quad (2.19)$$

This representation allows to compute bounds for  $s$  from upper and lower bounds for the energy norm of the linear combinations of the primal and dual errors  $\kappa e \pm \frac{1}{\kappa} \varepsilon$ , see equation (2.14).

The upper bounds for the energy norms  $\left\| \kappa e \pm \frac{1}{\kappa} \varepsilon \right\|$  are usually computed using implicit residual type error estimators. Most of these techniques provide estimates  $\hat{e}$  and  $\hat{\varepsilon}$  (see Lemma 3.2.1) verifying

$$\left\| \kappa e \pm \frac{1}{\kappa} \varepsilon \right\| \leq \left\| \kappa \hat{e} \pm \frac{1}{\kappa} \hat{\varepsilon} \right\|.$$

These estimates,  $\hat{e}$  and  $\hat{\varepsilon}$ , are discontinuous approximations of the error fields  $e$  and  $\varepsilon$  respectively. In order to obtain lower bounds for the energy norm  $\left\| \kappa e \pm \frac{1}{\kappa} \varepsilon \right\|$ , the discontinuous estimates  $\hat{e}$  and  $\hat{\varepsilon}$  may be post-processed to obtain continuous approximations of the error fields  $e$  and  $\varepsilon$  and thus, continuous approximations of  $\kappa e \pm \frac{1}{\kappa} \varepsilon$  denoted by  $\xi^{\pm}$  (see Section 3.3.2).

The lower bounds for the energy norm of the linear combinations of the primal and dual problems,  $\left\| \kappa e \pm \frac{1}{\kappa} \varepsilon \right\|_{\text{LB}}$ , may be then recovered from  $\xi^{\pm}$  using the dual characterization of the energy norm. Indeed, the energy norm of  $\left\| \kappa e \pm \frac{1}{\kappa} \varepsilon \right\|$  may be characterized (Oden and Prudhomme 1999, Theorem 4.3) using duality as

$$\left\| \kappa e \pm \frac{1}{\kappa} \varepsilon \right\|^2 = \sup_{v \in \mathcal{V}} \frac{R^{\pm}(v)^2}{\|v\|^2}, \quad (2.20)$$

where the residue  $R^\pm \in \mathcal{V}'$  is defined as

$$R^\pm(v) = \kappa R^P(v) \pm \frac{1}{\kappa} R^D(v). \quad (2.21)$$

As a consequence of (2.20), any continuous approximation of  $\kappa e \pm \frac{1}{\kappa} \varepsilon$ ,  $\xi^\pm \in \mathcal{V}$ , provides the lower bound

$$\|\kappa e \pm \frac{1}{\kappa} \varepsilon\|^2 \geq \frac{R^\pm(\xi^\pm)^2}{\|\xi^\pm\|^2},$$

where the previous inequality holds as an equality for  $\xi^\pm = \kappa e \pm \frac{1}{\kappa} \varepsilon$ . Therefore, accurate lower bounds for  $\|\kappa e \pm \frac{1}{\kappa} \varepsilon\|$  may be derived from good continuous approximations,  $\xi^\pm \in \mathcal{V}$ , of  $\kappa e \pm \frac{1}{\kappa} \varepsilon$ .

Once the upper and lower bounds for the quantities  $\|\kappa e \pm \frac{1}{\kappa} \varepsilon\|$  are computed, the upper and lower bounds for  $s$ ,  $s_u$  and  $s_l$  respectively, are recovered using equation (2.14), namely

$$s_l = \frac{1}{4} \frac{R^+(\xi^+)^2}{\|\xi^+\|^2} - \frac{1}{4} \|\kappa \hat{e} - \frac{1}{\kappa} \hat{\varepsilon}\|^2 \quad \text{and} \quad s_u = \frac{1}{4} \|\kappa \hat{e} + \frac{1}{\kappa} \hat{\varepsilon}\|^2 - \frac{1}{4} \frac{R^-(\xi^-)^2}{\|\xi^-\|^2}, \quad (2.22)$$

see (A.13).

The procedure to obtain bounds for  $s$  depending on the solution of a symmetric model problem is summarized in Figure 2.1.

## 2.2.2 Bounds for nonself-adjoint model problems

If the bilinear form  $a(\cdot, \cdot)$  is nonsymmetric, the exact representation of the quantity of interest  $s = a(e, \varepsilon)$  still holds. However, the parallelogram identity is no longer valid for  $a(\cdot, \cdot)$  and therefore it is not possible to find an analogous of equation (2.19) which allows to find accurate bounds for  $s$ .

Bounds for the output  $s$  may be deduced, as shown in (2.18), from

$$-\|e^s\| \|\varepsilon^s\| \leq s \leq \|e^s\| \|\varepsilon^s\|. \quad (2.23)$$

However, since these bounds are deduced using the Cauchy-Schwartz inequality, the resulting bounds are not usually sharp. Paraschivoiu et al. (1997) avoid the use of the Cauchy-Schwartz inequality reformulating the output  $s$  as the solution of a

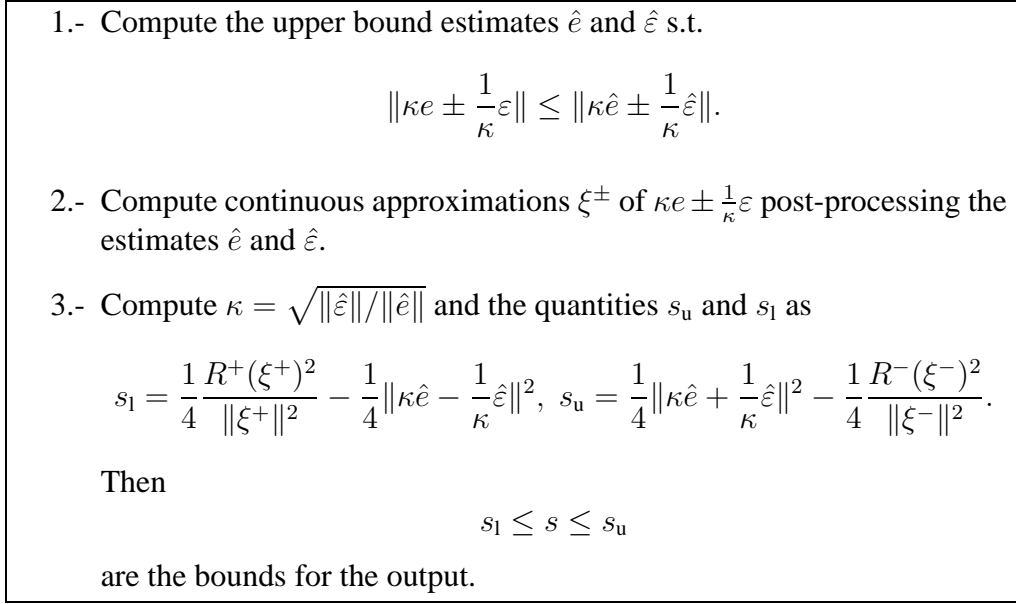


Figure 2.1: Main steps of the strategy used to obtain bounds for an output  $s$  depending on the solution of a symmetric boundary value problem.

minimization problem and making use of the saddle point theory. The proposed bounds are

$$-\frac{1}{4} \|\kappa e^s - \frac{1}{\kappa} \varepsilon^s\|^2 \leq s \leq \frac{1}{4} \|\kappa e^s + \frac{1}{\kappa} \varepsilon^s\|^2, \quad (2.24)$$

where  $\kappa \in \mathbb{R}$  is again a nonzero arbitrary scalar parameter (see Remark A.1.1). The parameter  $\kappa$  is selected so that the bounds are optimized. The same value of  $\kappa$  minimizes the upper bound and maximizes the lower bound and is given by  $\kappa = \sqrt{\|\varepsilon^s\|/\|e^s\|}$ . For the optimal parameter  $\kappa$ , the resulting bounds are

$$\frac{1}{2} a^s(e^s, \varepsilon^s) - \frac{1}{2} \|e^s\| \|\varepsilon^s\| \leq s \leq \frac{1}{2} a^s(e^s, \varepsilon^s) + \frac{1}{2} \|e^s\| \|\varepsilon^s\|,$$

improving the bounds given in (2.23) since  $|a^s(e^s, \varepsilon^s)| \leq \|e^s\| \|\varepsilon^s\|$ .

Equation (2.24) indicates that in order to bound the output of the error  $s$  it is sufficient to obtain upper bounds of the energy norm of the errors  $\kappa e^s \pm \frac{1}{\kappa} \varepsilon^s$ . Indeed, if  $\|\kappa e^s \pm \frac{1}{\kappa} \varepsilon^s\|_{\text{UB}}$  are upper bounds of  $\|\kappa e^s \pm \frac{1}{\kappa} \varepsilon^s\|$ , then the error in the output is readily bounded by

$$-\frac{1}{4} \|\kappa e^s - \frac{1}{\kappa} \varepsilon^s\|_{\text{UB}}^2 \leq s \leq \frac{1}{4} \|\kappa e^s + \frac{1}{\kappa} \varepsilon^s\|_{\text{UB}}^2. \quad (2.25)$$

It is worth noting the similarity of these bounds with the bounds given in equation (2.14) for symmetric model problems. The primal and dual errors appearing in (2.14) are replaced by its associated symmetric errors whereas the lower bounds of the energy norms are taken to be zero.

The error estimation procedure presented by Paraschivoiu et al. (1997) produces the bounds given in equation (2.25). The sharpness of these bounds is strongly dependent on the characteristics of the problem at hand and the desired final accuracy. In some cases, the provided bounds are too pessimistic (for instance, for convection-dominated convection-diffusion problems, one can choose an output for which the error in the output is practically zero and the obtained bounds using equation (2.25) degenerate as the convection parameter increases).

Motivated by this non-desired behavior of the resulting bounds, the deduction of the bounds proposed by Paraschivoiu et al. (1997) has been carefully examined (see Appendix A). The net result is the introduction of a relatively straightforward modification which allows to enhance the initial bounds provided by equation (2.25).

### Enhancement of the bounds for the output

The bounds given in equation (2.25) are computed using implicit residual type error estimators. Most of these techniques provide discontinuous estimates  $\hat{e}^s$  and  $\hat{\varepsilon}^s$  (Lemma 3.2.1) verifying

$$\|\kappa e^s \pm \frac{1}{\kappa} \varepsilon^s\| \leq \|\kappa \hat{e}^s \pm \frac{1}{\kappa} \hat{\varepsilon}^s\|$$

yielding the bounds for the output

$$-\frac{1}{4} \|\kappa \hat{e}^s - \frac{1}{\kappa} \hat{\varepsilon}^s\|^2 \leq s \leq \frac{1}{4} \|\kappa \hat{e}^s + \frac{1}{\kappa} \hat{\varepsilon}^s\|^2.$$

In order to enhance the bounds, continuous approximations,  $\xi^\pm \in \mathcal{V}$ , of  $\kappa e \pm \frac{1}{\kappa} \varepsilon$  may be computed, for instance post-processing  $\kappa \hat{e}^s \pm \frac{1}{\kappa} \hat{\varepsilon}^s$ . However, when dealing with nonsymmetric model problems it is not sufficient to have these approximations  $\xi^\pm$  to enhance the bounds. Once the continuous approximations  $\xi^\pm$  are computed, the error estimation procedure must be applied to the following symmetric problem

$$a^s(\xi^{s\pm}, v) = a(v, \xi^\pm) \quad \forall v \in \mathcal{V},$$

in order to compute the discontinuous estimate  $\hat{\xi}^{s\pm}$  such that  $\|\xi^{s\pm}\| \leq \|\hat{\xi}^{s\pm}\|$  (Theorem A.2.2).

The final upper and lower bounds  $s_u$  and  $s_l$  are, respectively, (A.24)

$$\begin{aligned} s_l &= \frac{1}{4} \frac{(2\kappa R^P(\xi^+) - a^s(\kappa\hat{e}^s - \frac{1}{\kappa}\hat{\varepsilon}^s, \hat{\xi}^{s+}))^2}{\|\hat{\xi}^{s+}\|^2} - \frac{1}{4} \|\kappa\hat{e}^s - \frac{1}{\kappa}\hat{\varepsilon}^s\|^2, \\ s_u &= \frac{1}{4} \|\kappa\hat{e}^s + \frac{1}{\kappa}\hat{\varepsilon}^s\|^2 - \frac{1}{4} \frac{(2\kappa R^P(\xi^-) - a^s(\kappa\hat{e}^s + \frac{1}{\kappa}\hat{\varepsilon}^s, \hat{\xi}^{s-}))^2}{\|\hat{\xi}^{s-}\|^2}. \end{aligned} \quad (2.26)$$

Moreover, if the bilinear form  $a(\cdot, \cdot)$  is symmetric, the bounds obtained for the symmetric model problem given in equation (2.22) are recovered.

*Remark 2.2.1.* In fact,  $s_u$  and  $s_l$  are upper and lower bounds for  $s$  if the estimates  $\hat{e}^s, \hat{\varepsilon}^s$  and  $\hat{\xi}^{s\pm}$  verify

$$a^s(\hat{e}^s, v) = a^s(e^s, v), \quad a^s(\hat{\varepsilon}^s, v) = a^s(\varepsilon^s, v) \quad \forall v \in \mathcal{V}$$

and

$$a^s(\hat{\xi}^{s\pm}, \kappa\hat{e}^s \pm \frac{1}{\kappa}\hat{\varepsilon}^s) = a^s(\xi^{s\pm}, \kappa\hat{e}^s \pm \frac{1}{\kappa}\hat{\varepsilon}^s),$$

see Appendix A, equations (A.19) and (A.23).

It is worth emphasizing, however, that most implicit residual type estimation techniques yield estimates  $\hat{e}^s, \hat{\varepsilon}^s$  and  $\hat{\xi}^{s\pm}$  verifying the previous conditions. In fact, the first conditions are sufficient conditions to proof that  $\|e^s\| \leq \|\hat{e}^s\|$  and  $\|\varepsilon^s\| \leq \|\hat{\varepsilon}^s\|$ , see Lemma 3.2.1.

The procedure is summarized in Figure 2.2. Note that the main difference between bounds for symmetric and nonsymmetric model problems is that in the nonsymmetric case it is not sufficient to be able to obtain a continuous approximation of  $\kappa e \pm \frac{1}{\kappa}\varepsilon$ , a new estimate depending on the continuous approximation must be evaluated. However, a careful examination of the procedure, shows that there is another added difficulty in the computation of the bounds. Is it easy to compute good continuous approximations of  $\kappa e \pm \frac{1}{\kappa}\varepsilon$  from the available approximations of  $e^s$  and  $\varepsilon^s$ ? If the resulting bounds are to be sharp, the selection of the continuous approximations must be exercised carefully.

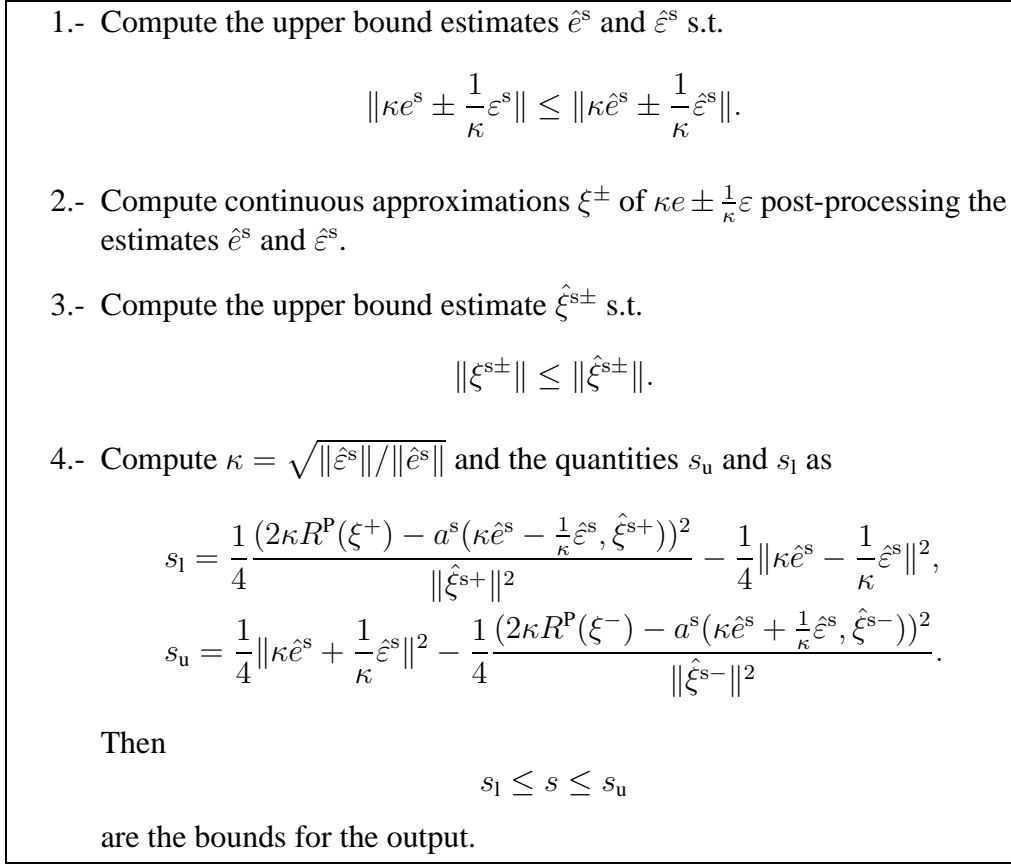


Figure 2.2: Main steps of the strategy used to obtain bounds for an output  $s$  depending on the solution of a nonsymmetric boundary value problem.

## 2.3 Summary

Obtaining bounds for quantities of interest is crucial in applications. Babuška and Miller (1984) proposed the use of the parallelogram identity in symmetric coercive model problems to obtain output bounds from upper and lower bounds measured in the energy norm. Lower bounds for the energy norm are computed using a simple post-processing technique. This has been extensively studied by Babuška, Strouboulis and Gangaraj (1999), Strouboulis, Babuška and Gangaraj (2000), Prudhomme, Oden, Westermann, Bass and Botkin (2003) and Díez et al. (2003).

Paraschivoiu et al. (1997) present a general framework to obtain bounds for the output of nonsymmetric coercive problems from upper bounds for the energy

---

norm. In this chapter, the same techniques have been used to sharpen the bounds. A slightly modification is introduced which allows to obtain more accurate bounds. The idea is analogous to the enhancement of the bounds in the symmetric model problem using post-processing techniques to obtain continuous approximations of the error field.

At this point, it is clear that techniques to obtain bounds for the output  $s$  may be developed from estimates of the error in the global energy norm of solutions of symmetric model problems.





## Chapter 3

# Bounds for the energy norm of solutions of self-adjoint model problems

The computation of bounds for linear outputs of interest can be accomplished invoking any error estimation procedure which allows to compute upper and lower bounds for the energy norm of the error in the finite element approximation of a *self-adjoint* or *symmetric* coercive variational problem.

This chapter is intended to provide the main ingredients present in any implicit residual type a posteriori error estimation technique which allows to compute upper and lower bounds for the energy norm of the error. In order to develop a general framework where the new techniques presented in this thesis can be appropriately defined, the main characteristics of the equilibrated residual method are described to introduce a common notation for the estimates. The choice of the equilibrated residual method as the representative to introduce the main notation is not casual. In addition to be one of the chief methods currently available to find upper bounds for the error measured in the energy norm, it is closely related with the two new estimates developed in Chapters 4 and 5. The estimation procedure introduced in Chapter 4 circumvents the need of flux-equilibration and results in a simple implementation that uses standard resources available in finite element codes, in contrast to the cumbersome implementation of the equilibrated residual method specially for 3D applications. The estimation procedure introduced in Chapter 5 currently entails the computation of equilibrated fluxes although it would be possible to use the ideas

introduced in Chapter 4 to side-step the computation of the equilibrated fluxes.

The chapter is structured as follows: after introducing the model problem, the derivation of upper bounds for the error measured in the energy norm is considered. In particular, the main steps of the equilibrated residual method are summarized. Then, the derivation of lower bounds is considered focusing primarily on the post-processing techniques recently developed. Finally, attention is paid to some practical implementation issues. The discussion of the derivation of upper bounds is valid as long as the local problems are solved exactly, which is of course infeasible. The approximation of the local problems over a reduced finite dimensional space and its consequences is discussed.

### 3.1 Model problem

Consider a general symmetric variational problem given in weak form as: find  $z \in \mathcal{V}$  such that

$$a(z, v) = R^*(v) \quad \forall v \in \mathcal{V}, \quad (3.1)$$

where  $R^*(\cdot) \in \mathcal{V}'$  is a continuous linear functional over  $\mathcal{V}$  and  $a(\cdot, \cdot)$  is a continuous coercive *symmetric* bilinear form.

The solution space is defined as  $\mathcal{V} = \{v \in [\mathcal{H}^1(\Omega)]^{n_{\text{sd}}}, v|_{\Gamma^{\text{D}}} = 0\}$ , where  $n_{\text{sd}}$  is the number of spatial dimensions,  $\mathcal{H}^1$  is the standard Sobolev space of square integrable functions and first derivatives, and the boundary  $\Gamma = \partial\Omega$  is divided into two complementary disjoint parts  $\Gamma^{\text{D}}$  and  $\Gamma^{\text{N}}$ , where essential and Neumann boundary conditions are imposed respectively.

To fix ideas, consider the scalar diffusion-reaction equation. The bilinear form and the r.h.s. term for this problem are of the form

$$a(w, v) = \int_{\Omega} \nu \nabla w \cdot \nabla v + \mu w v \, d\Omega, \quad R^*(v) = \int_{\Omega} f^* v \, d\Omega + \int_{\Gamma^{\text{N}}} g_N^* v \, d\Gamma - a(z_H, v),$$

for a strictly positive real coefficient  $\nu \in \mathcal{L}^\infty(\Omega)$ , a nonnegative real coefficient  $\mu \in \mathcal{L}^\infty(\Omega)$ , and where  $f^* \in \mathcal{H}^{-1}(\Omega)$ ,  $g_N^* \in \mathcal{H}^{-\frac{1}{2}}(\partial\Omega)$  and  $z_H \in \mathcal{V}^H$ . Similarly, the bilinear form and the r.h.s. term for the elasticity problem are of the form

$$a(\mathbf{w}, \mathbf{v}) = \int_{\Omega} \boldsymbol{\sigma}(\mathbf{w}) : \boldsymbol{\varepsilon}(\mathbf{v}) \, d\Omega, \quad R^*(\mathbf{v}) = \int_{\Omega} \mathbf{f}^* \cdot \mathbf{v} \, d\Omega + \int_{\Gamma^{\text{N}}} \mathbf{g}^* \cdot \mathbf{v} \, d\Gamma - a(\mathbf{z}_H, \mathbf{v}),$$

where  $\mathbf{f}^* \in [\mathcal{H}^{-1}(\Omega)]^2$ ,  $\mathbf{g}^* \in [\mathcal{H}^{-\frac{1}{2}}(\partial\Omega)]^2$  is the imposed traction distribution and  $\mathbf{z}_H \in \mathcal{V}^H$ .

It is clear that any linear combination of the primal and dual errors,  $\alpha e + \beta \varepsilon$ ,  $\alpha, \beta \in \mathbb{R}$  is the solution of problem (3.1) with  $R^*(v) = \alpha R^P(v) + \beta R^D(v)$ . In particular, the choice  $\alpha = \kappa, \beta = \pm \frac{1}{\kappa}$  is used to obtain the required upper and lower bounds for  $\|\kappa e \pm \frac{1}{\kappa} \varepsilon\|$ .

## 3.2 Upper bounds for the energy

The need of obtaining reliable upper bounds of the energy norm of the error motivates the use of implicit residual error estimators, which are currently the only type of estimators ensuring bounds for the error entailing only local computations. The underlying idea in any implicit estimate is to decompose the global residual problem (3.1) into a series of local independent boundary value problems posed on small subdomains. This leads to the twofold classification of the implicit estimates: the element residual methods and the subdomain residual methods depending on whether the local problems are posed over a single element or over a small patch of elements respectively.

In the present thesis the element residual methods providing upper bounds for the energy norm of the error are also named as equilibrated residual methods since they require the computation of boundary fluxes (or tractions) which are in equilibrium with the interior residual loads in each element of the mesh. The subdomain residual methods, on the other hand, are also named as *flux-free* residual methods as opposed to the equilibrated methods since they do not require a direct computation of equilibrated fluxes.

The structure of this section is as follows: first, the basic property whereby upper bounds for the energy norm may be derived is summarized. The property is a sufficient condition for an estimate to provide an upper bound for the error measured in the energy norm and does not depend on the method used to compute the estimate. In fact, this property is used to prove the upper bound property of either the estimate obtained using an equilibrated residual method and the flux-free estimate introduced in Chapter 4. Once this key property is introduced, the main features of

the equilibrated residual method are detailed emphasizing the characteristics which are needed to introduce the estimation procedure presented in Chapter 5.

### 3.2.1 Sufficient condition for the upper bound property

The following result summarizes a sufficient condition for an estimate to yield an upper bound for the error measured in the energy norm.

**Lemma 3.2.1.** *Any estimate  $\hat{z} \in \mathcal{W}$  of the function  $z \in \mathcal{V}$  verifying the weak error equation*

$$a(\hat{z}, v) = a(z, v) = R^*(v) \quad \forall v \in \mathcal{V}, \quad (3.2)$$

where  $\mathcal{W}$  is a suitable interpolation space for the estimate  $\hat{z}$ , is such that its norm is an upper bound of the energy norm of  $z$ , that is

$$\|z\| \leq \|\hat{z}\|.$$

*Proof.* Using equation (3.2) with  $v = z$ , the energy norm of  $z$  may be rewritten as

$$\|z\|^2 = a(z, z) = R^*(z) = a(\hat{z}, z).$$

Then, the upper bound property follows applying the Cauchy-Schwarz inequality, namely

$$\|z\|^2 = a(\hat{z}, z) \leq \|\hat{z}\| \|z\|.$$

□

In principle, one may try to find an estimate in  $\mathcal{V}$ , and consider  $\mathcal{W} = \mathcal{V}$ . Unfortunately, if  $\hat{z} \in \mathcal{V}$  verifies equation (3.2), then necessarily  $\hat{z} = z$ . Consequently, a space larger than  $\mathcal{V}$  is required to be able to find an estimate  $\hat{z}$  with an affordable computational effort.

The interpolation space  $\mathcal{W}$  is in most cases taken to be the *broken* space, that is, the space obtained from  $\mathcal{V}$  relaxing both the Dirichlet boundary conditions and the continuity of the functions across the edges of the mesh. It is worth noting that in this case the estimate  $\hat{z} \in \mathcal{W}$  is not uniquely determined by (3.2). Thus, different estimation techniques will provide different estimates  $\hat{z}$  living in the broken space and verifying equation (3.2).

### 3.2.2 The equilibrated residual method

The derivation of equilibrated fluxes or tractions on the element boundaries goes back to the work of Ladevèze (Ladevèze 1977, Ladevèze and Leguillon 1983) and the works developed by Kelly (1984), Bank and Weiser (1985), Ainsworth and Oden (1992) and Ainsworth and Oden (1993). The literature on the computation of equilibrated fluxes or tractions is copious but there are two works worth to be highlighted. The first work by Ladevèze and Maunder (1996) provides a geometrical interpretation of the computation of equilibrated tractions whereby a comparison of different equilibration procedures is done; in particular, the geometrical interpretation allows to relate the works of Ladevèze and Leguillon (1983), Bank and Weiser (1985) and Ainsworth and Oden (1992) amongst others. The second work is the book by Ainsworth and Oden (2000). Chapter 6 provides a clear and simple exposition of the method and the computation of the equilibrated tractions.

Here, the main features common to all the equilibrated residual methods are exposed, but a detailed explanation may be found in Ainsworth and Oden (2000). First, a domain decomposition strategy is used to transform the global residual problem (3.1) into a sequence of uncoupled local boundary value problems posed over the elements of the underlying mesh. The local boundary value problems being of Neumann type ensures the upper bound property, but care must be exercised in the choice of the data for the boundary conditions (imposed tractions) to ensure that the local problems are solvable. Then, the local problems are solved independently to obtain the upper bound for the energy norm.

#### Domain decomposition

Let  $\mathcal{T}_H = \{\Omega_k\}_{k=1}^{n_{e1}}$  be the partition of the computational domain  $\Omega$  associated with the finite element interpolation space  $\mathcal{V}^H$ , where  $\Omega_k$  denotes a generic element of the mesh. Let also  $\Gamma_H$  be the set of all edges in the mesh and  $\Lambda = \prod_{k=1}^{n_{e1}} [\mathcal{H}^{-\frac{1}{2}}(\partial\Omega_k)]^{n_{sd}}$  the space of integrable tractions in  $\Gamma_H$ . The set  $\Gamma_H$  is divided into two complementary disjoint sets, the boundary edges and the set of all interior edges of the mesh denoted by  $\Gamma^I$ . For each edge  $\gamma \in \Gamma_H$  a unit normal direction,  $\mathbf{n}^\gamma$ , is assigned such that, if  $\gamma$  is a boundary edge,  $\mathbf{n}^\gamma$  coincides with the outward unit normal to  $\Gamma$ .

Similarly, given an element  $\Omega_k$  and an edge of this element  $\gamma \in \partial\Omega_k$ , the outward normal to the element associated to  $\gamma$  is denoted by  $\mathbf{n}_k^\gamma$ . Then,  $\tau_k$  is defined as  $\tau_k|_\gamma = \mathbf{n}_k^\gamma \cdot \mathbf{n}^\gamma$ , that is:

$$\tau_k|_\gamma = \mathbf{n}_k^\gamma \cdot \mathbf{n}^\gamma = \begin{cases} 1 & \text{if } \mathbf{n}_k^\gamma = \mathbf{n}^\gamma \\ -1 & \text{if } \mathbf{n}_k^\gamma = -\mathbf{n}^\gamma. \end{cases}$$

Note that if  $\gamma = \partial\Omega_k \cap \partial\Omega_l$ , then  $\tau_k|_\gamma + \tau_l|_\gamma = 0$ .

The broken space  $\widehat{\mathcal{V}}$  is introduced by relaxing in  $\mathcal{V}$  both the Dirichlet homogeneous boundary conditions and the continuity of the functions across the edges of  $\Gamma_H$ , that is,

$$\widehat{\mathcal{V}} = \{ \hat{v} \in [\mathcal{L}^2(\Omega)]^{\text{n\text{sd}}}, \hat{v}|_{\Omega_k} \in [\mathcal{H}^1(\Omega_k)]^{\text{n\text{sd}}} \forall \Omega_k \in \Omega \}.$$

Given a function in the broken space  $\hat{v} \in \widehat{\mathcal{V}}$ , the jump of  $\hat{v}$  across the mesh edges is defined as

$$[[\hat{v}]]|_\gamma = \begin{cases} \hat{v}|_{\Omega_k} \tau_k|_\gamma + \hat{v}|_{\Omega_l} \tau_l|_\gamma, & \text{if } \gamma = \partial\Omega_k \cap \partial\Omega_l \in \Gamma^{\text{I}} \\ \hat{v}, & \text{if } \gamma \in \Gamma, \end{cases}$$

where the sign of the jump depends on the arbitrary choice of the edge normals. Note that if  $\hat{v}$  is a continuous function verifying the Dirichlet homogeneous boundary conditions,  $\hat{v} \in \mathcal{V}$ , then  $[[\hat{v}]] = 0$  in  $\Gamma^{\text{I}} \cup \Gamma^{\text{D}}$ . Then, given a broken function  $\hat{v} \in \widehat{\mathcal{V}}$ , the continuity at inter-elemental edges and Dirichlet homogeneous boundary conditions in  $\Gamma^{\text{D}}$  can be enforced weakly through the bilinear form  $b : \widehat{\mathcal{V}} \times \Lambda \rightarrow \mathbb{R}$

$$b(\hat{v}, \lambda) = \sum_{\gamma \in \Gamma^{\text{I}} \cup \Gamma^{\text{D}}} \int_\gamma \lambda [[\hat{v}]] \, d\Gamma = \sum_{k=1}^{n_{\text{e1}}} \int_{\partial\Omega_k \setminus \Gamma^{\text{N}}} \tau_k \lambda \hat{v}|_{\Omega_k} \, d\Gamma,$$

by imposing  $b(\hat{v}, \lambda) = 0$  for all  $\lambda \in \Lambda$ . Therefore, the space of test functions  $\mathcal{V}$  can be recovered as

$$\mathcal{V} = \{ \hat{v} \in \widehat{\mathcal{V}}, b(\hat{v}, \lambda) = 0 \quad \forall \lambda \in \Lambda \}.$$

The goal is now to obtain an estimate  $\hat{z} \in \widehat{\mathcal{V}}$  verifying equation (3.2) involving only the solution of local problems posed on the elements of the mesh. Let  $\hat{z} \in \widehat{\mathcal{V}}$  be the solution of the residual problem

$$a(\hat{z}, \hat{v}) = R^*(\hat{v}) + b(\hat{v}, \lambda) \quad \forall \hat{v} \in \widehat{\mathcal{V}}, \quad (3.3)$$

for a given function  $\lambda \in \Lambda$ , where the bilinear form  $a(\cdot, \cdot)$  and the residue  $R^*(\cdot)$  are generalized to accept *broken* functions in its arguments; that is, for  $\hat{v}, \hat{w} \in \hat{\mathcal{V}}$

$$a(\hat{w}, \hat{v}) = \sum_{k=1}^{n_{e1}} a_k(\hat{w}, \hat{v}), \quad R^*(\hat{v}) = \sum_{k=1}^{n_{e1}} R_k^*(\hat{v}),$$

where  $a_k(\cdot, \cdot)$  and  $R_k^*(\cdot)$  are the restrictions of the bilinear form  $a(\cdot, \cdot)$  and the residue  $R^*(\cdot)$  to the element  $\Omega_k$ .

In particular, since  $\mathcal{V} \subset \hat{\mathcal{V}}$ , equation (3.3) is valid for any  $v \in \mathcal{V}$ , namely

$$a(\hat{z}, v) = R^*(v) + b(v, \lambda) \quad \forall v \in \mathcal{V},$$

and since for any  $v \in \mathcal{V}$ ,  $b(v, \lambda) = 0 \quad \forall \lambda \in \Lambda$ , the previous equation is equivalent to the condition posed by equation (3.2). Therefore the energy norm of the estimate  $\hat{z} \in \hat{\mathcal{V}}$  computed from equation (3.3) provides an upper bound for  $\|z\|$ .

Clearly, from the definition of the broken space  $\hat{\mathcal{V}}$ , the equation for the estimate (3.3) decomposes into independent local problems posed on the elements of the mesh: find  $\hat{z}^k \in \mathcal{V}_k$  such that

$$a_k(\hat{z}^k, v) = R_k^*(v) + b_k(v, \lambda) \quad \forall v \in \mathcal{V}_k, \quad (3.4)$$

where  $\mathcal{V}_k$  is the restriction of the test space  $\mathcal{V}$  to the element  $\Omega_k$  and where  $b_k(\cdot, \cdot)$  is the local counterpart of the continuity form  $b(\cdot, \cdot)$ , namely

$$b_k(v, \lambda) = \int_{\partial\Omega_k \setminus \Gamma^N} \tau_k \lambda v \, d\Gamma.$$

It is worth noting that the additional terms  $b_k(v, \lambda)$  in equation (3.4) for the element  $\Omega_k$  are additional Neumann boundary conditions on the edges of the element (unless an edge belongs to  $\Gamma^N$ ). Thus, the function  $\lambda$  corresponds to additional Neumann boundary conditions applied on the edges of the mesh. This motivates that  $\lambda$  is named after equilibrated tractions.

The estimate  $\hat{z}$  is then obtained from the local estimates  $\hat{z}^k$  extending them to  $\Omega$  by setting the values outside  $\Omega_k$  to zero and defining  $\hat{z} = \sum_{k=1}^{n_{e1}} \hat{z}^k$ . In fact, the upper bound may be computed as

$$\|z\|^2 \leq \|\hat{z}\|^2 = \sum_{k=1}^{n_{e1}} \|\hat{z}^k\|_k^2,$$

where the local energy norm  $\|\cdot\|_k$  is defined as  $\|v\|_k^2 = a_k(v, v)$ .

### Choice of the equilibrated tractions $\lambda \in \Lambda$

In the definition of the local estimates  $\hat{z}^k$ , equation (3.4), it is tacitly assumed that the local problems admit a solution. However, since the local problems are subject to pure Neumann boundary conditions (unless the element  $\Omega_k$  abuts the Dirichlet boundary  $\Gamma^D$ ) in general a solution to these problems will not exist. This is due to the fact that in general the local bilinear form  $a_k(\cdot, \cdot)$  has a nontrivial kernel. For instance, in the scalar diffusion equation if the term  $\mu$  appearing in the bilinear form vanishes  $\mu = 0$ , the kernel of  $a_k(\cdot, \cdot)$  are the constant functions. This means that unless the local r.h.s.,  $R_k^*(\cdot) + b_k(\cdot, \lambda)$ , satisfies appropriate compatibility conditions, the problem will fail to possess solution, which will be the general case. For the scalar diffusion equation these compatibility conditions will reduce to ensure that  $R_k^*(1) + b_k(1, \lambda) = 0$ , where 1 stands for the unitary function in the element  $\Omega_k$ .

For an arbitrary choice of the function  $\lambda \in \Lambda$  the compatibility conditions will not be necessarily verified leading to unsolvable local problems. Thus, the function  $\lambda \in \Lambda$  has to be properly chosen so that the associated local problems are solvable.

The local problems given in equation (3.4) will be solvable as long as  $R_k^*(v) + b_k(v, \lambda) = 0$  for all the functions  $v$  in the kernel of the local bilinear form  $a_k(\cdot, \cdot)$  (see Parés, Díez and Huerta 2005, Theorem 5), that is, the compatibility condition reads

$$R_k^*(v) + b_k(v, \lambda) = 0 \quad \forall v \in \ker a_k.$$

The different equilibration techniques differ in the choice of the equilibrated tractions  $\lambda \in \Lambda$  satisfying the compatibility condition.

In Chapter 5 the approach proposed by Ladevèze and Leguillon (1983) is used. Since  $\ker(a_k) \subset \mathcal{V}_k^H$ , where  $\mathcal{V}_k^H$  is the restriction of the interpolation space  $\mathcal{V}^H$  to the element  $\Omega_k$ , the compatibility conditions may be ensured choosing equilibrated tractions verifying

$$R_k^*(v) + b_k(v, \lambda) = 0 \quad \forall v \in \mathcal{V}_k^H,$$

for all the elements of the mesh. These conditions may be rewritten in a compact form as: chose  $\lambda$  verifying

$$R^*(\hat{v}_H) + b(\hat{v}_H, \lambda) = 0 \quad \forall \hat{v}_H \in \hat{\mathcal{V}}^H, \quad (3.5)$$



where  $\widehat{\mathcal{V}}^H = \prod_{k=1}^{n_{e1}} \mathcal{V}_k^H$  is obtained from  $\mathcal{V}^H$  relaxing both the Dirichlet boundary conditions and the continuity across the edges of the mesh.

Equation (3.5) leads to a system of equations which do not uniquely determine  $\lambda$ . Moreover, the equations are coupled in the sense that the value of the flux or tractions in an edge is determined using the two neighboring elements containing it. Fortunately, the computation of the equilibrated fluxes  $\lambda$  does not require a global computation but it can be evaluated solving local problems on patches of elements. More precisely, for each node of the mesh, a local computation is required involving the elements containing this node.

The procedure to obtain the upper bounds for the energy norm of  $z$  is summarized in the box in Figure 3.1.

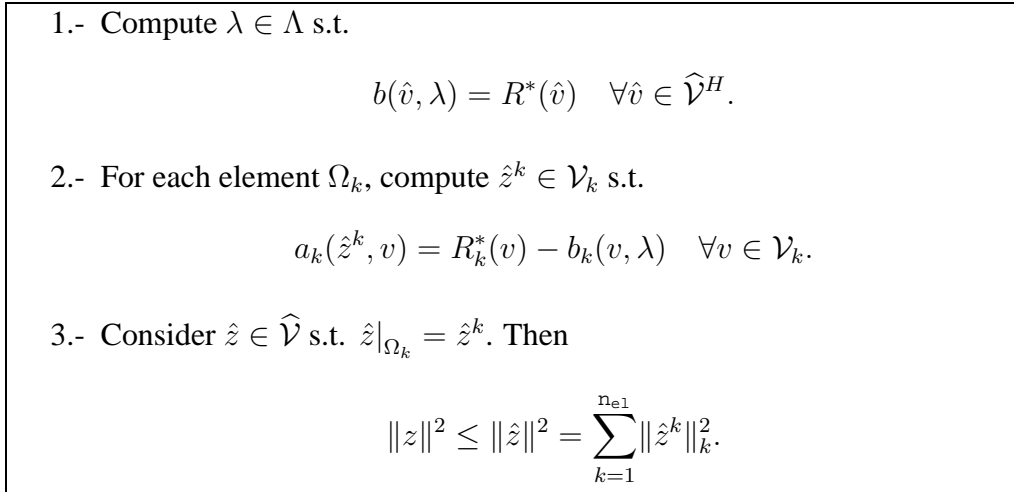


Figure 3.1: Main steps of the strategy used to obtain upper bounds for the energy norm of the solution of a symmetric boundary value problem.

### 3.3 Lower bounds for the energy

The growing interest in the computation of accurate bounds for quantities of interest has been accompanied by a surge of interest in the recovery of lower bounds for the energy norm of the error.

Although implicit residual type error estimation techniques providing lower bounds for the error measured in the energy norm as the ones presented by Díez, Egozcue and Huerta (1998) and Huerta and Díez (2000) may be used to obtain lower bounds for the energy norm, recently it has been proposed to obtain lower bounds using post-processing techniques (Babuška et al. 1999, Strouboulis et al. 2000, Prudhomme et al. 2003, Díez et al. 2003). In any goal-oriented algorithm yielding bounds for the output  $s$ , upper bounds for the energy norm of the error have to be computed. The idea is then to take advantage of the available upper bound estimates and with a simple post-process obtain lower bounds for the energy norm.

In any case, the derivation of the lower bounds for the energy norm is based in the following key property detailed by Díez et al. (2003, Theorem 2).

### 3.3.1 Sufficient condition for the lower bound property

**Lemma 3.3.1.** *Any continuous estimate  $\xi \in \mathcal{V}$  provides a parametric family of lower bounds for  $\|z\|$*

$$2\lambda R^*(\xi) - \lambda^2 \|\xi\|^2 \leq \|z\|^2, \quad (3.6)$$

depending on  $\lambda \in \mathbb{R}$ . The natural choice  $\lambda = 1$ , and the optimal choice for  $\lambda$  maximizing the lower bound,  $\lambda = \frac{R^*(\xi)}{\|\xi\|}$ , yield the bounds

$$2R^*(\xi) - \|\xi\|^2 \leq \frac{R^*(\xi)^2}{\|\xi\|^2} \leq \|z\|^2. \quad (3.7)$$

Moreover, the previous bounds are optimal whenever  $\xi = z$ .

*Proof.* Since  $\xi \in \mathcal{V}$ , it is possible to replace  $v$  by  $\xi$  in the residual equation for  $z$  (3.1). That is

$$a(z, \xi) = R^*(\xi). \quad (3.8)$$

Then, using equation (3.8), inequality (3.6) is proved considering the following algebraic manipulation:

$$0 \leq \|z - \lambda\xi\|^2 = \|z\|^2 + \lambda^2 \|\xi\|^2 - 2\lambda a(z, \xi) = \|z\|^2 + \lambda^2 \|\xi\|^2 - 2\lambda R^*(\xi).$$

Finally, the optimality of the bounds for  $\xi = z$  is proven noting that equation (3.1) with  $v = z$  yields to  $R^*(z) = \|z\|^2$ . Therefore, the bounds given in equation (3.7) are both optimal and equal to  $\|z\|^2$ .  $\square$

*Remark 3.3.1.* Given a continuous function  $\xi \in \mathcal{V}$ , the lower bound obtained in Lemma 3.3.1 by choosing the optimal parameter  $\lambda$  may be also derived from the dual characterization of the energy norm. Indeed, the energy norm of the function  $z$  may be computed as (Oden and Prudhomme 1999, Theorem 4.3)

$$\|z\| = \sup_{\xi \in \mathcal{V}} \frac{|R^*(\xi)|}{\|\xi\|},$$

and therefore,  $\forall \xi \in \mathcal{V}$

$$\frac{|R^*(\xi)|}{\|\xi\|} \leq \|z\|,$$

which is equivalent to the optimal lower bound provided by Lemma 3.3.1.

### 3.3.2 Lower bounds by post-processing

Implicit residual type error estimates providing lower bounds for the energy norm as the ones presented by Díez et al. (1998) and Huerta and Díez (2000) provide estimates  $\underline{z} \in \underline{\mathcal{V}} \subset \mathcal{V}$  solution of the optimization problem

$$\|\underline{z}\|^2 = \sup_{\xi \in \underline{\mathcal{V}}} 2R^*(\xi) - \|\xi\|^2 \leq \sup_{\xi \in \mathcal{V}} 2R^*(\xi) - \|\xi\|^2 = \|z\|^2. \quad (3.9)$$

Unfortunately, the natural choice  $\underline{\mathcal{V}} = \mathcal{V}^H$  yields to the trivial solution  $\underline{z} = 0$ . Consequently, a larger space than  $\mathcal{V}^H$  is required. Thus, care must be exercised in the choice of  $\underline{\mathcal{V}}$  since the expense of the error estimation procedure must be lower than the cost of directly computing a new approximation of the original problem. The interpolation space  $\underline{\mathcal{V}}$  is chosen so that the optimization problem (3.9) decomposes into local independent residual problems posed either in elements or in patches of elements, so that the computation of  $\underline{z}$  entails only local computations. It is worth noting that since  $\underline{\mathcal{V}} \subset \mathcal{V}$ , the estimate  $\underline{z}$  is a *continuous* approximation of the function  $z$ .

However, in goal-oriented error estimation techniques, upper bounds for the energy norm must also be computed and the idea is to use the upper bound estimates to compute lower bounds for the energy. Assume that  $\hat{z} \in \mathcal{W}$  is an upper bound estimate verifying equation (3.2) and thus yielding an upper bound of the energy norm of  $z$ , that is,  $\|z\| \leq \|\hat{z}\|$ . As mentioned before, if  $\hat{z} \in \mathcal{V}$  then necessarily  $\hat{z} = z$  and there is no need to compute a lower bound for the energy norm of the

error since the exact value is available. Obviously,  $\hat{z}$  will only be in  $\mathcal{V}$  for very particular problems and in general  $\hat{z} \notin \mathcal{V}$ . Therefore the approximation  $\hat{z}$  can not be directly taken to be the continuous function used to find a lower bound for the error. The idea is then to construct a smoothing operator from  $\mathcal{W}$  to  $\mathcal{V}$ ,  $\mathcal{S} : \mathcal{W} \rightarrow \mathcal{V}$ , and consider  $\xi = \mathcal{S}(\hat{z})$  and the associated lower bound

$$\frac{R^*(\mathcal{S}(\hat{z}))^2}{\|\mathcal{S}(\hat{z})\|^2} \leq \|z\|^2.$$

Evidently, the accuracy of the lower bound is directly related to the choice of the smoothing operator. Díez et al. (2003) and Prudhomme et al. (2003) study the choice of the smoothing operator for the thermic model problem in order to obtain good continuous approximations of the error, and thus yield accurate bounds.

The procedure to obtain lower bounds for the energy norm of  $z$  using a post-processing techniques is summarized in the box in Table 3.2.

|  |
|--|
| <p>1.- Define a smoothing operator <math>\mathcal{S} : \mathcal{W} \rightarrow \mathcal{V}</math></p> <p>2.- Consider the continuous function <math>\xi = \mathcal{S}(\hat{z})</math>.</p> <p>3.- Recover the lower bound for <math>\ z\ </math> as</p> $\frac{ R^*(\xi) }{\ \xi\ } \leq \ z\ .$ |
|--|

Figure 3.2: Main steps of the strategy used to obtain lower bounds for the energy norm of the solution of a symmetric boundary value problem.

In order to define a smoothing operator which yield a good approximation of the function  $z$ , Lemma 3.3.1 introduced by Díez et al. (2003) plays a key role. The chief purpose of the definition of the smoothing operator is to obtain an accurate lower bound for  $\|z\|$ . Consequently, the goal would be to determine a smoothing operator optimizing the lower bound, namely

$$\sup_{\mathcal{S}} \frac{R^*(\mathcal{S}(\hat{z}))^2}{\|\mathcal{S}(\hat{z})\|^2} \leq \|z\|^2.$$

Obviously, the optimal smoothing operator would lead to  $\mathcal{S}(\hat{z}) = z$ . However, one could consider a less ambitious goal and consider a parametric family of operators

$\mathcal{S}_{\text{par}}$  to determine the optimal smoothing operator in this subspace, which obviously would also provide a lower bound for  $\|z\|$

$$\sup_{\mathcal{S} \in \mathcal{S}_{\text{par}}} \frac{R^*(\mathcal{S}(\hat{z}))^2}{\|\mathcal{S}(\hat{z})\|^2} \leq \sup_{\mathcal{S}} \frac{R^*(\mathcal{S}(\hat{z}))^2}{\|\mathcal{S}(\hat{z})\|^2} \leq \|z\|^2.$$

At this point is when Lemma 3.3.1 comes in handy. The problem of determining the smoothing operator in  $\mathcal{S}_{\text{par}}$  optimizing the bound  $\frac{R^*(\mathcal{S}(\hat{z}))^2}{\|\mathcal{S}(\hat{z})\|^2}$  yields to a complex non-linear system. However, Lemma 3.3.1 suggests the alternative to optimize the natural lower bound  $2R^*(\mathcal{S}(\xi)) - \|\mathcal{S}(\xi)\|^2$  instead of the optimal lower bound  $\frac{R^*(\mathcal{S}(\hat{z}))^2}{\|\mathcal{S}(\hat{z})\|^2}$ . It is worth noting that the alternative optimization procedure,

$$\sup_{\mathcal{S} \in \mathcal{S}_{\text{par}}} 2R^*(\mathcal{S}(\xi)) - \|\mathcal{S}(\xi)\|^2 \leq \|z\|^2,$$

although still non-linear, depends quadratically on  $\mathcal{S}(\xi)$ .

The procedure to obtain the bounds consists of, first, determining the optimal smoothing operator  $\bar{\mathcal{S}}(\cdot)$  solution of

$$\bar{\mathcal{S}}(\cdot) = \arg \sup_{\mathcal{S} \in \mathcal{S}_{\text{par}}} 2R^*(\mathcal{S}(\xi)) - \|\mathcal{S}(\xi)\|^2.$$

Then, given the continuous function  $\xi = \bar{\mathcal{S}}(\hat{z})$ , the best computable lower bound for  $\|z\|$  is computed, namely

$$\frac{R^*(\bar{\mathcal{S}}(\hat{z}))^2}{\|\bar{\mathcal{S}}(\hat{z})\|^2} \leq \|z\|^2.$$

### 3.3.3 Optimization of the lower bounds

The purpose of this section is to briefly outline the main characteristics of three strategies to enhance the choice of the smoothing operator  $\mathcal{S}(\cdot)$  introduced by Díez et al. (2003) and Parés, Díez and Huerta (2005). Here it is assumed that the initial discontinuous function  $\hat{z}$  belongs to the broken space  $\widehat{\mathcal{V}}$ , that is,  $\hat{z}$  is continuous inside the elements of the mesh and only presents discontinuities on edges of the mesh. Thus, the goal is to construct a smoothing operator  $\mathcal{S} : \widehat{\mathcal{V}} \rightarrow \mathcal{V}$ . All three strategies start from the smoothing operator  $\mathcal{S}_{\text{ave}}(\cdot)$  which averages the discontinuities of a function across interelement edges, see Figure 3.3, and introduce different alternatives to improve the continuous approximations of  $z$ .

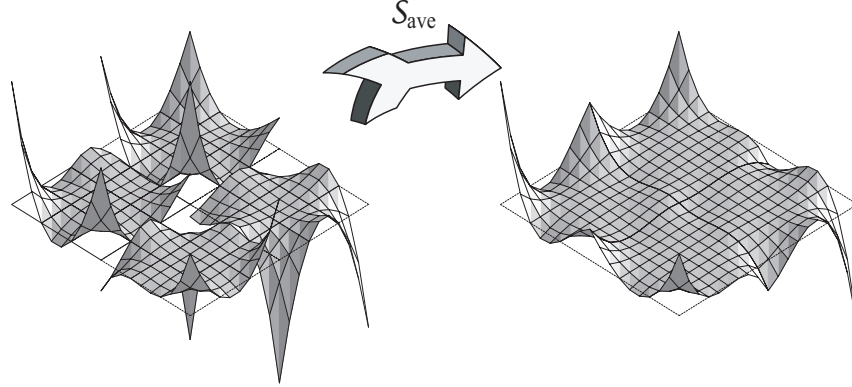


Figure 3.3: Smoothing operator  $\mathcal{S}_{\text{ave}}(\cdot)$  acting on a discontinuous function

### Interior enhancement

The averaging smoothing operator  $\mathcal{S}_{\text{ave}}(\cdot)$  acts on a function of the broken space modifying its values on the edges of the mesh to flatten the discontinuities across the edges. However, it is worth noting that the values at the nodes inside each element remain unchanged. Moreover, since the only restriction of the smoothing operator is to provide a continuous function, the values at the interior nodes in each element may be set arbitrarily since they do not affect to the continuity of the function. The goal is then to consider the optimal values of the continuous function inside each element.

Let  $\mathcal{V}_k^{\text{bub}}$  be the subspace of  $\mathcal{V}_k$  consisting of functions vanishing on the boundaries of the element  $\Omega_k$  (bubble functions, see Figure 3.4), and consider  $\mathcal{V}^{\text{bub}} = \bigoplus_{k=1}^{n_{e1}} \mathcal{V}_k^{\text{bub}}$ . Then, the continuous approximation of  $z$  is selected to be of the form  $\mathcal{S}_{\text{ave}}(\hat{z}) + z_{\text{bub}}$ , where  $z_{\text{bub}} \in \mathcal{V}^{\text{bub}}$  is found solving the optimization problem

$$z_{\text{bub}} = \arg \sup_{v \in \mathcal{V}^{\text{bub}}} 2R^*(\mathcal{S}_{\text{ave}}(\hat{z}) + v) - \|\mathcal{S}_{\text{ave}}(\hat{z}) + v\|^2.$$

Since the functions of the bubble space  $\mathcal{V}^{\text{bub}}$  vanish on the edges of each element the previous optimization problem decouples into local independent problems posed over each element of the mesh: find  $z_{\text{bub}}^k \in \mathcal{V}_k^{\text{bub}}$  such that

$$a_k(z_{\text{bub}}^k, v) = R_k^*(v) - a_k(\mathcal{S}_{\text{ave}}(\hat{z}), v) \quad \forall v \in \mathcal{V}_k^{\text{bub}},$$

and  $z_{\text{bub}}$  is then computed adding the local contributions  $z_{\text{bub}} = \sum_{k=1}^{n_{e1}} z_{\text{bub}}^k$ . Using this

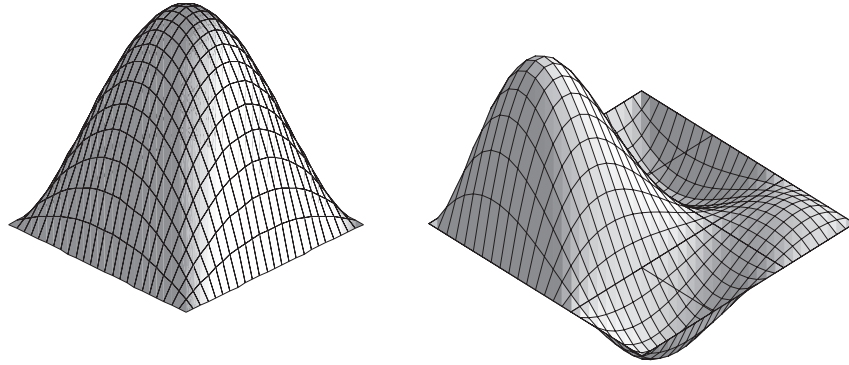


Figure 3.4: Examples of bubble functions in an element. Functions in  $\mathcal{V}_k^{\text{bub}}$ .

approach, the smoothing operator  $\mathcal{S}_{\text{int}}(\cdot)$  consists on first averaging the values of the function on the edges of the elements to recover a continuous function and then modify its values inside its element, that is,  $\mathcal{S}_{\text{int}}(\hat{z}) = \mathcal{S}_{\text{ave}}(\hat{z}) + z_{\text{bub}}$ .

### Constant fitting

In order to obtain an estimate  $\hat{z}$  verifying the condition given in equation (3.2), and thus yielding an upper bound for the energy norm, the estimate  $\hat{z}$  is computed solving local Neumann boundary value subproblems. The boundary conditions being of Neumann type means that, in general, the solution to the local problems is not unique. For instance, in the scalar diffusion equation if the term  $\mu$  appearing in the bilinear form vanishes  $\mu = 0$ , the local estimates are defined up to a constant, and in the mechanical setting, the local estimates are defined up to rigid body motions (translations and rotations). Although the value of the upper bound does not depend on the particular choice of the local estimates, this choice affects drastically to the lower bounds. To fix ideas, consider the scalar diffusion equation with no reaction term ( $\mu = 0$ ). Then, the local estimates computed using an equilibrated residual method,  $\hat{z}^k$  solution of (3.4), are determined up to a constant, that is,  $\hat{z}^k$  may be replaced by  $\hat{z}^k + c^k$  where  $c^k$  stands for a constant function on the element  $\Omega_k$ . The local norm of the estimates remains unchanged  $\|\hat{z}^k + c^k\| = \|\hat{z}^k\|$  and therefore the associated global upper bound does not depend on the choice for the constants  $c_k$ . However, in order to recover a good approximation of  $z$  from averaging the global estimate  $\hat{z} = \sum_{k=1}^{n_{\text{el}}} \hat{z}^k + c^k$ , the local constants can not be set arbitrarily. The work

presented by Díez et al. (2003) allows to optimally fit the local constants precluding the main drawback of the post-processing techniques.

The idea is really simple. Let  $\{\chi^1, \chi^2, \dots, \chi^{n_{e1}}\}$  be the space of piecewise constant functions (see Figure 3.5). Then, the global estimate provided by the equilibrated residual method is  $\hat{z} + \sum_{k=1}^{n_{e1}} c^k \chi^k$ ,  $\hat{z}$  being the global estimate computed from particular solutions of the local equations (3.4) and  $c^k \in \mathbb{R}$ . This estimate could be used to compute a continuous approximation of  $z$  using the smoothing operator described in the previous section, yielding the lower bound

$$\frac{R^*(\mathcal{S}_{\text{int}}(\hat{z} + \sum_{k=1}^{n_{e1}} c^k \chi^k))^2}{\left\| \mathcal{S}_{\text{int}}(\hat{z} + \sum_{k=1}^{n_{e1}} c^k \chi^k) \right\|^2} \leq \|z\|^2.$$

In order to fit the local constants, the following optimization procedure is considered

$$\sup_{c^k \in \mathbb{R}} 2R^*(\mathcal{S}_{\text{int}}(\hat{z} + \sum_{k=1}^{n_{e1}} c^k \chi^k)) - \left\| \mathcal{S}_{\text{int}}(\hat{z} + \sum_{k=1}^{n_{e1}} c^k \chi^k) \right\|^2.$$

In this case, the optimization with respect to the constants  $c^k$ , yields to a global linear  $n_{e1} \times n_{e1}$  system of equations (see Díez et al. 2003).

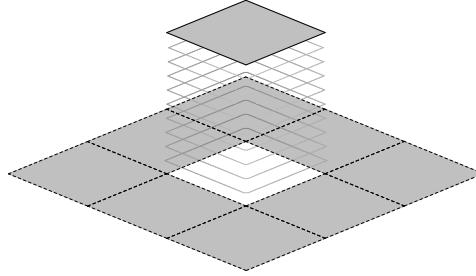


Figure 3.5: Local piecewise constant functions in an element,  $\chi^k$ .

### Global enhancement

Once the smoothing operator and thus the continuous approximation  $\xi$  are determined, the lower bound can be improved using a global computation in the coarse mesh. Note that if instead of using the approximation  $\xi$  in the expression  $\frac{R^*(\xi)^2}{\|\xi\|^2}$  one considers  $\xi + \xi_H$ ,  $\xi_H \in \mathcal{V}^H$ , then due to the Galerkin orthogonality, the only



difference in the expression is the replacement of the norm  $\|\xi\|^2$  by  $\|\xi + \xi_H\|^2$  in the denominator.

Thus, one can compute the best enhancement  $\xi_H$  maximizing the norm of  $\|\xi + \xi_H\|^2$ , that is

$$\xi_H = \arg \max_{v_H \in \mathcal{V}^H} \|\xi + v_H\|^2.$$

Hence,  $\xi_H \in \mathcal{V}^H$  is found solving the weak problem

$$a(\xi_H, v) = -a(\xi, v) \quad \forall v \in \mathcal{V}^H. \quad (3.10)$$

In this case,  $\|\xi + \xi_H\|^2 = \|\xi\|^2 - \|\xi_H\|^2$ , and thus the final lower bounds are obtained as

$$\frac{R^*(\xi)^2}{\|\xi\|^2 - \|\xi_H\|^2} \leq \|z\|^2.$$

This global enhancement in the coarse mesh coincides with the proposed pollution assessment of the error proposed by Huerta and Díez (2000). Obviously, since  $\xi$  is computed from  $\hat{z}$  which is obtained by performing only local computations, it does not account for pollution errors. The unestimated part of error,  $z - \xi$ , which includes the pollution effects, is denoted as global error and is the solution of the weak problem

$$a(z - \xi, v) = R^*(v) \quad \forall v \in \mathcal{V}. \quad (3.11)$$

Huerta and Díez (2000) assess the pollution error solving the equation for the global error (3.11) in the coarse mesh. This approach leads to the same enhancement  $\xi_H$  given in equation (3.10).

## 3.4 Implementation issues

### 3.4.1 Remarks on the derivation of upper bounds for the energy

The equilibrated residual method presented in Section 3.2.2 requires the true solution of the local problems (3.4) in order to obtain upper bounds for  $\|z\|$ . Of course, it is infeasible to require the exact solution of these problems since  $\mathcal{V}_k$  is an infinite dimensional space. In practice, the local problems are approximated using a finite dimensional subspace of  $\mathcal{V}_k$ . This leads to the loss of the upper bound property

with respect to the energy norm of the exact solution of problem (3.1),  $z$ , that is, the estimate  $\hat{z}$  does not necessarily provide an upper bound for  $\|z\|$ , it is possible to obtain an estimate  $\hat{z}$  such that  $\|\hat{z}\| \leq \|z\|$ .

However, in this case, the obtained bounds are still strict with respect to the energy norm of the reference solution  $z_h$ , where the reference solution  $z_h$  is the projection of  $z$  into the finite dimensional subspace of  $\mathcal{V}$  used to solve the local problems.

Let  $\mathcal{V}^h$  be the reference interpolation space obtained from  $\mathcal{V}^H$  either using an  $h$  or  $p$  refinement and consider the broken reference interpolation space,  $\widehat{\mathcal{V}}^h$ , obtained from  $\mathcal{V}^h$  relaxing both the Dirichlet boundary conditions and the continuity of the functions along the edges of the reference mesh. The local residual problems given by equation (3.4) are approximated by the weak problems: find  $\hat{z}_h^k \in \mathcal{V}_k^h$  such that

$$a_k(\hat{z}_h^k, v) = R_k^*(v) + b_k(v, \lambda) \quad \forall v \in \mathcal{V}_k^h, \quad (3.12)$$

where  $\mathcal{V}_k^h$  is the restriction of the reference space  $\mathcal{V}^h$  to the element  $\Omega_k$ , and in this case, the global estimate  $\hat{z}_h$  obtained adding the local estimates  $\hat{z}_h = \sum_{k=1}^{n_{e1}} \hat{z}_h^k$  is the solution of the weak problem:

$$a(\hat{z}_h, v) = R^*(v) + b(v, \lambda) \quad \forall v \in \widehat{\mathcal{V}}^h. \quad (3.13)$$

An immediate consequence of equation (3.13) is that the energy norm of  $\hat{z}_h$  is an upper bound of the energy norm of the reference solution  $z_h$ , as asserted above. Indeed, let  $z_h \in \mathcal{V}^h$  be the projection of  $z$  into  $\mathcal{V}^h$  solution of the weak problem

$$a(z_h, v) = R^*(v) \quad \forall v \in \mathcal{V}^h. \quad (3.14)$$

Since for any continuous function  $v \in \mathcal{V}^h$   $b(v, \lambda) = 0$ , replacing  $v = z_h$  in equation (3.13) yields

$$a(\hat{z}_h, z_h) = R^*(z_h),$$

and the upper bound property follows from a routine application of the Cauchy-Schwarz inequality

$$\|z_h\|^2 = R^*(z_h) = a(\hat{z}_h, z_h) \leq \|\hat{z}_h\| \|z_h\|.$$

Nonetheless, since the energy norm of the reference solution  $z_h$  underestimates the norm of  $z$ ,  $\|z_h\| \leq \|z\|$ , it is not possible, in general, to guarantee that  $\|\hat{z}_h\|$  provides an upper bound for  $\|z\|$ .

Actually, the loss of the upper bound property comes from the fact that in each element  $\|\hat{z}_h^k\|$  underestimates the norm of the local estimates  $\hat{z}^k$ , that is,  $\|\hat{z}_h^k\| \leq \|\hat{z}^k\|$ . Thus, although,  $\|z\|^2 \leq \|\hat{z}\|^2 = \sum_{k=1}^{n_{e1}} \|\hat{z}^k\|_k^2$ , it can not be guaranteed that

$$\|z\|^2 \leq \|\hat{z}_h\|^2 = \sum_{k=1}^{n_{e1}} \|\hat{z}_h^k\|_k^2.$$

Chapter 5 presents an alternative to the approximation of the local problems (3.4) using a finite dimensional subspace, which allows to obtain strict upper bounds of the energy norm of  $z$ . The outline of the method is as follows: upper bounds for the energy norm  $\|z\|$  may be obtained from an estimate  $\hat{z} \in \hat{\mathcal{V}}$  computed from the local estimates  $\hat{z}^k$  obtained solving the local problems given in equation (3.4) yielding

$$\|z\|^2 \leq \|\hat{z}\|^2 = \sum_{k=1}^{n_{e1}} \|\hat{z}^k\|_k^2.$$

If the local problems are solved using a finite dimensional subspace, that is, if instead of  $\hat{z}^k \in \mathcal{V}_k$  one computes  $\hat{z}_h^k \in \mathcal{V}_k^h$  verifying equation (3.12) then one has that  $\|\hat{z}_h^k\|_k \leq \|\hat{z}^k\|_k$  although  $\|\hat{z}_h^k\|_k$  asymptotically approaches  $\|\hat{z}^k\|_k$ . In order to obtain strict bounds for the energy norm, instead of approximating the local problems (3.12) using a submesh, the local problems are approximated using the complementary energy approach. The idea is instead of approximation the primal variables  $\hat{z}^k$ , the approximation is done in the dual variables (fluxes or stresses). The procedure provides for each element scalar quantities  $\nu^k$  such that  $\|\hat{z}^k\|_k^2 \leq \nu^k$ , and therefore the strict upper bound for the squared energy norm is found as

$$\|z\|^2 \leq \|\hat{z}\|^2 = \sum_{k=1}^{n_{e1}} \|\hat{z}^k\|_k^2 \leq \sum_{k=1}^{n_{e1}} \nu^k.$$

### 3.4.2 Remarks on the derivation of lower bounds for the energy

Standard a posteriori implicit error estimators, as the equilibrated residual method, provide upper bounds for the energy norm which are only strict upper bounds in the

asymptotic range since the local problems are usually approximated using a finite dimensional space. An estimate  $\hat{z}_h$  is obtained overestimating the energy norm of the reference solution,  $\|z_h\| \leq \|\hat{z}_h\|$ , but there is no guarantee that  $\hat{z}_h$  provides an upper bound for the exact solution  $z$ . In fact, in some numerical examples it happens that  $\|\hat{z}_h\| \leq \|z\|$ . Moreover, if  $\xi \in \mathcal{V}$  is a good continuous approximation of  $z$  providing a sharp lower bound for the energy norm of  $z$ , even the upper bound for  $\|z_h\|$  may be smaller than the lower bound for  $\|z\|$ , that is

$$\|\hat{z}_h\| \leq \frac{|R^*(\xi)|}{\|\xi\|}.$$

This result is somewhat incoherent and unsatisfactory especially when the upper and lower bounds for the energy norm are combined to assess an outputs of interest.

A consistent approach is to consider directly the goal of bounding the energy norm of the reference solution and state that these bounds are, in the asymptotic range, bounds for the norm of the exact solution. In this case, the lower bounds with respect to the norm of the reference solution may be found from Lemma 3.3.1 noting that if the continuous approximation  $\xi$  belongs to the reference interpolation space  $\mathcal{V}^h$ , then the lower bounds hold not only for the exact solution  $z$  but also for the reference solution  $z_h$ . That is, any continuous function  $\xi_h \in \mathcal{V}^h$  provides a parametric family of lower bounds for  $\|z_h\|$  and  $\|z\|$  at the same time,

$$2\lambda R^*(\xi_h) - \lambda^2 \|\xi_h\|^2 \leq \|z_h\|^2 \leq \|z\|^2.$$

Obviously, the particular bounds

$$2R^*(\xi_h) - \|\xi_h\|^2 \leq \frac{R^*(\xi_h)^2}{\|\xi_h\|^2} \leq \|z_h\|^2$$

also hold.

In some cases it may be of use to consider the nodal projection of a function in  $\mathcal{V}$  onto the reference space  $\mathcal{V}^h$ ,  $\pi^h : \mathcal{V} \rightarrow \mathcal{V}^h$ . Given a continuous approximation  $\xi \in \mathcal{V}$  of  $z$ ,  $\pi^h \xi \in \mathcal{V}^h$  provides an approximation of  $z_h$  and thus, a lower bound for the energy norm of the reference solution.

Regardless of the lower bounds being strict with respect to the energy norm of the exact or reference error, obtaining lower bounds for the energy norm is important in goal-oriented error estimation techniques.

Continuous approximations of the error fields are used in goal-oriented error estimation techniques to enhance the bounds for the output  $s$  given in equations (2.22) and (2.26) for a symmetric and nonsymmetric model problem respectively. In these expressions for the bounds, if the continuous approximations  $\xi^\pm$  of  $\kappa e \pm \frac{1}{\kappa} \varepsilon$  are taken to be zero, one would still get bounds for the error in the output  $s$ . Moreover, although the accuracy of the bounds will worsen, the rate of convergence will be the same. In fact, for regular problems where the finite element method has linear convergence, the bounds converge quadratically. In this case, although the initial bound gap  $s_u - s_l$  is large, with few effort it can be reduced to the desired tolerance using adaptive algorithms.

However, there are two important cases where the choice of the approximation function  $\xi^\pm$  is not a moot point. First, one can consider problems where the primal and dual errors are large, but nearly orthogonal, and thus, the output is nearly zero,  $s = a(e, \varepsilon) \approx 0$ . In this case, the use of  $\xi^\pm = 0$  leads to expressions for the bounds which do not take into account the orthogonality between the errors but the product of its norms yielding poor bounds. Another important case where the choice of  $\xi^\pm$  must be taken into consideration is for problems with really slow convergence. In this case, the reduction of the bound gap is costly, and thus it is important to start with good approximation of the bound gap. This appears for instance in the computation of outputs of interest in fracture mechanics where the convergence of the finite element method is slow (see Section (5.4.2)).



## Chapter 4

# Subdomain-based *flux-free* a posteriori error estimator

Parés, Díez and Huerta (2005) present a new subdomain-based flux-free error estimation technique to compute upper and lower bounds for linear-functional outputs of solutions of symmetric coercive model problems (such as the diffusion-reaction equation and the elasticity equations). The quantities of interest (functional outputs) are recovered combining upper and lower bounds of the energy norm for both the original problem (primal) and dual problem (associated with the selected functional output) using the technique detailed in Chapter 2.

The need of obtaining reliable and cost effective upper and lower bounds of the error measured in the energy norm has motivated the use of residual error estimators. Classical residual type estimators, which provide upper bounds for the energy norm of the error, require flux-equilibration procedures to properly set boundary conditions for local problems, see for instance, Ladevèze and Leguillon (1983) and Ainsworth and Oden (2000). Flux-equilibration is preformed by a non trivial algorithm, strongly dependent on the element type. Moreover, this procedure requires a data structure that is not natural in a standard finite element code; for instance, edges sides must be ordered and classified accounting for their nodes (improper order) and the elements they belong to.

The idea of using flux-free estimates, based on the partition-of-the-unity concept and using local subdomains different than elements, was first introduced by Babuška and Rheinboldt (1978a). Although their presented approach did not provide one-sided bounds of the error measured in the energy norm, it was the first to introduce

the basic property of subdomain error estimates: the decomposition of the residual into local contributions using the partition-of-the-unity property. Afterwards, Carsensen and Funken (1999/00), Machiels, Maday and Patera (2000), Morin, Nochetto and Siebert (2003) and Prudhomme, Nobile, Chamoin and Oden (2004) developed different flux-free estimates providing upper bounds for the energy norm of the error for solutions of the diffusion-reaction equation.

The main advantage of the flux-free approach is the simplicity of its implementation. Obviously, this has special relevance in 3D problems. Boundary conditions of local problems are trivial and the usual data structure of a finite element code is directly employed. Recently, Choi and Paraschivoiu (2004) compared flux-free estimates with standard *hybrid-flux* estimates in terms of both sharpness (effectivity) and computational efficiency. The main conclusion of this investigation is that, in most of test cases, the hybrid-flux estimates are more accurate while the overall computational cost is lower for flux-free estimates.

Parés, Díez and Huerta (2005) introduce a new flux-free error estimator improving the effectivity of previous approaches and with further implementation simplifications. The present chapter provides a brief description of this method. The idea is not to provide a complete description of the method, but rather to present the method in the same notation introduced in Chapter 3 so that the relation between the different estimation techniques becomes more apparent. After introducing the model problem, the methodology to obtain upper and lower bounds for the energy norm is presented. Although Parés, Díez and Huerta (2005) provide bounds for linear outputs, here, for clarity of exposition, only energy norm estimates are considered. The reader is referred to Chapter 2 to obtain bounds for the output from bounds for the energy norm. The new estimation technique is then compared with the one proposed by Carstensen and Funken (1999/00), Machiels et al. (2000), Morin et al. (2003) and Prudhomme et al. (2004). It is shown that the newly developed technique computes sharper bounds. The chapter concludes with the discussion of some computational aspects and the presentation of a numerical example.



## 4.1 Model problem

Consider a general symmetric coercive variational problem given in weak form as: find  $z \in \mathcal{V}$  such that

$$a(z, v) = R^*(v) \quad \forall v \in \mathcal{V}, \quad (4.1)$$

where the residue  $R^*(\cdot)$  is a linear functional orthogonal to the finite element space  $\mathcal{V}^H$ , that is,  $R^*(v_H) = 0 \quad \forall v_H \in \mathcal{V}^H$ . Note that this problem includes both the primal and dual residual problems. Indeed, taking  $R^*(v) = R^P(v)$ , equation (4.1) leads to the equation for the primal error (2.3) and thus,  $z = e$ . Similarly,  $R^*(v) = R^D(v)$  would lead to the dual residual problem of equation (2.11), that is,  $z = \varepsilon$ . In fact, any linear combination of the primal and dual errors,  $\alpha e + \beta \varepsilon$ ,  $\alpha, \beta \in \mathbb{R}$ , is the solution of a generalized problem with  $R^*(v) = \alpha R^P(v) + \beta R^D(v)$ . This is the reason why  $z$  will be often denoted as the error function.

## 4.2 Upper bound for the energy

In order to find upper bounds for the energy norm  $\|z\|$ , some notation must be introduced. Let  $x^i$ ,  $i = 1 \dots n_{\text{np}}$  denote the set of vertices of the finite mesh and  $\phi^i$  denote the corresponding first-order Lagrange shape functions which are characterized by  $\phi^i(x^j) = \delta_{ij}$  and by being a partition of unity, that is

$$\sum_{i=1}^{n_{\text{np}}} \phi^i(x) = 1 \quad \forall x \in \Omega. \quad (4.2)$$

The support of the nodal function  $\phi^i$ , is denoted by  $\omega^i$  and it is called the star centered in, or associated with, node  $x^i$ . It consists of the patch of elements containing vertex  $x^i$ .

The method is formulated starting with the residual equation (4.1). The underlying idea is to replace the global problem characterizing the exact solution  $z$ , by a sequence of independent problems posed on the stars. The basis functions,  $\phi^i$ , may be utilized in this purpose. With the aid of the partition of unity property (4.2) and the linearity of the residue,  $R^*(\cdot)$ , it follows that  $R^*(\cdot)$  may be decomposed into the

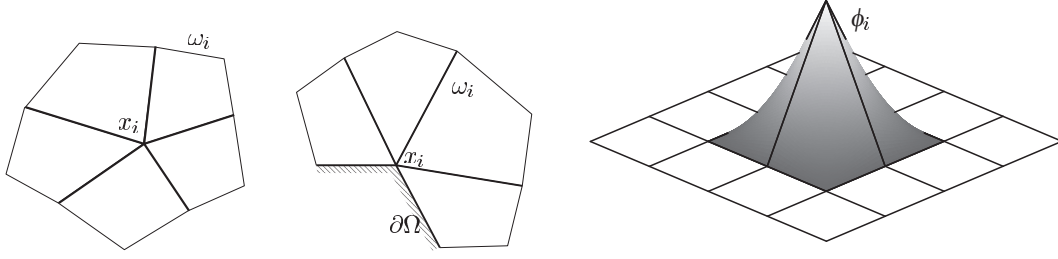


Figure 4.1: Representation of two stars centered in a node  $x_i$ ,  $\omega^i$  (left) and representation of a shape function  $\phi^i$ .

form

$$R^*(v) = R^*\left(\sum_{i=1}^{n_{np}} \phi^i v\right) = \sum_{i=1}^{n_{np}} R^*(\phi^i v), \quad (4.3)$$

for every  $v \in \mathcal{V}$ , and hence the global residual equation (4.1) may be rewritten as

$$a(z, v) = \sum_{i=1}^{n_{np}} R^*(\phi^i v) \quad \forall v \in \mathcal{V}.$$

An immediate consequence of this decomposition is that the error function  $z$  may be decomposed into  $z = \sum_{i=1}^{n_{np}} z^i$ , where  $z^i \in \mathcal{V}$  are the solutions of the global problems

$$a(z^i, v) = R^*(\phi^i v) \quad \forall v \in \mathcal{V}, \quad (4.4)$$

and hence, the energy norm of  $z$  can be recovered as

$$\|z\| = \left\| \sum_{i=1}^{n_{np}} z^i \right\|. \quad (4.5)$$

Of course, it is infeasible to solve for the exact solution of every global problem (4.4) which have the same complexity as the original one (4.1).

However, note that the function  $\phi^i v$  appearing on the r.h.s. of equation (4.4) is supported on  $\omega^i$ , thus,  $R^*(\phi^i v)$  is also supported on  $\omega^i$ , that is  $R^*(\phi^i v) = R^*(\phi^i v|_{\omega^i})$ , where  $v|_{\omega^i} \in \mathcal{V}(\omega^i)$ . The local space  $\mathcal{V}(\omega^i)$  will be denoted in the following by  $\mathcal{V}_{\omega^i}$ , and the local bilinear form associated with this space, which is the restriction of the bilinear form  $a(\cdot, \cdot)$  to the star  $\omega^i$ , is denoted by  $a_{\omega^i} : \mathcal{V}_{\omega^i} \times \mathcal{V}_{\omega^i} \rightarrow \mathbb{R}$ .

The subdomain residual problem is an approximation of the global problem (4.4) and consists of finding  $\hat{z}^i \in \mathcal{V}_{\omega^i}$  such that

$$a_{\omega^i}(\hat{z}^i, v) = R^*(\phi^i v) \quad \forall v \in \mathcal{V}_{\omega^i}. \quad (4.6)$$

Formally the functions  $\hat{z}^i$  are not defined in the whole domain  $\Omega$  but only in the star  $\omega^i$ . However, they can be naturally extended to  $\Omega$  by setting the values outside  $\omega^i$  to zero leading to a function which is generally discontinuous across the boundary of the star  $\mathcal{V}_{\omega^i}$ , that is,  $\hat{z}^i \in \widehat{\mathcal{V}}$ , where  $\widehat{\mathcal{V}}$  is the broken space introduced from  $\mathcal{V}$  relaxing the continuity of the functions across the edges of the mesh.

Motivated by equation (4.5), the global error estimate  $\hat{z}$  is obtained adding the local estimators  $\hat{z} = \sum_{i=1}^{n_{\text{np}}} \hat{z}^i$  and the upper bound for the energy norm is directly the norm of  $\hat{z}$ , that is,

$$\|z\| \leq \|\hat{z}\| = \left\| \sum_{i=1}^{n_{\text{np}}} \hat{z}^i \right\|.$$

Parés, Díez and Huerta (2005, Section 4.1) discuss the solvability of local problems (4.6). These local problems admit a solution if and only if the following compatibility condition holds

$$R^*(\phi^i v) = 0 \quad \forall v \in \ker a_{\omega^i},$$

where

$$\ker a_{\omega^i} = \{v \in \mathcal{V}_{\omega^i}, a_{\omega^i}(v, w) = 0 \forall w \in \mathcal{V}_{\omega^i}\},$$

see Parés, Díez and Huerta (2005, Theorem 5).

The solvability of the local variational problems depend on the verification of the compatibility condition for the functions in the kernel of the local bilinear operator  $a_{\omega^i}(\cdot, \cdot)$ . Thus, it depends on the model problem at hand.

The solvability in many cases (scalar diffusion-reaction equation, elasticity equations with higher-order elements, . . .) follows from the orthogonality of the residue  $R^*(\cdot)$  to the finite element space  $\mathcal{V}^H$ . However, local problems (4.6) are not solvable in all the cases. In those cases, Parés, Díez and Huerta (2005, Section 3.2) introduce a modification of the local residual problems (4.6) ensuring solvability and maintaining the upper bound property. The modified local equations are

$$a_{\omega^i}(\hat{z}^i, v) = R^*(\phi^i(v - \pi^H v)) \quad \forall v \in \mathcal{V}_{\omega^i}, \quad (4.7)$$

where  $\pi^H$  denote the nodal projection of a function in  $\mathcal{V}$  onto the finite element space  $\mathcal{V}^H$ , that is, for every vertex  $x^i$  of the finite element mesh,  $\pi^H : \mathcal{V} \longrightarrow \mathcal{V}^H$  is such that  $\pi^H(v)(x^i) = v(x^i)$ .

In this case, the compatibility condition reads

$$R^*(\phi^i(v - \pi^H v)) = 0 \quad \forall v \in \ker a_{\omega^i},$$

which always holds since  $\ker a_{\omega^i} \subset \mathcal{V}^H$ . Indeed, if  $v \in \ker a_{\omega^i}$ ,  $v \in \mathcal{V}^H$ , then  $v - \pi^H v = 0$ . Therefore the modified equation for the local estimates always admits a solution.

The global estimate  $\hat{z} = \sum_{i=1}^{n_{np}} \hat{z}^i \in \hat{\mathcal{V}}$  defined adding the local estimates  $\hat{z}^i$ , solution of either the residual problem (4.6) or (4.7), verifies the hypothesis of Lemma (3.2.1) and therefore, its energy norm provides an upper bound for  $\|z\|$  (Parés, Díez and Huerta 2005, Lemma 7).

Appendix B and Parés, Díez and Huerta (2005, Section 5) present a comparison of this new estimate with the estimates introduced by Carstensen and Funken (1999/00), Morin et al. (2003), Machiels et al. (2000) and Prudhomme et al. (2004) emphasizing the novelties of the presented approach.

### 4.3 Lower bound for the energy

The upper bound of the energy norm of the function  $z$ ,  $\|\hat{z}\|$ , is associated with the estimate  $\hat{z} \in \hat{\mathcal{V}}$  of the function  $z$  itself. The upper bound property is intrinsically related with the broken (discontinuous) nature of  $\hat{z}$ . On the contrary, a lower bound estimate is easily recovered from a continuous estimate of the function  $z$ , see Lemma 3.3.1. Thus, once  $\hat{z}$  is obtained, a continuous estimate,  $\xi \in \mathcal{V}$ , is computed from  $\hat{z}$ . Two different alternatives may be considered to compute  $\xi$  from  $\hat{z}$ . First, the strategy presented in detail by Díez et al. (2003) and outlined in Section 3.3.2 which is valid for any discontinuous estimate  $\hat{z}$  (discontinuous across inter-element edges or faces) can be readily implemented. Second, the weighting strategy, where the continuous estimate is obtained from

$$\xi = \sum_{i=1}^{n_{np}} \phi^i \hat{z}^i. \quad (4.8)$$

This approach uses the fact that the local estimates  $\hat{z}^i$  are continuous in each star. The discontinuities of  $\hat{z}^i$  on the boundary of each star  $\omega^i$  are smoothed by multi-

plying by  $\phi^i$ , which vanishes along the boundary of  $\omega^i$ . Consequently, this is the natural choice for the estimate  $\hat{z}$  presented in this chapter.

Moreover, in order to improve the quality of the estimate, any of the enhancements presented in Section 3.3.3 can be implemented.

## 4.4 Computational aspects

The subdomain-based flux-free residual method described above requires the true solution of the local problems (4.6) or (4.7). Of course, it is infeasible to require the exact solution of these problems since  $\mathcal{V}_{\omega^i}$  is an infinite dimensional space. In practice, the subdomain residual problems are approximated using a finite dimensional subspace of  $\mathcal{V}_{\omega^i}$ . This leads to the loss of the upper bound property with respect to the exact norm of the solution,  $\|z\|$ , but the bounds are still strict with respect to the energy norm of a reference solution. Strict upper bounds for the energy norm of the exact weak solution  $z$  using flux-free estimates could be obtained extending the ideas presented in Chapter 5. Instead of approximating the local problems (4.6) or (4.7) using a finite dimensional subspace, the complementary energy approach could be used to transform the local problems into computable feasibility problems without losing the upper bound property. The derivation of strict upper bounds for the energy norm of the error using flux-free estimates has not been yet exploited.

Let  $\mathcal{V}^h$  be a reference interpolation space obtained from  $\mathcal{V}^H$  either using an  $h$  or  $p$  refinement. Then, the subdomain residual problems given by equations (4.6) or (4.7) are approximated by: find  $\hat{z}_h^i \in \mathcal{V}_{\omega^i}^h$  such that

$$a_{\omega^i}(\hat{z}_h^i, v) = R^*(\phi^i v) \quad \text{or} \quad a_{\omega^i}(\hat{z}_h^i, v) = R^*(\phi^i(v - \pi^H v)) \quad \forall v \in \mathcal{V}_{\omega^i}^h, \quad (4.9)$$

where  $\mathcal{V}_{\omega^i}^h = \mathcal{V}^h(\omega^i)$  is the restriction of the reference space to the star. This leads to the estimate  $\hat{z}_h = \sum_{i=1}^{n_{\text{np}}} \hat{z}_h^i$  verifying

$$\|z_h\| \leq \|\hat{z}_h\|,$$

where  $z_h$  is the reference solution, that is, the projection of the solution  $z$  into the reference space  $\mathcal{V}^h$  (see equation (3.14)).

When a reference mesh is used to compute an upper bound for  $\|z\|$ , that is, whenever an estimate  $\hat{z}_h \in \mathcal{V}^h$  is computed such that  $\|z_h\| \leq \|\hat{z}_h\|$ , the upper bound property does not necessarily hold for  $\|z\|$ . In fact, if a lower bound is computed from  $\hat{z}_h$  using the weighting strategy (see equation (4.8)), it may occur that the lower bound is larger than the upper bound due to the fact that the continuous estimate  $\xi$  does not belong to  $\mathcal{V}^h$  but in a *larger* subspace of  $\mathcal{V}$ . In these cases, the projection of  $\xi$  into  $\mathcal{V}^h$ ,  $\pi^h(\xi)$ , can be taken as a continuous approximation of  $z_h$  (see Section 3.4.2), that is, the continuous estimate could be obtained from

$$\xi = \pi^h \left( \sum_{i=1}^{n_{mp}} \phi^i \hat{z}_h^i \right).$$

Parés, Díez and Huerta (2005, Section 7) comment some implementation issues which drastically simplify the implementation of flux-free estimates. The final algorithm results in a simple implementation specially for 3D applications.

## 4.5 Numerical examples

In this section, the behavior of the new estimate presented above is analyzed for a mechanical problem. A square thin plate with two holes proposed by Paraschivoiu et al. (1997) is considered. This is a plane-stress linear elastic problem loaded with a horizontal unit tension along the vertical edges  $\Gamma_0$ , see Figure 4.2. Note that the solution of this problem,  $\mathbf{u}$ , has corner singularities due to the interior rectangular cut-outs. Due to symmetry, only one fourth of the domain is analyzed. Values for Young's modulus and Poisson ratio are set to 1 and 0.3, respectively. Two meshes are considered, a coarse uniform mesh with 70 nodes and a finer one with 985 nodes, adapted heuristically. Upper and lower bounds for the energy norm of the reference error  $\|e_h\|$  ( $E_u$  and  $E_l$  respectively) are computed for both cases and the results are summarized in Table 4.1. The quality of the error estimates is measured with the index

$$\rho := \frac{\text{estimated error norm}}{\text{reference error norm}} - 1.$$

Index  $\rho$  is the usual effectivity index minus one. The accuracy of an error estimate is given by the absolute value of  $\rho$  and the sign indicates if the estimate is an

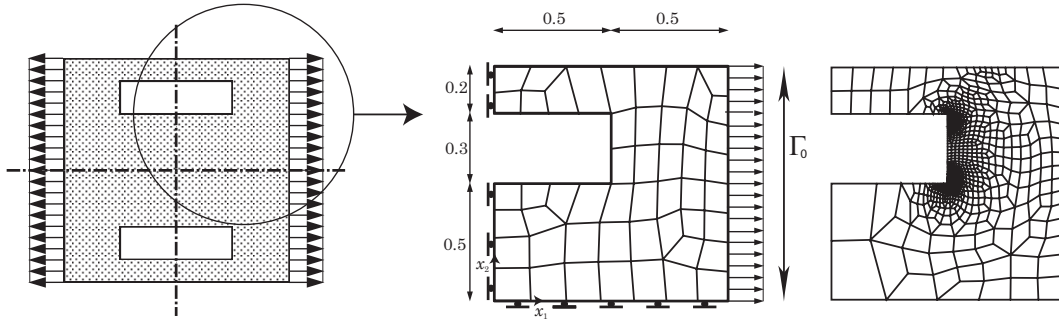


Figure 4.2: Thin plate model problem and meshes with 140 d.o.f. (center) and 1970 d.o.f.(right)

overestimation (positive  $\rho$ ) or an underestimation (negative  $\rho$ ) of the true error. For instance,  $\rho = 2\%$  indicates that the estimated error is larger than the reference error with a factor 1.02 and  $\rho = -3\%$  means that the reference error is underestimated by a factor 0.97. The effectivity index of the upper bound estimate is similar for the

Table 4.1: Upper and lower bounds for  $\|e_h\|$

| d.o.f. | $\ e_h\ $ | $\frac{\ e_h\ }{\ u_h\ }$ | $\rho(E_u)$ | $\rho(E_l)$ |
|--------|-----------|---------------------------|-------------|-------------|
| 140    | 0.146     | 12.8%                     | 17.9%       | -68.7%      |
| 1970   | 0.040     | 3.44%                     | 17.1%       | -70.1%      |

two meshes, and close to 1.17 ( $\rho \approx 17\%$ ). The lower bound effectivities are not as sharp, they are close to 0.3 ( $\rho \approx -70\%$ ). Spatial distributions of error  $E_u$  are displayed in Figures 4.3 and 4.4 for the uniform and adapted meshes, respectively.

It is worth noting that the error distributions for  $E_u$  are in good agreement with the reference error. The bad behavior of the local effectivity index in the first mesh, see Figure 4.3, is due to the fact that practically all the error is concentrated in a few relevant elements. The histogram in Figure 4.4 is narrow because the number of elements in the zones where the error is relevant is much higher for the second mesh.

Finally, Figure 4.5 shows a comparison between the proposed upper bound estimate,  $E_u$ , the flux-free techniques proposed by Machiels et al. (2000),  $E_u^\sigma$ , and by Carstensen and Funken (1999/00) and Morin et al. (2003),  $E_u^\phi$ , and a hybrid-flux upper bound estimate,  $E_u^{\text{hf}}$ , see (Ladevèze and Leguillon 1983, Ainsworth and Oden 2000). The estimates presented by Machiels et al. (2000), Carstensen and

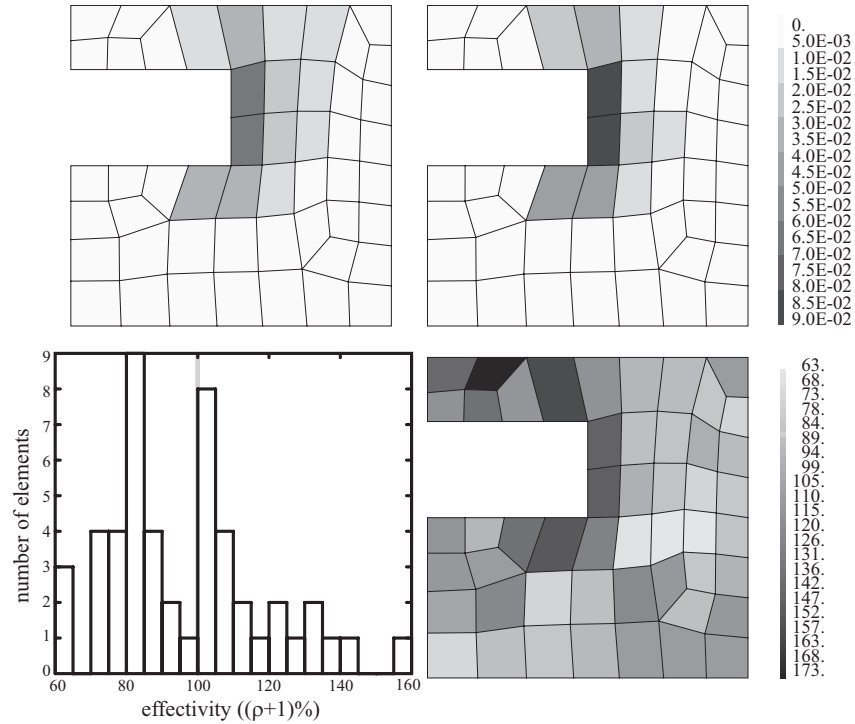


Figure 4.3: Spatial distribution of the reference error (top left), estimate  $E_u$  (top right), and local distribution of the effectivity indices  $(\rho + 1)\%$  (bottom) for the mesh with 140 d.o.f.

Funken (1999/00) and by Morin et al. (2003) are also compared with the proposed approach in Appendix B and in (Parés, Díez and Huerta 2005, Section 5).

The upper bound estimates are computed for a series of adapted triangular meshes. As expected all of them converge. Moreover, this is an example in which the hybrid-flux bound is sharper than the previously published flux-free upper bound estimates. In (Choi and Paraschivoiu 2004) the majority of the examples behave similarly. However, as already discussed the proposed flux-free bound is almost as sharp as the hybrid-flux one.



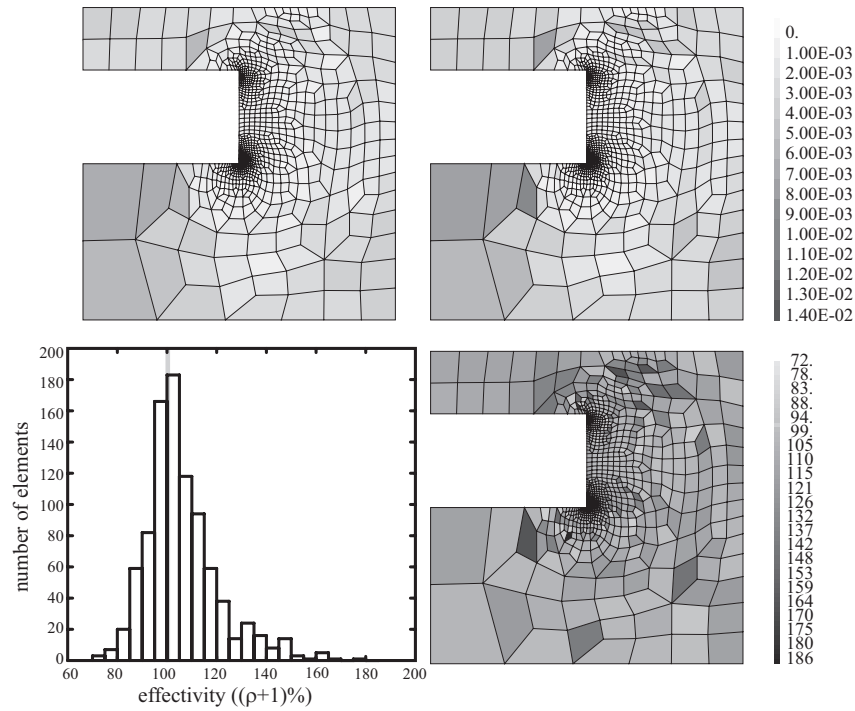


Figure 4.4: Spatial distribution of the reference error (top left), estimate  $E_u$  (top right), and local distribution of the effectivity indices  $(\rho + 1)\%$  (bottom) for the mesh with 1970 d.o.f.

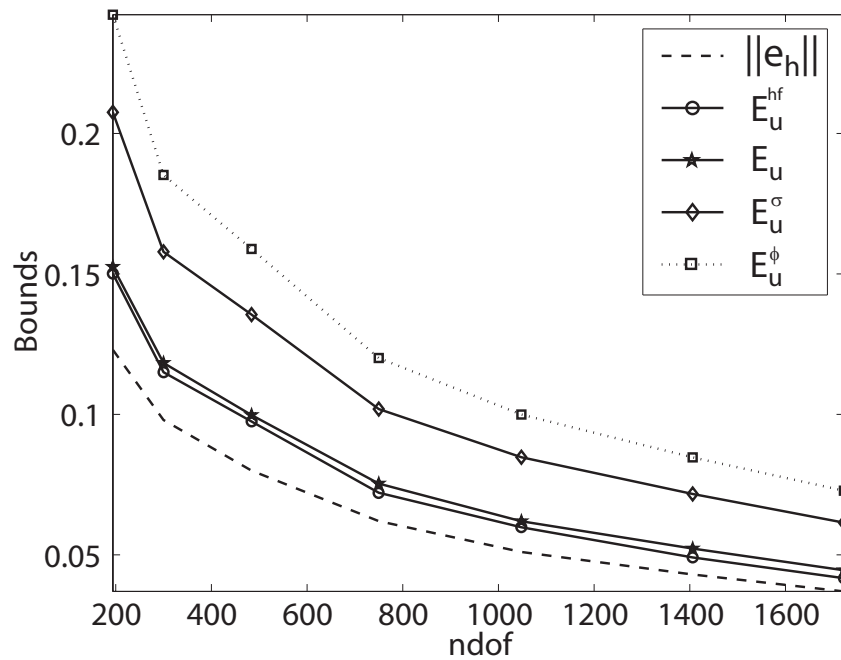


Figure 4.5: Comparison between flux-free and hybrid-flux estimates



## Chapter 5

# Strict bounds for the energy norm of weak solutions to the elasticity equations

Parés, Bonet, Huerta and Peraire (2005) present a method to compute upper and lower bounds for linear-functional outputs of the exact solutions of two dimensional elasticity equations. The method can be regarded as a generalization of the well known complementary energy principle. The desired output is cast as the supremum of a quadratic-linear convex functional over an infinite dimensional domain. Using duality the computation of an upper bound for the output of interest is reduced to a feasibility problem for the complementary, or dual, problem. In order to make the problem tractable, from a computational perspective, two additional relaxations that preserve the bounding property are introduced. First, the domain is triangulated and a domain decomposition strategy is used to generate a sequence of independent problems to be solved over each triangle. The Lagrange multipliers enforcing continuity are approximated using piecewise linear functions over the edges of the triangulation. Second, the solution of the adjoint problem is approximated over the triangulation using a standard Galerkin finite element approach. A lower bound for the output of interest is computed by repeating the process for the negative of the output. Reversing the sign of the computed upper bound for the negative of the output yields a lower bound for the actual output. The method can be easily generalized to three dimensions. However, a constructive proof for the existence of feasible solutions for the outputs of interest is only given in two dimensions. The

computed bound gap is found to converge optimally, that is, at the same rate as the finite element approximation. An attractive feature of the proposed approach is that it allows to generate a data set that can be used to certify and document the computed bounds. Using this data set and a simple algorithm, the correctness of the computed bounds can be established without recourse to the original code used to compute them.

The present chapter provides a brief description of the method presented by Parés, Bonet, Huerta and Peraire (2005). The idea is not to provide a complete description of the method, but rather to present the method in the same notation introduced in Chapter 3. After introducing the model problem, the methodology to obtain strict upper bounds for the energy norm of weak solutions of the elasticity equations is presented. Although Parés, Bonet, Huerta and Peraire (2005) provide strict bounds for linear outputs, here, for clarity of exposition, only energy norm estimates are considered. The reader is referred to Chapter 2 to obtain bounds for the output from bounds for the energy norm. Finally, the estimation procedure is used in two numerical examples: two mechanical test where the outputs vary from the average of the displacements in a part of the boundary to the  $J$ -integral (which is a non-linear output of the displacements).

## 5.1 Model problem

Consider the generalized elasticity problem with Neumann and homogeneous Dirichlet boundary conditions written in weak form as: find  $z \in \mathcal{V}$  such that

$$a(z, v) = R^*(v) \quad \forall v \in \mathcal{V}, \quad (5.1)$$

where  $\mathcal{V} = \{v \in [\mathcal{H}^1(\Omega)]^2, v|_{\Gamma_D} = \mathbf{0}\}$  acts both as the space of admissible displacement fields and the space of test functions. The linear forcing functional  $R^* \in \mathcal{V}'$

$$R^*(v) = \int_{\Omega} \mathbf{f}^* \cdot v \, d\Omega + \int_{\Gamma^N} \mathbf{g}^* \cdot v \, d\Gamma - a(z_H, v), \quad (5.2)$$

contains both the internal forces per unit volume  $\mathbf{f}^* \in [\mathcal{H}^{-1}(\Omega)]^2$  and the Neumann boundary tractions  $\mathbf{g}^* \in [\mathcal{H}^{-\frac{1}{2}}(\Gamma^N)]^2$  and  $a : \mathcal{V} \times \mathcal{V} \rightarrow \mathbb{R}$  is the symmetric coercive

bilinear form given by

$$a(\boldsymbol{w}, \boldsymbol{v}) = \int_{\Omega} \boldsymbol{\sigma}(\boldsymbol{w}) : \boldsymbol{\varepsilon}(\boldsymbol{v}) \, d\Omega.$$

Here,  $\boldsymbol{\varepsilon}(\boldsymbol{v})$  is the second order deformation tensor, which is defined as the symmetric part of the gradient tensor  $\nabla \boldsymbol{v}$ , that is,  $\boldsymbol{\varepsilon}(\boldsymbol{v}) = \frac{1}{2}(\nabla \boldsymbol{v} + (\nabla \boldsymbol{v})^T)$ . The stress tensor  $\boldsymbol{\sigma}(\boldsymbol{v})$ , is related to the deformation tensor through a linear constitutive relation of the form,  $\boldsymbol{\sigma}(\boldsymbol{v}) = \mathbb{C} : \boldsymbol{\varepsilon}(\boldsymbol{v})$ , where  $\mathbb{C}$  is the fourth-order elasticity tensor.

## 5.2 Strict upper bound for the energy

This section addresses the problem of computing upper bounds for the energy norm of the weak solution  $\boldsymbol{z}$ . Our point of departure are the local residual problems (3.4) introduced in Chapter 2: find  $\hat{\boldsymbol{z}}^k \in \mathcal{V}_k$  such that

$$a_k(\hat{\boldsymbol{z}}^k, \boldsymbol{v}) = R_k^*(\boldsymbol{v}) + b_k(\boldsymbol{v}, \boldsymbol{\lambda}) \quad \forall \boldsymbol{v} \in \mathcal{V}_k, \quad (5.3)$$

where  $\boldsymbol{\lambda}$  are the equilibrated local tractions ensuring the solvability of the local problems. The upper bound for the squared energy norm of  $\boldsymbol{z}$  is recovered from these local estimates as

$$\|\boldsymbol{z}\|^2 \leq \sum_{k=1}^{n_{e1}} \|\hat{\boldsymbol{z}}^k\|_k^2.$$

The local problems (5.3), although local, can not be solved exactly because  $\mathcal{V}_k$  is an infinite dimensional space. Moreover, if we replace  $\mathcal{V}_k$  with a finite dimensional subspace, the upper bound property is lost as shown in Section 3.4.1.

Since the local norms  $\|\hat{\boldsymbol{z}}^k\|_k^2$  can not be computed exactly, one faces the problem of finding computable upper bounds for these quantities, that is, find  $\nu^k \in \mathbb{R}$  such that

$$\|\hat{\boldsymbol{z}}^k\|_k^2 \leq \nu^k.$$

This would lead to the global upper bound

$$\|\boldsymbol{z}\|^2 \leq \sum_{k=1}^{n_{e1}} \|\hat{\boldsymbol{z}}^k\|_k^2 \leq \sum_{k=1}^{n_{e1}} \nu^k.$$

The upper bounds  $\nu^k$  are computed using a standard duality argument which transforms the problem of finding the solution of equation (5.3) over the infinite

dimensional space  $\mathcal{V}_k$  to a problem of finding a feasible solution in an appropriate finite dimensional space. Instead of approximating the primal variables (displacements)  $\hat{\mathbf{z}}_k$  over a finite dimensional space, which yields to a loss of the upper bound property, the approximation is done in the dual variables (stresses)  $\boldsymbol{\sigma}(\hat{\mathbf{z}}_k)$ , which allows to obtain the desired result.

Let  $\mathcal{S}_k$  denote the space of componentwise square-integrable stress fields in  $\Omega_k$ , that is,  $\mathcal{S}_k$  contains all the second-order tensors with  $\sigma_{ij} \in \mathcal{L}^2(\Omega_k) \quad \forall i, j$ . Then,  $\mathcal{S}_k^{eq}$  denotes the subset of  $\mathcal{S}_k$  which contains all the equilibrated stress fields with respect to  $R^*$  and  $\boldsymbol{\lambda}$ , that is,  $\hat{\boldsymbol{\sigma}}^k \in \mathcal{S}_k^{eq}$  verifies

$$\int_{\Omega_k} \hat{\boldsymbol{\sigma}}^k : \boldsymbol{\varepsilon}(\mathbf{v}) \, d\Omega = R_k^*(\mathbf{v}) + b_k(\mathbf{v}, \boldsymbol{\lambda}) \quad \forall \mathbf{v} \in \mathcal{V}_k. \quad (5.4)$$

The stress fields in  $\mathcal{S}_k^{eq}$  are usually referred to as being statically admissible. In addition, the complementary energy of a stress field  $\boldsymbol{\sigma}^k \in \mathcal{S}_k$  is defined as

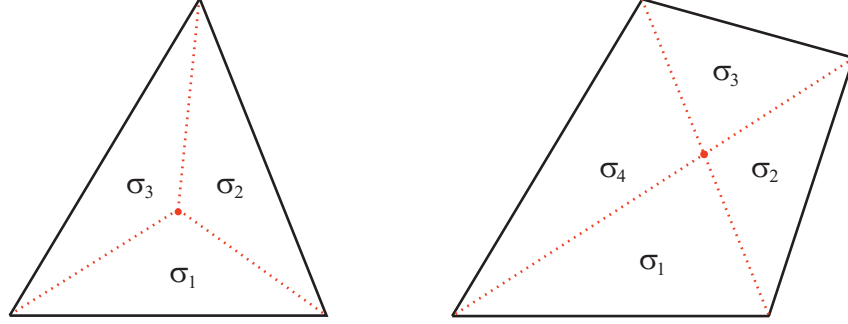
$$||| \boldsymbol{\sigma}^k |||_k^2 = \int_{\Omega_k} \boldsymbol{\sigma}^k : \mathbb{C}^{-1} : \boldsymbol{\sigma}^k \, d\Omega.$$

Parés, Bonet, Huerta and Peraire (2005, Lemma 1) provide the key to obtaining the local upper bounds  $\nu^k$ . It is sufficient to compute a statically admissible stress field  $\hat{\boldsymbol{\sigma}}^k \in \mathcal{S}_k^{eq}$ , and then evaluate its complementary energy. This follows from the fact that for any admissible stress field  $\hat{\boldsymbol{\sigma}}^k \in \mathcal{S}_k^{eq}$

$$\|\hat{\mathbf{z}}^k\|_k^2 \leq ||| \hat{\boldsymbol{\sigma}}^k |||_k^2.$$

Moreover Parés, Bonet, Huerta and Peraire (2005) show that one can chose the statically admissible stress field to be piecewise polynomial and provide a constructive proof of the existence of a piecewise polynomial equilibrated stress fields. The only requirement is that the forcing data  $\mathbf{f}^*$ ,  $\mathbf{g}^*$  and  $\boldsymbol{\lambda}$ , and the displacement field  $\mathbf{z}_H$  have to be piecewise polynomial functions.

The key point is to divide each element into three or four triangles (depending on the initial element being a triangular or a quadrilateral element respectively) and then consider the statically admissible stress field to be polynomial in each subtriangle (see Figure 5.1). The degree of the local polynomial fields in each subtriangle depends on the degrees of the forcing data  $\mathbf{f}^*$ ,  $\mathbf{g}^*$  and  $\boldsymbol{\lambda}$  and the displacement field  $\mathbf{z}_H$ . For instance, if  $\mathbf{f}^*$  is a constant distribution of internal forces in the element  $\Omega_k$ ,

Figure 5.1: Local subdivision of an element  $\Omega_k$  into subtriangles.

$\mathbf{g}^*$  and  $\boldsymbol{\lambda}$  are linear tractions imposed on the edges of the element and  $\mathbf{z}_H$  is a linear displacement field in  $\Omega_k$ , then, it is sufficient to consider the statically admissible stress field  $\hat{\boldsymbol{\sigma}}^k$  to be a linear polynomial in each subtriangle conforming the element. The reader is referred to (Parés, Bonet, Huerta and Peraire 2005, Appendix A) for a detailed construction of the statically admissible stress field.

The procedure to obtain strict upper bounds for the energy norm of  $\mathbf{z}$  is summarized in the box in Table 5.2.

|  |
|--|
| <p>1.- Compute <math>\boldsymbol{\lambda} \in \boldsymbol{\Lambda}</math> s.t.</p> $b(\hat{\mathbf{v}}, \boldsymbol{\lambda}) = R^*(\hat{\mathbf{v}}) \quad \forall \hat{\mathbf{v}} \in \hat{\mathcal{V}}^H.$ <p>2.- For each element <math>\Omega_k</math>, compute <math>\hat{\boldsymbol{\sigma}}^k \in \mathcal{V}_k</math> s.t.</p> $\int_{\Omega_k} \hat{\boldsymbol{\sigma}}^k : \boldsymbol{\varepsilon}(\mathbf{v}) \, d\Omega = R_k^*(\mathbf{v}) + b_k(\mathbf{v}, \boldsymbol{\lambda}) \quad \forall \mathbf{v} \in \mathcal{V}_k.$ <p>3.- Compute the upper bound as</p> $\ \mathbf{z}\ ^2 \leq \sum_{k=1}^{n_{el}} \ \hat{\boldsymbol{\sigma}}^k\ _k^2.$ |
|--|

Figure 5.2: Main steps of the strategy used to obtain strict upper bounds for the energy norm of the solution of a symmetric boundary value problem.

Let  $\boldsymbol{\mathcal{S}}$  denote the space of componentwise square-integrable stress fields in  $\Omega$ , that is,  $\boldsymbol{\mathcal{S}}$  contains all the second-order tensors with  $\sigma_{ij} \in \mathcal{L}^2(\Omega_k) \quad \forall i, j$ , with the

associated complementary energy

$$|||\boldsymbol{\sigma}|||^2 = \int_{\Omega} \boldsymbol{\sigma} : \mathbb{C}^{-1} : \boldsymbol{\sigma} \, d\Omega.$$

The local admissible stress fields  $\hat{\boldsymbol{\sigma}}_k \in \mathcal{S}_k$  are not defined in the whole domain but only in the element  $\Omega_k$ . However they can be naturally extended to  $\Omega$  setting the values outside  $\Omega_k$  to zero. Then, the upper bound for squared the energy norm of  $\mathbf{z}$  given by

$$\|\mathbf{z}\|^2 \leq \sum_{k=1}^{n_{e1}} |||\hat{\boldsymbol{\sigma}}_k|||_k^2,$$

may be computed also from the global stress field  $\hat{\boldsymbol{\sigma}} = \sum_{k=1}^{n_{e1}} \hat{\boldsymbol{\sigma}}_k \in \mathcal{S}$  as

$$\|\mathbf{z}\| \leq |||\hat{\boldsymbol{\sigma}}|||.$$

It is worth noting that the global stress field  $\hat{\boldsymbol{\sigma}}$  belongs to the space of globally admissible stress fields

$$\mathcal{S}^{eq} = \left\{ \boldsymbol{\sigma} \in \mathcal{S}, \int_{\Omega} \boldsymbol{\sigma} : \boldsymbol{\varepsilon}(\mathbf{v}) \, d\Omega = R^*(\mathbf{v}) \quad \forall \mathbf{v} \in \mathcal{V} \right\}.$$

### 5.2.1 Sufficient condition for the upper bound property

The following result summarize a sufficient condition for a global stress field to yield an upper bound for the error measured in the energy norm. In fact, the theorem states that every stress field globally admissible provides an upper bound for  $\|\mathbf{z}\|$ .

**Lemma 5.2.1.** *Any stress field  $\hat{\boldsymbol{\sigma}} \in \mathcal{S}^{eq}$ , that is, any stress field  $\hat{\boldsymbol{\sigma}} \in \mathcal{S}$  verifying the weak error equation*

$$\int_{\Omega} \hat{\boldsymbol{\sigma}} : \boldsymbol{\varepsilon}(\mathbf{v}) \, d\Omega = a(\mathbf{z}, \mathbf{v}) = R^*(\mathbf{v}) \quad \forall \mathbf{v} \in \mathcal{V}, \quad (5.5)$$

*is such that its complementary norm is an upper bound of the energy norm of  $\mathbf{z}$ , that is*

$$\|\mathbf{z}\| \leq |||\hat{\boldsymbol{\sigma}}|||.$$

*Proof.* Using equation (5.5) with  $\mathbf{v} = \mathbf{z}$ , the energy norm of  $\mathbf{z}$  may be rewritten as

$$\|\mathbf{z}\|^2 = a(\mathbf{z}, \mathbf{z}) = R^*(\mathbf{z}) = \int_{\Omega} \hat{\boldsymbol{\sigma}} : \boldsymbol{\varepsilon}(\mathbf{z}) \, d\Omega.$$



Then, a simple algebraic manipulation yields to

$$\begin{aligned} 0 &\leq \int_{\Omega} (\hat{\boldsymbol{\sigma}} - \boldsymbol{\sigma}(\mathbf{z})) : \mathbb{C}^{-1} : (\hat{\boldsymbol{\sigma}} - \boldsymbol{\sigma}(\mathbf{z})) \, d\Omega \\ &= \|\hat{\boldsymbol{\sigma}}\|^2 + \|\mathbf{z}\|^2 - 2 \int_{\Omega} \hat{\boldsymbol{\sigma}} : \boldsymbol{\varepsilon}(\mathbf{z}) \, d\Omega = \|\hat{\boldsymbol{\sigma}}\|^2 - \|\mathbf{z}\|^2, \end{aligned}$$

and the lemma is proved.  $\square$

### 5.3 Certification

An attractive feature of the proposed approach is that the piecewise polynomial equilibrated stress-like fields, which are computed as part of the bound process, can be used as certificates to guarantee the correctness of the computed bounds. It turns out that given a stress field it is easy to check whether this field corresponds to a valid certificate, and in the affirmative case it is straightforward to determine the value of the output that it can certify. In particular, the stress fields need to satisfy continuity of normal tractions,  $\boldsymbol{\sigma} \cdot \mathbf{n}$ , across elements, and membership of an appropriate space.

The idea of a certificate that is computed simultaneously with the solution has many attractive features. In particular, a certificate consisting of the data set necessary to describe the piecewise polynomial stress-like fields could be used to document the computed results. Exercising the certificate does not require access to the code used to compute it and can be done with a simple algorithm which does not require solving any system of equations. A very important point is that, if a certificate meets all the necessary conditions, which in turn are easy to verify, then there is no need to certify the code used to compute it. In practice, the size of these certificates depends on the required level of certainty. As expected, we shall find that high levels of certainty, i.e. small bound gaps, will often require longer certificates (larger data sets) than those required to certify less sharp claims.

In the case of computing strict upper bounds for the energy norm of  $\mathbf{z}$ , the certification procedure is summarized in Figure 5.3.

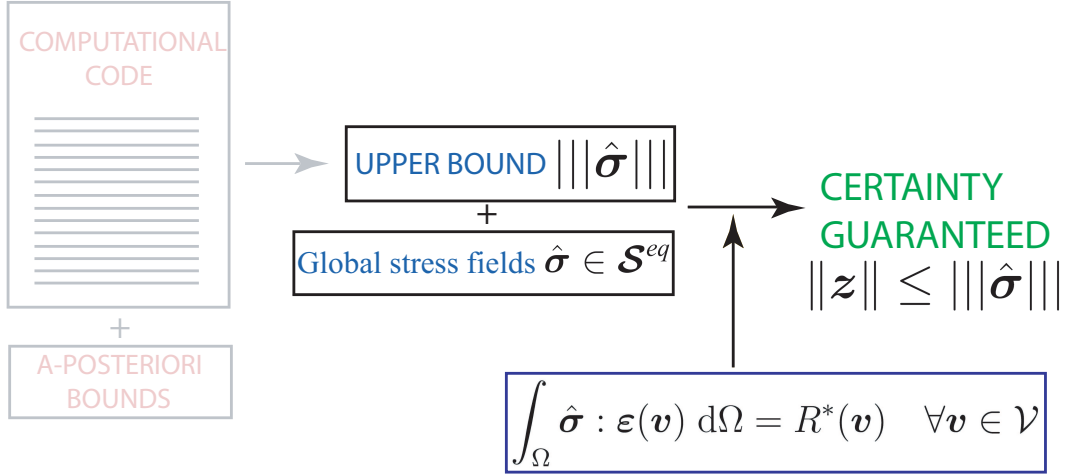


Figure 5.3: Summary of the certification procedure to obtain strict upper bounds for  $||z||$ .

## 5.4 Numerical examples

The presented method is illustrated with three numerical examples: a square plate with two interior rectangular cut-outs, the solution of which, has corner singularities, and two plates with cracks. The outputs of interest are displacements and reaction forces integrated over parts of the boundary in the first example, and the value of the  $J$ -integral at the crack tips in the second and third example.

The coarse mesh problems are solved using triangular linear finite elements, the hybrid fluxes are interpolated linearly over each edge of the mesh and the local equilibrated stress fields are taken to be piecewise linear in each triangle of the mesh. Three estimates of the error in the output  $s$  are considered: the upper and lower bounds ( $s_u$  and  $s_l$ , respectively) and their average,  $\underline{s} = (s_u + s_l)/2$ . This yields to the four estimates of the output  $\ell^{\mathcal{O}}(\mathbf{u})$  itself: the upper and lower bounds ( $\ell^{\mathcal{O}}(\mathbf{u}_H) + s_u$  and  $\ell^{\mathcal{O}}(\mathbf{u}_H) + s_l$ , respectively), their average ( $\ell^{\mathcal{O}}(\mathbf{u}_H) + \underline{s}$ ), and also the output given by the finite element approximation itself,  $\ell^{\mathcal{O}}(\mathbf{u}_H)$ .

In both examples, since the analytical solution is not known, the quality of the different estimates for the output  $\ell^{\mathcal{O}}(\mathbf{u})$  is measured in terms of the relative half bound gap,  $\rho_G$ , which is defined as half of the difference between the upper and

lower bounds for the output, divided by the average estimate

$$\rho_G = \frac{1}{2} \frac{(\ell^{\mathcal{O}}(\mathbf{u}_H) + s_u) - (\ell^{\mathcal{O}}(\mathbf{u}_H) + s_l)}{|\ell^{\mathcal{O}}(\mathbf{u}_H) + \underline{s}|} = \frac{1}{2} \frac{s_u - s_l}{|\ell^{\mathcal{O}}(\mathbf{u}_H) + \underline{s}|} \geq 0.$$

### 5.4.1 Square plate

A square thin plate with two rectangular holes is considered. Normal tractions are applied on the left and right sides of the plate (Paraschivoiu et al. 1997, Peraire and Patera 1997). Since the problem is symmetric, only one fourth of the plate is considered, as shown in Figure 5.4.

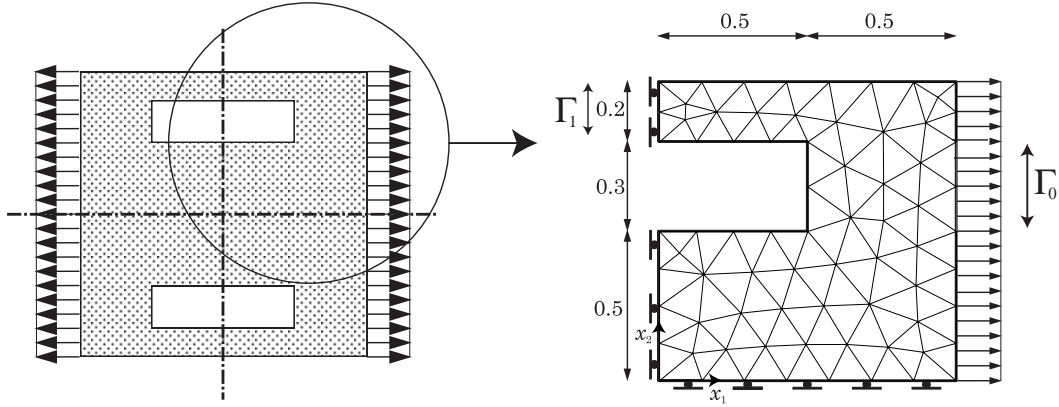


Figure 5.4: Model problem (left) and initial mesh (right)

Two outputs of interest are considered: the integral of the normal displacement over the boundary  $\Gamma_0$ , and the integrated normal component of the traction in  $\Gamma_1$ , that is,

$$\ell_0^{\mathcal{O}}(\mathbf{v}) = \int_{\Gamma_0} \mathbf{v} \cdot \mathbf{n} \, d\Gamma, \quad \ell_1^{\mathcal{O}}(\mathbf{v}) = \int_{\Gamma_1} \mathbf{n} \cdot \boldsymbol{\sigma}(\mathbf{v}) \cdot \mathbf{n} \, d\Gamma. \quad (5.6)$$

*Remark 5.4.1.* The dual residue associated with the first output is already in the form of equation (5.2) with  $\mathbf{g}^* = \mathbf{n}|_{\Gamma_0}$  and  $\mathbf{g}^* = 0$  elsewhere,  $\mathbf{f}^* = 0$  in  $\Omega$  and  $\mathbf{z}_H = \boldsymbol{\psi}_H$ , where  $\boldsymbol{\psi}_H$  is the finite element approximation of the dual problem (2.10) associated with the output  $\ell_0^{\mathcal{O}}(\cdot)$ , that is

$$R^{\mathcal{D}}(\mathbf{v}) = \int_{\Gamma_0} \mathbf{v} \cdot \mathbf{n} \, d\Gamma - a(\boldsymbol{\psi}_H, \mathbf{v}).$$

The residue associated with the second output, however, does not have the same form. In order to transform this residue into the form (5.2) considered here, an auxiliary function  $\chi$  is introduced. This function  $\chi$ , is such that  $\chi = 1$  on  $\Gamma_1$  and vanishes at all the other vertical boundaries. Then, if  $\mathbf{n}^1 = \mathbf{n}|_{\Gamma_1}$ ,

$$\ell_1^{\mathcal{O}}(\mathbf{u}) = \int_{\Gamma_1} \mathbf{n} \cdot \boldsymbol{\sigma}(\mathbf{u}) \cdot \mathbf{n} \, d\Gamma = a(\mathbf{u}, \chi \mathbf{n}^1) =: \tilde{\ell}_1^{\mathcal{O}}(\mathbf{u}),$$

and instead of working with the functional  $\ell_1^{\mathcal{O}}(\cdot)$ ,  $\tilde{\ell}_1^{\mathcal{O}}(\cdot)$  is considered. This is much easier since it corresponds to  $\mathbf{g}^* = 0$  on  $\Gamma^N$ ,  $\mathbf{f}^* = 0$  in  $\Omega$  and  $\mathbf{z}_H = \boldsymbol{\psi}_H - \chi \mathbf{n}^1$  in equation (5.2) where now  $\boldsymbol{\psi}_H$  denotes the finite element approximation of the dual problem associated with the output  $\tilde{\ell}_1^{\mathcal{O}}(\cdot)$ . That is,

$$R^D(\mathbf{v}) = \tilde{\ell}_1^{\mathcal{O}}(\mathbf{u}) - a(\mathbf{v}, \boldsymbol{\psi}_H) = a(\mathbf{v}, \boldsymbol{\psi}_H - \chi \mathbf{n}^1) = a(\boldsymbol{\psi}_H - \chi \mathbf{n}^1, \mathbf{v}),$$

due to the symmetry of the bilinear form  $a(\cdot, \cdot)$ .

Figure 5.5 and Tables 5.1 and 5.2 show the bounds obtained in this example. A nested sequence of meshes is considered. The initial mesh ( $h_{\text{ini}}$ ) is shown in

| $h$                 | displacement average               |  |  |  |          |
|---------------------|------------------------------------|--|--|--|----------|
|                     | $\ell^{\mathcal{O}}(\mathbf{u}_H)$ | $\ell^{\mathcal{O}}(\mathbf{u}_H) + s^-$ | $\ell^{\mathcal{O}}(\mathbf{u}_H) + s^+$ | $\ell^{\mathcal{O}}(\mathbf{u}_H) + \underline{s}$ | $\rho_G$ |
| $h_{\text{ini}}$    | .4060                              | .3794                                    | .5297                                    | .4546  | .1654    |
| $1/2h_{\text{ini}}$ | .4163                              | .4061                                    | .4706                                    | .4384  | .0736    |
| $1/4h_{\text{ini}}$ | .4207                              | .4172                                    | .4423                                    | .4298  | .0292    |
| $1/8h_{\text{ini}}$ | .4224                              | .4213                                    | .4309                                    | .4261  | .0113    |

Table 5.1: Bounds and relative bound gap in a series of uniformly refined  $h$ -meshes for  $\ell_0^{\mathcal{O}}(\mathbf{u})$

| $h$                 | reaction average                   |  |  |  |          |
|---------------------|------------------------------------|--|--|--|----------|
|                     | $\ell^{\mathcal{O}}(\mathbf{u}_H)$ | $\ell^{\mathcal{O}}(\mathbf{u}_H) + s^-$ | $\ell^{\mathcal{O}}(\mathbf{u}_H) + s^+$ | $\ell^{\mathcal{O}}(\mathbf{u}_H) + \underline{s}$ | $\rho_G$ |
| $h_{\text{ini}}$    | -.3199                             | -.3696                                   | -.2621                                   | -.3158   | .1702    |
| $1/2h_{\text{ini}}$ | -.3203                             | -.3438                                   | -.2982                                   | -.3210   | .0710    |
| $1/4h_{\text{ini}}$ | -.3211                             | -.3318                                   | -.3133                                   | -.3225   | .0286    |
| $1/8h_{\text{ini}}$ | -.3217                             | -.3265                                   | -.3189                                   | -.3227   | .0118    |

Table 5.2: Bounds and relative bound gap in a series of uniformly refined  $h$ -meshes for  $\ell_1^{\mathcal{O}}(\mathbf{u})$

Figure 5.4, and the refined meshes are obtained, as in the first example, dividing

each element into 4 new ones. The function  $\chi$  required in  $\tilde{\ell}_1^{\mathcal{O}}(\cdot)$ , is defined on the initial mesh by setting all the nodal values equal to zero except for those nodes on  $\Gamma_1$  which are given a value of unity.

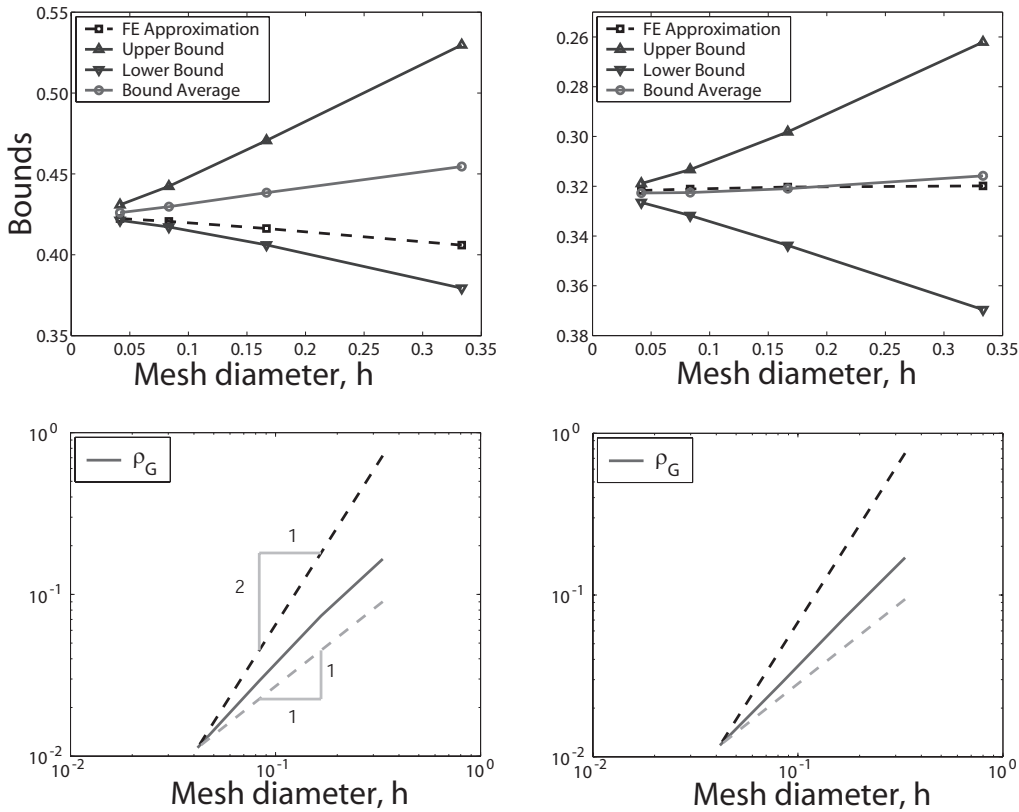


Figure 5.5: Bounds convergence for a uniform  $h$ -refinement (up) and for the displacement output  $\ell_0^{\mathcal{O}}(\mathbf{u})$  (left) and for the reaction output  $\ell_1^{\mathcal{O}}(\mathbf{u})$  (right)

This example shows that the bounds behave well even for problems with singularities. However, it is also observed that the convergence rate for the bounds, the finite element approximation and the bound average, is no longer  $\mathcal{O}(h^2)$ , although it is still faster than linear.

For the reaction output,  $\ell_1^{\mathcal{O}}(\mathbf{u})$ , an adaptive procedure has been employed starting with the mesh shown in Figure 5.4 where the bound gap  $\Delta_{\text{ini}}$  is 0.1075, and two target bound gaps have been considered  $\Delta_{\text{tol}} = \frac{1}{2}\Delta_{\text{ini}}$  and  $\Delta_{\text{tol}} = \frac{1}{10}\Delta_{\text{ini}}$ .

In order to achieve  $\Delta_{\text{tol}} = \frac{1}{2}\Delta_{\text{ini}}$  four new meshes are generated, where the bound gap for the last mesh is  $\Delta_f = 0.0471$ . The resulting sequence of meshes

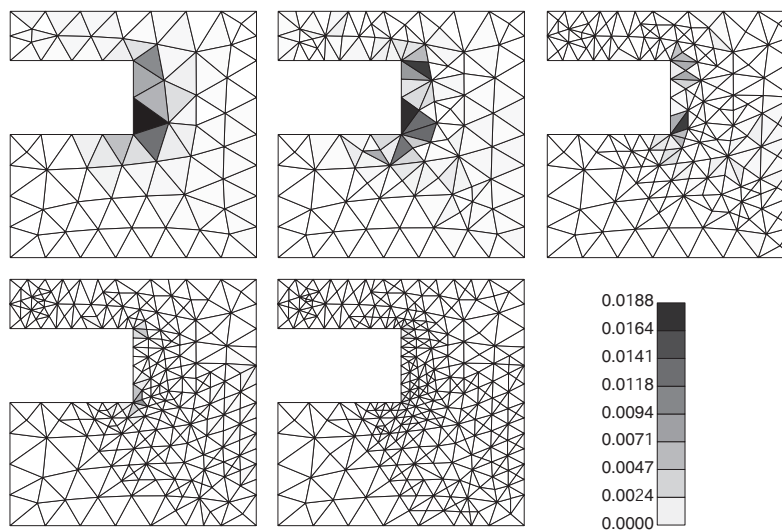


Figure 5.6: Sequence of adapted meshes for the output  $\ell_1^O(\mathbf{u})$  with desired final gap  $\Delta_{\text{tol}} = \frac{1}{2}\Delta_{\text{ini}}$  with  $n_{\text{el}} = 108, 165, 280, 405$  and  $538$

can be seen in Figure 5.6, where the local elementary contributions to the global bound gap are plotted in each element of the mesh. As can be seen not only the zone where the output is measured ( $\Gamma_1$ ) is refined, but also the corners where the solution is singular.

The values of the bounds for the adaptive procedure with the desired final gap  $\Delta_{\text{tol}} = \frac{1}{10}\Delta_{\text{ini}}$  are shown in Table 5.3.

| $n_{\text{el}}$ | $\Delta$ | $s_l$   | $s_u$   |
|-----------------|----------|---------|---------|
| 108             | .10749   | -.36957 | -.26208 |
| 222             | .18215   | -.38940 | -.20725 |
| 433             | .12171   | -.36880 | -.24709 |
| 811             | .07199   | -.35089 | -.27891 |
| 1387            | .03755   | -.33750 | -.29995 |
| 1966            | .02428   | -.33392 | -.30964 |
| 2532            | .01574   | -.32922 | -.31348 |
| 3069            | .01172   | -.32826 | -.31654 |
| 3564            | .00834   | -.32627 | -.31793 |

Table 5.3: Bounds in a series of adaptively  $h$ -refined meshes both for  $\ell_1^O(\mathbf{u})$  with desired final gap  $\Delta_{\text{tol}} = \frac{1}{10}\Delta_{\text{ini}}$

### 5.4.2 $J$ -integral

Here two plates with cracks are considered: the first plate with two edge cracks and the second plate with an inclined crack both subjected to a uniformly distributed tensile stress as shown in Figure 5.7. Both plates are assumed to be in plane strain and the value of the tensile force acting on the two ends of the plates is  $p = 1$ . The non-dimensionalized Young's modulus is 1.0 and the Poisson's ratio is 0.3. In the first example, due to the symmetry of the problem only one quarter of the plate is considered for the finite element analysis.

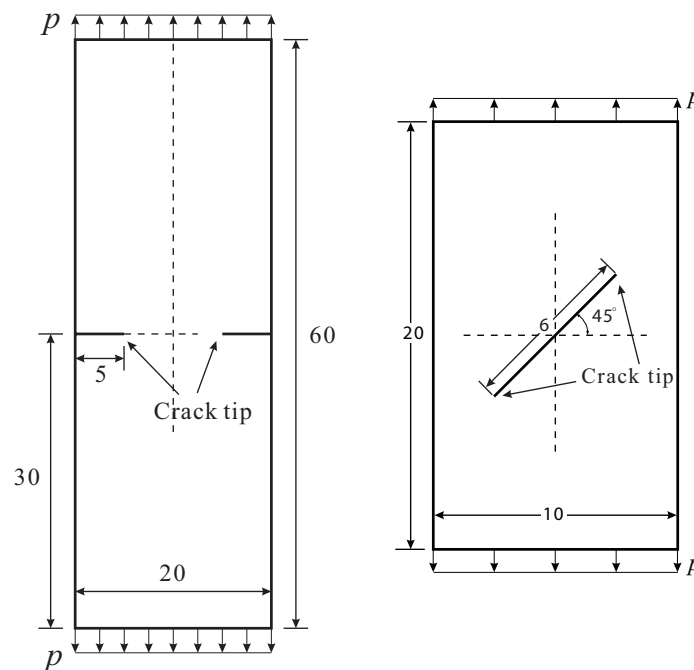


Figure 5.7:  $J$ -integral: double edge-cracked plate subjected to a uniform tensile stress (left) and plate with an inclined crack subjected to a uniform tensile stress (right)

In both examples the output of interest is the value of the  $J$ -integral which provides the energy release of the cracks. The  $J$ -integral is not a linear functional of the displacements, and thus the strategies presented in Chapter 2 do not directly provide bounds for this output. Xuan, Lee, Patera and Peraire (2004) present a method for computing upper and lower bounds for the value of the  $J$ -integral in two

dimensional linear fracture mechanics using a posteriori error estimates measuring the energy norm of a reference solution. Here, the same technique has been applied using the energy norm estimate described in this chapter. This allows to find strict bounds for the  $J$ -integral instead of bounds which are only strict with respect to a reference solution.

The obtention of the bounds proposed by Xuan et al. (2004) reformulates the  $J$ -integral as a bounded quadratic functional of the displacement and expands this quadratic functional into computable quantities plus additional linear and quadratic terms in the error. The linear terms are bounded using the strategy presented in this chapter along with the methodology introduced in Chapter 2 and the quadratic term is bounded with the energy norm of the error scaled by a suitable chosen continuity constant, which can be determined a priori. A detailed deduction of the bounds may be found in (Xuan et al. 2004) and in Appendix C.

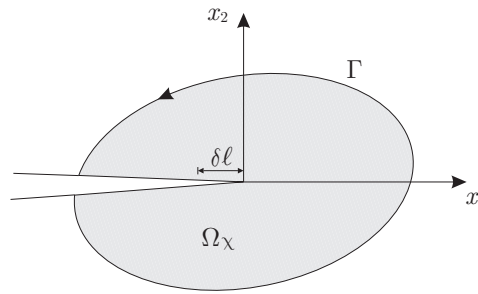


Figure 5.8: Crack geometry showing coordinate axes and the  $J$ -integral contour and domain of integration.

If one considers the geometry shown in Figure 5.8, the  $J$ -integral of the displacement field  $\mathbf{u}$  can be computed as

$$J(\mathbf{u}) = \int_{\Omega_\chi} \left( (\nabla\chi)^T \cdot \boldsymbol{\sigma}(\mathbf{u}) \frac{\partial \mathbf{u}}{\partial x_1} - \frac{\boldsymbol{\sigma}(\mathbf{u}) : \boldsymbol{\varepsilon}(\mathbf{u})}{2} \frac{\partial \chi}{\partial x_1} \right) d\Omega,$$

where the weighting function  $\chi$  is any function in  $H^1(\Omega_\chi)$  that is equal to one at the crack tip and vanishes on  $\Gamma$ . Here  $\Gamma$  denotes any path beginning at the bottom crack face and ending at the top crack face. Note that  $J(\mathbf{u})$  is a bounded quadratic functional of  $\mathbf{u}$ .



In the plate with two edge cracks, a 5 by 5 square area centered on the crack tip is taken as the support,  $\Omega_\chi$ , of the weighting function  $\chi$ , whereas in the plate with an inclined crack the support  $\Omega_\chi$  is a 3 by 3 square area centered on the crack tip (see Figure 5.9).

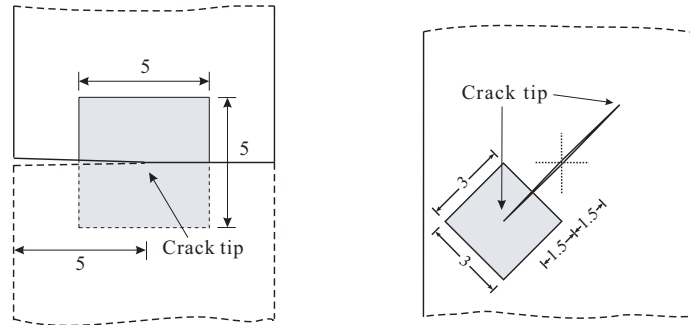


Figure 5.9: Support of weighting function  $\chi$  for the evaluation of the  $J$ -integral for the first example (left) and for the second example (right).

Four estimates of the  $J$ -integral are considered: the upper and lower bounds ( $J^+$  and  $J^-$ , respectively), their average  $J^{\text{ave}} = (J^+ + J^-)/2$ , and also the output given by the finite element approximation, denoted by  $J_H = J(\mathbf{u}_H)$ .

In the first example (plate with two edge cracks) an adaptive procedure has been used to reach a relative bound gap  $(J^+ - J^-)/(2J^{\text{ave}})$  of 5% and 2%. Table 5.4 shows the results for the output  $J_H$ , the computed upper and lower bounds,  $J^\pm$ , for  $J$ , and the relative bound gap for some of the steps of the adaptive procedure.

|                                     |          |         |         |         |         |         |         |
|-------------------------------------|----------|---------|---------|---------|---------|---------|---------|
| $n_{el}$                            | 416      | 525     | 759     | 1368    | 2962    | 10622   | 43733   |
| $J_H$                               | 17.4156  | 18.5208 | 19.1307 | 19.3498 | 19.4601 | 19.5196 | 19.5369 |
| $J^-$                               | -27.7619 | -2.7875 | 10.1176 | 15.1668 | 17.3981 | 18.6596 | 19.1712 |
| $J^+$                               | 86.0779  | 49.2769 | 31.5315 | 24.9273 | 22.0868 | 20.5178 | 19.9343 |
| $J^{\text{ave}}$                    | 29.158   | 23.245  | 20.825  | 20.047  | 19.742  | 19.589  | 19.553  |
| $\frac{J^+ - J^-}{2J^{\text{ave}}}$ | 1.952    | 1.110   | 0.514   | 0.243   | 0.119   | 0.047   | 0.0195  |

Table 5.4: Bound results for the plate with two edge cracks

Also the first three meshes of the adaptive procedure and the final mesh for the 5% relative bound gap are shown in Figure 5.10. It is worth noting that due to the

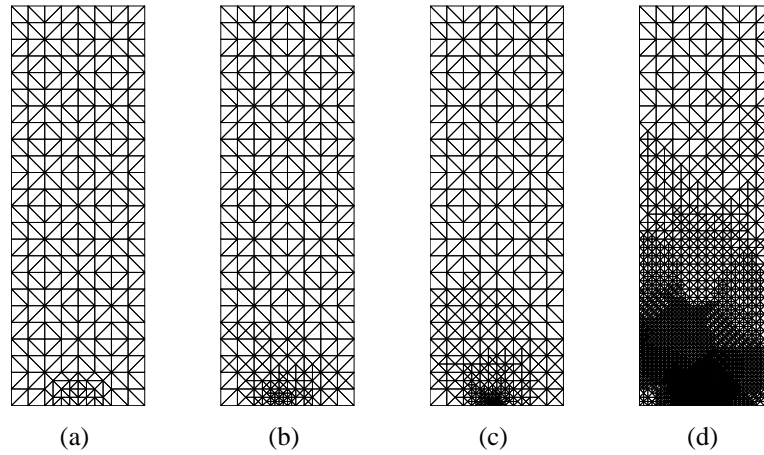


Figure 5.10: Finite element meshes: (a) coarse mesh  $n_{el} = 416$ , (b)  $n_{el} = 525$ , (c)  $n_{el} = 759$  and (d) final mesh for a relative bound gap of 5%,  $n_{el} = 10622$ .

slow convergence of the finite element solution for the problem at hand it is crucial to use adaptive strategies to yield accurate bounds for the output of interest  $J(\mathbf{u})$ .

As in the example of the plate with two edge cracks, an adaptive procedure has been used to reach the desired relative bound gap  $(J^+ - J^-)/(2J^{\text{ave}})$  for the plate with an inclined crack. Table 5.5 shows the results for the output  $J_H$ , the computed upper and lower bounds,  $J^\pm$ , for  $J(\mathbf{u})$ , and the relative bound gap for some of the steps of the adaptive procedure. Also the first four meshes of the adaptive procedure and the final mesh are shown in Figure 5.11. As in the previous example due to the slow convergence of the finite element solution for the problem at hand it is crucial to use adaptive strategies to yield accurate bounds for the output of interest  $J(\mathbf{u})$ .

|                                     |         |         |        |        |       |       |       |       |
|-------------------------------------|---------|---------|--------|--------|-------|-------|-------|-------|
| $n_{el}$                            | 164     | 302     | 632    | 1291   | 3231  | 8534  | 20217 | 41139 |
| $J_H$                               | 4.601   | 5.528   | 6.043  | 6.261  | 6.405 | 6.469 | 6.492 | 6.501 |
| $J^-$                               | -32.079 | -13.879 | -4.130 | 0.961  | 3.958 | 5.273 | 5.829 | 6.079 |
| $J^+$                               | 57.319  | 32.443  | 19.627 | 13.281 | 9.577 | 7.944 | 7.273 | 6.981 |
| $J^{\text{ave}}$                    | 12.620  | 9.282   | 7.746  | 7.121  | 6.766 | 6.609 | 6.551 | 6.530 |
| $\frac{J^+ - J^-}{2J^{\text{ave}}}$ | 3.542   | 2.495   | 1.533  | 0.865  | 0.415 | 0.202 | 0.110 | 0.069 |

Table 5.5: Bound results for the plate with an inclined cracks

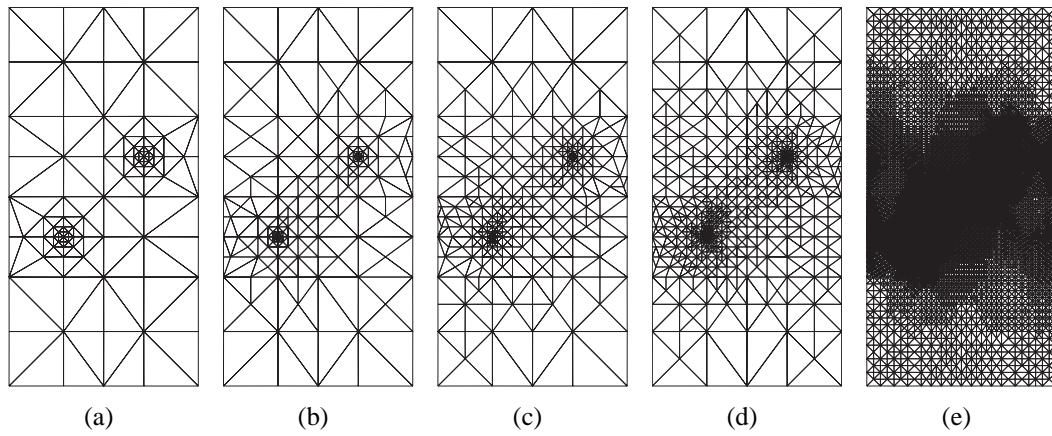


Figure 5.11: Finite element meshes: (a) coarse mesh  $n_{el} = 164$ , (b)  $n_{el} = 302$ , (c)  $n_{el} = 632$ , (d)  $n_{el} = 1291$  and (e) final mesh  $n_{el} = 41139$ .

The convergence of the obtained bounds for the output  $J(\mathbf{u})$  are also illustrated in Figure 5.12 where the behavior of the upper and lower bounds  $J^\pm$  and the bounds average  $J^{ave} = (J^+ + J^-)/2$  is shown.

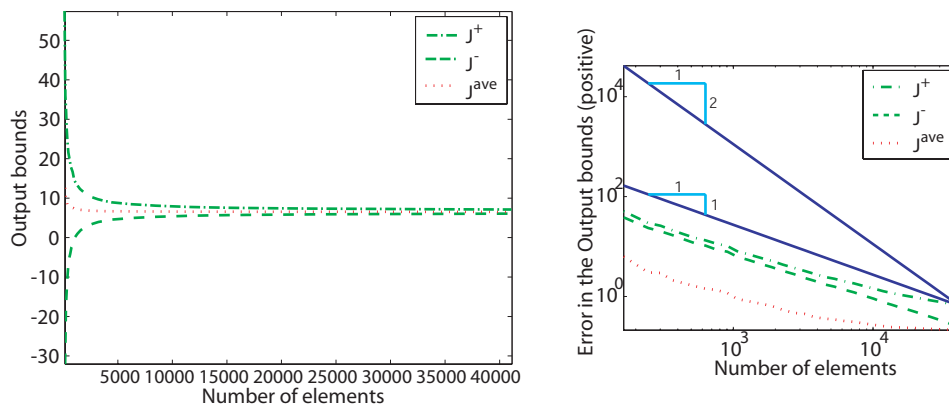


Figure 5.12: Convergence of the upper and lower bounds  $J^\pm$  and the bound average  $J^{ave}$ .



# Chapter 6

## Conclusions

The main contributions of this thesis are described in four items, all addressed to obtaining accurate bounds for outputs of interest.

The first contribution is a general framework to obtain upper and lower bounds for linear functional outputs of interest for both self-adjoint and nonself-adjoint operators. In the context of self-adjoint operators, the presentation shows a unified description of the two main existing strategies to obtain bounds for outputs from energy norm estimates (the “parallelogram identity based” and the “optimization-Lagrangian based”). For nonself-adjoint operators an enhancement of the existing bounds is proposed. Improving the effectivity in this case is necessary because existing error estimation strategies fail on obtaining sharp estimates if the skew-symmetric part is relevant.

Second, a simple postprocessing strategy has been presented to recover lower bounds for the energy norm from standard residual estimates producing upper bounds for the energy norm. The main idea is to smooth out the discontinuities of the upper bound estimate and obtain a continuous approximation to the error. This continuous approximation allows to recover a lower bound for the energy. For the pure diffusion problem (without reaction term), the discontinuous estimates yielding upper bounds for the energy are determined up to a local (element by element) constant. Although these constants do not affect the value of the upper bound, the choice of the constants results in very different smoothed continuous approximations, and therefore in lower bounds with very different accuracies. The presented strategy shows how to choose the local constants in order to maximize the lower bounds. Nu-

merical experiments demonstrate that the proposed strategy furnishes sharp lower estimates, of better quality than the original upper ones. The presented strategy may be used in the framework of error estimation for outputs of interest.

Third, a new subdomain-based flux-free error estimation technique is introduced. The implementation is much simpler compared to hybrid-flux estimation techniques because there is no need of using a flux equilibration algorithm. Moreover, the accuracy of the results (sharpness of the upper bounds) is drastically improved compared to other flux-free estimation techniques. In fact, it is at least comparable to hybrid-flux techniques. The resulting estimates yield guaranteed (and sharp) upper bounds of the energy norm of the reference error. A simple and painless postprocessing yields lower bounds of the energy norm of the error with a little extra computational cost. The distribution of the local contributions to the error are also accurately estimated, both for the energy norm of the error and for the error measured using some functional output. These estimates are therefore well suited to guide goal-oriented adaptive procedures.

Finally, a method to compute bounds for linear-functional outputs of weak solutions to the linear elasticity equations is presented. A distinctive feature of this method is that the computed bounds are strict with respect to the output of the exact solution, not with respect to some reference *truth* mesh. We believe this feature is of clear interest in real engineering practice. Numerical experiments show that the computed bounds are sharp and exhibit the proper converge. The method has been presented for the two dimensional elasticity equations, but the extension to three dimensions should not present additional difficulties. The major computational cost, additional to a standard finite element solution, is the solution of an adjoint problem for each output considered. The rest of the operations are local and result in a small computational overhead. The assumptions used in the presented approach restrict the applicability to problems with piecewise polynomial *forcing functions* (body forces, source terms ...) and with polygonal domains. Future work will focus on relaxing these constraints.

## 6.1 Future developments

The topics and methodologies analyzed in this thesis leave open research lines that are worthy to be studied in the next future.

First, future research should investigate extending the method to obtain strict bounds for outputs in elasticity (Parés, Bonet, Huerta and Peraire 2005) to non-polynomial forcing and non-polygonal domains, since many applications contain non-polynomial forcing and curved boundaries.

Also the new “flux-free” subdomain-residual error estimation technique (Parés, Díez and Huerta 2005) must be further investigated as an approach precluding the computation of equilibrated tractions. The presented method only provides bounds with respect to an enriched mesh and it would be interesting to obtain exact bounds with this methodology, independently of the underlying discretization of the domain.

Future research is needed to improve the effectiveness of the bounds obtained when dealing with nonsymmetric operators (such as the convection-diffusion-reaction equation). The enhancement of the bounds presented in Chapter 2 has been tested for a 1D convection-diffusion model problem yielding to much sharper bounds. Although for 1D model problems the upper bound estimation techniques provide the exact value of the energy norm of the errors, the previous bounds (non-enhanced) were not capable to capture the value of the output, resulting in a sequence of upper and lower bounds for the output degenerating as the convection parameter increase. The new approach allows to practically capture the exact value of the output. However, the particularity of the 1D model problem simplified the problem and the 2D implementation requires adopting new strategies.

Also, the derivation of strict bounds should be extended to different classes of problems. For instance, the Stokes equation and transient problems are worthy to be considered.

Finally, taken into account nonlinear functional outputs is also a promising research line.





# Appendix A

## Bounds for linear outputs of interest for a general variational problem

Chapter 2 presents a general framework to obtain bounds for linear functional output of solutions of coercive model problems. This appendix provides a detailed deduction of these bounds.

The method is based on obtaining an exact representation of  $s = \ell^{\mathcal{O}}(e)$ , different from  $s = a(e, \varepsilon)$ , see equation (2.12), which allows to easily recover bounds for  $s$ . This representation yields a systematic approach to calculate bounds for outputs of interest which are characterized by bounded linear functionals. Moreover, whenever  $a(\cdot, \cdot)$  is symmetric, it coincides with the alternative expression for  $s$  recovered from the parallelogram identity, equation (2.13).

This appendix is structured as follows: first, the output  $s$  is regarded as the solution of an optimization problem allowing to deduce bounds for the output. Then, the derivation of bounds for symmetric model problem is revised and finally, bounds for the nonsymmetric model problem are deduced in an analogous way.

### A.1 Energy reformulation

This section aims at finding an exact representation of the error in the output,  $s$ , allowing to compute bounds for  $s$  from available techniques for estimating the energy norm of the error in the finite element approximation of symmetric model problems.

The derivation of the alternative expression for  $s$  involves the following steps:

first, the output is regarded as the solution of a constrained minimization problem. Second, the constrained minimization is rewritten by means of a Lagrangian leading to an unconstrained minimization problem, and finally, the optimal Lagrange multipliers lead to the exact representation of  $s$  in terms of an unconstrained minimization problem.

### A.1.1 Minimization reformulation

Consider the energy-like functional  $\Psi : \mathcal{V} \rightarrow \mathbb{R}$  defined by

$$\Psi(v) = a^s(v, v) - R^P(v).$$

This functional has two essential properties. First, it is coercive on  $\mathcal{V}$ . Indeed, since  $a(\cdot, \cdot)$  is a coercive bilinear form,

$$a^s(v, v) = a(v, v) \geq c|v|^2 \quad \forall v \in \mathcal{V},$$

where  $|\cdot|$  is the norm of the Hilbert space  $\mathcal{V}$ . Second, it reduces to zero when  $v = e$  using equation (2.3) and the fact that  $a^s(e, e) = a(e, e)$ .

The energy-like functional  $\Psi(\cdot)$  allows to write the output as the minimum of the constrained minimization problem

$$\begin{aligned} \pm s &= \min_{v \in \mathcal{V}} \pm \ell^O(v) + \kappa^2 \Psi(v) \\ \text{s.t. } &a(v, \varphi) = R^P(\varphi) \quad \forall \varphi \in \mathcal{V}, \end{aligned}$$

for any arbitrary scalar parameter  $\kappa \in \mathbb{R}$ . Indeed, the constraint forces the solution to be  $e$  and the objective function for  $v = e$  is  $\pm s$  since  $\Psi(e) = 0$ . The signs  $\pm$  have been introduced in order to be able to recover both the upper and lower bounds at the same time (the  $+$  sign will lead to the lower bound,  $s_l$ , whereas the  $-$  sign will allow us to recover the upper bound,  $s_u$ ).

Now, introducing the quadratic-linear Lagrangian  $L^\pm(v, \varphi)$  given by

$$L^\pm(v, \varphi) = \pm \ell^O(v) + \kappa^2 \Psi(v) + R^P(\varphi) - a(v, \varphi),$$

the previous problem is equivalent to the unconstrained minimization

$$\pm s = \min_{v \in \mathcal{V}} \max_{\varphi \in \mathcal{V}} L^\pm(v, \varphi). \tag{A.1}$$

### A.1.2 Lagrange multiplier

The saddle point  $(\bar{v}, \bar{\varphi}^\pm)$  of the minimization problem (A.1) is found imposing that the variations of  $L^\pm(v, \varphi)$  with respect to  $v \in \mathcal{V}$  and  $\varphi \in \mathcal{V}$  must vanish, leading to  $\bar{v} = e$  and to the following weak problem for  $\bar{\varphi}^\pm \in \mathcal{V}$

$$a(v, \bar{\varphi}^\pm) = \pm \ell^\mathcal{O}(v) + \kappa^2(2a^s(e, v) - R^P(v)) \quad \forall v \in \mathcal{V}. \quad (\text{A.2})$$

Indeed, since  $L^\pm(v, \varphi) = \pm \ell^\mathcal{O}(v) + \kappa^2(a^s(v, v) - R^P(v)) + R^P(\varphi) - a(v, \varphi)$ , imposing that the variations with respect to  $\varphi \in \mathcal{V}$  must vanish leads to the condition: find  $\bar{v} \in \mathcal{V}$  such that

$$R^P(\varphi) - a(\bar{v}, \varphi) = 0 \quad \forall \varphi \in \mathcal{V},$$

which coincides with the residual problem for the primal error  $e$ , equation (2.3), and thus  $\bar{v} = e$ . Similarly, imposing that the variations with respect to  $v \in \mathcal{V}$  must vanish leads to the condition

$$\pm \ell^\mathcal{O}(v) + \kappa^2(2a^s(\bar{v}, v) - R^P(v)) - a(v, \bar{\varphi}^\pm) = 0 \quad \forall v \in \mathcal{V},$$

which yields to the weak problem (A.2) replacing  $\bar{v}$  for  $e$ .

Now, using the definition of the symmetric bilinear form  $a^s(\cdot, \cdot)$ , equation (2.5), and the residual equation (2.3), the equation determining  $\bar{\varphi}^\pm$  can be rewritten as

$$a(v, \bar{\varphi}^\pm) = \pm \ell^\mathcal{O}(v) + \kappa^2(a(e, v) + a(v, e) - R^P(v)) = \pm \ell^\mathcal{O}(v) + \kappa^2 a(v, e) \quad \forall v \in \mathcal{V}.$$

Hence  $\bar{\varphi}^\pm = \pm \psi + \kappa^2 e$ , where  $\psi \in \mathcal{V}$  is the solution of the dual or adjoint problem, see (2.9). Moreover, replacing  $\psi$  by the sum of its finite element approximation  $\psi_H$  and the associated dual error  $\varepsilon$ , that is,  $\psi = \psi_H + \varepsilon$ , the Lagrange multiplier  $\bar{\varphi}^\pm$  can be rewritten as  $\bar{\varphi}^\pm = \pm \psi_H + \kappa(\kappa e \pm \frac{1}{\kappa} \varepsilon)$ .

Finally, substituting the exact Lagrange multiplier  $\bar{\varphi}^\pm$  into equation (A.1) one obtains the exact representation for the output

$$\pm s = \min_{v \in \mathcal{V}} L^\pm(v, \pm \psi_H + \kappa(\kappa e \pm \frac{1}{\kappa} \varepsilon)). \quad (\text{A.3})$$

With the aid of the Galerkin orthogonality of the primal residual with respect to the finite element space  $\mathcal{V}^H$  and regrouping terms, the previous representation of the output  $s$  may be rewritten as

$$\pm s = \min_{v \in \mathcal{V}} \tilde{L}^\pm(v, \kappa e \pm \frac{1}{\kappa} \varepsilon), \quad (\text{A.4})$$

where

$$\tilde{L}^\pm(v, \varphi) = \kappa^2 a^s(v, v) - \kappa R^\mp(v) + \kappa(R^P(\varphi) - a(v, \varphi)),$$

$R^\mp(\cdot)$  being the residue defined in equation (2.21).

Indeed, given  $\varphi \in \mathcal{V}$ ,

$$\begin{aligned} L^\pm(v, \pm\psi_H + \kappa\varphi) &= \pm\ell^\mathcal{O}(v) + \kappa^2\Psi(v) + R^P(\pm\psi_H + \kappa\varphi) - a(v, \pm\psi_H + \kappa\varphi) \\ &= \pm\ell^\mathcal{O}(v) + \kappa^2\Psi(v) + \kappa R^P(\varphi) \mp a(v, \psi_H) - \kappa a(v, \varphi) \\ &= \kappa^2\Psi(v) \pm (\ell^\mathcal{O}(v) - a(v, \psi_H)) + \kappa(R^P(\varphi) - a(v, \varphi)) \\ &= \kappa^2\Psi(v) \pm R^D(v) + \kappa(R^P(\varphi) - a(v, \varphi)) \\ &= \kappa^2 a^s(v, v) - \kappa(kR^P(v) \mp \frac{1}{\kappa}R^D(v)) + \kappa(R^P(\varphi) - a(v, \varphi)) \\ &= \kappa^2 a^s(v, v) - \kappa R^\mp(v) + \kappa(R^P(\varphi) - a(v, \varphi)) = \tilde{L}^\pm(v, \varphi), \end{aligned}$$

from where,

$$\min_{v \in \mathcal{V}} L^\pm(v, \pm\psi_H + \kappa\varphi) = \min_{v \in \mathcal{V}} \tilde{L}^\pm(v, \varphi) \quad \forall \varphi \in \mathcal{V}.$$

In particular  $\varphi = \kappa e \pm \frac{1}{\kappa}\varepsilon$  leads to the equivalence between (A.3) and (A.4).

In fact, the equivalent exact representation of the error holds

$$\pm s = \min_{v \in \mathcal{V}} \max_{\varphi \in \mathcal{V}} \tilde{L}^\pm(v, \varphi), \quad (\text{A.5})$$

where now the saddle point is  $(\bar{v}, \bar{\varphi}^\pm) = (e, \kappa e \pm \frac{1}{\kappa}\varepsilon)$ . This exact representation of the error in the quantity of interest is the starting point of the derivation of the upper and lower bounds for  $s$ .

### A.1.3 Strong duality

From equation (A.5) lower bounds for the output  $\pm s$  may be deduced using strong duality of convex minimization and the saddle point property of the Lagrange multipliers as

$$\pm s = \min_{v \in \mathcal{V}} \max_{\varphi \in \mathcal{V}} \tilde{L}^\pm(v, \varphi) = \max_{\varphi \in \mathcal{V}} \min_{v \in \mathcal{V}} \tilde{L}^\pm(v, \varphi) \geq \min_{v \in \mathcal{V}} \tilde{L}^\pm(v, \varphi) \quad \forall \varphi \in \mathcal{V}. \quad (\text{A.6})$$

These bounds actually hold as an equality for  $\varphi = \bar{\varphi}^\pm = \kappa e \pm \frac{1}{\kappa}\varepsilon$  as shown in equation (A.4), recovering the exact representation for the output  $\pm s$ .

*Remark A.1.1.* If  $\varphi$  is taken to be 0, the resulting bounds are

$$\pm s \geq \min_{v \in \mathcal{V}} \tilde{L}^\pm(v, 0) = \min_{v \in \mathcal{V}} \kappa^2 a^s(v, v) - \kappa R^\mp(v). \quad (\text{A.7})$$

The saddle point of the previous minimization problem is  $\bar{v} \in \mathcal{V}$  verifying the following residual equation

$$a^s(\bar{v}, v) = \frac{1}{2\kappa} R^\mp(v) = \frac{1}{2\kappa} \left( \kappa R^P(v) \mp \frac{1}{\kappa} R^D(v) \right) \quad \forall v \in \mathcal{V}.$$

Therefore  $\bar{v} = \frac{1}{2\kappa} (\kappa e^s \mp \frac{1}{\kappa} \varepsilon^s)$  where  $e^s$  and  $\varepsilon^s$  are the symmetric primal and dual errors defined in equations (2.16) and (2.15) respectively. Moreover the value of the Lagrangian for  $v = \bar{v}$  is

$$\tilde{L}^\pm(\bar{v}, 0) = \kappa^2 a^s(\bar{v}, \bar{v}) - \kappa R^\mp(\bar{v}) = -\kappa^2 a^s(\bar{v}, \bar{v}) = -\frac{1}{4} \left\| \kappa e^s \mp \frac{1}{\kappa} \varepsilon^s \right\|^2.$$

The positive part of the inequality (A.7) yields the lower bound  $s \geq -\frac{1}{4} \left\| \kappa e^s - \frac{1}{\kappa} \varepsilon^s \right\|^2$ , whereas the negative part of the inequality yields  $-s \geq -\frac{1}{4} \left\| \kappa e^s + \frac{1}{\kappa} \varepsilon^s \right\|^2$ . Thus, multiplying this last inequality by minus 1, the final bounds for the output are

$$-\frac{1}{4} \left\| \kappa e^s - \frac{1}{\kappa} \varepsilon^s \right\|^2 \leq s \leq \frac{1}{4} \left\| \kappa e^s + \frac{1}{\kappa} \varepsilon^s \right\|^2. \quad (\text{A.8})$$

The strategy proposed by Paraschivoiu et al. (1997) yields to the previous bounds. The idea to enhance these bounds is simply the introduction of a non-zero approximation  $\varphi$  of the errors  $\kappa e \pm \frac{1}{\kappa} \varepsilon$ . The resulting bounds are at least of the same quality as the original ones (A.8).

## A.2 Bounds for the error in the quantity of interest

This section is devoted to find an analogous of equation (2.14) providing bounds for  $s$ , which are valid regardless of the bilinear form  $a(\cdot, \cdot)$  being symmetric or not.

The bounds for  $s$  are found starting with continuous approximations  $\xi^+$  and  $\xi^- \in \mathcal{V}$  of the errors  $\kappa e + \frac{1}{\kappa} \varepsilon$  and  $\kappa e - \frac{1}{\kappa} \varepsilon$  respectively. The characterization (A.6) of the output yields the bounds for  $s$

$$\min_{v \in \mathcal{V}} \tilde{L}^+(v, \xi^+) \leq s \leq -\min_{v \in \mathcal{V}} \tilde{L}^-(v, \xi^-). \quad (\text{A.9})$$

However, in order to be able to sharpen the bounds, two arbitrary scalar parameters  $\lambda^\pm \in \mathbb{R}$  are introduced as a scaling factor for the approximations  $\xi^\pm$ , that is, instead of the approximations  $\xi^\pm$ ,  $\lambda^\pm \xi^\pm$  are considered. The bounds associated to these new approximations are

$$\min_{v \in \mathcal{V}} \tilde{L}^+(v, \lambda^+ \xi^+) \leq s \leq - \min_{v \in \mathcal{V}} \tilde{L}^-(v, \lambda^- \xi^-). \quad (\text{A.10})$$

Note that these bounds depend on the parameters  $\lambda^\pm$ . Thus, they may be optimized to obtain sharper bounds resulting in

$$s_l = \max_{\lambda^+ \in \mathbb{R}} \min_{v \in \mathcal{V}} \tilde{L}^+(v, \lambda^+ \xi^+) \leq s \leq - \max_{\lambda^- \in \mathbb{R}} \min_{v \in \mathcal{V}} \tilde{L}^-(v, \lambda^- \xi^-) = s_u. \quad (\text{A.11})$$

Moreover, the bounds are optimal (they recover the output  $s$ ) if the continuous approximations  $\xi^\pm$  coincide with the errors  $\kappa e \pm \frac{1}{\kappa} \varepsilon$ .

The rest of the section is devoted to give the explicit expression for the bounds. First the symmetric model problem is considered to illustrate the procedure and to show that the derived bounds are equivalent to the bounds derived from the parallelogram identity (2.14). Then, the procedure is illustrated for a general nonsymmetric model problem.

### A.2.1 Bounds for self-adjoint model problems

In this section the bilinear form  $a(\cdot, \cdot)$  is assumed to be symmetric,  $a^s(v, w) = a(v, w)$ . Let  $\xi^\pm$  be continuous approximations of  $\kappa e \pm \frac{1}{\kappa} \varepsilon$  and consider the bounds given by equation (A.11). The explicit expressions of the upper and lower bounds for  $s$ ,  $s_u$  and  $s_l$  respectively, are found solving the minimization problems appearing in (A.11) with respect to  $v \in \mathcal{V}$  and then, the bounds are optimized with respect to  $\lambda^\pm \in \mathbb{R}$ . The following theorem gives the expression of the optimal upper and lower bounds given the approximations  $\xi^\pm$  of  $\kappa e \pm \frac{1}{\kappa} \varepsilon$ .

**Theorem A.2.1.** *Let  $\xi^+$  and  $\xi^- \in \mathcal{V}$  be two continuous functions. Then, the quantities  $s_l$  and  $s_u$  given by*

$$s_l = \frac{1}{4} \frac{R^+(\xi^+)^2}{\|\xi^+\|^2} - \frac{1}{4} \left\| \kappa e - \frac{1}{\kappa} \varepsilon \right\|^2, \quad s_u = \frac{1}{4} \left\| \kappa e + \frac{1}{\kappa} \varepsilon \right\|^2 - \frac{1}{4} \frac{R^-(\xi^-)^2}{\|\xi^-\|^2},$$

are a lower and an upper bound for the output  $s$  respectively, that is

$$s_l \leq s \leq s_u.$$

Moreover, the choice  $\xi^\pm = \kappa e \pm \frac{1}{\kappa} \varepsilon$ , leads to the optimal bounds

$$s = s_l = s_u = \frac{1}{4} \|\kappa e + \frac{1}{\kappa} \varepsilon\|^2 - \frac{1}{4} \|\kappa e - \frac{1}{\kappa} \varepsilon\|^2.$$

*Proof.* Given  $\xi^\pm \in \mathcal{V}$ , from equation (A.10)

$$\pm s \geq \min_{v \in \mathcal{V}} \tilde{L}^\pm(v, \lambda^\pm \xi^\pm),$$

where  $\lambda^\pm$  are two arbitrary scalar parameters. In particular,  $\pm s \geq \tilde{L}^\pm(\bar{v}^\pm, \lambda^\pm \xi^\pm)$ , where  $\bar{v}^\pm$  are the saddle points of the minimization of the Lagrangian, that is,  $\bar{v}^\pm = \arg \min_{v \in \mathcal{V}} \tilde{L}^\pm(v, \lambda^\pm \xi^\pm)$ .

Imposing that the variations of  $\tilde{L}^\pm(v, \lambda^\pm \xi^\pm)$  with respect to  $v \in \mathcal{V}$  must vanish, the following weak problem for  $\bar{v}^\pm \in \mathcal{V}$  is obtained

$$a^s(\bar{v}^\pm, v) = \frac{1}{2\kappa} R^\mp(v) + \frac{\lambda^\pm}{2\kappa} a(v, \xi^\pm) \quad \forall v \in \mathcal{V}. \quad (\text{A.12})$$

Now, the symmetry of  $a(\cdot, \cdot)$  involves that for any  $v \in \mathcal{V}$

$$R^\mp(v) = \kappa R^P(v) \mp \frac{1}{\kappa} R^D(v) = a(\kappa e \mp \frac{1}{\kappa} \varepsilon, v), \quad a(v, \xi^\pm) = a(\xi^\pm, v),$$

and the residual equation for  $\bar{v}^\pm$ , equation (A.12), transforms into

$$a(\bar{v}^\pm, v) = \frac{1}{2\kappa} \left( a(\kappa e \mp \frac{1}{\kappa} \varepsilon, v) + \lambda^\pm a(\xi^\pm, v) \right) \quad \forall v \in \mathcal{V},$$

leading to  $\bar{v}^\pm = \frac{1}{2\kappa} (\kappa e \mp \frac{1}{\kappa} \varepsilon + \lambda^\pm \xi^\pm)$ . Moreover,  $\tilde{L}^\pm(\bar{v}^\pm, \lambda^\pm \xi^\pm)$  may be rewritten using equation (A.12) with  $v = \bar{v}^\pm$  and rearranging the r.h.s. as

$$\tilde{L}^\pm(\bar{v}^\pm, \lambda^\pm \xi^\pm) = \kappa \lambda^\pm R^P(\xi^\pm) - \kappa^2 a(\bar{v}^\pm, \bar{v}^\pm).$$

Now, replacing  $\bar{v}^\pm$  by  $\frac{1}{2\kappa} (\kappa e \mp \frac{1}{\kappa} \varepsilon + \lambda^\pm \xi^\pm)$  in the previous equation it follows that

$$\begin{aligned} \tilde{L}^\pm(\bar{v}^\pm, \lambda^\pm \xi^\pm) &= \kappa \lambda^\pm R^P(\xi^\pm) - \frac{1}{4} \|\kappa e \mp \frac{1}{\kappa} \varepsilon + \lambda^\pm \xi^\pm\|^2 \\ &= \kappa \lambda^\pm R^P(\xi^\pm) - \frac{1}{4} \|\kappa e \mp \frac{1}{\kappa} \varepsilon\|^2 - \frac{1}{4} (\lambda^\pm)^2 \|\xi^\pm\|^2 - \frac{\lambda^\pm}{2} a(\kappa e \mp \frac{1}{\kappa} \varepsilon, \xi^\pm) \\ &= -\frac{1}{4} \|\kappa e \mp \frac{1}{\kappa} \varepsilon\|^2 - \frac{1}{4} (\lambda^\pm)^2 \|\xi^\pm\|^2 + \frac{\lambda^\pm}{2} \left( \kappa R^P(\xi^\pm) \pm \frac{1}{\kappa} R^D(\xi^\pm) \right) \\ &= -\frac{1}{4} \|\kappa e \mp \frac{1}{\kappa} \varepsilon\|^2 - \frac{1}{4} (\lambda^\pm)^2 \|\xi^\pm\|^2 + \frac{\lambda^\pm}{2} R^\pm(\xi^\pm). \end{aligned}$$

*Optimal parameter determination:* The optimal value for  $\lambda^\pm$ , that is, the value of the parameter which optimizes the bounds is

$$\lambda^\pm = \frac{R^\pm(\xi^\pm)}{\|\xi^\pm\|^2},$$

provided that  $\|\xi^\pm\|$  is nonzero. If  $\|\xi^\pm\|$  vanishes, the obtained bounds are the same as if  $\xi^\pm = 0$  and the selection of the parameter  $\lambda^\pm$  is a moot point. For the optimal selection of the parameter  $\lambda^\pm$ , the bounds are given by

$$\pm s \geq \frac{1}{4} \frac{R^\pm(\xi^\pm)^2}{\|\xi^\pm\|^2} - \frac{1}{4} \|\kappa e \mp \frac{1}{\kappa} \varepsilon\|^2.$$

The positive part of the previous equation leads to

$$s \geq \frac{1}{4} \frac{R^+(\xi^+)^2}{\|\xi^+\|^2} - \frac{1}{4} \|\kappa e - \frac{1}{\kappa} \varepsilon\|^2 = s_1$$

whereas the negative part multiplying the inequality by  $-1$  lead to,

$$s \leq \frac{1}{4} \|\kappa e + \frac{1}{\kappa} \varepsilon\|^2 - \frac{1}{4} \frac{R^-(\xi^-)^2}{\|\xi^-\|^2} = s_u.$$

Thus, the first part of the proof is concluded.

The optimality of the bounds for the choice  $\xi^\pm = \kappa e \pm \frac{1}{\kappa} \varepsilon$  follows directly from equation (A.4) and since for  $\xi^\pm = \kappa e \pm \frac{1}{\kappa} \varepsilon$ ,  $R^\pm(\xi^\pm) = \|\kappa e \pm \frac{1}{\kappa} \varepsilon\|^2 = \|\xi^\pm\|^2$ , both the upper and lower bounds lead to the parallelogram identity

$$s_1 = s_u = \frac{1}{4} \|\kappa e + \frac{1}{\kappa} \varepsilon\|^2 - \frac{1}{4} \|\kappa e - \frac{1}{\kappa} \varepsilon\|^2 = s.$$

□

### Practical computation of the bounds for self-adjoint model problems

Theorem A.2.1 states that in order to find bounds for the output  $s$ , it is sufficient to obtain continuous approximations  $\xi^\pm$  of  $\kappa e \pm \frac{1}{\kappa} \varepsilon$  and determine the energy norm  $\|\kappa e \pm \frac{1}{\kappa} \varepsilon\|$ . However, these bounds are of limited use as they stand in several aspects. First, they require knowledge of the exact error of the primal and dual problems. Second, for arbitrary choices of the functions  $\xi^\pm \in \mathcal{V}$  not taking into account that they must be approximations of  $\kappa e \pm \frac{1}{\kappa} \varepsilon$ , they will produce rather pessimistic



bounds. Fortunately, using standard a posteriori error estimation techniques, both shortcomings can be removed easily.

Standard a posteriori error estimation techniques allow to find estimates  $\hat{e}$  and  $\hat{\varepsilon}$  providing upper bounds for the energy norms  $\|\kappa e \pm \frac{1}{\kappa}\varepsilon\|$  (see Lemma 3.2.1)

$$\|\kappa e \pm \frac{1}{\kappa}\varepsilon\| \leq \|\kappa \hat{e} \pm \frac{1}{\kappa}\hat{\varepsilon}\|$$

and moreover using simple post-processing techniques, continuous approximations of the errors may be obtained from  $\hat{e}$  and  $\hat{\varepsilon}$ , thus providing the approximations  $\xi^\pm$  (see Section 3.3.2). The bounds are then recovered as

$$\frac{1}{4} \frac{R^+(\xi^+)^2}{\|\xi^+\|^2} - \frac{1}{4} \|\kappa \hat{e} - \frac{1}{\kappa}\hat{\varepsilon}\|^2 \leq s \leq \frac{1}{4} \|\kappa \hat{e} + \frac{1}{\kappa}\hat{\varepsilon}\|^2 - \frac{1}{4} \frac{R^-(\xi^-)^2}{\|\xi^-\|^2}. \quad (\text{A.13})$$

*Remark A.2.1.* The proposed bounds for the symmetric problem are equivalent to the bounds given in equation (2.14). The difference between the two characterization of the bounds, is that in equation (2.14) the lower bounds  $\|\kappa e \pm \frac{1}{\kappa}\varepsilon\|_{\text{LB}}^2$  have been replaced by the quantities

$$\frac{R^\pm(\xi^\pm)^2}{\|\xi^\pm\|^2},$$

using the dual characterization of the energy norm. Indeed, the energy norm of  $\|\kappa e \pm \frac{1}{\kappa}\varepsilon\|^2$  may be characterized using duality as

$$\|\kappa e \pm \frac{1}{\kappa}\varepsilon\|^2 = \sup_{v \in \mathcal{V}} \frac{R^\pm(v)^2}{\|v\|^2} \geq \frac{R^\pm(\xi^\pm)^2}{\|\xi^\pm\|^2} \quad \forall \xi^\pm \in \mathcal{V}.$$

From Theorem A.2.1 one can deduce that in order to find upper and lower bounds for the error in the output it is sufficient to find upper bounds for the energy norm of the linear combinations  $\kappa e \pm \frac{1}{\kappa}\varepsilon$  and set

$$-\|\kappa e - \frac{1}{\kappa}\varepsilon\|_{\text{UB}}^2 \leq s \leq \frac{1}{4} \|\kappa e + \frac{1}{\kappa}\varepsilon\|_{\text{UB}}^2.$$

Also to improve these bounds, continuous approximations,  $\xi^\pm$ , of  $\kappa e \pm \frac{1}{\kappa}\varepsilon$  may be used. These continuous approximations provide in fact, using the dual definition of the energy norm, lower bounds for the quantities  $\|\kappa e \pm \frac{1}{\kappa}\varepsilon\|$ . The derivation of the upper bounds for  $\|\kappa e + \frac{1}{\kappa}\varepsilon\|$  and the continuous approximations  $\xi^\pm$  of  $\kappa e + \frac{1}{\kappa}\varepsilon$  using a posteriori error estimation techniques is further discussed in Chapter 3.

## A.2.2 Bounds for nonself-adjoint model problems

As in the symmetric case, the bounds for the error in the output  $s$  are derived from (A.11). However, in this case, the saddle point  $\bar{v}^\pm$  can not be directly computed from  $e, \varepsilon$  and  $\xi^\pm$  but its computed from symmetric approximations of the errors and  $\xi^\pm$ . Once these approximations are introduced, the bounds are readily found as in the symmetric model problem with no additional difficulties.

**Theorem A.2.2.** *Let  $\xi^+$  and  $\xi^- \in \mathcal{V}$  be two continuous functions, and  $e^s, \varepsilon^s$  and  $\xi^{s\pm} \in \mathcal{V}$  be the solution of the global problems*

$$a^s(e^s, v) = R^P(v), \quad a^s(\varepsilon^s, v) = R^D(v) \quad \text{and} \quad a^s(\xi^{s\pm}, v) = a(v, \xi^\pm) \quad \forall v \in \mathcal{V}. \quad (\text{A.14})$$

Then, the quantities  $s_u$  and  $s_l$  given by

$$s_l = \frac{1}{4} \frac{(2\kappa R^P(\xi^+) - R^-(\xi^{s+}))^2}{\|\xi^{s+}\|^2} - \frac{1}{4} \left\| \kappa e^s - \frac{1}{\kappa} \varepsilon^s \right\|^2,$$

$$s_u = \frac{1}{4} \left\| \kappa e^s + \frac{1}{\kappa} \varepsilon^s \right\|^2 - \frac{1}{4} \frac{(2\kappa R^P(\xi^-) - R^+(\xi^{s-}))^2}{\|\xi^{s-}\|^2},$$

are a lower and an upper for the output  $s$  respectively, that is

$$s_l \leq s \leq s_u.$$

Moreover, the choice  $\xi^\pm = \kappa e \pm \frac{1}{\kappa} \varepsilon$ , lead to optimal bounds, that is,  $s = s_l = s_u$ .

*Proof.* The proof is analogous to the proof of Theorem A.2.1. The only difference is that now, since the bilinear form appearing in the Lagrangian  $\tilde{L}^\pm(\cdot, \cdot)$  is not  $a(\cdot, \cdot)$  but  $a^s(\cdot, \cdot)$ , in order to find the bounds the auxiliary functions  $e^s, \varepsilon^s$  and  $\xi^{s\pm} \in \mathcal{V}$  have to be introduced.

As in the symmetric case, given  $\xi^\pm \in \mathcal{V}$ , the bounds are recovered from equation (A.10) yielding

$$\pm s \geq \min_{v \in \mathcal{V}} \tilde{L}^\pm(v, \lambda^\pm \xi^\pm),$$

and hence,  $\pm s \geq \tilde{L}^\pm(\bar{v}^\pm, \lambda^\pm \xi^\pm)$ , for  $\bar{v}^\pm = \arg \min_{v \in \mathcal{V}} \tilde{L}^\pm(v, \lambda^\pm \xi^\pm)$ .

The weak problem for the saddle point  $\bar{v}^\pm \in \mathcal{V}$  is found imposing that the variations of  $\tilde{L}^\pm(v, \lambda^\pm \xi^\pm)$  with respect to  $v \in \mathcal{V}$  must vanish leading to

$$a^s(\bar{v}^\pm, v) = \frac{1}{2\kappa} R^\mp(v) + \frac{\lambda^\pm}{2\kappa} a(v, \xi^\pm) \quad \forall v \in \mathcal{V}, \quad (\text{A.15})$$

This residual problem may be rewritten taking into account that

$$R^\mp(v) = \kappa R^P(v) \mp \frac{1}{\kappa} R^D(v) = a^s(\kappa e^s \mp \frac{1}{\kappa} \varepsilon^s, v),$$

and the definition of the symmetric function  $\xi^{s\pm}$  (equation (A.14)) as

$$a^s(\bar{v}^\pm, v) = \frac{1}{2\kappa} a^s(\kappa e^s \mp \frac{1}{\kappa} \varepsilon^s, v) + \frac{\lambda^\pm}{2\kappa} a^s(\xi^{s\pm}, v) \quad \forall v \in \mathcal{V},$$

yielding  $\bar{v}^\pm = \frac{1}{2\kappa}(\kappa e^s \mp \frac{1}{\kappa} \varepsilon^s + \lambda^\pm \xi^{s\pm})$ . Then,  $\tilde{L}^\pm(\bar{v}^\pm, \lambda^\pm \xi^\pm)$  may be rewritten after rearranging terms as

$$\tilde{L}^\pm(\bar{v}^\pm, \tilde{\varphi}^\pm) = -\frac{1}{4} \|\kappa e^s \mp \frac{1}{\kappa} \varepsilon^s\|^2 - \frac{1}{4} (\lambda^\pm)^2 \|\xi^{s\pm}\|^2 + \frac{\lambda^\pm}{2} (2R^P(\xi^\pm) - R^\mp(\xi^{s\pm})).$$

Indeed

$$\begin{aligned} \tilde{L}^\pm(\bar{v}^\pm, \lambda^\pm \xi^\pm) &= \kappa^2 a^s(\bar{v}^\pm, \bar{v}^\pm) - \kappa R^\mp(\bar{v}^\pm) + \kappa (R^P(\lambda^\pm \xi^\pm) - a(\bar{v}^\pm, \lambda^\pm \xi^\pm)) \\ &= \kappa^2 a^s(\bar{v}^\pm, \bar{v}^\pm) - 2\kappa^2 a^s(\bar{v}^\pm, \bar{v}^\pm) + \kappa R^P(\lambda^\pm \xi^\pm) \\ &= \kappa R^P(\lambda^\pm \xi^\pm) - \kappa^2 a^s(\bar{v}^\pm, \bar{v}^\pm) = \kappa \lambda^\pm R^P(\xi^\pm) - \frac{1}{4} \|\kappa e^s \mp \frac{1}{\kappa} \varepsilon^s + \lambda^\pm \xi^{s\pm}\|^2 \\ &= \kappa \lambda^\pm R^P(\xi^\pm) - \frac{1}{4} \|\kappa e^s \mp \frac{1}{\kappa} \varepsilon^s\|^2 - \frac{(\lambda^\pm)^2}{4} \|\xi^{s\pm}\|^2 - \frac{\lambda^\pm}{2} a^s(\kappa e^s \mp \frac{1}{\kappa} \varepsilon^s, \xi^{s\pm}) \\ &= -\frac{1}{4} \|\kappa e^s \mp \frac{1}{\kappa} \varepsilon^s\|^2 - \frac{(\lambda^\pm)^2}{4} \|\xi^{s\pm}\|^2 + \frac{\lambda^\pm}{2} (2R^P(\xi^\pm) - a^s(\kappa e^s \mp \frac{1}{\kappa} \varepsilon^s, \xi^{s\pm})) \\ &= -\frac{1}{4} \|\kappa e^s \mp \frac{1}{\kappa} \varepsilon^s\|^2 - \frac{(\lambda^\pm)^2}{4} \|\xi^{s\pm}\|^2 + \frac{\lambda^\pm}{2} (2R^P(\xi^\pm) - R^\mp(\xi^{s\pm})). \end{aligned}$$

*Optimal parameter determination:* The optimal value for  $\lambda^\pm$  is

$$\lambda^\pm = \frac{2R^P(\xi^\pm) - R^\mp(\xi^{s\pm})}{\|\xi^{s\pm}\|^2},$$

leading to the bounds

$$\pm s \geq \frac{1}{4} \frac{(2\kappa R^P(\xi^\pm) - R^\mp(\xi^{s\pm}))^2}{\|\xi^{s\pm}\|^2} - \frac{1}{4} \|\kappa e^s \mp \frac{1}{\kappa} \varepsilon^s\|^2. \quad (\text{A.16})$$

Now taking the positive part of the previous equation leads to

$$s \geq \frac{1}{4} \frac{(2\kappa R^P(\xi^+) - R^-(\xi^{s+}))^2}{\|\xi^{s+}\|^2} - \frac{1}{4} \|\kappa e^s - \frac{1}{\kappa} \varepsilon^s\|^2 = s_1, \quad (\text{A.17})$$

and taking the negative part and multiplying the inequality by  $-1$ ,

$$s \leq \frac{1}{4} \left\| \kappa e^s + \frac{1}{\kappa} \varepsilon^s \right\|^2 - \frac{1}{4} \frac{(2\kappa R^P(\xi^-) - R^+(\xi^{s-}))^2}{\|\xi^{s-}\|^2} = s_u, \quad (\text{A.18})$$

and the first part of the proof is concluded.

Finally, the optimality of the bounds for the choice  $\xi^\pm = \kappa e \pm \frac{1}{\kappa} \varepsilon$  follows directly from equation (A.4).  $\square$

### Practical computation of the bounds for nonself-adjoint model problems

Theorem A.2.2 states that it is possible to obtain bounds for  $s$  given continuous approximations  $\xi^\pm$  of  $\kappa e \pm \frac{1}{\kappa} \varepsilon$  just computing the symmetric functions  $e^s$ ,  $\varepsilon^s$  and  $\xi^{s\pm}$ . Obviously, it is not possible in general to compute these functions exactly. As in the symmetric case, the energy norms  $\|\kappa e^s + \frac{1}{\kappa} \varepsilon^s\|$  may be replaced by an upper bound using techniques for estimating the energy norm of symmetric model problems. Thus, the only added difficulty between the bounds in the symmetric and nonsymmetric problems is the dependence of the bounds on the symmetric function  $\xi^{s\pm}$ . Fortunately, the same techniques to obtain upper bounds for the energy norm of solutions of symmetric model problems applied to  $\xi^{s\pm}$ , allow to replace the terms containing  $\xi^{s\pm}$  by approximations which still maintain the upper and lower bound property.

Let  $\hat{e}^s$  and  $\hat{\varepsilon}^s \in \widehat{\mathcal{V}}$  be two estimates verifying

$$a^s(\hat{e}^s, v) = R^P(v), \quad a^s(\hat{\varepsilon}^s, v) = R^D(v) \quad \forall v \in \mathcal{V} \quad (\text{A.19})$$

where  $\widehat{\mathcal{V}}$  is the *broken* space obtained from  $\mathcal{V}$  relaxing both the Dirichlet boundary conditions and the continuity of the functions across the edges of the mesh. The broken space is the most usual interpolation space for the estimates (see Chapter 3). It is worth noting that most residual implicit type error estimation techniques yield estimates verifying the previous condition.

From (A.19) it is easily shown that  $\hat{e}^s$  and  $\hat{\varepsilon}^s$  provide upper bounds for the energy norm of  $\kappa e^s \pm \frac{1}{\kappa} \varepsilon^s$ ,

$$\left\| \kappa e^s \pm \frac{1}{\kappa} \varepsilon^s \right\| \leq \left\| \kappa \hat{e}^s \pm \frac{1}{\kappa} \hat{\varepsilon}^s \right\|,$$

see Lemma 3.2.1. Thus the bounds

$$\pm s \geq \frac{1}{4} \frac{(2\kappa R^P(\xi^\pm) - R^\mp(\xi^{s\pm}))^2}{\|\xi^{s\pm}\|^2} - \frac{1}{4} \left\| \kappa e^s \mp \frac{1}{\kappa} \varepsilon^s \right\|^2$$

may be replaced by

$$\pm s \geq \frac{1}{4} \frac{(2\kappa R^P(\xi^\pm) - R^\mp(\xi^{s\pm}))^2}{\|\xi^{s\pm}\|^2} - \frac{1}{4} \left\| \kappa \hat{e}^s \mp \frac{1}{\kappa} \hat{\varepsilon}^s \right\|^2. \quad (\text{A.20})$$

Moreover, consider  $\xi^\pm$  to be a continuous approximation of  $\kappa e \pm \frac{1}{\kappa} \varepsilon$  obtained post-processing the estimates  $\hat{e}^s$  and  $\hat{\varepsilon}^s$ . Using again an estimation technique, one would compute  $\hat{\xi}^{s\pm} \in \hat{\mathcal{V}}$  such that

$$a^s(\hat{\xi}^{s\pm}, v) = a^s(\xi^{s\pm}, v) = a(v, \xi^\pm) \quad \forall v \in \mathcal{V}, \quad (\text{A.21})$$

yielding an upper bound for the energy norm of  $\xi^{s\pm}$ ,  $\|\xi^{s\pm}\| \leq \|\hat{\xi}^{s\pm}\|$ . Thus, the term  $\|\xi^{s\pm}\|$  in the bounds given by equation (A.20) may be replaced by  $\|\hat{\xi}^{s\pm}\|$  still maintaining the bounding property, namely

$$\pm s \geq \frac{1}{4} \frac{(2\kappa R^P(\xi^\pm) - R^\mp(\xi^{s\pm}))^2}{\|\hat{\xi}^{s\pm}\|^2} - \frac{1}{4} \left\| \kappa \hat{e}^s \mp \frac{1}{\kappa} \hat{\varepsilon}^s \right\|^2. \quad (\text{A.22})$$

Note that still remains an uncomputable term in the bound expressions  $R^\mp(\xi^{s\pm})$ . This term may be computed also using the estimates  $\hat{e}^s$ ,  $\hat{\varepsilon}^s$  and  $\hat{\xi}^{s\pm}$ . The computation of this term is described in the following.

From the definitions of  $\hat{e}^s$ ,  $\hat{\varepsilon}^s$ , since  $\xi^{s\pm} \in \mathcal{V}$

$$R^\mp(\xi^{s\pm}) = \kappa a^s(\hat{e}^s, \xi^{s\pm}) \mp \frac{1}{\kappa} a^s(\hat{\varepsilon}^s, \xi^{s\pm}) = a^s\left(\kappa \hat{e}^s \mp \frac{1}{\kappa} \hat{\varepsilon}^s, \xi^{s\pm}\right),$$

however, since  $\hat{e}^s$  and  $\hat{\varepsilon}^s$  are generally discontinuous functions ( $\hat{e}^s, \hat{\varepsilon}^s \notin \mathcal{V}$ ), equation (A.21) may not be used to replace  $\xi^{s\pm}$  by  $\hat{\xi}^{s\pm}$  in the previous equality.

However, the particular form of the residual problem (A.21) allows to compute an estimate  $\hat{\xi}^{s\pm} \in \hat{\mathcal{V}}$  for which equation (A.21) holds not only for any  $v \in \mathcal{V}$  but for any  $v \in \hat{\mathcal{V}}$ . Indeed, consider the local problems obtained restricting the global problem into an element  $\Omega_k \subset \Omega$ : find  $\hat{\xi}_k^{s\pm} \in \mathcal{V}_k$

$$a_k^s(\hat{\xi}_k^{s\pm}, v) = a_k(v, \xi^\pm) \quad \forall v \in \mathcal{V}_k.$$

These local problems do not have solvability problems, that is, there is no need to equilibrate the local problems, and thus  $\hat{\xi}^{s\pm} = \sum_{k=1}^{n_{el}} \hat{\xi}_k^{s\pm}$  verifies

$$a^s(\hat{\xi}^{s\pm}, v) = a^s(\xi^{s\pm}, v) = a(v, \xi^\pm) \quad \forall v \in \hat{\mathcal{V}}.$$

In particular for  $v = \kappa \hat{e}^s \mp \frac{1}{\kappa} \hat{\varepsilon}^s$ , yields

$$R^\mp(\xi^{s\pm}) = a^s(\kappa \hat{e}^s \mp \frac{1}{\kappa} \hat{\varepsilon}^s, \xi^{s\pm}) = a^s(\kappa \hat{e}^s \mp \frac{1}{\kappa} \hat{\varepsilon}^s, \hat{\xi}^{s\pm}), \quad (\text{A.23})$$

and the bounds for the output may be computed as

$$\pm s \geq \frac{1}{4} \frac{(2\kappa R^P(\xi^\pm) - a^s(\kappa \hat{e}^s \mp \frac{1}{\kappa} \hat{\varepsilon}^s, \hat{\xi}^{s\pm}))^2}{\|\hat{\xi}^{s\pm}\|^2} - \frac{1}{4} \|\kappa \hat{e}^s \mp \frac{1}{\kappa} \hat{\varepsilon}^s\|^2, \quad (\text{A.24})$$

as mentioned in (2.26).

A summary of the practical computation of the bounds for the output in a non-symmetric model problem is given in Figure A.1.

### Quality of the bounds using only upper bounds for the energy

From Theorem A.2.2 (taking  $\xi^\pm = 0$ ) one can deduce that in order to find upper and lower bounds for the error in the output it is sufficient to find upper bounds for the energy norm of the linear combinations  $\kappa e^s \pm \frac{1}{\kappa} \varepsilon^s$  and set

$$-\|\kappa e^s - \frac{1}{\kappa} \varepsilon^s\|_{\text{UB}}^2 \leq s \leq \frac{1}{4} \|\kappa e^s + \frac{1}{\kappa} \varepsilon^s\|_{\text{UB}}^2,$$

thus recovering the strategy proposed by Paraschivoiu et al. (1997).

The quality of these bounds depends on two factors: first, the quantities  $\|\kappa e^s \pm \frac{1}{\kappa} \varepsilon^s\|^2$  are replaced by an upper bound of them, and therefore the sharpness of the bounds is controlled by the accuracy of the upper bound error estimation techniques.

Second, the bounds do not take into account the terms

$$\frac{1}{4} \frac{(2\kappa R^P(\xi^\pm) - R^\mp(\xi^{s\pm}))^2}{\|\xi^{s\pm}\|^2}, \quad (\text{A.25})$$

for  $\xi^\pm = \kappa e \pm \frac{1}{\kappa} \varepsilon$ . If these terms are large compared to the value of the output  $s$ , then the quality of the proposed bounds will be poor.

The unestimated contribution (A.25) to the output  $s$ , can be rewritten with an algebraic manipulations as

$$\|\kappa e - \frac{1}{2}(\kappa e^s \mp \frac{1}{\kappa} \varepsilon^s)\|^2,$$

taking into account that for  $\xi^\pm = \kappa e \pm \frac{1}{\kappa} \varepsilon$ ,  $\xi^{s\pm} = 2(\kappa e - \frac{1}{2}(\kappa e^s \mp \frac{1}{\kappa} \varepsilon^s))$ . That is, in fact

$$\pm s = \|\kappa e - \frac{1}{2}(\kappa e^s \mp \frac{1}{\kappa} \varepsilon^s)\|^2 - \frac{1}{4} \|\kappa e^s \mp \frac{1}{\kappa} \varepsilon^s\|^2.$$

1.- Compute the upper bound estimates  $\hat{e}^s$  and  $\hat{\varepsilon}^s$  s.t.

$$a^s(\hat{e}^s, v) = R^P(v), \quad a^s(\hat{\varepsilon}^s, v) = R^D(v) \quad \forall v \in \mathcal{V},$$

yielding

$$\|\kappa e^s \pm \frac{1}{\kappa} \varepsilon^s\| \leq \|\kappa \hat{e}^s \pm \frac{1}{\kappa} \hat{\varepsilon}^s\|.$$

2.- Compute continuous approximations  $\xi^{\pm}$  of  $\kappa e \pm \frac{1}{\kappa} \varepsilon$  post-processing the estimates  $\hat{e}^s$  and  $\hat{\varepsilon}^s$ .

3.- Compute the upper bound estimates  $\hat{\xi}^{s\pm}$  s.t.

$$a^s(\hat{\xi}^{s\pm}, v) = a^s(v, \xi^{s\pm}) = a(v, \xi^{\pm}) \quad \forall v \in \hat{\mathcal{V}}$$

yielding

$$\|\xi^{s\pm}\| \leq \|\hat{\xi}^{s\pm}\|,$$

and

$$a^s(\kappa \hat{e}^s \mp \frac{1}{\kappa} \hat{\varepsilon}^s, \xi^{s\pm}) = a^s(\kappa \hat{e}^s \mp \frac{1}{\kappa} \hat{\varepsilon}^s, \hat{\xi}^{s\pm}).$$

4.- Compute  $\kappa = \sqrt{\|\hat{\varepsilon}^s\|/\|\hat{e}^s\|}$  and the quantities  $s_u$  and  $s_l$  as

$$s_l = \frac{1}{4} \frac{(2\kappa R^P(\xi^+) - a^s(\kappa \hat{e}^s - \frac{1}{\kappa} \hat{\varepsilon}^s, \hat{\xi}^{s+}))^2}{\|\hat{\xi}^{s+}\|^2} - \frac{1}{4} \|\kappa \hat{e}^s - \frac{1}{\kappa} \hat{\varepsilon}^s\|^2,$$

$$s_u = \frac{1}{4} \|\kappa \hat{e}^s + \frac{1}{\kappa} \hat{\varepsilon}^s\|^2 - \frac{1}{4} \frac{(2\kappa R^P(\xi^-) - a^s(\kappa \hat{e}^s + \frac{1}{\kappa} \hat{\varepsilon}^s, \hat{\xi}^{s-}))^2}{\|\hat{\xi}^{s-}\|^2}.$$

Then

$$s_l \leq s \leq s_u$$

are the bounds for the output.

Figure A.1: Main steps of the strategy used to obtain bounds for an output  $s$  depending on the solution of a nonsymmetric boundary value problem.





# Appendix B

## Comparison between subdomain-based *flux-free* residual methods

The subdomain residual method was first devised by Babuška and Rheinboldt (1978a) (see also Babuška and Rheinboldt 1978b, Babuška and Rheinboldt 1979). The presented approach proposes to solve the subdomain residual problems: find  $\eta^i \in \mathcal{V}_{\omega^i}^0$  such that

$$a_{\omega^i}(\eta^i, v) = R^*(v) \quad \forall v \in \mathcal{V}_{\omega^i}^0,$$

where  $\mathcal{V}_{\omega^i}^0$  is the local test space in the star  $\omega^i$  with Dirichlet homogeneous boundary conditions, and recover the global error estimator  $\eta$  by summing the norm contributions from the subdomains

$$\eta = \left( \sum_{i=1}^{n_{\text{np}}} \|\eta^i\|^2 \right)^{\frac{1}{2}}.$$

This approach leads to an estimate  $\eta$  that is no strict upper or lower bound for the energy norm of the error but however provides two-sided bounds for the error, that is, there exist two constants  $C_1$  and  $C_2$  depending only on the regularity of the elements of the mesh such that

$$C_1\eta \leq \|z\| \leq C_2\eta.$$

Carstensen and Funken (1999/00), Morin et al. (2003), Machiels et al. (2000) and Prudhomme et al. (2004) developed similar subdomain residual error estimation

techniques leading to upper bounds for the energy norm of the error in the context of scalar model problems. They present two-sided bounds leading to estimates  $\eta \in \mathbb{R}$  verifying

$$C\eta \leq \|z\| \leq \eta,$$

where  $C$  is a constant depending only on the regularity of the mesh.

This appendix presents a unified approach of the techniques presented by Carstensen and Funken (1999/00), Morin et al. (2003), Machiels et al. (2000) and Prudhomme et al. (2004) allowing to easily compare these techniques with the approach presented in Chapter 4. The rationale is to decompose the bilinear form  $a(\cdot, \cdot)$  in a sum of local contributions associated with each star. That is, weighted local bilinear forms  $a_{w^i}(\cdot, \cdot)$  are introduced such that

$$a(w, v) = \sum_{i=1}^{n_{np}} a_{w^i}(w, v). \quad (\text{B.1})$$

The weighted bilinear forms  $a_{w^i}(\cdot, \cdot)$  obviously depend on the problem at hand, but they can be defined in a general way introducing non-negative local integrable weights, supported in  $\omega^i$ , verifying the partition of the unity property, that is,  $w^i \in \mathcal{L}^2(\omega^i)$  such that

$$\sum_{i=1}^{n_{np}} w^i = 1 \quad \text{and} \quad w^i(x) \geq 0 \quad \forall x \in \omega^i.$$

Then, the weighted bilinear forms  $a_{w^i}(\cdot, \cdot)$  are obtained replacing the integrals appearing in  $a(\cdot, \cdot)$  by weighted integrals associated with the weights  $w^i$ . For instance, for the diffusion-reaction model problem, one would obtain

$$a_{w^i}(w, v) = \int_{\Omega} w^i (\nu \nabla w \cdot \nabla v + \mu w v) \, d\Omega,$$

whereas for the mechanical problem,

$$a_{w^i}(\mathbf{w}, \mathbf{v}) = \int_{\Omega} w^i \boldsymbol{\sigma}(\mathbf{w}) : \boldsymbol{\varepsilon}(\mathbf{v}) \, d\Omega.$$

Note that the weights  $w^i$  account for the overlapping of the stars.

Once the bilinear form is decomposed into local contributions, the local estimates  $\eta^i \in \mathcal{V}_{\omega^i}$  are computed solving the local equations

$$a_{w^i}(\eta^i, v) = R^P(\phi^i v) \quad \forall v \in \mathcal{V}_{\omega^i}. \quad (\text{B.2})$$

The global error estimator  $\eta$  is obtained by summing the weighted norms of the local estimates  $\eta^i$

$$\eta = \left( \sum_{i=1}^{n_{np}} \|\eta^i\|_{w^i}^2 \right)^{\frac{1}{2}} \geq \|z\|, \quad (\text{B.3})$$

where the local norm  $\|\cdot\|_{w^i}^2$  is the norm induced by the local weighted scalar product  $a_{w^i}(\cdot, \cdot)$ , that is  $\|v\|_{w^i}^2 = a_{w^i}(v, v)$ .

*Remark B.0.1.* Carstensen and Funken (1999/00), Morin et al. (2003), Machiels et al. (2000) and Prudhomme et al. (2004) are only concerned with the scalar (thermal) problem, however, the estimate  $\eta$  can be easily extended to the mechanical problem using the modification of the residue appearing in the r.h.s. of equation (B.3) as presented in (4.7).

The upper bound property,  $\|z\| \leq \eta$ , is proved in Parés, Díez and Huerta (2005, Theorem 12). The repeated use of the Cauchy-Schwarz inequality in the proof of the upper bound property suggests that the obtained upper bound is not as sharp as the upper bound associated with the estimate  $\hat{z}$  defined from the local estimates  $\hat{z}^i$  solution of (4.6). Theorem B.0.2 provides a comparison between the two estimates, but in order to prove the theorem, the following lemma, which is a particular case of the Chebyshev sum inequality, must be introduced.

**Lemma B.0.1.** *Let  $\langle \cdot, \cdot \rangle$  be a scalar product with associated norm  $|\cdot|$  acting on a space  $V$ , then given  $m$  functions in  $V$ ,  $a_i \in V$ , it holds that*

$$\left| \sum_{i=1}^m a_i \right|^2 \leq m \sum_{i=1}^m |a_i|^2.$$

*Proof.* Note that the previous inequality is equivalent to see that

$$m \sum_{i=1}^m |a_i|^2 - \left| \sum_{i=1}^m a_i \right|^2 \geq 0 \quad \forall a_i \in V.$$

The proof consists in showing that

$$m \sum_{i=1}^m |a_i|^2 - \left| \sum_{i=1}^m a_i \right|^2 = \sum_{i=1}^m \sum_{j=i+1}^m |a_i - a_j|^2 \geq 0, \quad (\text{B.4})$$

and this is done by induction.

Let  $m = 2$ , then  $2(a_1^2 + a_2^2) - (a_1 + a_2)^2 = (a_1 - a_2)^2$ , and the base case holds.

Assume now that equation (B.4) holds for  $m$ . The goal is to prove that equation (B.4) is also true for  $m + 1$ , that is

$$(m + 1) \sum_{i=1}^{m+1} |a_i|^2 - \left| \sum_{i=1}^{m+1} a_i \right|^2 = \sum_{i=1}^{m+1} \sum_{j=i+1}^{m+1} |a_i - a_j|^2,$$

which follows doing a little bit of algebra. The terms appearing in the r.h.s. of the previous equation may be rewritten as

$$\begin{aligned} (m + 1) \sum_{i=1}^{m+1} |a_i|^2 &= m \sum_{i=1}^m |a_i|^2 + \sum_{i=1}^m |a_i|^2 + (m + 1) |a_{m+1}|^2, \\ \left| \sum_{i=1}^{m+1} a_i \right|^2 &= \left| \sum_{i=1}^m a_i \right|^2 + |a_{m+1}|^2 + 2 \left\langle \sum_{i=1}^m a_i, a_{m+1} \right\rangle. \end{aligned}$$

Thus, denoting by  $L = (m + 1) \sum_{i=1}^{m+1} |a_i|^2 - \left| \sum_{i=1}^{m+1} a_i \right|^2$  and using the induction hypothesis

$$\begin{aligned} L &= m \sum_{i=1}^m |a_i|^2 - \left| \sum_{i=1}^m a_i \right|^2 + \sum_{i=1}^m |a_i|^2 + m |a_{m+1}|^2 - 2 \sum_{i=1}^m \langle a_i, a_{m+1} \rangle \\ &= \sum_{i=1}^m \sum_{j=i+1}^m |a_i - a_j|^2 + \sum_{i=1}^m |a_i|^2 + m |a_{m+1}|^2 - 2 \sum_{i=1}^m \langle a_i, a_{m+1} \rangle \\ &= \sum_{i=1}^m \sum_{j=i+1}^m |a_i - a_j|^2 + \sum_{i=1}^m (|a_i|^2 + |a_{m+1}|^2 - 2 \langle a_i, a_{m+1} \rangle) \\ &= \sum_{i=1}^m \sum_{j=i+1}^m |a_i - a_j|^2 + \sum_{i=1}^m |a_{m+1} - a_i|^2 = \sum_{i=1}^{m+1} \sum_{j=i+1}^{m+1} |a_i - a_j|^2, \end{aligned}$$

which concludes the proof.  $\square$

**Theorem B.0.2.** Let  $\hat{z} = \sum_{i=1}^{n_{np}} \hat{z}^i$  denote the error estimate obtained from the local estimates  $\hat{z}^i$  solution of the local problems given in equation (4.6) and  $\eta$  denote the estimate given in equation (B.3), then there exist a constant  $C$  depending only on the regularity of the mesh such that

$$\|\hat{z}\| \leq C\eta.$$

*Proof.* Let  $\mathcal{N}$  be the set of vertices in the mesh. It is possible to partition  $\mathcal{N}$  into the union of disjoint subsets  $\mathcal{N}_1, \mathcal{N}_2, \dots$  such that any pair of nodal basis functions in the same subset  $\mathcal{N}_r$  have nonoverlapping supports. Specifically, the condition that will be required is

$$\forall i, j \in \mathcal{N}_r : i \neq j \implies \text{int } \omega^i \cap \text{int } \omega^j \text{ is empty.}$$

It is easy to see that such a partitioning exists since one can simply choose each of the sets  $\mathcal{N}_r$  to consist of a single vertex. The smallest possible number of subsets is denoted by  $\nu$  and referred as the overlap index for the partition. The overlap index for a regular family of partitions may be bounded by a constant depending only on the regularity of the elements in the family. Let  $\mathcal{N}$  be a minimal partition, that is,  $\mathcal{N} = \{\mathcal{N}_r\}_{r=1 \dots \nu}$ .

Then, with the aid of Lemma B.0.1

$$\|\hat{z}\|^2 = \left\| \sum_{i=1}^{n_{\text{np}}} \hat{z}^i \right\|^2 = \left\| \sum_{i \in \mathcal{N}} \hat{z}^i \right\|^2 = \left\| \sum_{r=1}^{\nu} \sum_{i \in \mathcal{N}_r} \hat{z}^i \right\|^2 \leq \nu \sum_{r=1}^{\nu} \left\| \sum_{i \in \mathcal{N}_r} \hat{z}^i \right\|^2 = \nu \sum_{r=1}^{\nu} \sum_{i \in \mathcal{N}_r} \|\hat{z}^i\|^2,$$

where the last equality follows from the fact that the stars involved in  $\mathcal{N}_r$  are nonoverlapping. Thus,

$$\|\hat{z}\|^2 \leq \nu \sum_{i=1}^{n_{\text{np}}} \|\hat{z}^i\|^2. \quad (\text{B.5})$$

The proofs ends showing that  $\sum_{i=1}^{n_{\text{np}}} \|\hat{z}^i\|^2 \leq \eta^2$ . Indeed, taking  $v = \hat{z}^i$  in equation (4.6) and (B.2) and applying the Cauchy-Schwarz inequality

$$\|\hat{z}^i\|^2 = a_{\omega^i}(\hat{z}^i, \hat{z}^i) = R^*(\phi^i \hat{z}^i) = a_{\mathbf{w}^i}(\eta^i, \hat{z}^i) \leq \|\eta^i\|_{\mathbf{w}^i} \|\hat{z}^i\|_{\mathbf{w}^i}.$$

Then, summing over all the stars and noticing that  $\|\hat{z}^i\|_{\mathbf{w}^i} \leq \|\hat{z}^i\|$  (since  $\|\mathbf{w}^i\|_{\infty} \leq 1$ ) it follows that

$$\sum_{i=1}^{n_{\text{np}}} \|\hat{z}^i\|^2 \leq \sum_{i=1}^{n_{\text{np}}} \|\eta^i\|_{\mathbf{w}^i} \|\hat{z}^i\|_{\mathbf{w}^i} \leq \sum_{i=1}^{n_{\text{np}}} \|\eta^i\|_{\mathbf{w}^i} \|\hat{z}^i\|,$$

where applying the Cauchy-Schwarz inequality again lead to the desired expression

$$\sum_{i=1}^{n_{\text{np}}} \|\hat{z}^i\|^2 \leq \left( \sum_{i=1}^{n_{\text{np}}} \|\eta^i\|_{\mathbf{w}^i}^2 \right)^{\frac{1}{2}} \left( \sum_{i=1}^{n_{\text{np}}} \|\hat{z}^i\|^2 \right)^{\frac{1}{2}} = \eta \left( \sum_{i=1}^{n_{\text{np}}} \|\hat{z}^i\|^2 \right)^{\frac{1}{2}}. \quad (\text{B.6})$$

Finally, joining equations (B.5) and (B.6), the estimate  $\|\hat{z}\|$  can be bounded by

$$\|\hat{z}\| \leq \sqrt{\nu}\eta.$$

□

Theorem B.0.2 does not state that the estimate  $\|\hat{z}\|$  provides a better upper bound for the energy norm of the error, because of the factor  $\sqrt{\nu} > 1$ . However, the bound  $\|\hat{z}\| \leq \sqrt{\nu}\eta$  is obtained using repeatedly the Cauchy-Schwarz inequality, the property  $\|v\|_{w^i} \leq \|v\|$  and also Lemma B.0.1. Thus, the constant  $\sqrt{\nu}$  is dictated by the most pathological functions  $\hat{z}^i$  leading to

$$\left\| \sum_{i=1}^{n_{np}} \hat{z}^i \right\|^2 = \nu \sum_{i=1}^{n_{np}} \|\hat{z}^i\|^2,$$

where there is no cancellation between the different estimates  $\hat{z}^i$ . Then, in practice one would expect to find that  $\|\hat{z}\| \leq C\eta$  for smaller values of  $C$  which in the worst-case scenario would at most lead to  $C = \sqrt{\nu}$ . In fact, if one considers a regular triangular mesh the theory of graphs informs us that  $\nu = 4$  and taking then the local weighted functions  $w^i = \frac{1}{3}$ , one has that  $\|\hat{z}\| \leq \sqrt{\frac{4}{3}}\eta$ .

Numerical examples confirm this impression: the estimate  $\|\hat{z}\|$  provides sharper bounds for the energy norm of the error than the estimate  $\eta$  leading to  $C \ll 1$  (see Section 4.5).

# Appendix C

## Strict bounds for outputs of interest: bounds for the $J$ -integral

This appendix presents an *a-posteriori* method for computing strict upper and lower bounds of the  $J$ -integral in two dimensional linear elasticity. The  $J$ -integral, which is typically expressed as a contour integral, is recast as a surface integral which yields a quadratic continuous functional of the displacement. By expanding the quadratic output about an approximate finite element solution, the output is expressed as a known computable quantity plus linear and quadratic functionals of the solution error. The quadratic component is bounded by the energy norm of the error scaled by a continuity constant, which is determined explicitly. The linear component is expressed as an inner product of the errors in the displacement and in a computed adjoint solution, and bounded using standard *a-posteriori* error estimation techniques. The method is illustrated with two fracture problems in plane strain elasticity. An important feature of the method presented is that the computed bounds are strict with respect to the weak solution of the elasticity equation.

Xuan et al. (2004) present a method for computing bounds for the  $J$ -integral which is a quadratic functional of the solution field. However, the bounds for the  $J$ -integral are strict only with respect to a reference solution. Xuan et al. (2005) continues the work of Xuan et al. (2004) presenting bounds which are strict with respect to the weak solution of the elasticity equations and not for a reference solution. The bounds for the linear component of the output are computed using the error estimation technique presented in Chapter 5 and in (Parés, Bonet, Huerta and Peraire 2005).

This appendix is structured as follows: first the expansion of the  $J$ -integral as a known computable term plus linear and quadratic functionals is provided. Then, upper and lower bounds for the quadratic term are provided and finally the upper and lower bounds for the linear term are treated.

## Reformulation of the $J$ -integral

Consider the elasticity problem with Neumann and homogeneous Dirichlet boundary conditions written in weak form as: find  $\mathbf{u} \in \mathcal{V}$  such that

$$a(\mathbf{u}, \mathbf{v}) = \ell(\mathbf{v}) \quad \forall \mathbf{v} \in \mathcal{V}, \quad (\text{C.1})$$

where  $\mathcal{V} = \{\mathbf{v} \in [\mathcal{H}^1(\Omega)]^2, \mathbf{v}|_{\Gamma_D} = \mathbf{0}\}$ . The linear forcing functional  $\ell \in \mathcal{V}'$

$$\ell(\mathbf{v}) = \int_{\Omega} \mathbf{f} \cdot \mathbf{v} \, d\Omega + \int_{\Gamma^N} \mathbf{g} \cdot \mathbf{v} \, d\Gamma, \quad (\text{C.2})$$

contains both the internal forces per unit volume  $\mathbf{f} \in [\mathcal{H}^{-1}(\Omega)]^2$  and the Neumann boundary tractions  $\mathbf{g} \in [\mathcal{H}^{-\frac{1}{2}}(\Gamma^N)]^2$  and  $a : \mathcal{V} \times \mathcal{V} \rightarrow \mathbb{R}$  is the symmetric coercive bilinear form given by

$$a(\mathbf{w}, \mathbf{v}) = \int_{\Omega} \boldsymbol{\sigma}(\mathbf{w}) : \boldsymbol{\varepsilon}(\mathbf{v}) \, d\Omega.$$

For a two-dimensional linear elastic body the energy release rate,  $J(\mathbf{u})$ , can be calculated as a path independent line integral known as the  $J$ -integral (Rice 1968). If the geometry shown in Figure C.1 is considered, the  $J$ -integral has the following expression,

$$J(\mathbf{u}) = \int_{\Gamma} \left( W^e n_1 - \mathbf{T} \cdot \frac{\partial \mathbf{u}}{\partial x_1} \right) d\Gamma,$$

where  $\Gamma$  is any path beginning at the bottom crack face and ending at the top crack face,  $W^e = (\boldsymbol{\sigma} : \boldsymbol{\varepsilon})/2$  is the strain energy density,  $\mathbf{T}$  is the traction given as  $\mathbf{T} = \boldsymbol{\sigma} \cdot \mathbf{n}$ , and  $\mathbf{n} = (n_1, n_2)$  is the outward unit normal to  $\Gamma$ .

An alternative expression for  $J(\mathbf{u})$  was proposed by Li, Shih and Needleman (1985), where the contour integral is transformed into the following area integral expression,

$$J(\mathbf{u}) = \int_{\Omega_\chi} \left( (\nabla \chi)^T \cdot \boldsymbol{\sigma} \frac{\partial \mathbf{u}}{\partial x_1} - W^e \frac{\partial \chi}{\partial x_1} \right) d\Omega. \quad (\text{C.3})$$



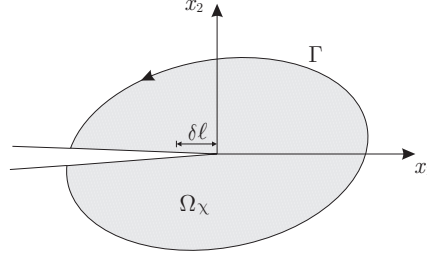


Figure C.1: Crack geometry showing coordinate axes and the  $J$ -integral contour and domain of integration.

Here, the weighting function  $\chi$  is any function in  $H^1(\Omega_\chi)$  that is equal to one at the crack tip and vanishes on  $\Gamma$ .

For a given  $\chi$ ,  $J(\mathbf{u})$  is a bounded quadratic functional of  $\mathbf{u}$ . To be able to obtain bounds for  $J(\mathbf{u})$  it is convenient to make the quadratic dependence of  $J(\mathbf{u})$  more explicit. To this end, define the bilinear form  $\bar{q}(\mathbf{w}, \mathbf{v}) : \mathcal{V} \times \mathcal{V} \rightarrow \mathbb{R}$  as,

$$\bar{q}(\mathbf{w}, \mathbf{v}) = \int_{\Omega_\chi} (\nabla \chi)^T \cdot \boldsymbol{\sigma}(\mathbf{w}) \frac{\partial v}{\partial x_1} d\Omega - \int_{\Omega_\chi} \frac{1}{2} \boldsymbol{\sigma}(\mathbf{w}) : \boldsymbol{\varepsilon}(\mathbf{v}) \frac{\partial \chi}{\partial x_1} d\Omega,$$

and its symmetric part  $q(\mathbf{w}, \mathbf{v}) : \mathcal{V} \times \mathcal{V} \rightarrow \mathbb{R}$ ,  $q(\mathbf{w}, \mathbf{v}) = \frac{1}{2}(\bar{q}(\mathbf{w}, \mathbf{v}) + \bar{q}(\mathbf{v}, \mathbf{w}))$ . It is clear from these definitions that,  $J(\mathbf{u}) = q(\mathbf{u}, \mathbf{u})$ , and that there exists  $\eta < \infty$  such that,

$$|q(\mathbf{v}, \mathbf{v})| \leq \eta \|\mathbf{v}\|^2 \quad \forall \mathbf{v} \in \mathcal{V}. \quad (\text{C.4})$$

The goal is to compute upper and lower bounds, for  $J(\mathbf{u})$ , where  $\mathbf{u}$  satisfies problem (C.1). Let  $\mathbf{u}_H \in \mathcal{V}_H$  be the finite element approximation of  $\mathbf{u}$  lying in the finite dimensional subspace  $\mathcal{V}_H \subset \mathcal{V}$ . For simplicity, we shall assume that  $\mathcal{V}_H$  is the space of piecewise linear continuous functions defined over a triangulation,  $\mathcal{T}_H$ , of  $\Omega$  which satisfies the Dirichlet boundary conditions. An approximation to  $J(\mathbf{u})$ ,  $J_H$ , can be obtained as  $J_H = q(\mathbf{u}_H, \mathbf{u}_H)$ , where, for convenience,  $\chi$  in (C.3) is chosen to be piecewise linear over the elements  $T_H \in \mathcal{T}_H$ . Exploiting the bilinearity of  $q(\mathbf{w}, \mathbf{v})$ , we can write

$$\begin{aligned} J(\mathbf{u}) - J_H &= q(\mathbf{u}, \mathbf{u}) - q(\mathbf{u}_H, \mathbf{u}_H) = q(\mathbf{u} - \mathbf{u}_H, \mathbf{u} - \mathbf{u}_H) + 2q(\mathbf{u}, \mathbf{u}_H) - 2q(\mathbf{u}_H, \mathbf{u}_H) \\ &= q(\mathbf{e}, \mathbf{e}) + 2q(\mathbf{e}, \mathbf{u}_H), \end{aligned}$$

where  $\mathbf{e} = \mathbf{u} - \mathbf{u}_H$  is the error in the approximation  $\mathbf{u}_H$ . It is clear that if we are able to compute bounds  $Q$  and  $L^\pm$  for the quadratic and linear error terms,

$$|q(\mathbf{e}, \mathbf{e})| \leq Q \quad \text{and} \quad L^- \leq q(\mathbf{e}, \mathbf{u}_H) \leq L^+,$$

then, the bounds for  $J(\mathbf{u})$ ,  $J^\pm$ , follow as,

$$J^- \equiv J_H - Q + 2L^- \leq J(\mathbf{u}) \leq J_H + Q + 2L^+ \equiv J^+.$$

### Bounds for the quadratic error term

Xuan et al. (2004) show that for two dimensional linear elasticity, a suitable value for the continuity constant in expression (C.4) is given by

$$\eta_\chi = \max_{T_H \in \mathcal{T}_H} \frac{(3K + 4\mu)|\nabla\chi|^2}{4\sqrt{(3K + \mu)\left(3\mu\left(\frac{\partial\chi}{\partial x_1}\right)^2 + (3K + 4\mu)\left(\frac{\partial\chi}{\partial x_2}\right)^2\right)}}, \quad (\text{C.5})$$

where  $\mu = E/(2(1+\nu))$  is the elastic shear modulus,  $K$  is the elastic bulk modulus which is given by  $K = E/(1+2\nu)/(3(1-\nu^2))$  for plane stress, and  $K = E/(3(1-2\nu))$  for plain strain. In these expressions,  $E$  is Young's elastic modulus and  $\nu$  is the Poisson's ratio. Therefore,

$$|q(\mathbf{e}, \mathbf{e})| \leq \eta_\chi \|\mathbf{e}\|^2.$$

The computation of a bound for  $q(\mathbf{e}, \mathbf{e})$  is straightforward once a bound for the error in the energy norm  $\|\mathbf{e}\|$  has been obtained.

The error  $\mathbf{e} \in \mathcal{V}$  is the solution of the residual equation

$$a(\mathbf{e}, \mathbf{v}) = R^P(\mathbf{v}) \quad \forall \mathbf{v} \in \mathcal{V}, \quad (\text{C.6})$$

where

$$R^P(\mathbf{v}) = \int_{\Omega} \mathbf{f} \cdot \mathbf{v} \, d\Omega + \int_{\Gamma^N} \mathbf{g} \cdot \mathbf{v} \, d\Gamma - a(\mathbf{u}_H, \mathbf{v}).$$

The residual equation for the error (C.6) corresponds to  $\mathbf{f}^* = \mathbf{f}$ ,  $\mathbf{g}^* = \mathbf{g}$  and  $\mathbf{z}_H = \mathbf{u}_H$  in (5.1). Thus, one may apply the error estimation procedure presented in Chapter 5 to obtain a statically admissible stress field  $\boldsymbol{\sigma}^e \in \mathcal{S}$  verifying

$$\int_{\Omega} \boldsymbol{\sigma}^e : \boldsymbol{\varepsilon}(\mathbf{v}) \, d\Omega = R^P(\mathbf{v}) \quad \forall \mathbf{v} \in \mathcal{V}, \quad (\text{C.7})$$

and yielding the upper bound

$$\|\mathbf{e}\| \leq \|\|\sigma^e\|\|,$$

and therefore the bound for the quadratic term  $q(\mathbf{e}, \mathbf{e})$

$$Q = \eta_\chi \|\|\sigma^e\|\|^2.$$

### Bounds for the linear error term

In order to derive upper and lower bounds for the linear term  $q(\mathbf{e}, \mathbf{u}_H)$ , the procedure to obtain bounds for linear outputs of interest for symmetric model problems detailed in Chapter 2 may be considered.

Denote by  $s = q(\mathbf{e}, \mathbf{u}_H)$ . In order to find bounds for  $s$ , the following dual or adjoint problem is introduced: find  $\psi \in \mathcal{V}$  such that

$$a(\mathbf{v}, \psi) = q(\mathbf{v}, \mathbf{u}_H) \quad \forall \mathbf{v} \in \mathcal{V}. \quad (\text{C.8})$$

The finite element approximation of the dual problem is denoted by  $\psi_H \in \mathcal{V}_H$  and its associated error is  $\varepsilon = \psi - \psi_H \in \mathcal{V}$  solution of the residual problem

$$a(\mathbf{v}, \varepsilon) = q(\mathbf{v}, \mathbf{u}_H) - q(\psi, \mathbf{u}_H) = R^D(\mathbf{v}) \quad \forall \mathbf{v} \in \mathcal{V}.$$

Then, bounds for  $s$  may be obtained from equation (2.14) as

$$-\frac{1}{4} \|\kappa \mathbf{e} - \frac{1}{\kappa} \varepsilon\|_{\text{UB}}^2 \leq s \leq \frac{1}{4} \|\kappa \mathbf{e} + \frac{1}{\kappa} \varepsilon\|_{\text{UB}}^2,$$

and the problem reduces to find upper bounds for the linear combinations  $\kappa \mathbf{e} \pm \frac{1}{\kappa} \varepsilon$ .

Let  $\sigma^e \in \mathcal{S}$  be a statically admissible stress field for the primal residual problem verifying equation (C.7). Consider in an analogous form the statically admissible stress field for the dual residual problem  $\sigma^\varepsilon \in \mathcal{S}$  verifying

$$\int_{\Omega} \sigma^\varepsilon : \varepsilon(\mathbf{v}) \, d\Omega = R^D(\mathbf{v}) \quad \forall \mathbf{v} \in \mathcal{V}.$$

Then, the linear combination  $\kappa \sigma^e \pm \frac{1}{\kappa} \sigma^\varepsilon$  is a statically admissible stress field for the combined residual problem

$$a(\kappa \mathbf{e} \pm \frac{1}{\kappa} \varepsilon, \mathbf{v}) = \kappa R^P(\mathbf{v}) \pm \frac{1}{\kappa} R^D(\mathbf{v}) \quad \forall \mathbf{v} \in \mathcal{V},$$

and therefore

$$\|\kappa e \pm \frac{1}{\kappa} \varepsilon\| \leq \|\kappa \sigma^e \pm \frac{1}{\kappa} \sigma^e\|.$$

The bounds for the linear term are

$$L^- = -\frac{1}{4} \|\kappa \sigma^e - \frac{1}{\kappa} \sigma^e\|^2, \quad L^+ = \frac{1}{4} \|\kappa \sigma^e + \frac{1}{\kappa} \sigma^e\|^2,$$

and the value of the arbitrary parameter  $\kappa$  which optimizes the bounds is  $\kappa = \sqrt{\|\sigma^e\|/\|\sigma^e\|}$ .

The procedure to obtain the bounds for the  $J$ -integral is summarized in the box in Table C.2.

|   |
|---|
| <p>1.- Compute <math>\mathbf{u}_H</math> and <math>\psi_H \in \mathcal{V}_H</math> s.t.</p> $a(\mathbf{u}_H, \mathbf{v}) = \ell(\mathbf{v}) \quad \forall \mathbf{v} \in \mathcal{V}_H,$ $a(\mathbf{v}, \psi_H) = q(\mathbf{v}, \mathbf{u}_H) \quad \forall \mathbf{v} \in \mathcal{V}_H.$ <p>2.- Compute <math>\sigma^e</math> and <math>\sigma^\varepsilon \in \mathcal{S}</math> s.t.</p> $\int_{\Omega} \sigma^e : \varepsilon(\mathbf{v}) \, d\Omega = R^P(\mathbf{v}) \quad \forall \mathbf{v} \in \mathcal{V},$ $\int_{\Omega} \sigma^\varepsilon : \varepsilon(\mathbf{v}) \, d\Omega = R^D(\mathbf{v}) \quad \forall \mathbf{v} \in \mathcal{V}.$ <p>3.- Compute <math>\eta_\chi</math> as</p> $\eta_\chi = \max_{T_H \in \mathcal{T}_H} \frac{(3K + 4\mu) \nabla \chi ^2}{4\sqrt{(3K + \mu)\left(3\mu\left(\frac{\partial \chi}{\partial x_1}\right)^2 + (3K + 4\mu)\left(\frac{\partial \chi}{\partial x_2}\right)^2\right)}}.$ <p>4.- Compute <math>\kappa = \sqrt{\ \sigma^\varepsilon\ /\ \sigma^e\ }</math> and <math>J^-</math> and <math>J^+</math> as</p> $J^- = J(\mathbf{u}_H) - \eta_\chi \ \sigma^e\ ^2 - \frac{1}{2} \ \kappa \sigma^e - \frac{1}{\kappa} \sigma^\varepsilon\ ^2,$ $J^+ = J(\mathbf{u}_H) + \eta_\chi \ \sigma^e\ ^2 + \frac{1}{2} \ \kappa \sigma^e + \frac{1}{\kappa} \sigma^\varepsilon\ ^2.$ |
|---|

Figure C.2: Main steps of the strategy used to obtain upper bounds for the energy norm of the solution of a symmetric boundary value problem.

# Bibliography

- Ainsworth, M. and Oden, J. (1992), ‘A procedure for a posteriori error estimation for h-p finite element methods’, *Reliability in computational mechanics (Kraków, 1991). Computer Methods in Applied Mechanics and Engineering* **101**(1-3), 73–96.
- Ainsworth, M. and Oden, J. (1993), ‘A unified approach to a posteriori error estimation using element residual methods’, *Numerische Mathematic* **65**(1), 23–50.
- Ainsworth, M. and Oden, J. T. (2000), *A Posteriori Error Estimation in Finite Element Analysis*, Pure and Applied Mathematics. A Wiley-Interscience Series of Texts, Monographs, and Tracts, John Wiley and Sons, Inc.
- Babuška, I. and Miller, A. (1984), ‘The post-processing approach in the finite element method. Part 3. A posteriori error estimates and adaptive mesh selection’, *International Journal for Numerical Methods in Engineering* **20**(12), 2311–2324.
- Babuška, I. and Rheinboldt, W. C. (1978a), ‘Error estimates for adaptive finite element computations’, *SIAM Journal on Numerical Analysis* **15**(4), 736–754.
- Babuška, I. and Rheinboldt, W. C. (1978b), ‘A posteriori error estimates for the finite element method’, *International Journal for Numerical Methods in Engineering* **12**, 1597–1615.
- Babuška, I. and Rheinboldt, W. C. (1979), ‘Adaptive approaches and reliability estimations in finite element analysis’, *Computer Methods in Applied Mechanics and Engineering* **17/18**(part 3), 519–540.
- Babuška, I., Strouboulis, T. and Gangaraj, S. K. (1999), ‘Guaranteed computable bounds for the exact error in the finite element solution. Part I: One-dimensional model problem’, *Computer Methods in Applied Mechanics and Engineering* **176**(1-4), 51–79. New advances in computational methods (Cachan, 1997).

- 
- Bank, R. E. and Weiser, A. (1985), ‘Some a posteriori error estimators for elliptic partial differential equations’, *Mathematics of Computation* **44**(170), 283–301.
- Brenner, S. C. and Scott, L. R. (1994), *The Mathematical Theory of Finite Element Methods*, number 15 in ‘Texts in Applied Mathematics’, Springer-Verlag.
- Carstensen, C. and Funken, S. A. (1999/00), ‘Fully reliable localized error control in the FEM’, *SIAM Journal on Scientific Computing* **21**(4), 1465–1484 (electronic).
- Choi, H. W. and Paraschivoiu, M. (2004), ‘Adaptive computations of a posteriori finite element output bounds: a comparison of the ”hybrid-flux” approach and the ”flux-free” approach’, *Computer Methods in Applied Mechanics and Engineering* **193**(36-38), 4001–4033.
- Díez, P., Egozcue, J. J. and Huerta, A. (1998), ‘A posteriori error estimation for standard finite element analysis’, *Computer Methods in Applied Mechanics and Engineering* .
- Díez, P., Parés, N. and Huerta, A. (2003), ‘Recovering lower bounds of the error by postprocessing implicit residual a posteriori error estimates’, *International Journal for Numerical Methods in Engineering* **56**(10), 1465–1488.
- Huerta, A. and Díez, P. (2000), ‘Error estimation including pollution assessment for nonlinear finite element analysis’, *Computer Methods in Applied Mechanics and Engineering* **181**(1-3), 21–41.
- Kelly, D. (1984), ‘The self-equilibration of residuals and complementary a posteriori error estimates in the finite element method’, *International Journal for Numerical Methods in Engineering* **20**, 1491–1506.
- Ladevèze, P. (1977), Nouvelle procédure d’estimation d’erreur relative à la méthode des éléments finis et applications, in ‘Actes des Journées Eléments Finis’. Rennes, France.
- Ladevèze, P. and Leguillon, D. (1983), ‘Error estimate procedure in the finite element method and applications’, *SIAM Journal on Numerical Analysis* **20**(3), 485–509.
- Ladevèze, P. and Maunder, E. A. W. (1996), ‘A general method for recovering equilibrating element tractions’, *Computer Methods in Applied Mechanics and Engineering* **137**(12), 111–151.

- Larsson, F., Hansbo, P. and Runesson, K. (2002), ‘Strategies for computing goal-oriented a posteriori error measures in non-linear elasticity’, *International Journal for Numerical Methods in Engineering* **55**(12), 879–894.
- Li, F. Z., Shih, C. F. and Needleman, A. (1985), ‘A comparison of methods for calculating energy release rates’, *Engineering Fracture Mechanics* **21**(2), 405–421.
- Machiels, L., Maday, Y. and Patera, A. T. (2000), ‘A “flux-free” nodal Neumann subproblem approach to output bounds for partial differential equations’, *Comptes Rendus des Séances de l’Académie des Sciences. Série I. Mathématique* **330**(3), 249–254.
- Maday, Y., Patera, A. T. and Peraire, J. (1999), ‘A general formulation for a posteriori bounds for output functionals of partial differential equations; application to the eigenvalue problem’, *Comptes Rendus des Séances de l’Académie des Sciences. Série I. Mathématique* **328**(9), 823–828.
- Morin, P., Nochetto, R. H. and Siebert, K. G. (2003), ‘Local problems on stars: a posteriori error estimators, convergence, and performance’, *Mathematics of Computation* **72**(243), 1067–1097.
- Oden, J. T. and Prudhomme, S. (1999), Goal-oriented error estimation and adaptivity for the finite element method, report 99-15, TICAM.
- Oden, J. T. and Prudhomme, S. (2001), ‘Goal-oriented error estimation and adaptivity for the finite element method’, *Computers & Mathematics with Applications. An International Journal* **41**(5-6), 735–756.
- Paraschivoiu, M., Peraire, J. and Patera, A. T. (1997), ‘A posteriori finite element bounds for linear-functional outputs of elliptic partial differential equations’, *Computer Methods in Applied Mechanics and Engineering* **150**(1-4), 289–312. Symposium on Advances in Computational Mechanics, Vol. 2 (Austin, TX, 1997).
- Parés, N., Bonet, J., Huerta, A. and Peraire, J. (2005), ‘The computation of bounds for linear-functional outputs of weak solutions to the two-dimensional elasticity equations’, *Computer Methods in Applied Mechanics and Engineering* **194**. In press.
- Parés, N., Díez, P. and Huerta, A. (2005), ‘Subdomain-based flux-free a posteriori error estimators’, *Computer Methods in Applied Mechanics and Engineering* **194**. In press.

- Patera, A. T. and Peraire, J. (2003), A general Lagrangian formulation for the computation of a posteriori finite element bounds, in 'Error estimation and adaptive discretization methods in computational fluid dynamics', Vol. 25 of *Lecture Notes in Computational Science and Engineering*, Springer, Berlin, pp. 159–206.
- Peraire, J. and Patera, A. T. (1997), Bounds for linear-functional outputs of coercive partial differential equations: Local indicators and adaptive refinement, in 'Proceedings of the Workshop On New Advances in Adaptive Computational Methods in Mechanics', P. Ladeveze and J. Oden, eds., Elsevier, Cachan, September 17-19, 1997.
- Prudhomme, S., Nobile, F., Chamoin, L. and Oden, J. T. (2004), 'Analysis of a subdomain-based error estimator for finite element approximations of elliptic problems', *Numerical Methods for Partial Differential Equations* **20**(2), 165–192.
- Prudhomme, S. and Oden, J. T. (1999), 'On goal-oriented error estimation for elliptic problems: application to the control of pointwise errors', *Computer Methods in Applied Mechanics and Engineering* **176**(1-4), 313–331. New advances in computational methods (Cachan, 1997).
- Prudhomme, S., Oden, J. T., Westermann, T., Bass, J. and Botkin, M. E. (2003), 'Practical methods for a posteriori error estimation in engineering applications', *International Journal for Numerical Methods in Engineering* **56**(8), 1193–1224.
- Rice, J. (1968), 'A path independent integral and approximate analysis of strain concentration by notches and cracks', *Journal of Applied Mechanics* **35**, 379–386.
- Sarrate, J., Peraire, J. and Patera, A. T. (1999), 'A posteriori finite element error bounds for non-linear outputs of the helmholtz equation', *International Journal for Numerical Methods in Fluids* **31**(1), 17–36.
- Sauer-Budge, A. M., Bonet, J., Huerta, A. and Peraire, J. (2004), 'Computing bounds for linear functionals of exact weak solutions to poisson's equation', *SIAM Journal on Numerical Analysis* **42**(4), 1610–1630.
- Sauer-Budge, A. M. and Peraire, J. (2004), 'Computing bounds for linear functionals of exact weak solutions to the advection-diffusion-reaction equation', *SIAM Journal on Scientific Computing* **26**(2), 636–652.
- Strouboulis, T., Babuška, I. and Gangaraj, S. K. (2000), 'Guaranteed computable bounds for the exact error in the finite element solution. Part II: Bounds for



---

the energy norm of the error in two dimensions', *International Journal for Numerical Methods in Engineering* **47**(1-3), 427–475. Richard H. Gallagher Memorial Issue.

Xuan, Z. C., Lee, K. H., Patera, A. T. and Peraire, J. (2004), Computing upper and lower bounds for the  $j$ -integral in two-dimensional linear elasticity, in 'Proceedings of the SMA symposium. Singapore'.

Xuan, Z. C., Parés, N. and Peraire, J. (2005), 'Computing upper and lower bounds for the  $j$ -integral in two-dimensional linear elasticity', *Computer Methods in Applied Mechanics and Engineering* **194**. In press.



# Recovering lower bounds of the error by postprocessing implicit residual a posteriori error estimates

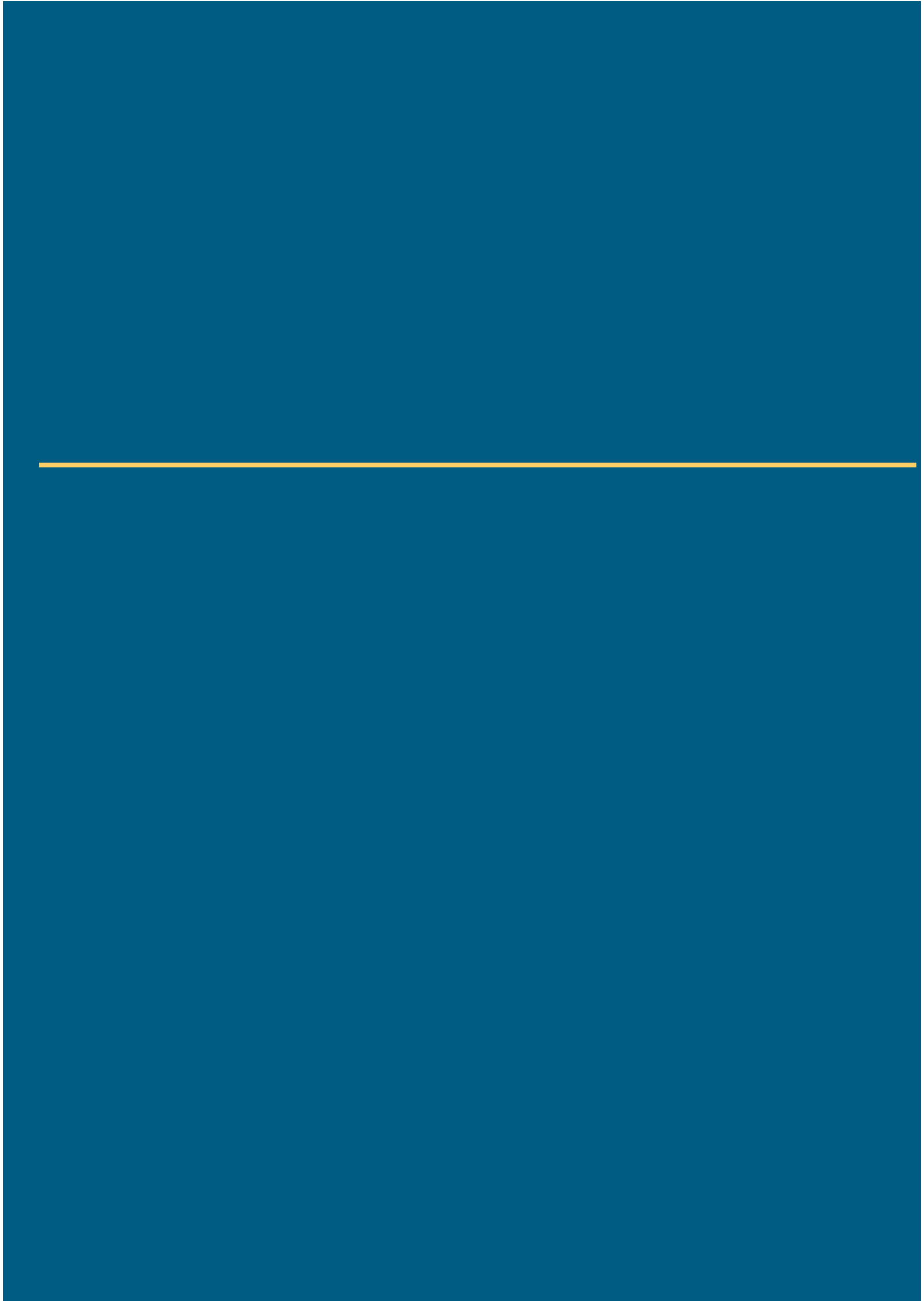
Díez P., Parés N. and Huerta A.

---

*International Journal for Numerical*

*Methods in Engineering*

**56** (10), 1465-1488



## Recovering lower bounds of the error by postprocessing implicit residual *a posteriori* error estimates

Pedro Díez, Núria Parés and Antonio Huerta<sup>\*,†</sup>

*Departament de Matemàtica Aplicada III, ETSECCPB, Universidad Politècnica de Catalunya, Barcelona, Spain*

### SUMMARY

Classical residual type error estimators approximate the error flux around the elements and yield upper bounds of the exact (or reference) error. Lower bounds of the error are also needed in goal oriented adaptivity and for bounds on functional outputs. This work introduces a simple and cheap strategy to recover a lower bound estimate from standard upper bound estimates. This lower bound may also be used to assess the effectivity of the former estimate and to improve it. Copyright © 2003 John Wiley & Sons, Ltd.

KEY WORDS: implicit residual type error estimator; upper and lower bounds; quality assessment

### 1. INTRODUCTION

Implicit residual-type error estimators require to set proper boundary conditions for the local (usually element by element) error equations. If these boundary conditions are of Neumann type [1, 2] the obtained estimates are upper bounds of the error. The error estimators based on the error in the constitutive relation introduced by Ladevèze [3, 4] may also be classified in this group and also overestimate the error. The selection of the flux on the interelement edges may use either a trivial flux averaging [1] or a more sophisticated recovering technique yielding equilibrated residuals [2, 3]. The equilibrated residual strategies are expected to furnish more realistic boundary conditions for the local problems and, consequently, to yield better error estimates.

On the other hand, residual-type error estimators using Dirichlet boundary conditions in the local error equations [5, 6] yield lower bounds of the error. Basically, the lower bound property is induced by the continuity of the obtained estimate.

---

\*Correspondence to: Antonio Huerta, Departament de Matemàtica Aplicada III, ETSECCPB, Universidad Politècnica de Catalunya, Campus Norte UPC, Mòdulo C-2, Jordi Giruna 1, E-08034 Barcelona, Spain

†E-mail: antonio.huerta@upc.es

Contract/grant sponsor: Ministerio de Educación y Cultura; contract/grant number: TAP98-0421

*Received 5 November 2001*

*Revised 6 May 2002*

*Accepted 6 June 2002*

The comparison of these two approaches suggest the idea of postprocessing residual-type error estimators yielding upper bound, enforcing continuity and obtaining a lower bound of the error with a small supplementary effort.

The idea of obtaining a couple of upper and lower bound estimates at the same time is also suggested by the goal oriented adaptive strategies [7, 8]. Indeed, in the context of symmetric (self-adjoint) problems, these strategies require both a lower and an upper bound of the error in the standard energy norm to assess the error in an output of interest. However, the approach introduced in Reference [8] allows also to obtain upper and lower bounds for functional outputs of non-symmetric problems.

The approach presented here is based on the postprocessing of the upper bound estimate  $e_{\text{est}}$ , which is discontinuous. The postprocessing introduces a correction  $e_{\text{cor}}$  such that the corrected error distribution,  $e_{\text{cont}} := e_{\text{est}} + e_{\text{cor}}$ , is continuous. Thus, the correction  $e_{\text{cor}}$  must compensate the discontinuities of  $e_{\text{est}}$ . Then, a lower bound is computed straightforward using  $e_{\text{est}}$  and  $e_{\text{cor}}$ .

The remainder of the paper is structured as follows. The model problem is stated in Section 2. Section 3 is devoted to introduce the local and global versions of error equation, and the reference error. In Section 4, the residual-type error estimators approximating the local flux are described. The upper bound property of this kind of estimators is easily proved. Attention is paid to the solvability problems of the pure diffusion case. Then, in Section 5, the estimate  $e_{\text{est}}$  yielding an upper bound is corrected to enforce its continuity and a lower bound is recovered. Also at this point, some additional effort must be done to deal with the pure diffusion case, where the original estimate is locally determined up to a constant. These local constants do not affect the norm of  $e_{\text{est}}$  but do condition  $e_{\text{cor}}$  and, consequently, in order to have an optimal correction, it is worthy to select them properly. Numerical examples demonstrating the good behaviour of the proposed strategy are shown in Section 6.

## 2. STATEMENT OF THE PROBLEM

### 2.1. Model problem

Let us consider the following linear Neumann boundary value problem in an open, bounded domain  $\Omega \subset \mathbb{R}^2$

$$\left. \begin{aligned} -\nabla \cdot (a \nabla u) + bu &= s && \text{in } \Omega \\ a \nabla u \cdot n &= g_N && \text{on } \partial \Omega \end{aligned} \right\} \quad (1)$$

In order to simplify the presentation, the boundary conditions are assumed to be only of Neumann type. Accounting for Dirichlet or mixed boundary conditions does not introduce any additional conceptual difficulty. Moreover, in order to ensure ellipticity, it is assumed that

$$\begin{aligned} 0 < \underline{a} &\leq a(x) \leq \bar{a} \\ 0 &\leq \underline{b} \leq b(x) \leq \bar{b} \end{aligned}$$

for some  $\underline{a}$ ,  $\bar{a}$ ,  $\underline{b}$  and  $\bar{b}$ .

The weak form of this problem reads: find  $u \in H^1(\Omega)$  such that

$$a(u, v) = \int_{\Omega} sv \, d\Omega + \int_{\partial\Omega} g_N v \, d\Gamma \quad \forall v \in H^1(\Omega) \quad (2)$$

where

$$a(u, v) := \int_{\Omega} (a \nabla u \cdot \nabla v + buv) \, d\Omega$$

and  $H^1(\Omega)$  stands for the standard Sobolev space.

The Galerkin finite element method provides an approximation  $u_h$  to  $u$ , lying in a finite-dimensional space  $V_h \subset H^1(\Omega)$  and verifying

$$a(u_h, v) = \int_{\Omega} sv \, d\Omega + \int_{\partial\Omega} g_N v \, d\Gamma \quad \forall v \in V_h \quad (3)$$

The finite-dimensional space  $V_h$  is associated with a finite element mesh of characteristic size  $h$ . The degree of the complete polynomials used in the interpolation of  $V_h$  is denoted by  $p$ . The geometric support of the elements of this mesh are open subdomains denoted by  $\Omega_k$ ,  $k = 1, \dots, n_{\text{elem}}$ . It is assumed that  $\bar{\Omega} = \bigcup_k \bar{\Omega}_k$  (the mesh covers the whole domain) and  $\Omega_k \cap \Omega_l = \emptyset$  for  $k \neq l$  (different elements have in common, at most, part of their boundary). The derivation of *a priori* estimates requires further regularity conditions for the mesh. The precise assumptions on the meshes may be found in Reference [9, Section 1.3.3].

The goal of *a posteriori* error estimation is to assess the accuracy of the approximate solution  $u_h$ , that is, to evaluate and measure the error,  $e := u - u_h$ , or an approximation to it. The error is measured using some functional norm. One of the most popular options is the energy norm induced by  $a(\cdot, \cdot)$ :

$$\|e\| := [a(e, e)]^{1/2} \quad (4)$$

Local restrictions of the norm are needed to describe the spatial distribution of the error. In the following, the restriction of  $a(\cdot, \cdot)$  to the element  $\Omega_k$  ( $k = 1, \dots, n_{\text{elem}}$ ) is denoted by  $a_k(\cdot, \cdot)$ . Thus, the restriction of  $\|\cdot\|$  to  $\Omega_k$ ,  $\|\cdot\|_k$ , is induced by  $a_k(\cdot, \cdot)$ . In order to describe the spatial distribution of the error, the value of  $\|e\|_k$  in each element is estimated.

## 2.2. Error equations and reference error

The global equation for the error is recovered from Equation (2), replacing  $u$  by  $u_h + e$ :

$$a(e, v) = \int_{\Omega} sv \, d\Omega + \int_{\partial\Omega} g_N v \, d\Gamma - a(u_h, v) =: \mathbf{R}(v) \quad \forall v \in H^1(\Omega) \quad (5)$$

The r.h.s. term of Equation (5),  $\mathbf{R}(v)$ , is the weak residual associated with the approximate solution  $u_h$ .

The local counterpart of Equation (5) is derived integrating the weighted residual of the strong form, Equation (1), in  $\Omega_k$ . It reads,

$$a_k(e, v) = \mathbf{R}_k(v) + \int_{\partial\Omega_k \cap \Omega} a \nabla u \cdot n v \, d\Gamma \quad \forall v \in H^1(\Omega_k) \quad (6)$$

where  $\mathbf{R}_k(v)$  is the restriction of  $\mathbf{R}(v)$  to  $\Omega_k$ :

$$\mathbf{R}_k(v) := \int_{\Omega_k} sv \, d\Omega + \int_{\partial\Omega_k \cap \partial\Omega} g_N v \, d\Gamma - a_k(u_h, v) \quad (7)$$

Note that the last term of the r.h.s. of Equation (6) accounts for the unknown flux on the interelement edges. In other words, the boundary conditions of the local problem are not known.

The error is estimated approximating the solution of the local error equation (6). The characterization of any residual-type error estimator requires to select both:

- the finite-dimensional space where the local error equation is solved (local  $h$ - or  $p$ -refinement) and
- the unknown boundary conditions for the local problems.

The first point is related with the concept of reference error. Residual *a posteriori* error estimation techniques are based on assessing and bounding the reference error and not the error itself. For all practical purposes, the exact value of the error,  $e$ , is replaced by a reference (or ‘truth’) error,  $e_{\text{ref}}$ , lying in a finite-dimensional space much refined with respect to the computational space  $V_h$ . Let us denote by  $V^{\text{ref}}$  this refined space.  $V^{\text{ref}}$  is generated either as a  $h$  or  $p$ -refinement of  $V_h$ . That is, denoting by  $\tilde{h}$  and  $\tilde{p}$  the characteristic element size and the degree of interpolation of the elements generating  $V^{\text{ref}}$ , either  $\tilde{h} \ll h$  or  $\tilde{p} \gg p$  holds.

Thus, the reference error,  $e_{\text{ref}} \in V^{\text{ref}}$ , verifies the discrete form of Equation (5), that is

$$a(e_{\text{ref}}, v) = \mathbf{R}(v) \quad \forall v \in V^{\text{ref}} \quad (8)$$

The direct computation of  $e_{\text{ref}}$  is computationally unaffordable because it requires to solve a system of equations with the number of degrees of freedom equal to the dimension of  $V^{\text{ref}}$ .

The fact of using a reference error (that is, replacing the continuous space  $H^1(\Omega)$  by the refined space  $V^{\text{ref}}$ , and the exact error  $e$  by the reference error  $e_{\text{ref}}$ ) does not introduce a significant loss of accuracy in the error estimation procedure. Consequently, the quality of a residual-type error estimation procedure depends essentially on the approximation of the local boundary conditions.

### 3. STANDARD RESIDUAL-TYPE ERROR ESTIMATES

Standard residual-type error estimators [1–3] solve the local error equation (6) using approximated Neumann boundary conditions. The values of the flux  $a\nabla u \cdot n|_{\partial\Omega_k \cap \Omega}$ , see Equation (6), are determined or approximated along the boundary of each element  $\Omega_k$ . This section is devoted to briefly describe this kind of estimators and to recall the proof of their upper bound property.

#### 3.1. Approximation of fluxes

The approximation of the flux is based on smoothing the approximate flux  $a\nabla u_h \cdot n$ , which is discontinuous. The basic idea due to Bank and Weiser [1] is to average the approximate flux on every interelement edge. Let  $\Gamma_m$ , for  $m = 1, \dots, n_{\text{int}}$ , be the interelement edges of the mesh.



That is, for every  $m \in \{1, \dots, n_{\text{int}}\}$  they exist  $k, l \in \{1, \dots, n_{\text{elem}}\}$ ,  $k \neq l$ , such that  $\Gamma_m = \overline{\Omega_k} \cap \overline{\Omega_l}$ . Then

$$a \nabla u|_{\Gamma_m} \simeq [a \nabla u_h]_A := \frac{1}{2}(a \nabla u_h|_{\partial \Omega_l} + a \nabla u_h|_{\partial \Omega_k}) \quad \text{for } m = 1, \dots, n_{\text{int}} \tag{9}$$

where  $[\cdot]_A$  stands for the average on  $\Gamma_m$ . The approximation given in Equation (9) is used in Equation (6).

More sophisticated flux averaging procedures are used by other authors [2, 3] in order to obtain equilibrated local problems. They improve the efficiency of the estimator. Here, the simplest averaging is used for illustration purposes. In fact, the following developments are also valid for these approaches: it suffices to use a more complicated definition for the average  $[a \nabla u_h]_A$ .

### 3.2. Discrete local residual equation

Thus, the error estimate  $e_{\text{est}}$  is computed locally by solving the following problem: find  $e_{\text{est}} \in V_k^{\text{ref}}$  such that

$$a_k(e_{\text{est}}, v) = \mathbf{R}_k(v) + \int_{\partial \Omega_k \cap \Omega} [a \nabla u_h]_A \cdot n v \, d\Gamma \quad \forall v \in V_k^{\text{ref}} \tag{10}$$

where  $V_k^{\text{ref}}$  is the restriction of  $V^{\text{ref}}$  to  $\Omega_k$ , that is

$$V_k^{\text{ref}} := \{v \in H^1(\Omega_k) / \exists \tilde{v} \in V^{\text{ref}}, v = \tilde{v}|_{\Omega_k}\} \tag{11}$$

Equation (10) is the discrete version of Equation (6) using the approximation given by Equation (9).

Note that the sum of the spaces  $V_k^{\text{ref}}$  is not equal to  $V^{\text{ref}}$ . In fact,  $V_{\text{brok}}^{\text{ref}} := \bigoplus_k V_k^{\text{ref}}$  is a space of ‘broken’ functions. In order to recover  $V^{\text{ref}}$  it is necessary to restrict the space forcing the continuity:  $V^{\text{ref}} = V_{\text{brok}}^{\text{ref}} \cap \mathcal{C}^0$ .

A global equation for the error estimate  $e_{\text{est}}$  is found summing up Equation (10) for all  $k$  ( $k = 1, \dots, n_{\text{elem}}$ ),

$$a(e_{\text{est}}, v) = \mathbf{R}(v) + \sum_{m=1}^{n_{\text{int}}} \int_{\Gamma_m} [a \nabla u_h]_A \cdot [vn]_J \, d\Gamma \quad \forall v \in V_{\text{brok}}^{\text{ref}} \tag{12}$$

where  $[vn]_J$  stands for the jump of  $vn$  across  $\Gamma_m = \overline{\Omega_k} \cap \overline{\Omega_l}$ , that is,

$$[vn]_J := (v|_{\Omega_k})n_k + (v|_{\Omega_l})n_l \tag{13}$$

being  $n_k = -n_l$  the corresponding outward normal unit vectors. The recovered flux, see Section 3.1, is said to be consistent if the approximation of the flux is continuous, i.e. if the approximation of  $a \nabla u|_{\Gamma_m}$  is the same viewed from  $\Omega_k$  and from  $\Omega_l$ . In order to derive Equation (12) it is necessary that the recovered fluxes are consistent.

Furthermore, if the test functions are continuous, i.e. if  $v$  is in  $V^{\text{ref}} \subset V_{\text{brok}}^{\text{ref}}$ , then  $[vn]_J = 0$  and from Equation (12) one gets

$$a(e_{\text{est}}, v) = \mathbf{R}(v) \quad \forall v \in V^{\text{ref}}, \text{ where still } e_{\text{est}} \in V_{\text{brok}}^{\text{ref}} \tag{14}$$

In other words, if the consistency condition is satisfied, the interelement edges are not a source of flux in the global error equation (for  $v$  continuous). In the following, some properties of the estimate  $e_{\text{est}}$  are derived replacing  $v$  in Equation (14) by particular functions in  $V^{\text{ref}}$ .

*Remark 1*

In Equation (12), the definition of  $a(\cdot, \cdot)$  must be generalized to accept ‘broken’ functions in the arguments. Thus, for  $v, w \in V_{\text{brok}}^{\text{ref}}$ ,

$$a(w, v) := \sum_{k=1}^{n_{\text{elem}}} a_k(w, v) \quad (15)$$

Of course, this generalized definition coincides with the standard one when the arguments are in  $H^1(\Omega)$ .

*3.3. Upper bound property*

The consistency condition implies that the error estimates computed using Equation (10) are upper bounds of the reference error. Although this is a well-known property of this kind of estimators, the corresponding theorem is revisited and proved here because it is important in the following.

*Theorem 1*

The error estimate  $e_{\text{est}}$  computed solving Equation (10) yields an upper bound of the error, that is

$$\varepsilon_{\text{upp}} := \|e_{\text{est}}\|^2 \geq \|e_{\text{ref}}\|^2 \quad (16)$$

*Proof*

Taking  $v = e_{\text{ref}}$  in Equations (14) and (8) it follows that

$$a(e_{\text{est}}, e_{\text{ref}}) = a(e_{\text{ref}}, e_{\text{ref}}) \quad (17)$$

Then, the proof is completed by the following algebraic manipulation.

$$\begin{aligned} 0 \leq a(e_{\text{ref}} - e_{\text{est}}, e_{\text{ref}} - e_{\text{est}}) &= a(e_{\text{ref}}, e_{\text{ref}}) + a(e_{\text{est}}, e_{\text{est}}) - 2 \overbrace{a(e_{\text{est}}, e_{\text{ref}})}^{=a(e_{\text{ref}}, e_{\text{ref}})} \\ &= a(e_{\text{est}}, e_{\text{est}}) - a(e_{\text{ref}}, e_{\text{ref}}) \quad \square \end{aligned}$$

*Remark 2*

It is worth noting that the upper bound  $\varepsilon_{\text{upp}}$  is defined in Equation (16) as the squared norm of the error estimate. This is because the use of squared norms simplifies the presentation. Thus, in the following, the estimates of the squared error norms, approximations of  $\|e_{\text{ref}}\|^2$ , are denoted by  $\varepsilon_{\star}$ .

*Remark 3*

In the general case,  $e_{\text{est}}$  is not continuous (it is in  $V_{\text{brok}}^{\text{ref}}$  but not in  $V^{\text{ref}}$ ). Thus, in general, it is not possible to take  $v = e_{\text{est}}$  in Equation (14). However, if a particular choice of the boundary conditions of the local problems leads to a continuous estimate  $e_{\text{est}}$ , then it can be

easily shown that  $a(e_{\text{est}}, e_{\text{est}}) \leq a(e_{\text{ref}}, e_{\text{ref}})$  and, consequently,  $a(e_{\text{est}}, e_{\text{est}}) = a(e_{\text{ref}}, e_{\text{ref}})$ . That is, the choice of the Neumann boundary conditions giving a continuous estimate is optimal.

### 3.4. Solvability problems when $b = 0$

If the reaction term vanishes in Equation (1) ( $b = 0$ ), the solvability of the local Neumann problem, Equation (10), requires proper data ensuring equilibrium. It is well known that if the source term  $s$  (body load) is not equilibrated by the prescribed boundary flux, the Neumann problem does not have any solution. Locally (in element  $\Omega_k$ ), the equilibrium condition reads

$$\int_{\Omega_k} s \, d\Omega + \int_{\partial\Omega_k \cap \partial\Omega} g_N \, d\Gamma + \int_{\partial\Omega_k \cap \Omega} [a\nabla u_h]_A \cdot n \, d\Gamma = 0 \quad (18)$$

The simple averaging described in Equation (9) does not enforce the equilibrium condition.

Two different strategies may be used in order to ensure the solvability of the local problems. A first option is to use approximation of fluxes yielding equilibrated local problems.

The second strategy is to restrict the set of admissible functions in the local problem eliminating from the local interpolation space the kernel of the l.h.s. of Equation (10). In fact, the second and third estimators introduced by Bank and Weiser in Reference [1] use this strategy. These estimators are used in the numerical examples and are they denoted by  $e_2$  and  $e_3$ , respectively.

#### Remark 4

The description of these estimators requires to introduce the hierarchical decomposition of  $V^{\text{ref}}$ ,  $V^{\text{ref}} = V_h \oplus V^{\text{com}}$ , where  $V^{\text{com}}$  is the hierarchical complement of  $V_h$  in  $V^{\text{ref}}$ . The space  $V^{\text{com}}$  contains the functions  $v$  of  $V^{\text{ref}}$  such that the degrees of freedom (nodal values) of  $v$  corresponding to  $V_h$  are null. Typically, for  $p$ -refinement, the functions of  $V^{\text{com}}$  are of the bubble type. Then, for all  $v \in V^{\text{ref}}$ ,  $\exists! v_h \in V_h$  and  $\exists! v_{\text{com}} \in V^{\text{com}}$  such that  $v = v_h + v_{\text{com}}$ . Thus, the nodal projection from  $V^{\text{ref}}$  to  $V_h$ ,  $\mathcal{I} : V^{\text{ref}} \rightarrow V_h$  is defined such that  $\mathcal{I}(v) = v_h$ .

The second estimator,  $e_2$  is then computed as the solution of the following local problem:

$$a_k(e_2, v) = \mathbf{R}_k(v - \mathcal{I}(v)) + \int_{\partial\Omega_k \cap \Omega} [a\nabla u_h]_A \cdot n(v - \mathcal{I}(v)) \, d\Gamma \quad \forall v \in V_k^{\text{ref}} \quad (19)$$

where the restriction of  $e_2$  to  $\Omega_k$  is in  $V_k^{\text{ref}}$  and, therefore, the global  $e_2$  is in  $V_{\text{brok}}^{\text{ref}}$ .

The third estimator,  $e_3$ , is locally computed as the solution of

$$a_k(e_3, v) = \mathbf{R}_k(v) + \int_{\partial\Omega_k \cap \Omega} [a\nabla u_h]_A \cdot n v \, d\Gamma \quad \forall v \in V_k^{\text{com}} \quad (20)$$

where the local restriction of  $V^{\text{com}}$ ,  $V_k^{\text{com}}$ , must be understood in the same sense as in Equation (11).

It is worth noting that  $e_2$  is an upper bound for the reference error but  $e_3$  is not. Indeed, summing up the local Equation (19) on  $k$  one gets a global equation for  $e_2$  where  $v$  ranges on  $V_{\text{brok}}^{\text{ref}}$  and the same rationale given for  $e_{\text{est}}$ , see Theorem 1, can be followed to deduce that  $\|e_2\| \geq \|e_{\text{ref}}\|$ . On the contrary, in the global equation corresponding to Equation (20),  $v$  ranges on  $V_{\text{brok}}^{\text{com}}$ . The upper bound property cannot be deduced in this case because  $V^{\text{ref}} \not\subset V_{\text{brok}}^{\text{com}}$ . However, in the asymptotic range, that is for  $h$  small enough, numerical evidence shows that  $e_3$  behaves also as an upper bound.

## 4. CORRECTION AND LOWER BOUND RECOVERING

In the previous section, see Remark 3, it has been noted that the overestimation of the error is associated with the continuity defaults of the estimate  $e_{\text{est}}$ . In fact, it has been observed that if the flux splitting is such that  $e_{\text{est}}$  is continuous, then the estimate  $e_{\text{est}}$  is optimal. Thus, the idea developed in this section is to introduce a correction of the error estimate in order to enforce its continuity. This correction allows to deduce a lower bound of the reference (and exact) error and, hence, to assess the effectivity of the original error estimate.

## 4.1. Correction and lower bound

Babuška and co-workers originally proposed to obtain a lower bound  $\varepsilon_{\text{low}}$  from a continuous corrected estimate [10, 11]. Here, the evaluation of the lower bound is improved by defining a scalar parametric family  $\varepsilon_{\text{low}}(\lambda)$ . Moreover, it is proved in this section that an optimal value of  $\lambda$ ,  $\lambda_{\text{opt}}$  exists and that it can be easily evaluated. Note that the optimal estimate,  $\varepsilon_{\text{low}}(\lambda_{\text{opt}})$  corresponds to the expression proposed in Reference [12], where the optimality of this choice is not mentioned.

Recall that  $e_{\text{est}} \in V_{\text{brok}}^{\text{ref}}$ , that is  $e_{\text{est}}$  is, in general, not continuous. Let  $e_{\text{cor}} \in V_{\text{brok}}^{\text{ref}}$  be a correction of  $e_{\text{est}}$  such that

$$e_{\text{cont}} := e_{\text{est}} + e_{\text{cor}} \in V^{\text{ref}} \quad (21)$$

that is, such that the corrected error  $e_{\text{cont}}$  is continuous.

Given a corrected estimate  $e_{\text{cont}}$ , a parametric family of lower bound estimates is found.

*Theorem 2*

Let  $e_{\text{est}}$  be an error estimate verifying the hypothesis of Theorem 1 and, therefore, being an upper bound of the reference error. Let  $e_{\text{cont}}$  be a corrected estimate as described in Equation (21). Then, for any scalar  $\lambda \in \mathbb{R}$ , the expression

$$\varepsilon_{\text{low}}(\lambda) := 2\lambda a(e_{\text{est}}, e_{\text{cont}}) - \lambda^2 \|e_{\text{cont}}\|^2 \quad (22)$$

is a lower bound of the reference error norm, that is,

$$\varepsilon_{\text{low}}(\lambda) \leq \|e_{\text{ref}}\|^2 \quad (23)$$

*Proof*

Since  $e_{\text{cont}}$  is continuous, it is possible to replace  $v$  by  $e_{\text{cont}}$  in Equations (14) and (8). That is,

$$a(e_{\text{est}}, e_{\text{cont}}) = a(e_{\text{ref}}, e_{\text{cont}}) \quad (24)$$

Then, using Equation (24), inequality (23) is proved considering the following algebraic manipulation:

$$\begin{aligned} 0 \leq a(e_{\text{ref}} - \lambda e_{\text{cont}}, e_{\text{ref}} - \lambda e_{\text{cont}}) &= a(e_{\text{ref}}, e_{\text{ref}}) + \lambda^2 a(e_{\text{cont}}, e_{\text{cont}}) - 2\lambda a(e_{\text{ref}}, e_{\text{cont}}) \\ &= \|e_{\text{ref}}\|^2 + \lambda^2 \|e_{\text{cont}}\|^2 - 2\lambda a(e_{\text{est}}, e_{\text{cont}}) \\ &= \|e_{\text{ref}}\|^2 - \varepsilon_{\text{low}}(\lambda) \quad \square \end{aligned}$$

Thus, once the corrected estimate  $e_{\text{cont}}$  is obtained, a lower bound of the error is recovered computing  $\varepsilon_{\text{low}}(\lambda)$ , for any value of  $\lambda$ . The natural choice,  $\lambda = 1$ , see References [10, 11, 13, 15], results in

$$\varepsilon_{\text{low}}(1) = 2a(e_{\text{est}}, e_{\text{cont}}) - \|e_{\text{cont}}\|^2 = \|e_{\text{est}}\|^2 - \|e_{\text{cor}}\|^2 \quad (25)$$

which in practice only requires the extra computation of  $\|e_{\text{cor}}\|$ .

However, the optimal choice for  $\lambda$  is the value that maximizes the lower bound  $\varepsilon_{\text{low}}(\lambda)$ . It is obvious from Equation (22) that this optimal value is

$$\lambda_{\text{opt}} = \frac{a(e_{\text{est}}, e_{\text{cont}})}{\|e_{\text{cont}}\|^2} \quad (26)$$

Consequently, given an upper bound estimate  $e_{\text{est}}$ , the optimal lower bound associated with a corrected estimate  $e_{\text{cont}}$  is

$$\varepsilon_{\text{low}}^{\text{opt}} := \varepsilon_{\text{low}}(\lambda_{\text{opt}}) = \frac{a(e_{\text{est}}, e_{\text{cont}})^2}{\|e_{\text{cont}}\|^2} \quad (27)$$

This is, in fact, the expression adopted in Reference [12].

*Remark 5*

Both  $\varepsilon_{\text{low}}^{\text{opt}}$  and  $\varepsilon_{\text{low}}(1)$  are exact if the recovering technique to obtain the corrected estimate  $e_{\text{cont}}$  is optimal. Indeed, if the corrected estimate coincides with the reference error, that is  $e_{\text{cont}} = e_{\text{ref}}$ , then

$$\varepsilon_{\text{low}}^{\text{opt}} = \varepsilon_{\text{low}}(1) = \|e_{\text{ref}}\|^2$$

Thus, both the lower bounds given by Equations (25) and (27) are sharp provided that the determination of the corrected estimate  $e_{\text{cont}}$  is accurate. In fact, the strategy used to obtain  $e_{\text{cont}}$  is oriented to enforce  $e_{\text{cont}} \approx e_{\text{ref}}$ .

Obviously, given  $e_{\text{cont}}$ , the estimate  $\varepsilon_{\text{low}}^{\text{opt}}$  is sharper than  $\varepsilon_{\text{low}}(1)$ . Consequently, once  $e_{\text{cont}}$  is determined,  $\varepsilon_{\text{low}}^{\text{opt}}$  is used to evaluate the lower bound. Nevertheless, in order to set a criterion for the determination of  $e_{\text{cont}}$ , the expression of  $\varepsilon_{\text{low}}(1)$ , Equation (25), is preferred to the expression of  $\varepsilon_{\text{low}}^{\text{opt}}$ , Equation (27). This is detailed in the next section.

#### 4.2. Determination of the corrected estimate $e_{\text{cont}}$

This section describes the smoothing process that builds up the corrected estimate  $e_{\text{cont}}$ . The degrees of freedom of the original estimate,  $e_{\text{est}}$ , affecting the continuity (associated with edges and corners) are simply averaged. This part of the smoothing process is standard [10, 11]. Here, the remaining degrees of freedom affecting the interior of the elements (bubble functions inside the elements) are set using an objective optimality criterion. The presentation is based in the formulation of the parametric family of scalar lower bounds,  $\varepsilon_{\text{low}}(\lambda)$  introduced in Section 4.1.

The correction  $e_{\text{cor}}$  and, consequently, the corrected estimate  $e_{\text{cont}}$  and the corresponding lower bound  $\varepsilon_{\text{low}}^{\text{opt}}$  are not unique. Any function  $e_{\text{cont}} \in \mathcal{V}^{\text{ref}}$  produces a lower bound  $\varepsilon_{\text{low}}^{\text{opt}}$ . However, as noted in Remark 5, in order to obtain a sharp lower bound  $e_{\text{cont}}$  must be selected in order to fairly approximate  $e_{\text{ref}}$ . Assuming that  $e_{\text{est}}$  is a proper approximation of  $e_{\text{ref}}$  but

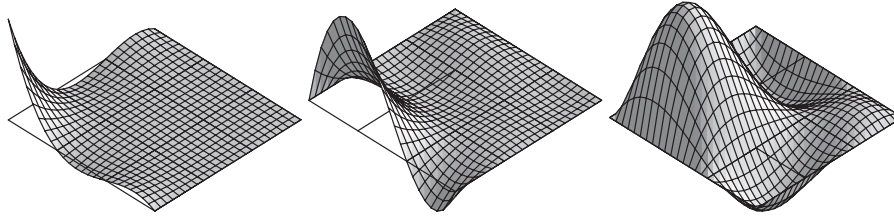


Figure 1. Classification and denomination of the interpolation functions in  $V_{\text{brok}}^{\text{ref}}$ . The functions affecting the boundaries, both associated with corners (left) and edges (centre) are responsible of the continuity. The interior bubble functions (right) do not affect the continuity and they are set in order to obtain the sharper lower bound  $\varepsilon_{\text{low}}$ .

in a broken space, a natural choice is to take the average of the estimated error along the interelement edges.

In order to formalize this averaging, the following decomposition of the local reference interpolation space  $V_k^{\text{ref}}$  is considered:

$$V_k^{\text{ref}} = V_k^{\text{corner}} \oplus V_k^{\text{edge}} \oplus V_k^{\text{bubble}} \tag{28}$$

where  $V_k^{\text{bubble}}$  is the subspace containing the bubble functions (vanishing on  $\partial\Omega_k$ ),  $V_k^{\text{edge}}$  contains the functions having non-zero values in the boundary and vanishing in the corner nodes of element  $\Omega_k$  and  $V_k^{\text{corner}}$  accounts for the degrees of freedom associated with the corner nodes, see Figure 1 for an illustration. This local decomposition induces the definition of the following global spaces:

$$\begin{aligned} V_{\text{brok}}^{\text{corner}} &:= \bigoplus_k V_k^{\text{corner}} & V^{\text{corner}} &:= V_{\text{brok}}^{\text{corner}} \cap V^{\text{ref}} \\ V_{\text{brok}}^{\text{edge}} &:= \bigoplus_k V_k^{\text{edge}} & V^{\text{edge}} &:= V_{\text{brok}}^{\text{edge}} \cap V^{\text{ref}} \\ V^{\text{bubble}} &:= \bigoplus_k V_k^{\text{bubble}} \end{aligned}$$

Note that  $V^{\text{bubble}}$  does not have a ‘broken’ version because the bubble functions do not introduce discontinuities along the edges. Thus,  $V_{\text{brok}}^{\text{ref}}$  and  $V^{\text{ref}}$  are decomposed as

$$V_{\text{brok}}^{\text{ref}} = V_{\text{brok}}^{\text{corner}} \oplus V_{\text{brok}}^{\text{edge}} \oplus V^{\text{bubble}} \quad \text{and} \quad V^{\text{ref}} = V^{\text{corner}} \oplus V^{\text{edge}} \oplus V^{\text{bubble}} \tag{29}$$

Consequently, the estimate  $e_{\text{est}}$  is uniquely represented by the following decomposition:

$$e_{\text{est}} = e_{\text{est}}^{\text{corner}} + e_{\text{est}}^{\text{edge}} + e_{\text{est}}^{\text{bubble}} \tag{30}$$

where  $e_{\text{est}}^{\text{corner}} \in V_{\text{brok}}^{\text{corner}}$ ,  $e_{\text{est}}^{\text{edge}} \in V_{\text{brok}}^{\text{edge}}$  and  $e_{\text{est}}^{\text{bubble}} \in V^{\text{bubble}}$ , and  $e_{\text{cont}} \in V^{\text{ref}}$  is uniquely decomposed as

$$e_{\text{cont}} = e_{\text{cont}}^{\text{corner}} + e_{\text{cont}}^{\text{edge}} + e_{\text{cont}}^{\text{bubble}} \tag{31}$$

where  $e_{\text{cont}}^{\text{corner}} \in V^{\text{corner}}$ ,  $e_{\text{cont}}^{\text{edge}} \in V^{\text{edge}}$  and  $e_{\text{cont}}^{\text{bubble}} \in V^{\text{bubble}}$ . The determination of  $e_{\text{cont}}$  requires to set the proper values for  $e_{\text{cont}}^{\text{corner}}$ ,  $e_{\text{cont}}^{\text{edge}}$  and  $e_{\text{cont}}^{\text{bubble}}$ .

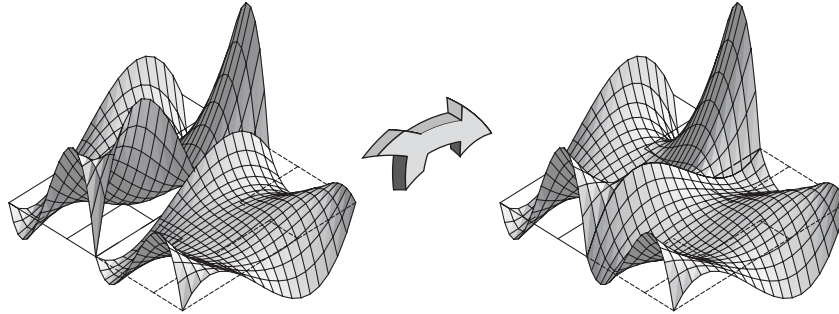


Figure 2. Averaging of the degrees of freedom associated with the edges.

Following Remark 5,  $e_{\text{cont}}$  is determined starting from  $e_{\text{est}}$  and such that  $e_{\text{cont}}$  is likely a good approximation to  $e_{\text{ref}}$ . The application transforming  $e_{\text{est}}$  in  $e_{\text{cont}}$  is denoted by  $\mathcal{M}$ :

$$\begin{aligned} \mathcal{M} : V_{\text{brok}}^{\text{ref}} &\rightarrow V^{\text{ref}} \\ e_{\text{est}} &\mapsto e_{\text{cont}} \end{aligned}$$

Thus, to characterize the smoothing operator  $\mathcal{M}$  it is sufficient to describe  $e_{\text{cont}}$  as a function of  $e_{\text{est}}$ , that is  $e_{\text{cont}}^{\text{corner}}$ ,  $e_{\text{cont}}^{\text{edge}}$  and  $e_{\text{cont}}^{\text{bubble}}$  as functions of  $e_{\text{est}}^{\text{corner}}$ ,  $e_{\text{est}}^{\text{edge}}$  and  $e_{\text{est}}^{\text{bubble}}$ . Indeed,  $\mathcal{M}$  is described by the way it maps  $e_{\text{est}}$  into  $e_{\text{cont}}$ . Thus, in order to characterize  $\mathcal{M}$  it suffices to define the decomposition of the  $e_{\text{cont}} = \mathcal{M}(e_{\text{est}})$ , that is  $e_{\text{cont}}^{\text{corner}}$ ,  $e_{\text{cont}}^{\text{edge}}$  and  $e_{\text{cont}}^{\text{bubble}}$ , in terms of the original estimate  $e_{\text{est}}$  or its decomposition.

In order to enforce continuity, the ‘corner’ and ‘edge’ components are smoothed independently, that is  $e_{\text{cont}}^{\text{corner}} = \mathcal{M}(e_{\text{est}}^{\text{corner}})$  and  $e_{\text{cont}}^{\text{edge}} = \mathcal{M}(e_{\text{est}}^{\text{edge}})$ . As already mentioned, the simplest option is to average the discontinuous values. In a 2-D framework, every interelement edge  $\Gamma_m$  ( $m = 1, \dots, n_{\text{int}}$ ) is shared by two elements, say  $\Gamma_m = \overline{\Omega_k} \cap \overline{\Omega_l}$  and, therefore

$$e_{\text{cont}}^{\text{edge}}|_{\Gamma_m} := \frac{1}{2}(e_{\text{est}}^{\text{edge}}|_{\Omega_k} + e_{\text{est}}^{\text{edge}}|_{\Omega_l}) \tag{32}$$

see Figure 2 for illustration. The same strategy is adopted for the corner points. The contribution of the interpolation functions associated with the corner points,  $e_{\text{cont}}^{\text{corner}}$  is computed averaging the values of the discontinuous function  $e_{\text{est}}^{\text{corner}}$  in each corner point. That results in an expression similar to Equation (32) where, for every corner point, the number of values to average is equal to the number of elements to which the corner point belongs. This is illustrated in Figure 3.

Once  $e_{\text{cont}}^{\text{corner}}$  and  $e_{\text{cont}}^{\text{edge}}$  are set it is necessary to find the value of  $e_{\text{cont}}^{\text{bubble}}$ . It is worth noting that the choice for  $e_{\text{cont}}^{\text{bubble}}$  does not affect the continuity of  $e_{\text{cont}}$ . The value of  $e_{\text{cont}}^{\text{bubble}}$  is therefore selected such that the obtained estimate is as sharp as possible.

Recall that, once  $e_{\text{cont}}$  is determined, the sharper lower bound is  $\varepsilon_{\text{low}}^{\text{opt}}$ , see Equation (27). Then, the first idea is to select  $e_{\text{cont}}^{\text{bubble}}$  such that, given  $e_{\text{cont}}^{\text{corner}}$  and  $e_{\text{cont}}^{\text{edge}}$ , it maximizes  $\varepsilon_{\text{low}}^{\text{opt}}$ . However, this criterion leads to a non-linear global (referred to the whole domain) equation which is difficult to solve. On the contrary, finding  $e_{\text{cont}}^{\text{bubble}}$  such that  $\varepsilon_{\text{low}}(1)$ , see Equation (25), is maximum leads to a simple linear local (element by element) equation. This is stated in the following theorem:

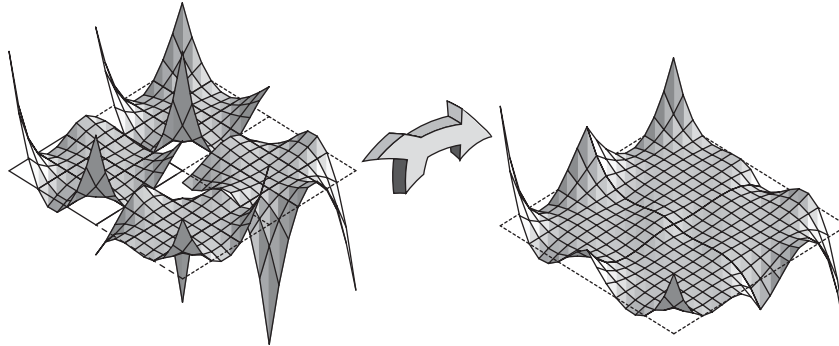


Figure 3. Averaging of the degrees of freedom associated with the corners.

*Theorem 3*

Let  $e_{est}$  be an error estimate verifying the hypothesis of Theorem 1 and, therefore, being an upper bound of the reference error. Let  $e_{cont} = e_{cont}^{corner} + e_{cont}^{edge} + e_{cont}^{bubble}$  be a corrected estimate. Assume that  $e_{cont}^{corner}$  and  $e_{cont}^{edge}$  are obtained by averaging. Then, the value of  $e_{cont}^{bubble}$  maximizing  $\varepsilon_{low}(1)$  is such that

$$a(e_{cont}^{bubble}, v) = a(e_{est} - e_{cont}^{corner} - e_{cont}^{edge}, v) \quad \forall v \in V^{bubble} \tag{33}$$

*Proof*

Recall that  $\varepsilon_{low}(1) = \|e_{est}\|^2 - \|e_{cont} - e_{est}\|^2$ , therefore maximize  $\varepsilon_{low}(1)$  is equivalent to minimize

$$\|e_{cont} - e_{est}\| = \|e_{cont}^{bubble} - (e_{est} - e_{cont}^{corner} - e_{cont}^{edge})\|$$

The problem is reformulated as: find  $e_{cont}^{bubble} \in V^{bubble}$  such that  $\|e_{cont}^{bubble} - (e_{est} - e_{cont}^{corner} - e_{cont}^{edge})\|$  is minimum. Obviously, the solution of this problem is the projection of  $e_{est} - e_{cont}^{corner} - e_{cont}^{edge}$  on  $V^{bubble}$  which satisfies Equation (33).  $\square$

Thus, taking  $e_{cont}^{bubble}$  as the solution of Equation (33) completes the determination of  $\mathcal{M}$ . Note that, in this case,  $e_{cont}$  depends on the ‘corner’ and ‘edge’ components of  $e_{est}$ .

*Remark 6*

The smoothing operator  $\mathcal{M}$  is linear because  $e_{cont}^{corner}$ ,  $e_{cont}^{edge}$  and  $e_{cont}^{bubble}$  are linear functions of  $e_{est}^{corner}$ ,  $e_{est}^{edge}$  and  $e_{est}^{bubble}$ . Moreover, the quality of the lower bound  $e_{low}^{opt}$  depends on the ability of  $\mathcal{M}$  to approximate the reference error  $e_{ref}$ . Note this quality depends only on the averaging on the boundaries. It suffices that  $e_{cont}$  coincides with  $e_{ref}$  on the interelement boundaries (i.e. for  $e_{cont}^{corner} + e_{cont}^{edge}$ ) to obtain an exact error assessment. That is if

$$e_{cont}|_{\Gamma_m} = e_{ref}|_{\Gamma_m} \quad \text{for every } m = 1, \dots, n_{int}$$

then  $e_{cont} = e_{ref}$  and, consequently (see Remark 5),

$$\varepsilon_{low}^{opt} = \varepsilon_{low}(1) = \|e_{ref}\|^2$$



4.3. Computational aspects

The selection of the optimal value of  $e_{\text{cont}}^{\text{bubble}}$  is performed solving Equation (33). These computations can be done locally, element by element, because the bubble spaces are orthogonal (the supports of the bubbles are disjoint). Thus, once  $e_{\text{cont}}^{\text{corner}}$  and  $e_{\text{cont}}^{\text{edge}}$  are computed by simple averaging, the restriction of  $e_{\text{cont}}^{\text{bubble}}$  to  $\Omega_k$ ,  $e_{\text{cont}}^{\text{bubble}}|_{\Omega_k}$  is computed solving the local version of Equation (33):

$$a_k(e_{\text{cont}}^{\text{bubble}}|_{\Omega_k}, v) = a_k(e_{\text{est}} - e_{\text{cont}}^{\text{corner}} - e_{\text{cont}}^{\text{edge}}, v) \quad \forall v \in V_k^{\text{bubble}} \tag{34}$$

Equation (34) results in a small system of linear equations that must be solved to compute  $e_{\text{cont}}^{\text{bubble}}|_{\Omega_k}$ . The number of equations for each local problem is equal to the number of ‘bubble’ degrees of freedom in the reference discretization. For example, for lagrangian quadrilateral elements, this number is equal to  $(1 - \tilde{p})^2$ , being  $\tilde{p}$  the degree of the polynomials used to generate  $V^{\text{ref}}$ .

4.4. Assessment of the effectivity index and average estimate

Once the lower bound of the error is computed, the effectivity index of the original estimate  $\|e_{\text{est}}\|$  may be easily assessed. Let  $\eta_{\text{est}}$  be the effectivity index associated with  $e_{\text{est}}$ ,

$$\eta_{\text{est}} := \frac{\|e_{\text{est}}\|}{\|e_{\text{ref}}\|} \tag{35}$$

The upper bound property ensures  $\eta_{\text{est}} \geq 1$ . Nevertheless  $\eta_{\text{est}}$  may be very large and it is not possible, in the general case, to assess the quality of the estimate. Using the lower bound  $\varepsilon_{\text{low}}$  of the error, an upper bound of the effectivity index  $\eta^+$  is easily computed:

$$\eta^+ := \frac{\|e_{\text{est}}\|}{\sqrt{\varepsilon_{\text{low}}}} = \geq \eta_{\text{est}} \tag{36}$$

This pessimistic value of the effectivity index is sharp when the lower bound error estimate  $\varepsilon_{\text{low}}$  is sharp.

Once the upper and the lower bounds of the error,  $\varepsilon_{\text{upp}} = \|e_{\text{est}}\|^2$  and  $\varepsilon_{\text{low}}$ , are available the average estimate is introduced

$$\varepsilon_{\text{ave}} := \frac{1}{2}(\varepsilon_{\text{upp}} + \varepsilon_{\text{low}}) \tag{37}$$

Remark 7

As noted in Remark 2, the estimates  $\varepsilon_{\star}$  represent approximations to the squared norms of the error. The average of the squared norms is larger than the simple averaging of the norms, that is,

$$\frac{1}{2}(\varepsilon_{\text{upp}} + \varepsilon_{\text{low}}) \geq \left[ \frac{1}{2}(\sqrt{\varepsilon_{\text{upp}}} + \sqrt{\varepsilon_{\text{low}}}) \right]^2$$

The behaviour of this average estimate is analysed in the examples presented in Section 6.

5. FITTING LOCAL ARBITRARY CONSTANTS FOR  $b = 0$

In problems without a reaction term, the lower bounds of the error obtained with the previously discussed techniques have a poor (very low) effectivity index. In this section, a strategy to

preclude this deficiency is introduced. If  $b=0$  in Equation (1) (pure diffusion, no reaction)  $e_{\text{est}}$  is locally determined up to a constant because

$$\|e_{\text{est}}\|_k = \|e_{\text{est}} + c_k\|_k, \quad k = 1, \dots, n_{\text{elem}} \quad (38)$$

Then, the estimate  $e_{\text{est}}$  may be replaced by  $e_{\text{est}} + \sum_{k=1}^{n_{\text{elem}}} c_k \phi_k$  without changing the upper bound  $\varepsilon_{\text{upp}}$ , being  $\{\phi_1, \phi_2, \dots, \phi_{n_{\text{elem}}}\}$  the basis of the space of piecewise constant functions. That is, for  $k = 1, \dots, n_{\text{elem}}$ ,

$$\phi_k(\mathbf{x}) = \begin{cases} 1 & \text{if } \mathbf{x} \in \Omega_k \\ 0 & \text{if } \mathbf{x} \notin \Omega_k \end{cases} \quad (39)$$

The upper bound estimate  $\varepsilon_{\text{upp}}$  is independent of the constants  $c_k$ . Nevertheless, the choice of the constants  $c_k$  affects drastically the value of the corrected error,  $e_{\text{cont}}$ . Moreover, the correction strategy is expected to work properly only if the average values of  $e_{\text{est}}$  are close to  $e_{\text{ref}}$ , see Remark 5. If the constants are set arbitrarily, the value of the correction cannot be expected to be optimal.

Consequently, the constants  $c_k$ ,  $k = 1, \dots, n_{\text{elem}}$ , are taken as unknowns and they are determined such that the resulting lower bound is somehow optimal. Let  $\mathbf{c} = [c_1 \dots c_{n_{\text{elem}}}]$  be the vector of unknown constants. The corrected estimate  $e_{\text{cont}}$  may be seen as a function of  $\mathbf{c}$ :

$$e_{\text{cont}}(\mathbf{c}) := \mathcal{M}\left(e_{\text{est}} + \sum_{k=1}^{n_{\text{elem}}} c_k \phi_k\right) = \mathcal{M}(e_{\text{est}}) + \sum_{k=1}^{n_{\text{elem}}} c_k \mathcal{M}(\phi_k) \quad (40)$$

It is clear from Equation (40) that, due to the linearity of  $\mathcal{M}$ ,  $e_{\text{cont}}(\mathbf{c})$  is linear. Both the lower bounds  $\varepsilon_{\text{low}}(1)$  and  $\varepsilon_{\text{low}}^{\text{opt}}$  depend on  $\mathbf{c}$  through  $e_{\text{cont}}$ . The criterion used to select  $\mathbf{c}$  is obviously to maximize the lower bound. The maximization of  $\varepsilon_{\text{low}}^{\text{opt}}$  is the more natural option because  $\varepsilon_{\text{low}}^{\text{opt}}$  is the sharper error bound. Nevertheless, similarly to the previous section, finding  $\mathbf{c}$  that optimizes  $\varepsilon_{\text{low}}^{\text{opt}}$  requires to solve a non-linear problem. On the contrary, to find  $\mathbf{c}$  such that  $\varepsilon_{\text{low}}(1)$  is maximum leads to a simple linear problem. Thus, the criterion for determining  $\mathbf{c}$  is based on maximizing  $\varepsilon_{\text{low}}(1)$  rather than  $\varepsilon_{\text{low}}^{\text{opt}}$ .

The dependence of  $\varepsilon_{\text{low}}(1)$  on  $\mathbf{c}$  is written by introducing Equation (40) in Equation (25) and replacing  $e_{\text{est}}$  by  $e_{\text{est}} + \sum_{k=1}^{n_{\text{elem}}} c_k \phi_k$ :

$$\begin{aligned} \varepsilon_{\text{low}}(1) &= \left\| e_{\text{est}} + \sum_{k=1}^{n_{\text{elem}}} c_k \phi_k \right\|^2 - \left\| e_{\text{est}} + \sum_{k=1}^{n_{\text{elem}}} c_k \phi_k - \mathcal{M}(e_{\text{est}}) - \sum_{k=1}^{n_{\text{elem}}} \mathcal{M}(\phi_k) c_k \right\|^2 \\ &= \|e_{\text{est}}\|^2 - \left\| e_{\text{est}} - \mathcal{M}(e_{\text{est}}) - \sum_{k=1}^{n_{\text{elem}}} \mathcal{M}(\phi_k) c_k \right\|^2 \end{aligned} \quad (41)$$

Then, to maximize  $\varepsilon_{\text{low}}(1)$  is equivalent to minimize the function  $F(\mathbf{c})$  defined by

$$F(\mathbf{c}) := \left\| e_{\text{est}} - \mathcal{M}(e_{\text{est}}) - \sum_{k=1}^{n_{\text{elem}}} \mathcal{M}(\phi_k) c_k \right\|^2$$

The coefficients  $c_k$  that minimize  $F(\mathbf{c})$  are obtained imposing that  $\sum_{k=1}^{n_{\text{elem}}} \mathcal{M}(\phi_k) c_k$  is the projection of  $e_{\text{est}} - \mathcal{M}(e_{\text{est}})$  on the space generated by the functions  $\mathcal{M}(\phi_k)$ , for  $k = 1 \dots n_{\text{elem}}$

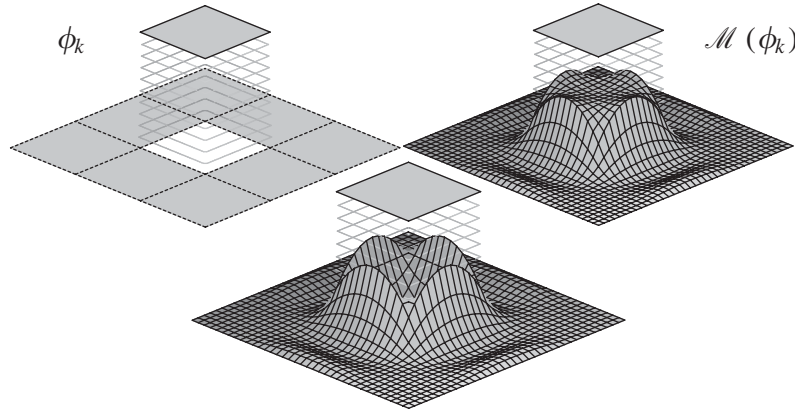


Figure 4. Construction of  $\mathcal{M}(\phi_k)$  (right) from  $\phi_k$  (left). The function in the centre accounts only for the 'corner' and 'edge' terms, before adding the 'bubble' term that affects only the interior of the elements. Note that the influence of using the proper 'bubble' contribution is very important.

(that is, the image by  $\mathcal{M}$  of the space of piecewise constant functions). Figure 4 illustrates the shape of the functions  $\mathcal{M}(\phi_k)$  and their construction from  $\phi_k$ .

Thus, the equation to be satisfied by the coefficients  $c_k$  is

$$\sum_{k=1}^{n_{\text{elem}}} c_k a(\mathcal{M}(\phi_k), \mathcal{M}(\phi_l)) = a(e_{\text{est}} - \mathcal{M}(e_{\text{est}}), \mathcal{M}(\phi_l)) \quad \text{for } l = 1, \dots, n_{\text{elem}} \quad (42)$$

That is,  $\mathbf{c}$  is computed as the solution of a linear  $n_{\text{elem}} \times n_{\text{elem}}$  system of equations.

Once the coefficients  $c_k$  are computed, the corresponding corrected estimate  $e_{\text{cont}}$  is introduced in the expression of  $\varepsilon_{\text{low}}^{\text{opt}}$  to obtain the sharper error lower bound.

Numerical experiments demonstrate that the correction obtained with this strategy yields sharp lower bound estimates because the obtained correction  $e_{\text{cont}}$  is a much better approximation to  $e_{\text{ref}}$ , see Figure 5. On the contrary, the correction for the standard estimate (i.e. with arbitrary constants) yields lower bound estimates of poor quality.

It is worth noting that the constants  $c_k$  are determined solving the global system of equations (42). Thus, adding these constants to the original estimate  $e_{\text{est}}$  accounts for the influence of the whole domain in the local errors. Consequently, the estimate  $e_{\text{cont}}$  using this information may be used to assess the pollution errors, that is, the errors affecting each zone of the domain coming from far from its close neighbourhood.

## 6. NUMERICAL EXAMPLES

We study in this section the behaviour of the postprocessing estimate presented above. The examples selected are such that the analytical exact solution is known and they have been used by other authors to assess the performance of similar techniques [1, 2]. The quality of the error estimates is measured using the index  $\rho$

$$\rho = \frac{\text{estimated error}}{\text{exact (or reference) error}} - 1$$

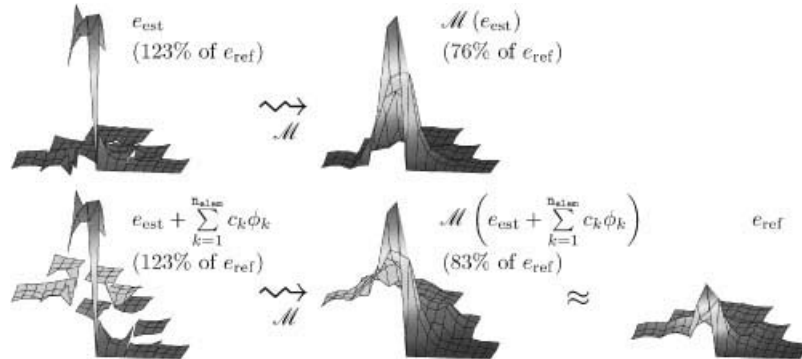


Figure 5. Illustration of the constant fitting process: the raw estimate  $e_{est}$  with arbitrary constants is smoothed into  $\mathcal{M}(e_{est})$  (top), the smoothed version of the estimate corrected with the optimal constants is much more similar to the reference error (bottom): in the example the underestimation is improved from 76% (without constant fitting) to 83%.

that is, the effectivity index minus one. The use of  $\rho$  is preferred because the sign of  $\rho$  indicates if the estimate is an upper or a lower bound (positive if upper, negative if lower) and the absolute value indicates the quality of the estimate (good quality if  $|\rho|$  small). In the following, the value of  $\rho$  corresponding to every estimate is denoted with the same subscript, that is,

$$\rho_{\star} = \frac{\sqrt{\varepsilon_{\star}}}{\|e\|} - 1$$

where the subscript  $\star$  takes the values ‘upp’, ‘low’ and ‘ave’. The superscript C for  $\rho_{low}$ ,  $\rho_{low}^C$ , is used to denote the correction obtained with the determination of elementwise constants introduced in Section 5. Moreover, we also use the version  $\rho^+$  corresponding to the assessed effectivity index  $\eta^+$  ( $\rho^+ := \eta^+ - 1$ ), see Equation (36).

As noted in Section 3.4, the second and third estimators introduced in Reference [1], denoted by  $e_2$  and  $e_3$ , respectively, are used as the original upper bound estimates  $e_{est}$ . In the examples, the performance of these estimates is analysed throughout the values of  $\rho_{upp}$ .

### 6.1. Example 1

In the first example the reaction–diffusion equation is solved,  $a = 1$  and  $b = 1$  in Equation (1). The problem is defined in the squared domain  $\Omega = (0, 1) \times (0, 1)$ . The boundary conditions are set to be Dirichlet homogeneous (that is  $u = 0$ ) on  $\Gamma_D := \{(x, 0); 0 \leq x \leq 1\}$  and Neumann homogeneous (that is  $\partial u / \partial n = 0$ ) elsewhere on  $\partial\Omega$ . The source term  $s$  is taken such that the exact solution has the following analytical expression:

$$u(x, y) = \frac{1}{2000} x^2(1 - x)^2 e^{10x^2} y^2(1 - y)^2 e^{10y} \tag{43}$$

see Figure 6 for a representation. The second example described in this section is stated such that the solution  $u$  is exactly the same.

The approximate solution  $u_h$  is computed using a bilinear interpolation ( $p = 1$ ) whereas the error estimates  $e_2$  and  $e_3$  are computed using a bicubic interpolation ( $\tilde{p} = 3$ ).

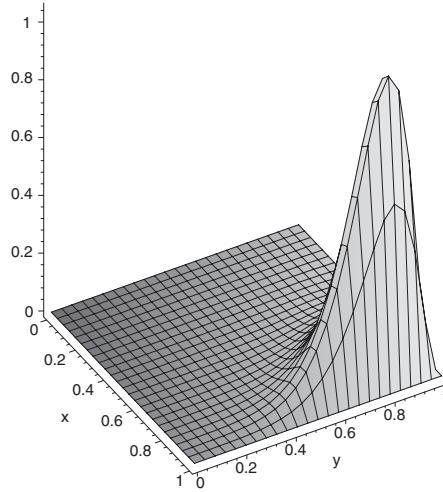


Figure 6. Examples 1 and 2: exact solution.

Table I. Example 1: results in a series of uniformly  $h$ -refined meshes.

| Dof  | $\ e\ /\ u\ $ | $\ e_{\text{ref}}\ /\ u\ $ | Estimate $e_2$ |                     |                     |                     | Estimate $e_3$ |                     |                     |                     |
|------|---------------|----------------------------|----------------|---------------------|---------------------|---------------------|----------------|---------------------|---------------------|---------------------|
|      |               |                            | $\rho^+$       | $\rho_{\text{upp}}$ | $\rho_{\text{low}}$ | $\rho_{\text{ave}}$ | $\rho^+$       | $\rho_{\text{upp}}$ | $\rho_{\text{low}}$ | $\rho_{\text{ave}}$ |
| 36   | 0.8469        | 0.7726                     | 0.3453         | 0.1589              | -0.1386             | 0.0210              | 0.2713         | 0.0544              | -0.1706             | -0.0514             |
| 121  | 0.4331        | 0.4036                     | 0.2428         | 0.1221              | -0.0971             | 0.0184              | 0.2116         | 0.0569              | -0.1277             | -0.0310             |
| 441  | 0.3083        | 0.3064                     | 0.3258         | 0.2132              | -0.0849             | 0.0745              | 0.2737         | 0.1706              | -0.0809             | 0.0524              |
| 1681 | 0.2093        | 0.2092                     | 0.2578         | 0.1831              | -0.0594             | 0.0688              | 0.1843         | 0.1263              | -0.0489             | 0.0424              |
| 6561 | 0.1144        | 0.1144                     | 0.1129         | 0.0845              | -0.0255             | 0.0310              | 0.0691         | 0.0498              | -0.0181             | 0.0164              |

The proposed approach is used to recover new estimates in two sequences of increasingly refined meshes. In the first series of meshes the refinement is uniform, in the second one the refinement follows an adaptive strategy based on the error assessment [14].

The results concerning the uniformly refined meshes are summarized in Table I and Figure 7.

In a similar manner, the results concerning the adaptively refined meshes are summarized in Table II and Figure 8. The sequence of adapted meshes is shown in Figure 9.

It is worth noting in Tables I and II that the difference between the exact error (in this case is known) and the reference error is negligible for accurate enough meshes. As expected, the values of  $\rho_{\text{upp}}$  are indeed positive and the values of  $\rho_{\text{low}}$  negative. The value of  $\rho^+$  is greater than  $\rho_{\text{upp}}$ . Note that  $\rho^+$  is computed without any information on the exact (or reference) solution but it furnishes a good approximation of the exact effectivity index. Moreover, for most of the meshes (except for the coarsest) the value of the corrected estimate  $\varepsilon_{\text{low}}$  is better than the original estimate  $\varepsilon_{\text{upp}}$  ( $|\rho_{\text{low}}| < |\rho_{\text{upp}}|$ ), that results on  $\rho_{\text{ave}} > 0$ .

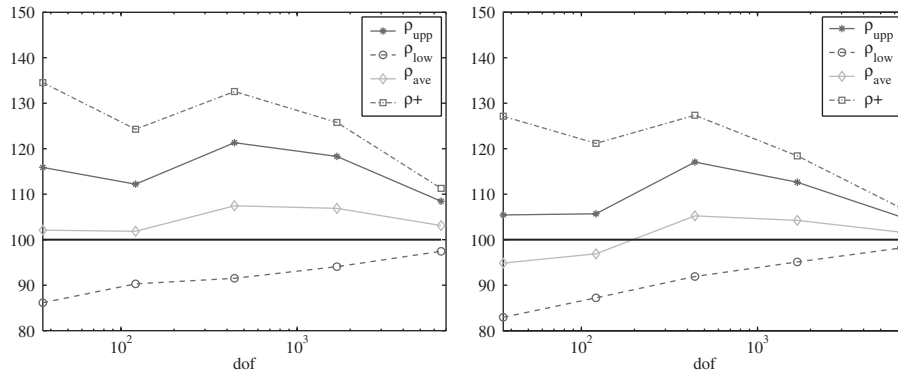


Figure 7. Example 1: performance of the estimates following a uniform  $h$ -refinement process for the estimates  $e_2$  (left) and  $e_3$  (right).

Table II. Example 1: results in a series of adaptively  $h$ -refined meshes.

| Dof  | Error         |                     | Estimate $e_2$ |              |              |              | Estimate $e_3$ |              |              |              |
|------|---------------|---------------------|----------------|--------------|--------------|--------------|----------------|--------------|--------------|--------------|
|      | $\ e\ /\ u\ $ | $\ e_{ref}\ /\ u\ $ | $\rho^+$       | $\rho_{upp}$ | $\rho_{low}$ | $\rho_{ave}$ | $\rho^+$       | $\rho_{upp}$ | $\rho_{low}$ | $\rho_{ave}$ |
| 36   | 0.8469        | 0.7726              | 0.3453         | 0.1589       | -0.1386      | 0.0210       | 0.2713         | 0.0544       | -0.1706      | -0.0514      |
| 2550 | 0.0798        | 0.0798              | 0.0822         | 0.0645       | -0.0164      | 0.0248       | 0.0517         | 0.0354       | -0.0155      | 0.0103       |
| 2905 | 0.0478        | 0.0478              | 0.1263         | 0.1136       | -0.0113      | 0.0530       | 0.1129         | 0.0622       | -0.0456      | 0.0098       |
| 3574 | 0.0433        | 0.0433              | 0.1279         | 0.1152       | -0.0113      | 0.0539       | 0.1108         | 0.0614       | -0.0445      | 0.0098       |

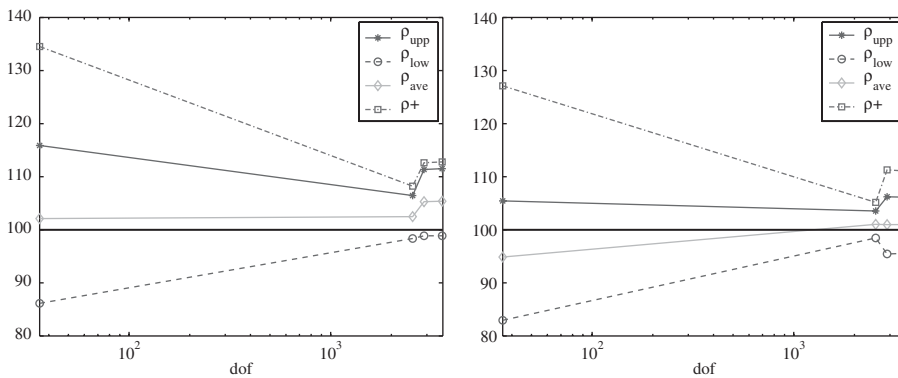


Figure 8. Example 1: performance of the estimates following an adaptive  $h$ -refinement process for the estimates  $e_2$  (left) and  $e_3$  (right)

*Remark 8*

As expected, the adaptive procedure optimizes the computational resources and yields lower error with less degrees of freedom. However, the adapted meshes have distorted elements, see Figure 9, and the quality of the estimates  $e_2$  and  $e_3$  is slightly degraded in adapted meshes, see Figure 8. This phenomenon produces a peculiar end effect in the plots describing the evolution of the effectivity along the adaptive process. This effect does not appear in the

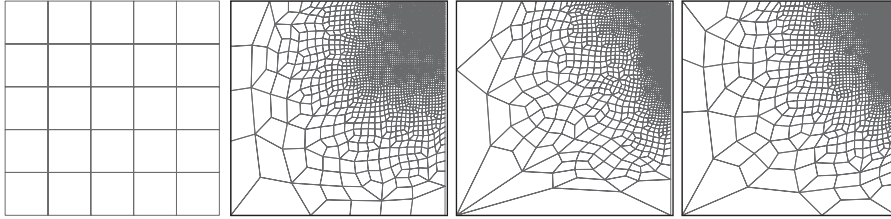


Figure 9. Example 1: Sequence of adapted meshes with 36, 2550, 2950 and 3574 dof.

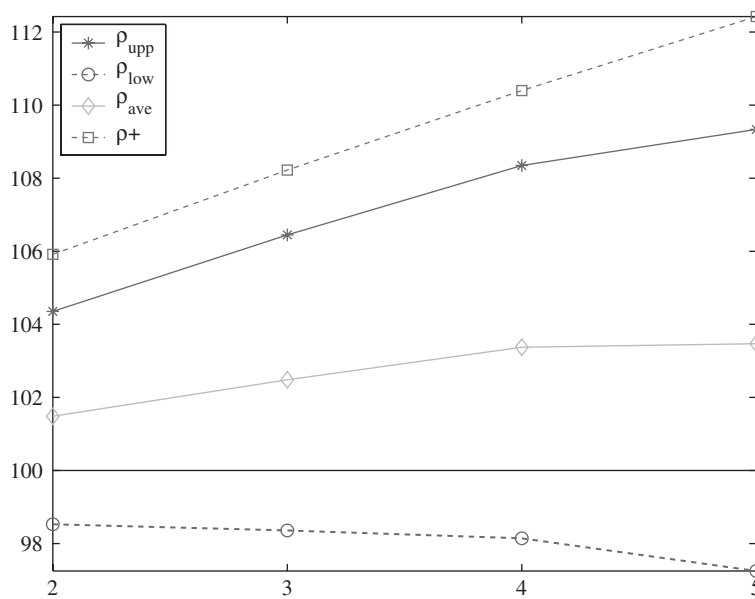


Figure 10. Example 1: performance of the estimators and using different degrees of interpolation in the reference space ( $\tilde{p}$ ).

uniform  $h$ -refinement process where all the meshes are structured, see Figure 7. The proposed lower bound corrects this behaviour in the case of the estimate  $e_2$  but not in the case of  $e_3$ . In this example, the average  $\varepsilon_{\text{ave}}$  performs very well in the sense that behaves as a new estimate, mostly a new upper bound, much more reliable than the original one. The same effect is observed in the next example, see Figure 11.

The effect of varying the degree of interpolation in the reference space ( $\tilde{p}$ ) is investigated for one of the meshes (the second mesh of the adaptive process, with 2550 dof) and for the estimate  $e_2$ . We are interested in assessing the influence of  $\tilde{p}$  in the error estimate and the corresponding corrections. The results are shown in Figure 10. Note that the effectivity of the original estimate,  $e_{\text{est}}$  is not improved by using a larger  $\tilde{p}$ . On the contrary, the larger values of  $\tilde{p}$  are associated with the poorer quality estimates. Nevertheless, the quality of the postprocessed lower bounds is not so sensitive to the variations of  $\tilde{p}$  and their quality does not depend on  $\tilde{p}$ .

Table III. Example 2: results in a series of uniformly  $h$ -refined meshes.

| Dof  | $\ e\ /\ u\ $ | $\ e_{\text{ref}}\ /\ u\ $ | Estimate $e_2$ |                     |                     |                       |                     |
|------|---------------|----------------------------|----------------|---------------------|---------------------|-----------------------|---------------------|
|      |               |                            | $\rho^+$       | $\rho_{\text{upp}}$ | $\rho_{\text{low}}$ | $\rho_{\text{low}}^C$ | $\rho_{\text{ave}}$ |
| 36   | 0.8483        | 0.7737                     | 0.2729         | 0.1571              | -0.1177             | -0.0909               | 0.0405              |
| 121  | 0.4342        | 0.4046                     | 0.2059         | 0.1217              | -0.0838             | -0.0698               | 0.0304              |
| 441  | 0.3091        | 0.3072                     | 0.2220         | 0.2131              | -0.0461             | -0.0073               | 0.1084              |
| 1681 | 0.2099        | 0.2098                     | 0.1844         | 0.1831              | -0.0321             | -0.0011               | 0.0949              |
| 6561 | 0.1148        | 0.1148                     | 0.0849         | 0.0845              | -0.0148             | -0.0003               | 0.0430              |

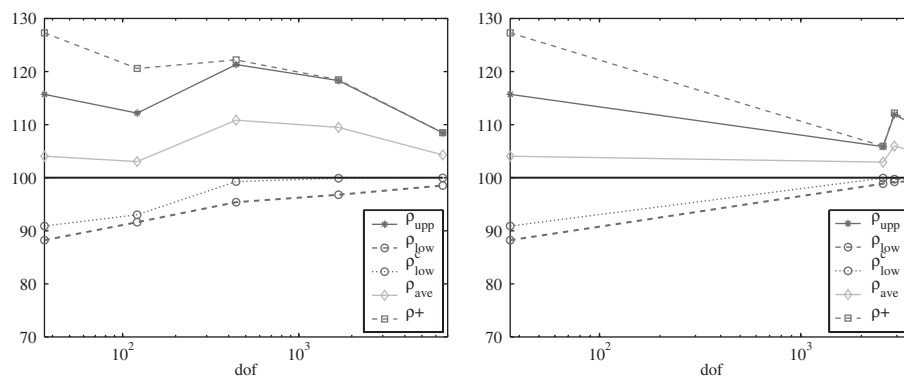


Figure 11. Example 2: performance of the estimates following a uniform (left) and an adaptive (right)  $h$ -refinement process for the estimate  $e_2$ .

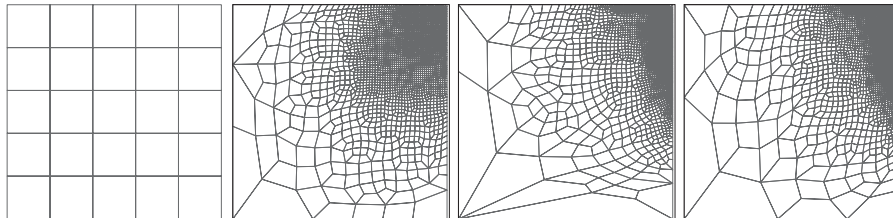


Figure 12. Example 2: sequence of adapted meshes with 36, 2561, 2918 and 3628 dof.

6.2. Example 2

Now, we consider the Poisson equation,  $a = 1$  and  $b = 0$  in Equation (1). The domain and the boundary conditions are exactly the same as in the previous example. The source term  $s$  is taken such that the exact solution is also the same, see Equation (43). In this example, we only study the application of the developed postprocessing strategy to the  $e_2$  estimate.

Again, the proposed strategy is used in a series of uniformly and adaptively  $h$ -refined meshes. The results for the uniformly refined meshes are summarized in Table III and Figure 11. Figure 12 shows a sequence of adapted meshes and Table IV with Figure 11 describe the behaviour of the different estimates. The notation  $\rho_{\text{low}}^C$  is introduced to denote the



Table IV. Example 2: results in a series of adaptively  $h$ -refined meshes.

| Dof  | $\ e\ /\ u\ $ | $\ e_{\text{ref}}\ /\ u\ $ | Estimate $e_2$ |                     |                     |                                |                     |
|------|---------------|----------------------------|----------------|---------------------|---------------------|--------------------------------|---------------------|
|      |               |                            | $\rho^+$       | $\rho_{\text{upp}}$ | $\rho_{\text{low}}$ | $\rho_{\text{low}}^{\text{C}}$ | $\rho_{\text{ave}}$ |
| 36   | 0.8483        | 0.7737                     | 0.2729         | 0.1571              | -0.1177             | -0.0909                        | 0.0405              |
| 2561 | 0.0785        | 0.0785                     | 0.0593         | 0.0586              | -0.0112             | -0.0007                        | 0.0294              |
| 2918 | 0.0482        | 0.0482                     | 0.1216         | 0.1186              | -0.0077             | -0.0027                        | 0.0596              |
| 3628 | 0.0432        | 0.0432                     | 0.1038         | 0.1008              | -0.0070             | -0.0027                        | 0.0503              |

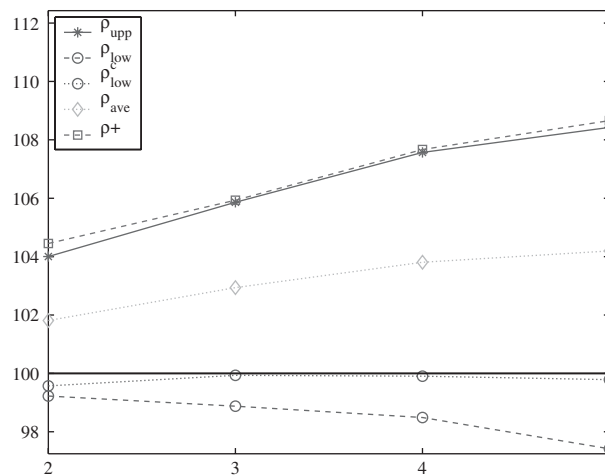


Figure 13. Example 2: performance of the estimators and using different degrees of interpolation in the reference space ( $\tilde{p}$ ).

correction introduced in Section 5. As expected, the value of  $\rho_{\text{low}}^{\text{C}}$  is much better than the value of  $\rho_{\text{low}}$ .

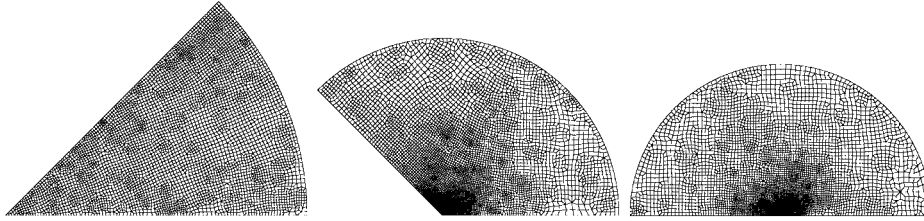
The influence of  $\tilde{p}$  in the different estimates is shown in Figure 13. These results correspond to the second mesh of the adaptive process, with 2561 dof. Once again, due to the phenomenon described in the previous example, increasing  $\tilde{p}$  does not result in a better effectivity index for the upper bound estimate. Nevertheless, the lower bound estimate  $e_{\text{cont}}$  with the constant element by element correction (measured by  $\rho_{\text{low}}^{\text{C}}$ ) is roughly independent of  $\tilde{p}$  and much better compared to the original estimate.

### 6.3. Example 3

This example was introduced in Reference [1]. We consider the Laplace equation,  $a=1$ ,  $b=0$  and  $s=0$  in Equation (1). As in the previous example, only the  $e_2$  estimate is used with the proposed postprocessing strategy.

The domain  $\Omega$  is defined by  $\Omega = \{(r, \theta) : 0 < r < 1, 0 < \theta < k\pi/4\}$  where  $r$  and  $\theta$  are the polar coordinates and the analytical solution is

$$u(r, \theta) = r^{2/k} \sin\left(\frac{2\theta}{k}\right) \tag{44}$$

Figure 14. Example 3: adapted meshes for  $k = 1$  (left)  $k = 3$  (centre) and  $k = 4$  (right).Table V. Example 3,  $k = 1$ : results in a series of adaptively  $h$ -refined meshes.

| Dof  | $\ e\ /\ u\ $ | $\ e_{\text{ref}}\ /\ u\ $ | Estimate $e_2$ |                     |                     |                       |                     |
|------|---------------|----------------------------|----------------|---------------------|---------------------|-----------------------|---------------------|
|      |               |                            | $\rho^+$       | $\rho_{\text{upp}}$ | $\rho_{\text{low}}$ | $\rho_{\text{low}}^C$ | $\rho_{\text{ave}}$ |
| 69   | 0.0397        | 0.0397                     | 0.3788         | 0.3730              | -0.0109             | -0.0042               | 0.1993              |
| 1637 | 0.0069        | 0.0069                     | 0.1250         | 0.1224              | -0.0052             | -0.0022               | 0.0619              |
| 3938 | 0.0044        | 0.0044                     | 0.1925         | 0.1849              | -0.0109             | -0.0064               | 0.0934              |
| 4668 | 0.0040        | 0.0040                     | 0.2051         | 0.1992              | -0.0092             | -0.0048               | 0.1019              |

Table VI. Example 3,  $k = 3$ : results in a series of adaptively  $h$ -refined meshes.

| Dof  | $\ e\ /\ u\ $ | $\ e_{\text{ref}}\ /\ u\ $ | Estimate $e_2$ |                     |                     |                       |                     |
|------|---------------|----------------------------|----------------|---------------------|---------------------|-----------------------|---------------------|
|      |               |                            | $\rho^+$       | $\rho_{\text{upp}}$ | $\rho_{\text{low}}$ | $\rho_{\text{low}}^C$ | $\rho_{\text{ave}}$ |
| 169  | 0.0298        | 0.0294                     | 1.0167         | 0.6153              | -0.2412             | -0.1991               | 0.2749              |
| 580  | 0.0139        | 0.0138                     | 0.7023         | 0.3618              | -0.2468             | -0.2001               | 0.1168              |
| 1436 | 0.0078        | 0.0077                     | 0.4700         | 0.3375              | -0.1211             | -0.0901               | 0.1438              |
| 3795 | 0.0047        | 0.0047                     | 0.3860         | 0.3242              | -0.0626             | -0.0446               | 0.1546              |
| 6585 | 0.0036        | 0.0035                     | 0.3407         | 0.2861              | -0.0615             | -0.0407               | 0.1345              |

That is,  $\Omega$  is a circular sector and  $k$  is a parameter that sets both the size of the domain and the regularity of the solution. In the following, we consider the cases  $k = 1, 3$  and  $4$ . Dirichlet boundary conditions are imposed along  $\theta = 0$  and Neumann boundary conditions are forced on the rest of the boundary. The boundary conditions are such that the exact solution is the analytical expression given in Equation (44).

For each one of the values of  $k$ , the error assessment is performed for a sequence of adapted meshes. Figure 14 shows examples of adapted meshes for each value of  $k$ .

The results are shown in Tables V, VI and VII for  $k = 1, 3$  and  $4$ , respectively, and also in Figure 15. It is worth noting that using the constant fitting (the difference between  $\rho_{\text{low}}^C$  and  $\rho_{\text{low}}$ , see Figure 15) is relevant specially for  $k = 4$ , that is, when the singularity pollutes the error estimate based only on local computations.

In order to analyse the spatial distribution of the estimated error, Figure 16 shows the histograms describing the occurrences of the values of local (element by element) effectivity indices for both the estimated error and the lower estimate. The example corresponds to the second mesh obtained for  $k = 1$  (with 1637 dof). An almost uniform distribution is obtained

Table VII. Example 3,  $k = 4$ : results in a series of adaptively  $h$ -refined meshes.

| Dof  | $\ e\ /\ u\ $ | $\ e_{\text{ref}}\ /\ u\ $ | Estimate $e_2$ |                     |                     |                                |                     |
|------|---------------|----------------------------|----------------|---------------------|---------------------|--------------------------------|---------------------|
|      |               |                            | $\rho^+$       | $\rho_{\text{upp}}$ | $\rho_{\text{low}}$ | $\rho_{\text{low}}^{\text{C}}$ | $\rho_{\text{ave}}$ |
| 220  | 0.1310        | 0.1200                     | 0.4449         | 0.1384              | -0.2960             | -0.2121                        | -0.0211             |
| 372  | 0.0587        | 0.0548                     | 0.4858         | 0.2090              | -0.2626             | -0.1863                        | 0.0304              |
| 723  | 0.0312        | 0.0297                     | 0.5364         | 0.2418              | -0.2603             | -0.1917                        | 0.0477              |
| 3297 | 0.0126        | 0.0122                     | 0.4800         | 0.2293              | -0.2346             | -0.1694                        | 0.0491              |
| 6859 | 0.0077        | 0.0076                     | 0.4154         | 0.2440              | -0.1790             | -0.1211                        | 0.0770              |

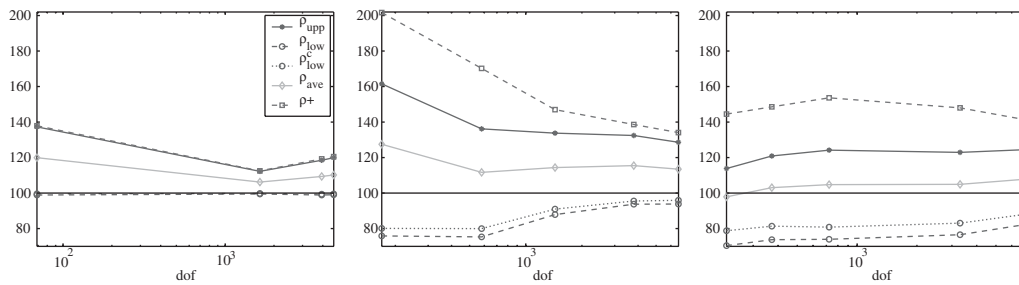


Figure 15. Example 3: performance of the estimates following an adaptive  $h$ -refinement for  $k = 1$  (left),  $k = 3$  (centre) and  $k = 4$  (right).

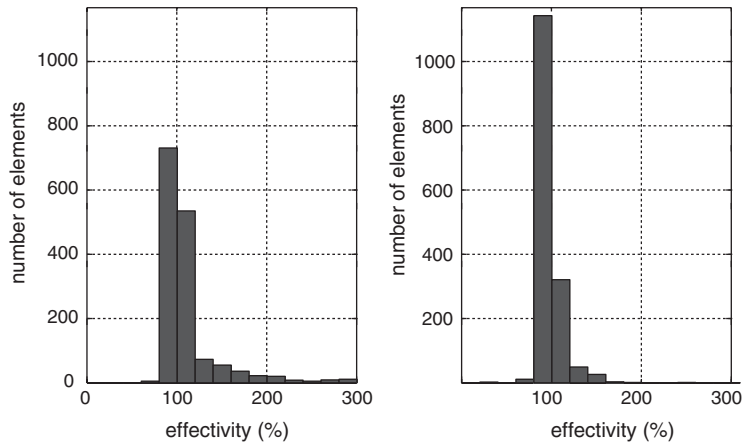


Figure 16. Example 3: Histogram representing the occurrences of the local effectivity index for  $e_2$  (left) and for the proposed strategy (right).

since the values are close to 100%. As expected, the second Bank and Weiser estimator  $e_2$  produces local estimates which overestimate almost everywhere the exact error. The local corrected estimates, as expected, underestimate the exact error. The bound property for the global error is then reproduced locally in most elements.

## 7. CONCLUDING REMARKS

A simple postprocessing strategy has been presented to recover lower bound estimates from standard residual estimators producing upper bounds of the error. The main idea is to smooth the discontinuous estimate  $e_{\text{est}}$  and to obtain a continuous approximation  $e_{\text{cont}}$  to the reference error  $e_{\text{ref}}$ . A lower bound of the error is computed using  $e_{\text{cont}}$ .

For the pure diffusion problem (when the reaction term in the PDE vanishes) the estimate  $e_{\text{est}}$  is determined up to a local (element by element) constant. In order to improve the postprocessing in this situation the local arbitrary constants are found such that the sharpest lower bound is obtained.

Numerical experiments show that the proposed strategy furnishes sharp lower estimates, of better quality than the original upper ones.

The presented strategy may be used in the framework of error estimation for outputs of interest, where upper and lower bounds of the energy error measure are required.

## REFERENCES

1. Bank RE, Weiser A. Some a posteriori error estimators for elliptic partial differential equations. *Mathematics of Computation* 1985; **44**:283–301.
2. Ainsworth M, Oden JT. A unified approach to a posteriori error estimation using element residual methods. *Numerische Mathematik* 1993; **65**:25–30.
3. Ladevèze P, Leguillon D. Error estimation procedures in the finite element method and applications. *SIAM Journal on Numerical Analysis* 1983; **20**:485–509.
4. Ladevèze P, Pelle JP, Rougeot PH. Error estimation and mesh optimization for classical finite elements. *Engineering Computations* 1991; **8**:69–80.
5. Díez P, Egozcue JJ, Huerta A. A posteriori error estimation for standard finite element analysis. *Computer Methods in Applied Mechanics and Engineering* 1998; **163**:141–157.
6. Huerta A, Díez P. Error estimation including pollution assessment for nonlinear finite element analysis. *Computer Methods in Applied Mechanics and Engineering* 2000; **181**:21–41.
7. Oden JT, Prudhomme S, Westermann TA, Bass JM. Progress on practical methods of error estimation for engineering calculations. In *Proceedings of the 2nd European Conference on Computational Mechanics*, ECCOMAS & IACM, Cracow, 2001.
8. Paraschivoiu M, Peraire J, Patera A. A posteriori finite element bounds for linear-functional outputs of elliptic partial differential equations. *Computer Methods in Applied Mechanics and Engineering* 1997; **150**:289–312.
9. Ainsworth M, Oden JT. *A Posteriori Error Estimation in Finite Element Analysis* (1st edn). Wiley: Chichester, 2000.
10. Babuška I, Strouboulis T, Gangaraj SK. Guaranteed computable bounds for the exact error in the finite element solution. Part I: One dimensional problem. *Computer Methods in Applied Mechanics and Engineering* 1999; **176**:51–79.
11. Strouboulis T, Babuška I, Gangaraj SK. Guaranteed computable bounds for the exact error in the finite element solution. Part II: bounds for the energy norm of the error in two dimensions. *International Journal for Numerical Methods in Engineering* 2000; **47**:427–475.
12. Prudhomme S, Oden JT. Simple techniques to improve the reliability of a posteriori error estimates for finite element approximations. In *Proceedings of the 2nd European Conference on Computational Mechanics*, ECCOMAS & IACM, Cracow, 2001.
13. Díez P, Parés N, Huerta A. Simple assessment of the effectivity index of residual a posteriori error estimators: recovering lower bounds. In *Proceedings of the 2nd European Conference on Computational Mechanics*, ECCOMAS & IACM, Cracow, 2001.
14. Díez P, Huerta A. A unified approach to remeshing strategies for finite element  $h$ -adaptivity. *Computer Methods in Applied Mechanics and Engineering* 1999; **176**:215–229.
15. Díez P, Parés N, Huerta A. Obtention de bornes supérieures et inférieures de l'erreur avec des estimateurs résiduels. In *Proceedings of the 5th 'Colloque National en Calcul des Structures'* CSMA, Giens, 2001.

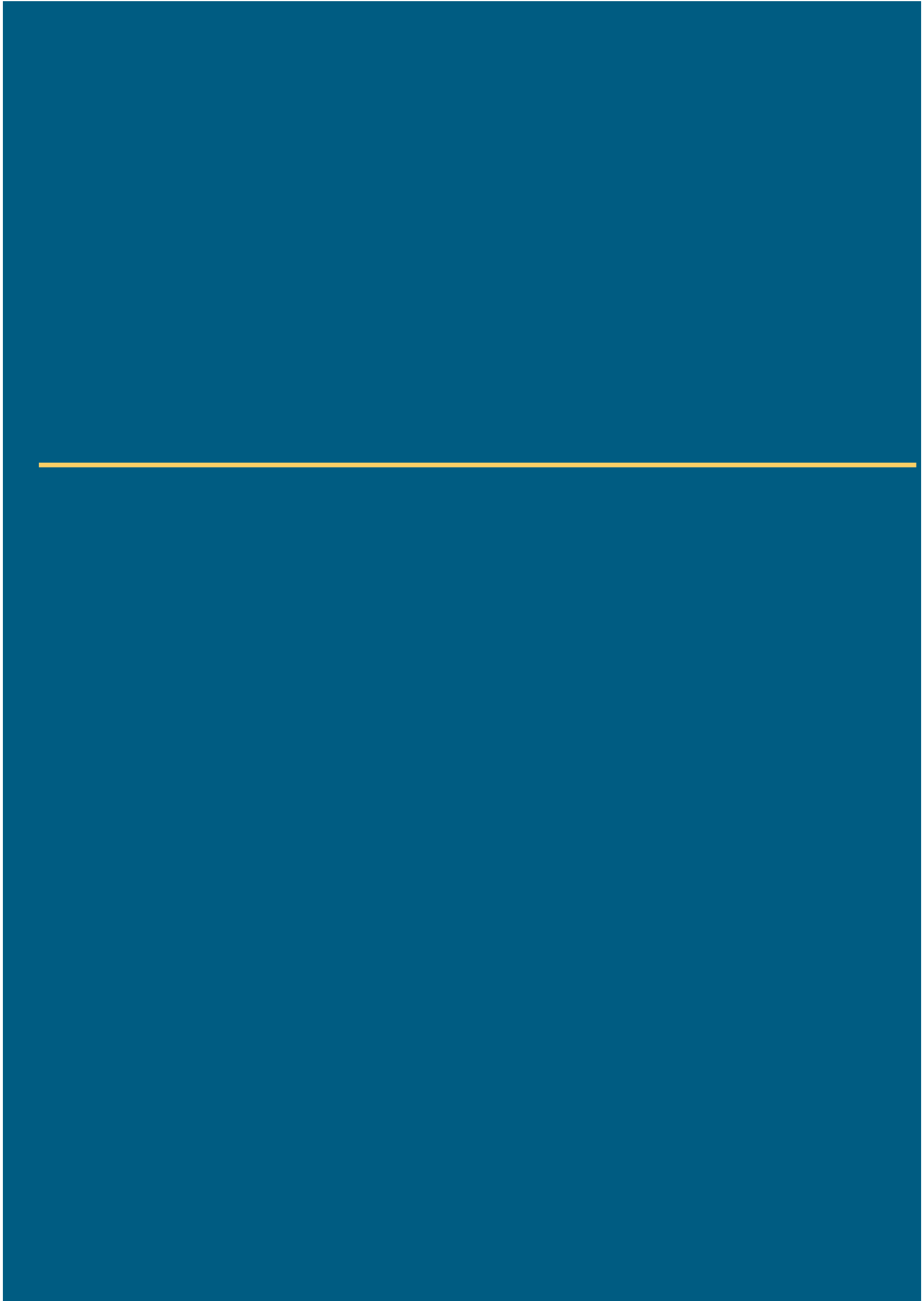
# Subdomain-based flux-free a posteriori error estimators

Parés N., Díez P. and Huerta A.

---

*Computer Methods in Applied  
Mechanics and Engineering*

**194.** In press.





Available online at [www.sciencedirect.com](http://www.sciencedirect.com)

SCIENCE @ DIRECT®

Comput. Methods Appl. Mech. Engrg. xxx (2005) xxx–xxx

**Computer methods  
in applied  
mechanics and  
engineering**

[www.elsevier.com/locate/cma](http://www.elsevier.com/locate/cma)

# Subdomain-based flux-free a posteriori error estimators

Núria Parés, Pedro Díez, Antonio Huerta \*

*Laboratori de Càlcul Numèric, Departament de Matemàtica Aplicada III, Universitat Politècnica de Catalunya,  
Modul C2, Jordi Girona 1-3, Barcelona E-08034, Spain*

Received 15 March 2004; received in revised form 20 June 2004; accepted 27 June 2004

---

## Abstract

A new residual type flux-free error estimator is presented. It estimates upper and lower bounds of the error in energy norm. The proposed approach precludes the main drawbacks of standard residual type estimators, circumvents the need of flux-equilibration and results in a simple implementation that uses standard resources available in finite element codes. This is specially interesting for 3D applications where the implementation of this technique is as simple as in 2D. Recall that on the contrary, the complexity of the flux-equilibration techniques increases drastically in the 3D case. The bounds for the energy norm of the error are used to produce upper and lower bounds of linear functional outputs, representing quantities of engineering interest. The presented estimators demonstrate their efficiency in numerical tests producing sharp estimates both for the energy and the quantities of interest.

© 2005 Elsevier B.V. All rights reserved.

*Keywords:* Error estimation; Error bounds; Functional outputs; Engineering outputs; Goal-oriented error estimation; Residual-based estimators

---

## 1. Introduction

Assessment of functional outputs of the solution (goal-oriented error estimation) in *computational mechanics* problems is a real need in standard engineering practice. In particular, end-users of finite element codes are interested in obtaining bounds for quantities of engineering interest. Techniques providing these bounds require using error estimators in the energy norm of the solution. Bounds for quantities of interest (functional outputs) are recovered combining upper and lower bounds of the energy error for both the original problem (primal) and a dual problem (associated with the selected functional output) [1–3].

---

\* Corresponding author.

*E-mail address:* [antonio.huerta@upc.es](mailto:antonio.huerta@upc.es) (A. Huerta).

*URL:* <http://www-lacan.upc.es> (A. Huerta).

The need of obtaining reliable upper and lower bounds of the error has motivated the use of residual error estimators, which are currently the only type of estimators ensuring bounds for the error. Classical residual type estimators, which provide upper bounds of the error, require flux-equilibration procedures to properly set boundary conditions for local problems [4,2]. Flux-equilibration is performed by a complex algorithm, strongly dependent on the element type and requiring a data structure that is not natural in a standard finite element code.

The idea of using flux-free estimates, based on the partition-of-the-unity concept and using local subdomains different than elements, has been already proposed in [5–7]. The main advantage of this approach is the simplicity of the implementation. Obviously, this is specially important in the 3D case. The boundary conditions of the local problems are trivial and the usual data structure of a finite element code is directly employed. Recently, in [8], the flux-free estimates have been compared with the standard *hybrid-flux* estimates in terms of both their sharpness (effectivity) and their computational efficiency. The main conclusion of this investigation is that in most of the test cases the hybrid-flux estimates are more accurate while the overall computational cost is lower for the flux-free estimates.

This paper introduces a new flux-free error estimator improving the effectivity of previous approaches and with a further simplification in the implementation. The remainder of the paper is structured as follows. In Section 2, the model problem is described. The development of this technique is motivated by the need of assessing and bounding the error of the functional outputs of the solution. Then, in Section 3, a procedure to obtain upper and lower bounds of the energy norm is presented. Section 4 is devoted to analyze the features of the proposed estimates, including proofs of the main properties. In Section 5 the estimates introduced here are compared with previously published flux-free techniques. The energy norm estimates are used in Section 6 to assess the error in quantities of interest. Computational aspects of the proposed methodology and some implementation details are discussed in Section 7. Finally, in Section 8, the different estimators are used in four numerical examples, from a simple 2D thermal problem to a 3D mechanical test.

## 2. Statement of the problem

### 2.1. Model problem

Let  $\Omega \subset \mathbb{R}^{n_{\text{sd}}}$  be an open, bounded domain with piecewise linear boundary and  $n_{\text{sd}}$  the number of spatial dimensions. Moreover,  $\partial\Omega$  is divided in two disjoint parts  $\Gamma_D$  and  $\Gamma_N$  such that  $\bar{\Gamma}_N \cup \bar{\Gamma}_D = \partial\Omega$ ,  $\bar{\Gamma}_N \cap \bar{\Gamma}_D = \emptyset$  and  $\Gamma_D$  is a non-empty set. Let  $\mathbf{u}$  be the solution of the linear elasticity problem,

$$\begin{cases} -\nabla \cdot \boldsymbol{\sigma}(\mathbf{u}) = \mathbf{s} & \text{in } \Omega, \\ \boldsymbol{\sigma}(\mathbf{u}) \cdot \mathbf{n} = \mathbf{t} & \text{on } \Gamma_N, \\ \mathbf{u} = \mathbf{u}_D & \text{on } \Gamma_D, \end{cases} \quad (1)$$

where  $\mathbf{t}$  and  $\mathbf{u}_D$  are the imposed traction and boundary displacements, respectively.

The weak solution of this problem is  $\mathbf{u} \in \mathcal{U}$  verifying

$$a(\mathbf{u}, \mathbf{v}) = l(\mathbf{v}) \quad \forall \mathbf{v} \in \mathcal{V}, \quad (2)$$

where

$$a(\mathbf{u}, \mathbf{v}) = \int_{\Omega} \boldsymbol{\sigma}(\mathbf{u}) : \boldsymbol{\varepsilon}(\mathbf{v}) \, d\Omega, \quad l(\mathbf{v}) = \int_{\Omega} \mathbf{s} \cdot \mathbf{v} \, d\Omega + \int_{\Gamma_N} \mathbf{t} \cdot \mathbf{v} \, d\Gamma. \quad (3)$$



The usual solution and test spaces are defined  $\mathcal{U} = \{\mathbf{u} \in [\mathcal{H}^1(\Omega)]^{\text{nsd}}, \mathbf{u}|_{\Gamma_D} = \mathbf{u}_D\}$  and  $\mathcal{V} = \{\mathbf{v} \in [\mathcal{H}^1(\Omega)]^{\text{nsd}}, \mathbf{v}|_{\Gamma_D} = 0\}$ , where  $\mathcal{H}^1$  is the standard Sobolev space of square integrable functions and first derivatives. The bilinear form  $a(\cdot, \cdot)$  induces the energy norm, which is denoted by  $\|\cdot\|$ , that is,  $\|\mathbf{v}\|^2 = a(\mathbf{v}, \mathbf{v})$ .

The finite element interpolation spaces  $\mathcal{U}^H \subset \mathcal{U}$  and  $\mathcal{V}^H \subset \mathcal{V}$  are associated with a finite element mesh of characteristic size  $H$  and degree  $p$  for the complete interpolation polynomial base. The geometric support of the elements for a given mesh are open subdomains denoted by  $\Omega_k$ ,  $k = 1, \dots, n_{e1}$ , where  $\bar{\Omega} = \cup_k \bar{\Omega}_k$ . It is also assumed that different elements do not overlap, that is,  $\Omega_k \cap \Omega_l = \emptyset$  for  $k \neq l$ .

Then, the finite element solution  $\mathbf{u}_H$  which is an approximation to  $\mathbf{u}$ , lies in the finite dimensional space  $\mathcal{U}^H$  and verifies

$$a(\mathbf{u}_H, \mathbf{v}) = l(\mathbf{v}) \quad \forall \mathbf{v} \in \mathcal{V}^H. \quad (4)$$

## 2.2. Error equations and reference error

The goal of a posteriori error estimation is to assess the accuracy of the finite element solution  $\mathbf{u}_H$ , that is, to evaluate and measure the error,  $\mathbf{e} := \mathbf{u} - \mathbf{u}_H$ , which belongs to  $\mathcal{V}$ , either in the energy norm  $\|\mathbf{e}\|$ , or in a quantity of interest  $l^0(\mathbf{e})$ .

The global equation for the error is recovered from (2) replacing the exact solution  $\mathbf{u}$  by  $\mathbf{u}_H + \mathbf{e}$  and using the linearity of the first argument of  $a(\cdot, \cdot)$

$$a(\mathbf{e}, \mathbf{v}) = l(\mathbf{v}) - a(\mathbf{u}_H, \mathbf{v}) =: R^p(\mathbf{v}) \quad \forall \mathbf{v} \in \mathcal{V}, \quad (5)$$

where  $R^p(\cdot)$  stands for the weak residual associated to the finite element approximation  $\mathbf{u}_H$ .

In practice, the exact error  $\mathbf{e}$  is replaced by a reference error,  $\mathbf{e}_h$ , lying in a finite dimensional space  $\mathcal{V}^h$  much richer than the original finite element space  $\mathcal{V}^H$ . That is, the exact solution  $\mathbf{u}$  is replaced by the reference (or truth) solution  $\mathbf{u}_h$ ; consequently,  $\mathbf{u} \approx \mathbf{u}_h = \mathbf{u}_H + \mathbf{e}_h$ . The reference error is the projection of the exact error into the reference space, that is,  $\mathbf{e}_h \in \mathcal{V}^h$  is the solution of the problem

$$a(\mathbf{e}_h, \mathbf{v}) = R^p(\mathbf{v}) \quad \forall \mathbf{v} \in \mathcal{V}^h. \quad (6)$$

The direct computation of  $\mathbf{e}_h$  is computationally unaffordable because the size of the system of equations is the dimension of  $\mathcal{V}^h$ . The idea behind any implicit residual type error estimator is to solve a set of local problems instead of the global problem (6). In each of these local problems, boundary conditions must be properly defined in order to obtain a good approximation of the error and to ensure solvability.

## 2.3. Estimation of outputs of interest

Attention is usually centered in bounding output quantities  $l^0(\mathbf{u})$ , where  $l^0(\cdot)$  is a linear functional, see for instance [1,9,10,3,11]. These strategies introduce a dual (or adjoint) problem with respect to the selected output. The weak form of the dual problem reads: find  $\boldsymbol{\psi} \in \mathcal{V}$  such that

$$a(\mathbf{v}, \boldsymbol{\psi}) = l^0(\mathbf{v}) \quad \forall \mathbf{v} \in \mathcal{V}.$$

The finite element approximation of the dual problem is  $\boldsymbol{\psi}_H \in \mathcal{V}^H$  such that

$$a(\mathbf{v}, \boldsymbol{\psi}_H) = l^0(\mathbf{v}) \quad \forall \mathbf{v} \in \mathcal{V}^H. \quad (7)$$

Finally, the dual reference error is  $\boldsymbol{\epsilon}_h \in \mathcal{V}^h$ , such that

$$a(\mathbf{v}, \boldsymbol{\epsilon}_h) = l^0(\mathbf{v}) - a(\mathbf{v}, \boldsymbol{\psi}_H) =: R^D(\mathbf{v}) \quad \forall \mathbf{v} \in \mathcal{V}^h, \quad (8)$$

where  $R^D$  is the weak residual associated with  $\boldsymbol{\psi}_H$ .

If  $\mathbf{v}$  is replaced by  $\mathbf{e}^h$  in (8), then using Galerkin orthogonality and the parallelogram identity, the following representation of  $l^\mathcal{O}(\mathbf{e}_h)$  can be obtained

$$l^\mathcal{O}(\mathbf{e}_h) = a(\mathbf{e}_h, \boldsymbol{\epsilon}_h) = \frac{1}{4} \left\| \kappa \mathbf{e}_h + \frac{1}{\kappa} \boldsymbol{\epsilon}_h \right\|^2 - \frac{1}{4} \left\| \kappa \mathbf{e}_h - \frac{1}{\kappa} \boldsymbol{\epsilon}_h \right\|^2 \quad (9)$$

for any arbitrary scalar parameter  $\kappa$ . To simplify the notation the arguments in the squared norms of the r.h.s. in (9) are denoted by  $\mathbf{z}_h^\pm = \kappa \mathbf{e}_h \pm \frac{1}{\kappa} \boldsymbol{\epsilon}_h$ .

In fact, in order to bound the output of the error,  $l^\mathcal{O}(\mathbf{e}_h)$ , the r.h.s. of (9) indicates that it is sufficient to bound the energy norm of  $\mathbf{z}_h^+$  and  $\mathbf{z}_h^-$ , (i.e. the energy norm of linear combinations of  $\mathbf{e}_h$  and  $\boldsymbol{\epsilon}_h$ ).

Define  $E_u[\mathbf{v}]$  and  $E_l[\mathbf{v}]$  as the upper and lower bound of  $\|\mathbf{v}\|^2$ , respectively. Note that  $E_u[\mathbf{v}]$  and  $E_l[\mathbf{v}]$  are not functions; instead, it is a convenient notation of the upper and lower bounds of  $\|\mathbf{v}\|^2$ . Thus, once the bounds for  $\|\mathbf{z}_h^\pm\|^2$  are computed, namely

$$E_l[\mathbf{z}_h^\pm] \leq \|\mathbf{z}_h^\pm\|^2 \leq E_u[\mathbf{z}_h^\pm],$$

the output of the error is readily bounded as

$$\frac{1}{4} E_l[\mathbf{z}_h^+] - \frac{1}{4} E_u[\mathbf{z}_h^-] \leq l^\mathcal{O}(\mathbf{e}_h) \leq \frac{1}{4} E_u[\mathbf{z}_h^+] - \frac{1}{4} E_l[\mathbf{z}_h^-]. \quad (10)$$

This procedure is summarized in Fig. 1 where bounds for the output of interest of the reference approximation,  $l^\mathcal{O}(\mathbf{u}_h)$ , are also shown:  $l^\mathcal{O}(\mathbf{u}_H)$  is added to each term of inequality (10). Next section introduces a methodology to obtain both upper and lower bound error estimates in energy norm. This approach is then used to compute  $E_u[\mathbf{z}_h^+]$ ,  $E_u[\mathbf{z}_h^-]$ ,  $E_l[\mathbf{z}_h^+]$  and  $E_l[\mathbf{z}_h^-]$ .

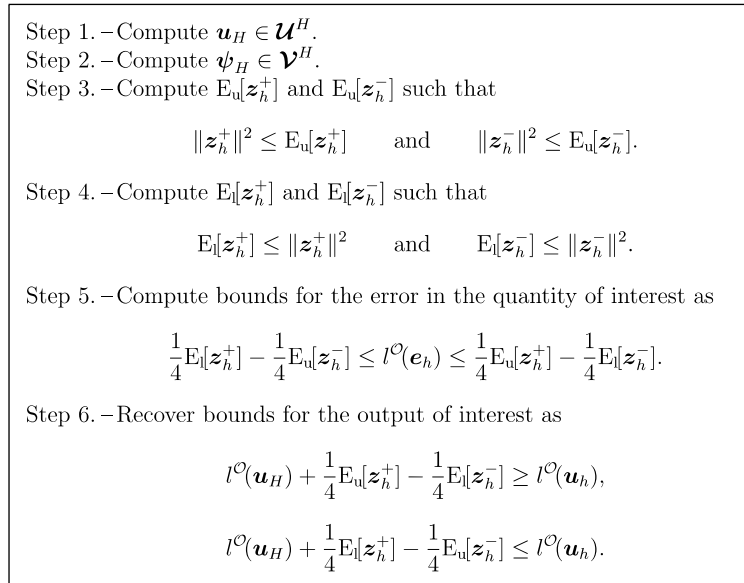


Fig. 1. Strategy to obtain bounds for the quantity of interest  $l^\mathcal{O}(\mathbf{u}_h)$ .

### 3. Estimation of the energy norm of the error

In this section, error estimates yielding upper and lower bounds of the energy norm are presented. For the sake of simplicity, the presentation concerns only the primal problem. The methodology is general and it is also applicable to the dual problem or to linear combinations of both.

#### 3.1. Definitions and preliminaries

Let  $\mathbf{x}^i$ ,  $i = 1, \dots, n_{np}$  denote the vertices of the elements in the computational mesh (thus linked to  $\mathcal{U}^H$ ) and  $\phi^i$  the corresponding linear (or bilinear or trilinear) shape functions, which are such that  $\phi^i(\mathbf{x}^j) = \delta_{ij}$ . The support of  $\phi^i$  is denoted by  $\omega^i$  and it is called the star centered in, or associated with, vertex  $\mathbf{x}^i$ .

It is important to recall that the linear shape functions based on the vertices are a *partition of unity*. Using this essential property and the linearity of the weak residue  $R^P(\cdot)$ , defined in (5), for every  $\mathbf{v} \in [\mathcal{H}^1(\Omega)]^{n_{sd}}$  the following equality holds

$$R^P(\mathbf{v}) = R^P\left(\sum_{i=1}^{n_{np}} \phi^i \mathbf{v}\right) = \sum_{i=1}^{n_{np}} R^P(\phi^i \mathbf{v}). \quad (11)$$

Note that  $R^P(\phi^i \mathbf{v})$  vanishes if  $\text{supp } \mathbf{v} \cap \omega^i = \emptyset$ , because  $\omega^i$  is the support of  $\phi^i$ . Therefore, the residue is decomposed into local contributions over each star. This basic property is the key idea to define residual estimators based in stars. Similar approaches have been used in Refs. [12,5–7].

Let  $\mathcal{V}_{\omega^i}^h$ , and  $\mathcal{V}_{\omega^i}^H$  denote the local restrictions of the reference and finite element spaces to the star  $\omega^i$ , that is,

$$\mathcal{V}_{\omega^i}^h := \mathcal{V}^h \cap [\mathcal{H}^1(\omega^i)]^{n_{sd}} \quad \text{and} \quad \mathcal{V}_{\omega^i}^H := \mathcal{V}^H \cap [\mathcal{H}^1(\omega^i)]^{n_{sd}}.$$

Formally any function  $\mathbf{v} \in \mathcal{V}_{\omega^i}^h$  (in particular,  $\mathbf{v} \in \mathcal{V}_{\omega^i}^H \subset \mathcal{V}_{\omega^i}^h$ ) is not defined in the whole domain  $\Omega$  but only in the star  $\omega^i$ . However, here any  $\mathbf{v} \in \mathcal{V}_{\omega^i}^h$  is naturally extended to  $\Omega$  by setting the values outside  $\omega^i$  to zero. Thus, functions in  $\mathcal{V}_{\omega^i}^h$  are continuous in  $\omega^i$  but generally discontinuous across the boundary of the star  $\omega^i$ .

The local restriction  $\mathcal{V}^h$  to the element  $\Omega_k$ ,  $\mathcal{V}_{\Omega_k}^h := \mathcal{V}^h \cap [\mathcal{H}^1(\Omega_k)]^{n_{sd}}$ , is also extended to  $\Omega$  in the same way. This induces the *broken space*, namely

$$\mathcal{V}_{\text{brok}}^h := \bigoplus_{k=1}^{n_{el}} \mathcal{V}_{\Omega_k}^h.$$

Note that functions in  $\mathcal{V}_{\text{brok}}^h$  may present discontinuities across the inter-element edges (or faces) and that  $\mathcal{V}_{\omega^i}^h \subset \mathcal{V}_{\text{brok}}^h$ .

The bilinear form  $a(\cdot, \cdot)$  and the energy norm are generalized to accept *broken* functions in its arguments; that is, for  $\mathbf{v}$  and  $\mathbf{w} \in \mathcal{V}_{\text{brok}}^h$ ,

$$a(\mathbf{v}, \mathbf{w}) := \sum_{k=1}^{n_{el}} a_{\Omega_k}(\mathbf{v}, \mathbf{w}) \quad \text{and} \quad \|\mathbf{v}\|^2 := \sum_{k=1}^{n_{el}} \|\mathbf{v}\|_k^2,$$

where  $a_{\Omega_k}(\cdot, \cdot)$  is the restriction of the bilinear form  $a(\cdot, \cdot)$  to the element  $\Omega_k$  and  $\|\mathbf{v}\|_k^2 = a_{\Omega_k}(\mathbf{v}, \mathbf{v})$ .

For further developments it is also necessary to introduce the nodal projections of any function in  $\mathcal{V}$  onto the finite element space,  $\mathcal{V}^H$ , and the reference space,  $\mathcal{V}^h$ . That is,  $\pi^H : \mathcal{V} \rightarrow \mathcal{V}^H$  such that  $\pi^H \mathbf{v}(\hat{\mathbf{x}}^i) = \mathbf{v}(\hat{\mathbf{x}}^i)$  where  $\hat{\mathbf{x}}^i$  denote every node on the finite element mesh, and  $\pi^h : \mathcal{V} \rightarrow \mathcal{V}^h$  such that  $\pi^h \mathbf{v}(\hat{\mathbf{x}}^i) = \mathbf{v}(\hat{\mathbf{x}}^i)$  where  $\hat{\mathbf{x}}^i$  denote now the nodal points of the reference mesh.

### 3.2. Upper bound estimate of the reference error

The strategy to compute upper bound estimates of the reference error,  $E_u[e_h]$ , consist in, first, the evaluation of the finite element solution  $\mathbf{u}_H$ , which is necessary to compute the residue  $R^P$  and, second, the appraisal of the local estimates  $\tilde{\mathbf{e}}^{\omega^i} \in \mathcal{V}_{\omega^i}^h$  solving problems in each star  $\omega^i$

$$a_{\omega^i}(\tilde{\mathbf{e}}^{\omega^i}, \mathbf{v}) = R^P(\phi^i(\mathbf{v} - \pi^H \mathbf{v})) \quad \forall \mathbf{v} \in \mathcal{V}_{\omega^i}^h, \quad (12)$$

where  $a_{\omega^i}(\cdot, \cdot)$  is the restriction of the bilinear form  $a(\cdot, \cdot)$  to the star  $\omega^i$ . Then, adding the local estimates, which have been extended into  $\mathcal{V}_{\text{brok}}^h$ , a global estimate  $\tilde{\mathbf{e}} \in \mathcal{V}_{\text{brok}}^h$  is obtained,

$$\tilde{\mathbf{e}} := \sum_{i=1}^{n_{\text{np}}} \tilde{\mathbf{e}}^{\omega^i} \quad (13)$$

and the upper bound of the energy norm of the reference error is recovered computing the norm of the estimate  $\tilde{\mathbf{e}}$ , that is,  $E_u[e_h] := \|\tilde{\mathbf{e}}\|^2 \geq \|e_h\|^2$ . Fig. 2 describes this strategy in four steps.

Note that the error estimator described above does not require any computation of fluxes (stresses) along the boundary of the elements (it is *flux-free*).

**Remark 1.** In the r.h.s. of (12) the projection  $\pi^H$  has been introduced in order to equilibrate the local problem and ensure its solvability. This is analyzed in Section 4.1. However, for scalar problems and mechanical problems with high-order elements (at least quadratic) the r.h.s. does not require the projection. That is, Eq. (12) reduces to

$$a_{\omega^i}(\tilde{\mathbf{e}}^{\omega^i}, \mathbf{v}) = R^P(\phi^i \mathbf{v}) \quad \forall \mathbf{v} \in \mathcal{V}_{\omega^i}^h.$$

**Remark 2.** In Section 7 another expression for the r.h.s. of (12) is proposed to drastically simplify the practical implementation of this estimator.

### 3.3. Lower bound estimates

The upper bound estimate of the squared energy norm,  $E_u[e_h]$ , is associated with the estimate  $\tilde{\mathbf{e}}$  of the error function. The upper bound property is intrinsically related with the broken (discontinuous) nature

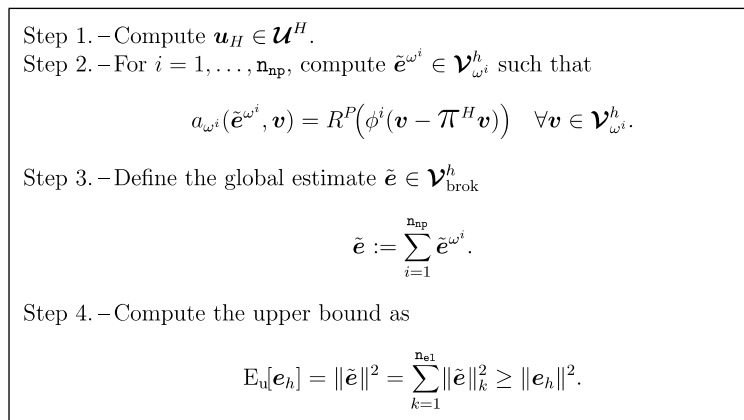


Fig. 2. Upper bound for the squared energy norm of the reference error.

of  $\tilde{e}$ . On the contrary, a lower bound estimate is easily recovered from a continuous estimate of the error function, see [13]. Thus, once  $\tilde{e}$  is obtained, a continuous estimate of the error function,  $\tilde{e}_{\text{cont}}$ , is computed by simple post-processing. Two different alternatives can be considered to compute  $\tilde{e}_{\text{cont}}$  from  $\tilde{e}$ . First, the strategy presented in detail in [13] and valid for any discontinuous estimate  $\tilde{e}$  (discontinuous across inter-element edges or faces) can be readily implemented. It averages the discontinuities of the function across inter-element edges/faces and produces a continuous function that belongs to  $\mathcal{V}^h$ . Second, the weighting strategy, where the continuous estimate is obtained from

$$\tilde{e}_{\text{cont}} := \pi^h \left( \sum_{i=1}^{n_{\text{sp}}} \phi^i \tilde{e}^{\omega^i} \right). \quad (14)$$

This approach uses the fact that local estimates  $\tilde{e}^{\omega^i}$  are continuous in each star. The discontinuities of  $\tilde{e}^{\omega^i}$  on the boundary of each star  $\omega^i$  are smoothed by multiplying by  $\phi^i$ , which vanishes along the boundary of  $\omega^i$ . Consequently, this is the natural choice for the estimates presented in this paper. The projection into the reference mesh  $\mathcal{V}^h$  ensures that the evaluation of  $R^P(\tilde{e}_{\text{cont}})$  is easily performed. Note that  $\phi^i \tilde{e}^{\omega^i}$  may not belong to  $\mathcal{V}^h$ .

For both averaging strategies, a lower bound,  $E_1[e_h]$ , of the energy norm of the reference error is obtained from  $\tilde{e}_{\text{cont}}$  as

$$E_1[e_h] := \frac{(R^P(\tilde{e}_{\text{cont}}))^2}{\|\tilde{e}_{\text{cont}}\|^2} \leq \|e_h\|^2. \quad (15)$$

Moreover, in order to improve the quality of the estimate the global enhancement strategy proposed in [14] can be implemented. First,  $\tilde{e}^G \in \mathcal{V}^H$  is computed solving

$$a(\tilde{e}^G, \mathbf{v}) = -a(\tilde{e}_{\text{cont}}, \mathbf{v}) \quad \forall \mathbf{v} \in \mathcal{V}^H \quad (16)$$

and then, the lower bound given in (15) is improved using  $\|\tilde{e}^G\|^2$  as

$$E_1^G[e_h] := \frac{(R^P(\tilde{e}_{\text{cont}}))^2}{\|\tilde{e}_{\text{cont}}\|^2 - \|\tilde{e}^G\|^2} \leq \|e_h\|^2. \quad (17)$$

This strategy is summarized in Fig. 3.

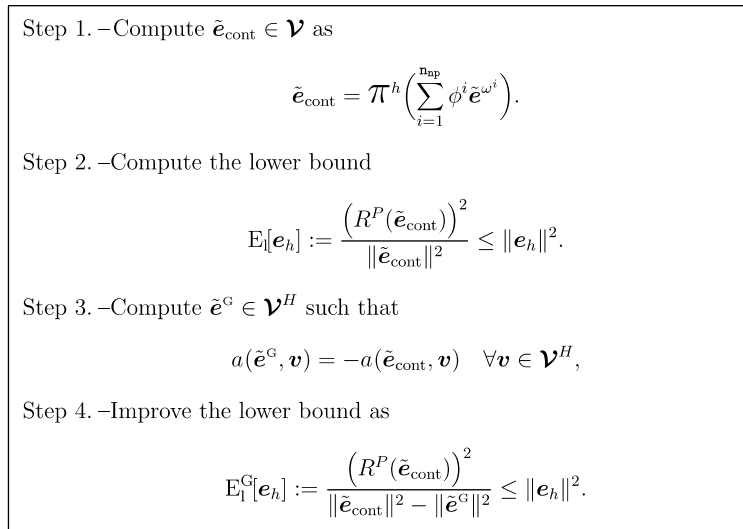


Fig. 3. Lower bounds for the squared energy norm of the error.

**Remark 3.** The evaluation of  $\tilde{z}^G$  from Eq. (16) is equivalent to the resolution of a system of equations with the same matrix used to compute  $\mathbf{u}_H$ , see Eq. (4).

#### 4. Analysis and properties of the proposed estimates

The upper bound estimate  $E_u[e_h]$  is obtained without any flux recovery or flux splitting technique. The effect of the flux jumps across each edge of the mesh is implicitly taken into account because the support of the local problems are the stars, which include the inter-element edges/faces. There is no need to compute and postprocess fluxes of the finite element solution,  $\mathbf{u}_H$ , along the inter-element edges/faces. Thus, the proposed estimate has two very attractive features:

- (1) there is no need to compute fluxes and flux jumps along the element boundaries, and
- (2) there is no need to perform any flux equilibration.

Consequently, it is especially well suited to assess the error in a 3D framework, where the cost of computing the boundary fluxes and their equilibration is usually extremely large. Moreover, it is also important to notice that flux-free estimators only require the standard data structure already present in any standard finite element code. In particular, there is no need to have structured the information on edges/faces for the evaluation of the fluxes. The remainder of this section is devoted to analyze the main properties of the estimates introduced above.

##### 4.1. Solvability of the local error equation

The local Eq. (12) is solved in each star  $\omega^i$  in order to compute the local estimate  $\tilde{e}^{\omega^i}$ . Note that the r.h.s. term of (12),  $R^P(\phi^i(\mathbf{v} - \pi^H \mathbf{v}))$ , does not coincide with the obvious decomposition of the residue given in Eq. (11), namely

$$a_{\omega^i}(\tilde{e}^{\omega^i}, \mathbf{v}) = R^P(\phi^i \mathbf{v}) \quad \forall \mathbf{v} \in \mathcal{V}_{\omega^i}^h. \quad (18)$$

The term  $R^P(\phi^i \mathbf{v})$  has been replaced in (12) by  $R^P(\phi^i(\mathbf{v} - \pi^H \mathbf{v}))$ . This is done to ensure the solvability of the local equation.

**Theorem 4.** *The local problem for the estimate  $\tilde{e}^{\omega^i}$ ,*

$$a_{\omega^i}(\tilde{e}^{\omega^i}, \mathbf{v}) = R^P(\phi^i(\mathbf{v} - \pi^H \mathbf{v})) \quad \forall \mathbf{v} \in \mathcal{V}_{\omega^i}^h$$

*is solvable.*

This Theorem is based on the following one, which can be found in [15, Thm. 9.2.30],

**Theorem 5.** *Let  $\mathcal{V}$  be a Hilbert space and  $a(\cdot, \cdot)$  be a bilinear form acting on  $\mathcal{V} \times \mathcal{V}$ . Let also  $\mathcal{V} = \ker a \oplus \widehat{\mathcal{V}}$  be a decomposition of  $\mathcal{V}$ , that is, for any given  $\mathbf{v} \in \mathcal{V}$ , there exists a unique pair  $(\mathbf{v}_a, \hat{\mathbf{v}}) \in \ker a \times \widehat{\mathcal{V}}$  such that  $\mathbf{v} = \mathbf{v}_a + \hat{\mathbf{v}}$ , where*

$$\ker a := \{\mathbf{v}_a \in \mathcal{V} | a(\mathbf{v}_a, \mathbf{w}) = 0 \quad \forall \mathbf{w} \in \mathcal{V}\}.$$

*Assume also that the bilinear form  $a(\cdot, \cdot)$  is coercive on  $\widehat{\mathcal{V}}$ , that is*

$$\exists \gamma > 0 \text{ such that } a(\hat{\mathbf{v}}, \hat{\mathbf{v}}) \geq \gamma \|\hat{\mathbf{v}}\|^2 \quad \forall \hat{\mathbf{v}} \in \widehat{\mathcal{V}}.$$

Then, the variational problem: find  $u \in \mathcal{V}$  such that

$$a(\mathbf{u}, \mathbf{v}) = l(\mathbf{v}) \quad \forall \mathbf{v} \in \mathcal{V}$$

is solvable if and only if the following compatibility condition holds:

$$l(\mathbf{v}) = 0 \quad \forall \mathbf{v} \in \ker a.$$

Solvability of a variational problem depends on the verification of the compatibility condition for the functions in the kernel of the bilinear operator. Thus, solvability of Eq. (12) depends on the model problem at hand. Here the mechanical problem is discussed but the following remark is concerned with scalar equations.

**Remark 6.** Consider the scalar diffusion–reaction equation. The bilinear form for this problem is

$$a(u, v) = \int_{\Omega} v \nabla u \cdot \nabla v + \mu u v \, d\Omega$$

for a strictly positive real coefficient  $v \in \mathcal{L}^{\infty}(\Omega)$  and a non-negative real coefficient  $\mu \in \mathcal{L}^{\infty}(\Omega)$ , and its restriction to a star  $\omega^i$  is, as previously, denoted by  $a_{\omega^i}(u, v)$ . A strictly positive reaction term in  $a_{\omega^i}(u, v)$  ensures the solvability of local Eq. (12) since the  $\ker a_{\omega^i} = \emptyset$ . For  $\mu|_{\omega^i} = 0$ , the kernel of the bilinear operator  $a_{\omega^i}(\cdot, \cdot)$  is the one dimensional space of constants,  $\mathbb{P}^0(\omega^i)$ . Then, Eq. (12) is solvable if and only if the compatibility condition holds, namely

$$R^P(\phi^i c) = c R^P(\phi^i) = 0 \quad \forall c \in \mathbb{P}^0(\omega^i),$$

which follows from the orthogonality of the primal residual to the finite element space  $\mathcal{V}^H$ , since  $\phi^i \in \mathcal{V}^H$ .

The bilinear form for the elasticity problem is defined in (3) and the kernel of its restriction to  $\omega^i$ ,  $a_{\omega^i}(\cdot, \cdot)$ , is defined by the solid rigid motions, that is, the zero energy modes. In 1D, the rigid body motions are only translations, that is, the one dimensional space of constants,  $\mathbb{P}^0(\omega^i)$ . In this case, as in the scalar (thermal) problem, for  $c \in \mathbb{P}^0(\omega^i)$ ,  $R^P(\phi^i c) = c R^P(\phi^i) = 0$  due to the Galerkin orthogonality and therefore the compatibility equation holds and Eq. (12) is solvable.

However for 2D and 3D mechanical problems, the solid rigid motions include also rotations. For instance in a 2D setup, the kernel of  $a_{\omega^i}(\cdot, \cdot)$  is a space of three dimensions generated by two translations  $\mathbf{t}_x$  and  $\mathbf{t}_y$  and one rotation  $\theta$ . The rotation is a linear function and consequently  $\phi^i \theta$  does not always belong to  $\mathcal{V}^H$  (for instance, for linear triangular elements,  $\phi^i \theta \notin \mathcal{V}^H$ ), and hence  $R^P(\phi^i \theta)$  is not necessarily zero. Thus, since the compatibility condition does not hold in general, it cannot be guaranteed that Eq. (12) is solvable. From a mechanical viewpoint, the forces associated with  $R^P(\phi^i \theta)$  are not equilibrated (the sum of forces is zero but the sum of moments does not vanish).

For domains with piecewise linear boundaries, a natural way to circumvent this problem is to consider higher-order Lagrange elements. If  $\phi^i$  are first-order Lagrange (linear, bilinear or trilinear) shape functions associated to the set of vertices of the finite mesh, since  $\theta$  is also a linear combination of first-order Lagrange shape functions then  $\phi^i \theta \in \mathcal{V}^H$  and thus the compatibility condition is verified. Thus, Eq. (12) is solvable.

However, this is no longer valid for first-order Lagrange elements nor for domains with, curved boundaries where rotation,  $\theta$ , is not characterized by a linear combination of first-order shape functions (recall that an isoparametric transformation is used for curved elements). In this case, an alternative is to correct the r.h.s. of (12) to ensure that the compatibility condition is verified also for rotations.

Since  $R^P(\mathbf{v}) = R^P(\mathbf{v} - \pi^H \mathbf{v})$  for all  $\mathbf{v} \in \mathcal{V}$  by Galerkin orthogonality, the partition defined by (11) can be redefined as

$$R^P(\mathbf{v}) = \sum_{i=1}^{n_{np}} R^P(\phi^i(\mathbf{v} - \boldsymbol{\pi}^H \mathbf{v})),$$

which leads to Eq. (12). In this case, any rigid solid motion  $\mathbf{v}_{rm}$  belongs to  $\mathcal{V}^H$ , thus  $\mathbf{v}_{rm} - \boldsymbol{\pi}^H \mathbf{v}_{rm} = 0$  and consequently problem (12) is solvable.

#### 4.2. The upper bound property

The following results summarize the basic property that most residual type error estimators based in a flux-equilibration technique verify, see Refs. [16,17], and also proves the upper bound property of this type of estimates.

**Lemma 7.** Any estimate  $\tilde{\mathbf{e}} \in \mathcal{V}_{\text{brok}}^h$  verifying the weak error equation

$$a(\tilde{\mathbf{e}}, \mathbf{v}) = R^P(\mathbf{v}) \quad \forall \mathbf{v} \in \mathcal{V}^h \quad (19)$$

is such that the norm of  $\tilde{\mathbf{e}}$  is an upper bound of the energy norm of the reference error, that is

$$\|\tilde{\mathbf{e}}\|^2 \geq \|e_h\|^2.$$

**Proof.** First the following trivial expansion is performed

$$0 \leq \|e_h - \tilde{\mathbf{e}}\|^2 = \|e_h\|^2 + \|\tilde{\mathbf{e}}\|^2 - 2a(\tilde{\mathbf{e}}, e_h).$$

Now, replacing  $\mathbf{v}$  by  $e_h$  in (19) and using equality  $R^P(e_h) = \|e_h\|^2$ , see (6), the upper bound property is obtained as follows

$$0 \leq \|e_h - \tilde{\mathbf{e}}\|^2 = \|e_h\|^2 + \|\tilde{\mathbf{e}}\|^2 - 2R^P(e_h) = \|\tilde{\mathbf{e}}\|^2 - \|e_h\|^2. \quad \square$$

Thus, to prove that  $E_u[e_h]$  is an upper bound of the energy norm of the error, it is only necessary to check that the global estimate  $\tilde{\mathbf{e}}$ , defined in (13), verifies Eq. (19).

**Theorem 8.** The estimate  $\tilde{\mathbf{e}} = \sum_{i=1}^{n_{np}} \tilde{\mathbf{e}}^{\omega^i}$ , where  $\tilde{\mathbf{e}}^{\omega^i}$  is the solution of the local problem given in (12), is such that

$$E_u[e_h] = \|\tilde{\mathbf{e}}\|^2 \geq \|e_h\|^2.$$

**Proof.** Using Eqs. (12) and (13) together with Galerkin orthogonality

$$\begin{aligned} a(\tilde{\mathbf{e}}, \mathbf{v}) &= \sum_{i=1}^{n_{np}} a(\tilde{\mathbf{e}}^{\omega^i}, \mathbf{v}) = \sum_{i=1}^{n_{np}} a_{\omega^i}(\tilde{\mathbf{e}}^{\omega^i}, \mathbf{v}) = \sum_{i=1}^{n_{np}} R^P(\phi^i(\mathbf{v} - \boldsymbol{\pi}^H \mathbf{v})) = R^P\left(\sum_{i=1}^{n_{np}} \phi^i(\mathbf{v} - \boldsymbol{\pi}^H \mathbf{v})\right) \\ &= R^P(\mathbf{v}) - R^P(\boldsymbol{\pi}^H \mathbf{v}) = R^P(\mathbf{v}) \quad \forall \mathbf{v} \in \mathcal{V}^h. \end{aligned}$$

And the proof is concluded using Lemma 7.  $\square$

#### 4.3. Lower bound by post-processing

The following theorem, see [13] for a detailed proof, shows that every continuous function yields a lower bound of the energy norm of the error. In particular those obtained by post-processing as indicated in Section 3.3. Obviously, for an arbitrary estimate  $\tilde{\mathbf{e}}_{\text{cont}}$ , the corresponding lower bound may have a very poor quality. The best choice for  $\tilde{\mathbf{e}}_{\text{cont}}$  is either  $e$  or  $e_h$ , in order to obtain  $E_l$  equal to  $\|e\|^2$  or  $\|e_h\|^2$ . Therefore, to



obtain sharp lower bounds, the estimate  $\tilde{e}_{\text{cont}}$  must be a good approximation of the actual error (either exact or reference).

**Theorem 9.** For any  $\tilde{e}_{\text{cont}} \in \mathcal{V}$ , a lower bound of the energy norm of the exact error is recovered as

$$0 \leq E_1[\mathbf{e}] := \frac{(R^P(\tilde{e}_{\text{cont}}))^2}{\|\tilde{e}_{\text{cont}}\|^2} \leq \|\mathbf{e}\|^2.$$

Moreover, if  $\tilde{e}_{\text{cont}} \in \mathcal{V}^h \subset \mathcal{V}$ , the lower bound is also a lower bound with respect to the energy norm of the reference error, that is,

$$0 \leq E_1[\mathbf{e}_h] := \frac{(R^P(\tilde{e}_{\text{cont}}))^2}{\|\tilde{e}_{\text{cont}}\|^2} \leq \|\mathbf{e}_h\|^2 \leq \|\mathbf{e}\|^2.$$

#### 4.4. Enhancing the lower bound

The continuous function  $\tilde{e}_{\text{cont}}$  is obtained by performing only local computations, consequently the corresponding estimate  $E_1$  does not account for pollution errors. The unestimated part of error,  $\mathbf{e} - \tilde{e}_{\text{cont}}$  includes the pollution effects and it is denoted as global error. In order to assess pollution, the equation for the global error is solved on the coarse mesh following the methodology proposed in [14]. Thus,  $\tilde{e}^G \in \mathcal{V}^H$  is computed using Eq. (16) and the enhanced lower bound estimate,  $E_1^G[\mathbf{e}_h]$ , is obtained using (17). The following theorem states that  $E_1^G[\mathbf{e}_h]$  is also a lower bound of the squared error energy norm.

**Theorem 10.** Let  $\tilde{e}^G \in \mathcal{V}^H$  be the solution of

$$a(\tilde{e}^G, \mathbf{v}) = -a(\tilde{e}_{\text{cont}}, \mathbf{v}) \quad \forall \mathbf{v} \in \mathcal{V}^H,$$

where  $\tilde{e}_{\text{cont}} \in \mathcal{V}^h$  is any continuous estimate. Then

$$E_1^G[\mathbf{e}_h] := \frac{(R^P(\tilde{e}_{\text{cont}}))^2}{\|\tilde{e}_{\text{cont}}\|^2 - \|\tilde{e}^G\|^2} \leq \|\mathbf{e}_h\|^2.$$

**Proof.** Let  $\tilde{e}_{\text{cont}}^G := \tilde{e}_{\text{cont}} + \tilde{e}^G$ , thus using Theorem 9,

$$\frac{(R^P(\tilde{e}_{\text{cont}}^G))^2}{\|\tilde{e}_{\text{cont}}^G\|^2} \leq \|\mathbf{e}_h\|^2.$$

First, the residue is modified as

$$R^P(\tilde{e}_{\text{cont}}^G) = R^P(\tilde{e}_{\text{cont}}) + R^P(\tilde{e}^G) = R^P(\tilde{e}_{\text{cont}}),$$

because the weak residue vanishes for every function in the finite element space  $\mathcal{V}^H$  (Galerkin orthogonality). And second, the proof is completed replacing the denominator by

$$\|\tilde{e}_{\text{cont}}^G\|^2 = \|\tilde{e}_{\text{cont}}\|^2 + \|\tilde{e}^G\|^2 + 2a(\tilde{e}_{\text{cont}}, \tilde{e}^G) = \|\tilde{e}_{\text{cont}}\|^2 + \|\tilde{e}^G\|^2 - 2a(\tilde{e}^G, \tilde{e}^G) = \|\tilde{e}_{\text{cont}}\|^2 - \|\tilde{e}^G\|^2,$$

where Eq. (16) is used replacing  $\mathbf{v}$  by  $\tilde{e}^G$ .  $\square$

It is worth noting that the non-enhanced lower bound  $E_1[\mathbf{e}_h]$  is also an error estimate. The computation of  $\tilde{e}^G$  and  $E_1^G[\mathbf{e}_h]$  is only performed to improve the quality of the error assessment: the value of the enhanced estimate is larger and the lower bound property is conserved. Therefore, the enhanced estimate,  $E_1^G[\mathbf{e}_h]$ , is sharper.

## 5. Comparison with other existing methods

Refs. [5–7] use apparently similar techniques to obtain upper bounds of the error in the context of a scalar model problem. This section is devoted to compare these techniques with the approach proposed in this paper. The rationale in [5–7] is to decompose the bilinear form  $a(\cdot; \cdot)$  in a sum of local contributions associated with each star. That is, weighted local bilinear forms  $a^{w^i}(\cdot, \cdot)$  are introduced such that

$$a(\mathbf{u}, \mathbf{v}) = \sum_{i=1}^{n_{np}} \int_{\omega^i} w^i \boldsymbol{\sigma}(\mathbf{u}) : \boldsymbol{\varepsilon}(\mathbf{v}) \, d\Omega =: \sum_{i=1}^{n_{np}} a^{w^i}(\mathbf{u}, \mathbf{v}), \quad (20)$$

where the weights  $w^i$  account for the overlapping of the stars verifying the partition of the unity property:

$$\sum_{i=1}^{n_{np}} w^i = 1.$$

The local norm induced by  $a^{w^i}(\cdot, \cdot)$  is denoted by  $\|\cdot\|_{w^i}$ , that is,  $\|\mathbf{v}\|_{w^i}^2 := a^{w^i}(\mathbf{v}, \mathbf{v})$ .

Two different choices for  $w^i$  have been considered. In [6], for each element  $\Omega_k$  of the star  $\omega^i$ , the proposed weight is  $w^i|_{\Omega_k} = (1/\sigma_k)$  where  $\sigma_k$  is the number of vertices of the element  $\Omega_k$ . In [5,7], the local weights are the shape functions,  $w^i = \phi^i$ .

Once the bilinear form is decomposed into local contributions, the local estimates  $\hat{\mathbf{e}}^{\omega^i} \in \mathcal{V}_{\omega^i}^h$  are computed solving the local equation

$$a^{w^i}(\hat{\mathbf{e}}^{\omega^i}, \mathbf{v}) = R^P(\phi^i(\mathbf{v} - \pi_v^H \mathbf{v})) \quad \forall \mathbf{v} \in \mathcal{V}_{\omega^i}^h. \quad (21)$$

**Remark 11.** In fact, in [5–7] the r.h.s. of (21) does not include the projection  $\pi_v^H$ . This is because these papers are only concerned with scalar (thermal) problems and, consequently, the solvability issues discussed in Section 4.1 are not relevant.

The upper bound of  $\|\mathbf{e}_h\|^2$  is obtained adding the local weighted norms of  $\hat{\mathbf{e}}^{\omega^i}$ , that is

$$\hat{E}_u[\mathbf{e}_h] := \sum_{i=1}^{n_{np}} \|\hat{\mathbf{e}}^{\omega^i}\|_{w^i}^2 \geq \|\mathbf{e}_h\|^2.$$

The strategy to obtain the upper bound estimate is summarized in Fig. 4.

Note that  $E_u[\mathbf{e}_h]$  and  $\hat{E}_u[\mathbf{e}_h]$  are computed with completely different expressions. The former is the norm of a sum and the latter is the sum of local norms.

The only difference between Eqs. (12) and (21) is the bilinear form in the l.h.s. term. However, the upper bounds  $E_u[\mathbf{e}_h]$  and  $\hat{E}_u[\mathbf{e}_h]$  have a different expression and, consequently, the analysis of the properties of the estimates follows a different strategy.

The following theorem states that  $\hat{E}_u[\mathbf{e}_h]$  is indeed an upper bound of the squared energy norm of the error.

**Theorem 12.** Let  $\hat{\mathbf{e}}^{\omega^i}$  be the solution of the local Eq. (21) where  $w^i$  are a partition of unity. Then,

$$\hat{E}_u[\mathbf{e}_h] = \sum_{i=1}^{n_{np}} \|\hat{\mathbf{e}}^{\omega^i}\|_{w^i}^2$$

is an upper bound of the squared energy norm of the reference error, namely,  $\hat{E}_u[\mathbf{e}_h] \geq \|\mathbf{e}_h\|^2$ .

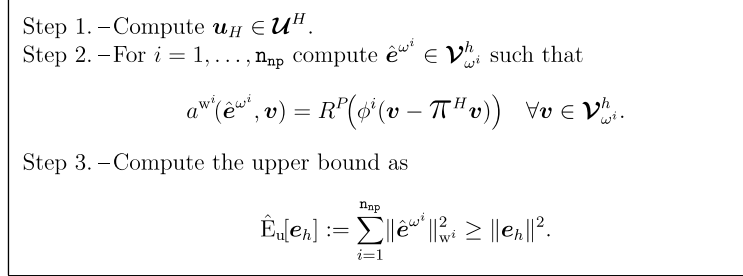


Fig. 4. Alternative upper bound for the squared energy norm of the reference error.

**Proof.** The decomposition of the bilinear form  $a(\cdot, \cdot)$  defined in (20) leads to the following decomposition of the energy norm

$$\|\mathbf{v}\| = a(\mathbf{v}, \mathbf{v})^{\frac{1}{2}} = \left( \sum_{i=1}^{n_{np}} a^{w^i}(\mathbf{v}, \mathbf{v}) \right)^{\frac{1}{2}} = \left( \sum_{i=1}^{n_{np}} \|\mathbf{v}\|_{w^i}^2 \right)^{\frac{1}{2}}.$$

Moreover, using Eqs. (11) and (21) together with Galerkin orthogonality, the squared energy norm of the reference error is rewritten as

$$\|\mathbf{e}_h\|^2 = a(\mathbf{e}_h, \mathbf{e}_h) = R^P(\mathbf{e}_h) = \sum_{i=1}^{n_{np}} R^P(\phi^i(\mathbf{e}_h - \pi^H \mathbf{e}_h)) = \sum_{i=1}^{n_{np}} a^{w^i}(\hat{\mathbf{e}}^{\omega^i}, \mathbf{e}_h).$$

Combining these two decompositions and with repeated use of the Cauchy–Schwartz inequality the proof is completed

$$\begin{aligned} \|\mathbf{e}_h\|^2 &= \left| \sum_{i=1}^{n_{np}} a^{w^i}(\hat{\mathbf{e}}^{\omega^i}, \mathbf{e}_h) \right| \leq \sum_{i=1}^{n_{np}} \left| a^{w^i}(\hat{\mathbf{e}}^{\omega^i}, \mathbf{e}_h) \right| \leq \sum_{i=1}^{n_{np}} \|\hat{\mathbf{e}}^{\omega^i}\|_{w^i} \|\mathbf{e}_h\|_{w^i} \\ &\leq \left( \sum_{i=1}^{n_{np}} \|\hat{\mathbf{e}}^{\omega^i}\|_{w^i}^2 \right)^{\frac{1}{2}} \left( \sum_{i=1}^{n_{np}} \|\mathbf{e}_h\|_{w^i}^2 \right)^{\frac{1}{2}} \leq \hat{E}_u[\mathbf{e}_h]^{\frac{1}{2}} \|\mathbf{e}_h\|. \quad \square \end{aligned}$$

**Remark 13.** The repeated use of the Cauchy–Schwartz inequality in the proof of Theorem 12 suggests that the obtained upper bound is not as sharp as the upper bound associated with the estimate  $\tilde{\epsilon}$ . The numerical examples confirm this impression: the estimate  $E_u[\mathbf{e}_h]$  is usually sharper than  $\hat{E}_u[\mathbf{e}_h]$ .

## 6. Bounds of the error in outputs of interest

As shown in Section 2.3, in order to estimate bounds of the error in the output of interest  $l^\theta(\mathbf{e}_h)$ , upper and lower bounds of  $\|\mathbf{z}_h^\pm\|$  are necessary instead of bounds of  $\|\mathbf{e}_h\|$ . Recall that  $\mathbf{z}_h^\pm = \kappa \mathbf{e}_h \pm \frac{1}{\kappa} \epsilon_h$  where the error of the primal and dual problems are involved. This section presents and discusses the particular evaluation of  $E_u[\mathbf{z}_h^\pm]$  and  $E_l[\mathbf{z}_h^\pm]$ . These values, as indicated by Eq. (10), allow to bound  $\|\mathbf{z}_h^\pm\|$ . Note also that bounds for the output of interest,  $l^\theta(\mathbf{u}_h)$ , can be computed adding  $l^\theta(\mathbf{u}_H)$  to each term of inequality (10).

### 6.1. Upper bound computation of $\|\mathbf{z}_h^\pm\|$

In order to determine  $E_u[\mathbf{z}_h^\pm]$  the error estimate of both the primal and dual problem are necessary. Section 3 describes the evaluation of the primal error estimate. The same methodology is used to estimate the dual error,  $\tilde{\epsilon}$ , by adding local estimates  $\tilde{\epsilon}^{\omega^i} \in \mathcal{V}_{\omega^i}^h$ , computed from

$$a_{\omega^i}(\mathbf{v}, \tilde{\epsilon}^{\omega^i}) = R^D(\phi^i(\mathbf{v} - \boldsymbol{\pi}^H \mathbf{v})) \quad \forall \mathbf{v} \in \mathcal{V}_{\omega^i}^j. \quad (22)$$

Then, the upper bound for  $\|\mathbf{z}_h^\pm\|^2$ ,  $E_u[\mathbf{z}_h^\pm]$ , is obtained as summarized in Fig. 5 and based on the following Lemma.

**Lemma 14.** *The estimate  $E_u[\mathbf{z}_h^\pm] := 2\|\tilde{\epsilon}\| \|\tilde{\epsilon}\| \pm 2a(\tilde{\epsilon}, \tilde{\epsilon})$  is such that*

$$E_u[\mathbf{z}_h^\pm] \geq \|\mathbf{z}_h^\pm\|^2.$$

**Proof.** Since  $a(\cdot, \cdot)$  is a symmetric bilinear form, the following equation for  $\mathbf{z}_h^\pm$  holds,

$$a(\mathbf{z}_h^\pm, \mathbf{v}) = \kappa R^P(\mathbf{v}) \pm \frac{1}{\kappa} R^D(\mathbf{v}) =: R^\pm(\mathbf{v}) \quad \forall \mathbf{v} \in \mathcal{V}^h.$$

Then, according to Lemma 7, an estimate  $\tilde{\mathbf{z}}^\pm \in \mathcal{V}_{\text{brok}}^h$  yields an upper bound of the energy norm of  $\mathbf{z}_h^\pm$  if

$$a(\tilde{\mathbf{z}}^\pm, \mathbf{v}) = R^\pm(\mathbf{v}) \quad \forall \mathbf{v} \in \mathcal{V}^h. \quad (23)$$

Recall now that the primal estimate  $\tilde{\epsilon}$  verifies Eq. (19), see Lemma 7. Similarly,  $\tilde{\epsilon}$  verifies

$$a(\tilde{\epsilon}, \mathbf{v}) = R^D(\mathbf{v}) \quad \forall \mathbf{v} \in \mathcal{V}^h.$$

Thus, introducing  $\tilde{\mathbf{z}}^\pm := \kappa \tilde{\epsilon} \pm \tilde{\epsilon}/\kappa \in \mathcal{V}_{\text{brok}}^h$  Eq. (23) holds true. The proof is completed taking  $\kappa^2 = \|\tilde{\epsilon}\|/\|\tilde{\epsilon}\|$  in

$$E_u[\mathbf{z}_h^\pm] = \|\tilde{\mathbf{z}}^\pm\|^2 = \kappa^2 \|\tilde{\epsilon}\|^2 + \frac{1}{\kappa^2} \|\tilde{\epsilon}\|^2 \pm 2a(\tilde{\epsilon}, \tilde{\epsilon}) = 2\|\tilde{\epsilon}\| \|\tilde{\epsilon}\| \pm 2a(\tilde{\epsilon}, \tilde{\epsilon}). \quad \square$$

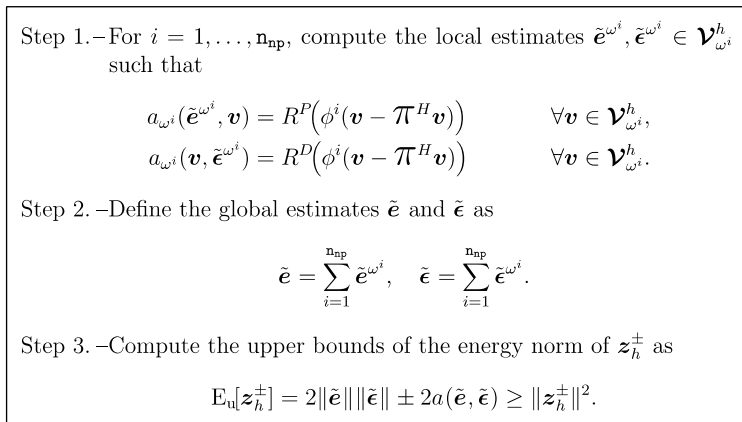


Fig. 5. Upper bounds for the squared energy norm of  $\mathbf{z}_h^\pm$ .

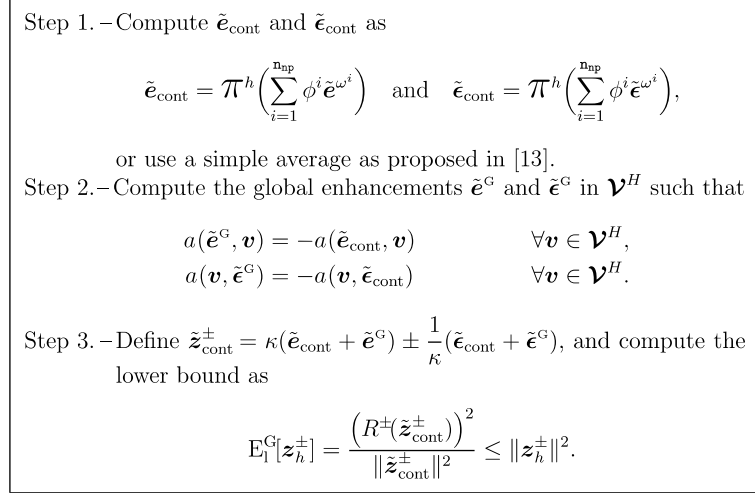


Fig. 6. Lower bounds for the squared energy norm of  $\mathbf{z}_h^\pm$ .

## 6.2. Lower bound computation of $\|\mathbf{z}_h^\pm\|$

Upper bound estimates of  $\tilde{z}^\pm$  are also post-processed to obtain  $\tilde{z}_{\text{cont}}^\pm$  as described in Section 3.3. The strategy used to obtain lower bounds of  $\mathbf{z}_h^\pm$  is summarized in Fig. 6.

**Theorem 15.** Let  $\tilde{e}^G \in \mathcal{V}^H$  and  $\tilde{\epsilon}^G \in \mathcal{V}^H$  be the global enhancements computed from  $\tilde{e}_{\text{cont}} \in \mathcal{V}^h$  and  $\tilde{\epsilon}_{\text{cont}} \in \mathcal{V}^h$ , respectively. Then, the estimate  $\tilde{z}_{\text{cont}}^\pm := \kappa(\tilde{e}_{\text{cont}} + \tilde{e}^G) \pm (\tilde{\epsilon}_{\text{cont}} + \tilde{\epsilon}^G)/\kappa$  provides a lower bound of the energy norm of  $\mathbf{z}_h^\pm$  that is,

$$E_1^G[\mathbf{z}_h^\pm] := \frac{\left( R^\pm(\tilde{z}_{\text{cont}}^\pm) \right)^2}{\|\tilde{z}_{\text{cont}}^\pm\|^2} \leq \|\mathbf{z}_h^\pm\|^2.$$

**Proof.** The proof is a direct consequence of Theorem 9 if both  $\|\tilde{e}\|$  and  $\|\tilde{\epsilon}\|$  are non-zero. The case  $\|\tilde{e}\| = 0$  or  $\|\tilde{\epsilon}\| = 0$  is trivial because it implies that either  $\mathbf{e} = \mathbf{0}$  or  $\boldsymbol{\varepsilon} = \mathbf{0}$ , and therefore  $l^0(\mathbf{e}) = 0$ . In this case, the obvious lower bound  $E_1^G[\mathbf{z}_h^\pm]$  is 0.  $\square$

## 7. Computational aspects

### 7.1. Simplified computation of the weak residual

Eq. (12) is in fact the fundamental equation that is solved repeatedly. The weak residual in its r.h.s. is not trivial to compute because for  $\mathbf{v} \in \mathcal{V}^h$ , in general,  $\phi^i \mathbf{v}$  does not belong to the finite element reference space,  $\mathcal{V}^h$ . Note that this is also the case for  $\phi^i(\mathbf{v} - \Pi^h \mathbf{v})$ . However, the evaluation of the weak residual is drastically simplified if its argument is projected into  $\mathcal{V}^h$ .

For every  $\mathbf{v} \in \mathcal{V}^h$  the following quality holds

$$R^P(\mathbf{v}) = R^P(\Pi^h \mathbf{v}) = R^P \left( \Pi^h \left( \sum_{i=1}^{n_{\text{np}}} \phi^i \mathbf{v} \right) \right) = \sum_{i=1}^{n_{\text{np}}} R^P(\Pi^h(\phi^i \mathbf{v})).$$

Thus the same partition proposed in (11) can be performed with the residual acting on the projection and consequently, Eq. (12) can be rewritten as find  $\tilde{\mathbf{e}}^{\omega^j} \in \mathcal{V}_{\omega^j}^h$  such that

$$a_{\omega^j}(\tilde{\mathbf{e}}^{\omega^j}, \mathbf{v}) = R^P(\pi^h(\phi^i(\mathbf{v} - \pi^H \mathbf{v}))) \quad \forall \mathbf{v} \in \mathcal{V}_{\omega^j}^h.$$

The behavior of the estimates obtained either introducing the projection,  $\pi^h$ , or not is similar as shown in the numerical examples. However the implementation of the r.h.s. term described in the previous equation is much simpler. This is because the argument of  $R^P(\cdot)$  is reinjected in the reference space, which is a standard finite element space. Moreover,  $\phi^i \in \mathcal{V}^h$  and  $\mathbf{v} \in \mathcal{V}^h$ , thus  $\pi^h(\phi^i \mathbf{v})$  is computed by a simple product of nodal values of  $\phi^i$  and  $\mathbf{v}$  (or  $\mathbf{v} - \pi^H \mathbf{v}$  when necessary).

## 7.2. Spatial distribution of upper bound estimates

The upper bound estimate  $E_u[\mathbf{e}_h]$  presented in Section 3.2 can be decomposed into positive contributions of each element of the mesh, thus providing local indicators of the value of the local energy norm of the reference error  $\|\mathbf{e}_h\|_k$ , that is,

$$E_u[\mathbf{e}_h] = \|\tilde{\mathbf{e}}\|^2 = \sum_{k=1}^{n_{e1}} \|\tilde{\mathbf{e}}\|_k^2$$

and  $\|\tilde{\mathbf{e}}\|_k$  is the local indicator for  $\|\mathbf{e}_h\|_k$ .

## 8. Numerical examples

In this section, the behavior of the estimates presented above is analyzed both for thermal and mechanical model problems. Some of the selected examples have been used by other authors to assess performance of similar techniques [1,3,18]. The quality of the error estimates is measured with the index

$$\rho := \frac{\text{estimated error norm}}{\text{true error norm}} - 1,$$

where the “true” error is either the exact error (if available) or the reference error. Index  $\rho$  is the usual effectivity index minus one. The accuracy of the error estimate is given by the absolute value of  $\rho$  and the sign indicates if the estimate is an overestimation (positive  $\rho$ ) or an underestimation (negative  $\rho$ ) of the true error. For instance,  $\rho = 2\%$  indicates that estimated error is larger than the “true” error with a factor 1.02 and  $\rho = -3\%$  means that the “true” error is underestimated by a factor 0.97.

In the remainder of the section  $\rho$  is used to assess the quality of the different estimates, for instance of  $E_u[\mathbf{e}_h]$ . Note however that  $E_u[\mathbf{e}_h]$  is an estimate of the squared energy norm, but the corresponding  $\rho$  index is computed using directly the approximation of the error norm (not squared). Moreover, when the exact error,  $\mathbf{e}$ , is known it is always used to compute  $\rho$ . Thus,  $\rho(E_u[\mathbf{e}_h])$  is defined either as

$$\rho(E_u[\mathbf{e}_h]) := \frac{\sqrt{E_u[\mathbf{e}_h]}}{\|\mathbf{e}\|} - 1 \quad \text{or as} \quad \rho(E_u[\mathbf{e}_h]) := \frac{\sqrt{E_u[\mathbf{e}_h]}}{\|\mathbf{e}_h\|},$$

depending on the availability of  $\mathbf{e}$ . This definition is extended to the other studied error estimates, for instance  $\rho(\hat{E}_u[\mathbf{e}_h])$ .

8.1. Thermal problem with energy norm assessment

First, the scalar benchmark is solved, see [2,3,13]. A squared domain,  $\Omega = ]0, 1[ \times ]0, 1[$ , with Dirichlet homogeneous boundary conditions on  $\delta\Omega$  and a source term are chosen such that the exact solution, given in Fig. 7, has the following analytical expression

$$u(x, y) = x^2(1 - x)^2(e^{10x^2} - 1)y^2(1 - y)^2(e^{10y^2} - 1)/2000.$$

The behavior of the energy norm estimates is analyzed comparing the estimates with the exact energy error norm  $\|e\|$ . Two different non-structured and non-uniform meshes have been considered, see Fig. 7. In both cases, the approximate solution  $u_H$  is computed using quadrilateral meshes with bilinear interpolation ( $p = 1$ ), and the reference space is associated with a mesh of size  $h = H/4$  (i.e. each element of the  $H$ -mesh is divided into 16 new elements).

Table 1 presents the  $\rho$  indices for estimates of upper bounds. Two versions of  $E_u[e_h]$  are shown, one using the projection  $\pi_h$  in the r.h.s. term of the local Eq. (12), as described in Section 7.1 to simplify computations, and another without the proposed projection. Two versions of  $\rho(\hat{E}_u[e_h])$  are also evaluated. They correspond to the different weighting functions  $w^i$  for the bilinear form in the l.h.s. of Eq. (21) described in Section 5, and formerly proposed in [6] ( $w^i = 1/\sigma$ ) and in [5,7] ( $w^i = \phi^i$ ).

Paradoxically, in this example, the two upper bounds estimates proposed in this paper (and associated with  $E_u[e_h]$ ) provide negative values of  $\rho$ . This is because all the presented estimates are upper bounds with respect to the reference error,  $e_h$ , that is, they are larger than the reference error but not necessarily larger than the exact error  $e$  which is used to compute  $\rho$ . Then, even if  $\|e_h\|^2 \leq E_u[e_h]$  stands, in this case the estimates are very sharp and we have  $\|e_h\|^2 < E_u[e_h] < \|e\|^2$ , see Table 1. The estimates corresponding to  $\hat{E}_u[e_h]$  are far from being sharp, they yield an overestimation of more than 60% (for  $w^i = 1/\sigma$ ) and 20% (for  $w^i = \phi^i$ ).

As expected, the effectivity indices of  $E_u[e_h]$  are better than the effectivity indices of  $\hat{E}_u[e_h]$  (for both versions  $w^i = 1/\sigma$  and  $w^i = \phi^i$ ). Table 1 also shows that the proposed projection  $\pi^h$  in the r.h.s. term of local Eq. (12), as described in Section 7.1, does not modify substantially the values of effectivity indices. Recall that

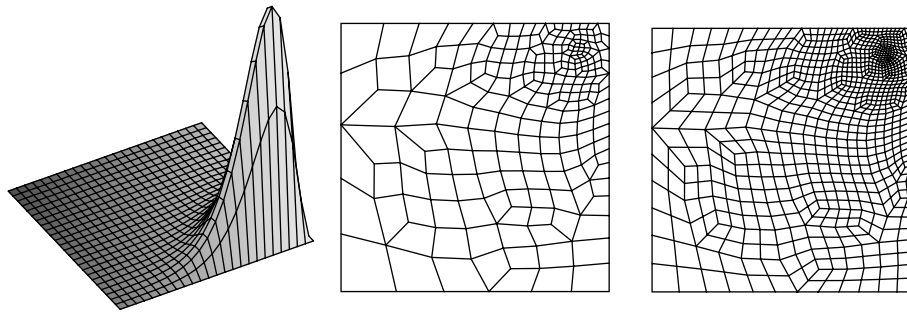


Fig. 7. Thermal problem: exact solution (left) and meshes with 240 d.o.f. (center) and 913 d.o.f. (right).

Table 1  
Thermal problem:  $\rho$  indices for upper bound energy norm estimates

| d.o.f. | $\frac{\ e\ }{\ u\ }$ (%) | $\frac{\ e_h\ }{\ e\ }$ (%) | $\rho(E_u[e_h])$ without $\pi^h$ (%) | $\rho(E_u[e_h])$ with $\pi^h$ (%) | $\rho(\hat{E}_u[e_h])_{w^i = 1/\sigma}$ (%) | $\rho(\hat{E}_u[e_h])_{w^i = \phi^i}$ (%) |
|--------|---------------------------|-----------------------------|--------------------------------------|-----------------------------------|---|---|
| 240    | 24.9                      | 95.9                        | -3.26                                | -3.34                             | 63.9  | 23.4                                      |
| 913    | 15.1                      | 96.6                        | -2.59                                | -2.65                             | 67.8  | 25.8                                      |

Table 2  
Thermal problem:  $\rho$  indices for lower bound energy norm estimates

| d.o.f. | $\rho(E_1[e_h])$ without $\pi^h$ (%) | $\rho(E_1[e_h])$ with $\pi^h$ (%) | $\rho(\hat{E}_1[e_h])w^j = 1/\sigma$ (%) | $\rho(\hat{E}_1[e_h])w^j = \phi^j$ (%) |
|--------|--------------------------------------|-----------------------------------|--|--|
| 240    | -19.9                                | -19.6                             | -19.9                                    | -31.0                                  |
| 913    | -18.9                                | -18.6                             | -18.9                                    | -29.6                                  |

Table 3  
Thermal problem:  $\rho$  indices for lower bound energy norm estimates with global enhancement

| d.o.f. | $\rho(E_1^G[e_h])$ without $\pi^h$ (%) | $\rho(E_1^G[e_h])$ with $\pi^h$ (%) | $\rho(\hat{E}_1^G[e_h])w^j = 1/\sigma$ (%) | $\rho(\hat{E}_1^G[e_h])w^j = \phi^j$ (%) |
|--------|--|-------------------------------------|--|--|
| 240    | -6.58                                  | -7.30                               | -6.58                                      | -4.67                                    |
| 913    | -5.79                                  | -6.50                               | -5.79                                      | -3.97                                    |

the use of  $\pi^h$  simplifies considerably the implementation of the estimator and it is therefore strongly recommended.

Effectivity indices for lower bound estimates are displayed in Tables 2 and 3. Estimates  $\hat{E}_1[e_h]$  and  $\hat{E}_1^G[e_h]$  are computed in the same fashion as  $E_1[e_h]$  but using the continuous function resulting of smoothing  $\hat{e} = \sum_i \hat{e}^{w^i}$  instead of  $\tilde{e}$ , see Eq. (13). Table 3 shows the results obtained applying the global enhancement discussed in Section 4.4. These results indicate that lower bounds are not sensitive to the original (discontinuous) estimate, which provides the upper bound. All estimates in Table 2 perform similarly. The effect of the global enhancement is however very important: the effectivity indices improve drastically from Table 2 to Table 3. Recall that in all these tables  $\rho$  is computed with respect to the exact error because the exact solution is known.

From a qualitative viewpoint, it is worth noting that the estimated error distribution is in good agreement with the exact error distribution, both for the estimate proposed here ( $E_u[e_h]$ ) and for the estimates proposed in [5–7] (the two versions of  $\hat{E}_u[e_h]$ ).

Fig. 8 shows the spatial distribution of the local effectivity index and the histogram representing the occurrences of local effectivity indices. The histogram shows the number of elements with local effectivity

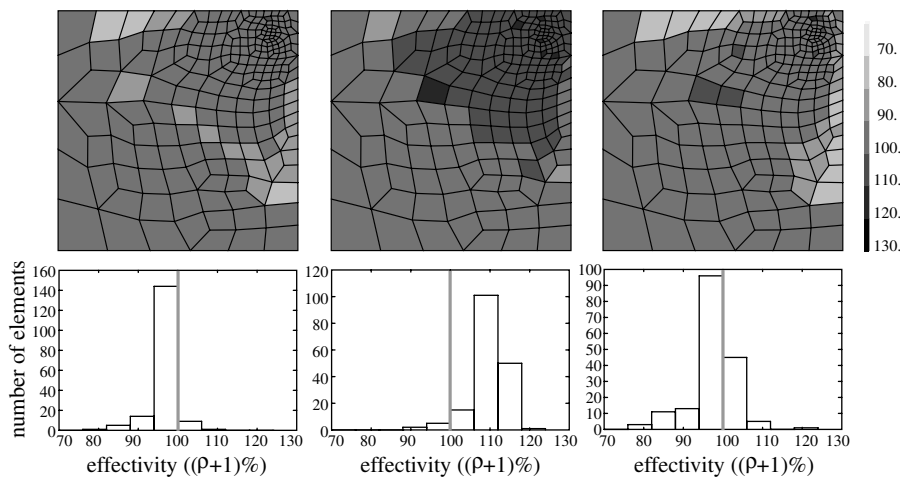


Fig. 8. Thermal problem: spatial distribution of the local effectivity (top) and histograms for local effectivity (bottom). The results correspond to  $E_u[e_h]$  (left),  $\hat{E}_u[e_h]$  with  $w^j = \phi^j$  (center) and  $E_1^G[e_h]$  (right).



in a given range. The histograms show a good behavior of the estimate if they display a narrow distribution (all elements have similar local effectivity indices) concentrated around 100%. Observe that the local values associated with the estimate  $E_u[e_h]$  proposed here are much more accurate than the values corresponding to  $E_u[e_h]$ .

**Remark 16.** Elements with a small local error are not taken into account. Because the areas where error is small are not interesting from, an adaptive viewpoint. Moreover, in these areas, small defaults in the error assessment lead to very bad effectivity indices (small absolute error but large relative error). Here, the criterion used is to suppress in the histograms elements such that the local error norm is lower than  $\|e_h\|/4n_{e1}$  (being  $n_{e1}$  the number of elements). That results on neglecting 20% of the elements approximately.

8.2. Thin plate energy error assessment

A square thin plate with two holes proposed in [19] is considered next. This is a plane-stress linear elastic problem loaded with a horizontal unit tension along the vertical edges  $\Gamma_0$ , see Fig. 9. Note that the solution of this problem has corner singularities due to the interior rectangular cut-outs. Due to symmetry, only one fourth of the domain is analyzed. Values for Young’s modulus and Poisson ratio are set to 1 and 0.3, respectively.

Two meshes are considered, a coarse uniform mesh with 70 nodes and a finer one with 850 nodes, adapted heuristically. Error estimates  $E_u[e_h]$  and  $E_1^G[e_h]$  are computed for both cases and results are summarized in Table 4. The effectivity index of the upper bound estimate is similar for the two meshes, and close to 1.17 ( $\rho \approx 17\%$ ). The lower bound effectivity are not as sharp, they are close to 0.3 ( $\rho \approx -70\%$ ).

Spatial distributions of error  $E_u[e_h]$  are displayed in Figs. 10 and 11 for the uniform and adapted meshes, respectively. Note that they are computed using  $\pi^h$  in the r.h.s of Eq. (12). It is worth noting that the error distributions for  $E_u[e_h]$  are in good agreement with the reference error. The bad behavior of the local effectivity index in the first mesh, see Fig. 10, is due to the fact that practically all the error is concentrated in a

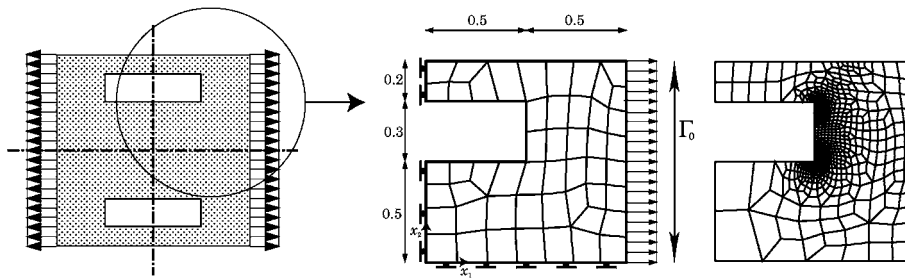


Fig. 9. Thin plate model problem and meshes with 140 d.o.f. (center) and 1970 d.o.f. (right).

Table 4

Thin plate: upper and lower bounds for  $\|e_h\|$

| d.o.f. | $\ e_h\ $ | $\frac{\ e_h\ }{\ u_h\ }$ (%) | $\rho(E_u[e_h])$ with $\pi^h$ (%) | $\rho(E_1^G[e_h])$ with $\pi^h$ (%) |
|--------|-----------|-------------------------------|-----------------------------------|-------------------------------------|
| 140    | 0.146     | 12.8                          | 17.9                              | -68.7                               |
| 1970   | 0.040     | 3.44                          | 17.1                              | -70.1                               |

| d.o.f. | $\ e_h\ $ | $\frac{\ e_h\ }{\ u_h\ }$ | $\rho(E_u[e_h])$<br>with $\mathcal{T}^h$ | $\rho(E_1^G[e_h])$<br>with $\mathcal{T}^h$ |
|--------|-----------|---------------------------|--|--|
| 140    | 0.146     | 12.8%                     | 17.9%                                    | -68.7%                                     |
| 1970   | 0.040     | 3.44%                     | 17.1%                                    | -70.1%                                     |

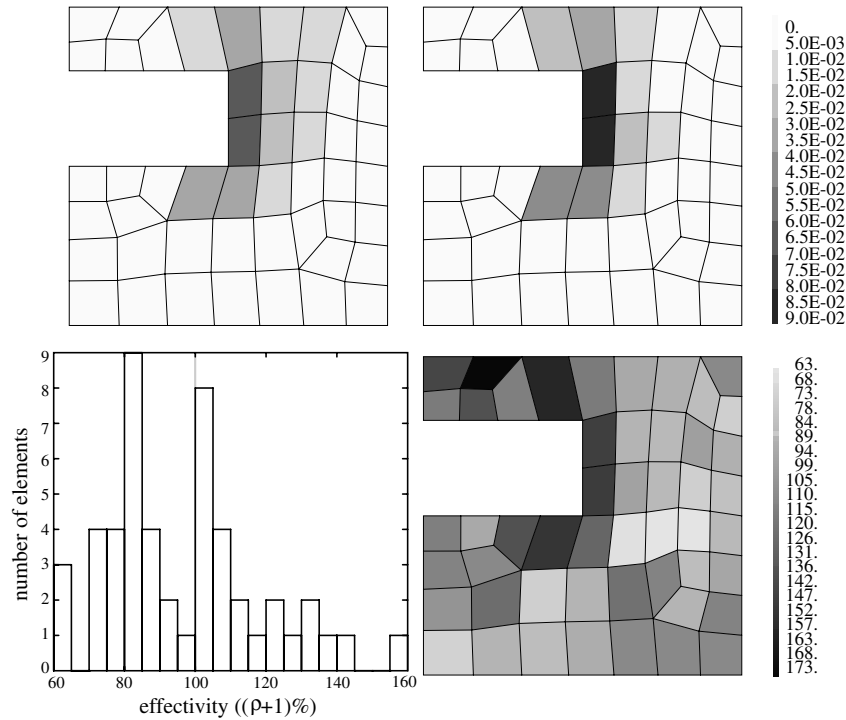


Fig. 10. Thin plate: spatial distribution of the reference error (top left), estimate  $E_u[e_h]$  (top right), and local distribution of the effectivity indices  $(\rho + 1)\%$  (bottom) for the mesh with 140 d.o.f.

few relevant elements. The histogram in Fig. 11 is narrow because the number of elements in the zones where the error is relevant is much higher for the second mesh.

Finally, Fig. 12 shows a comparison between the proposed upper bound estimate,  $\sqrt{E_u[e_h]}$ , the flux-free techniques proposed in [6] ( $w^i = 1/\sigma$ ) and in [5,7] ( $w^i = \phi^i$ ), and a hybrid-flux upper bound estimate, see [4,2]. The upper bound estimates are computed for a series of adapted triangular meshes. As expected all of them converge. Moreover, this is an example in which the hybrid-flux bound is sharper than the previously published flux-free upper bound estimates. In [8] the majority of the examples behave similarly. However, as already discussed the proposed flux-free bound is as sharp as the hybrid-flux one.

### 8.3. Assessment of outputs of interest for a crack opening problem

The error estimator presented in this paper is applied to the crack opening problem proposed in [18]. The specimen is described in Fig. 13. Loads are a uniform pressure in the upper round cavity and a uniform normal traction pulling the left upper part of the specimen. Displacements are set to zero along  $\Gamma_D$ , around

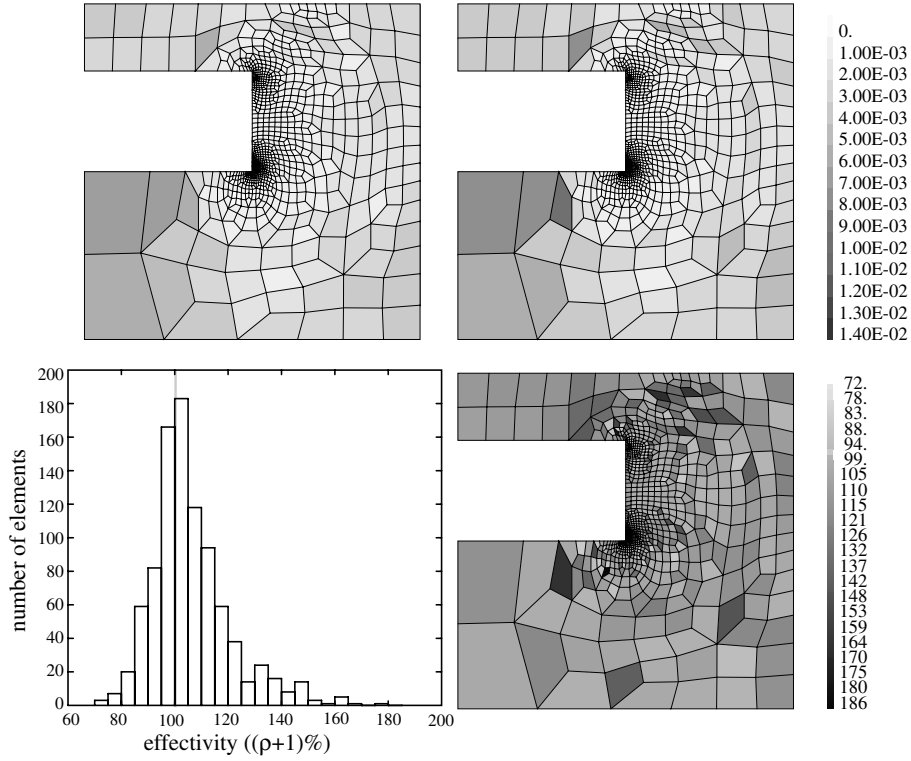


Fig. 11. Thin plate: spatial distribution of the reference error (top left), estimate  $E_u[e_h]$  (top right), and local distribution of the effectivity indices  $(\rho + 1)\%$  (bottom) for the mesh with 2588 d.o.f.

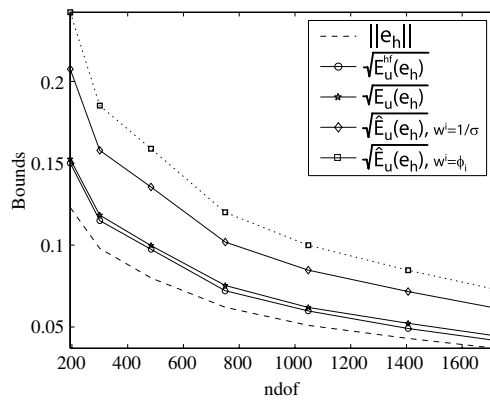


Fig. 12. Thin plate: comparison between flux-free and hybrid-flux estimates.

the centered round cavity. The edges of the crack are denoted by  $\Gamma_1$  (right) and  $\Gamma_2$  (left). The quantity of interest is taken as the average opening along the crack, that is,

$$I^0(\mathbf{u}) = - \int_{\Gamma_1} \mathbf{u} \cdot \mathbf{n} d\Gamma - \int_{\Gamma_2} \mathbf{u} \cdot \mathbf{n} d\Gamma.$$

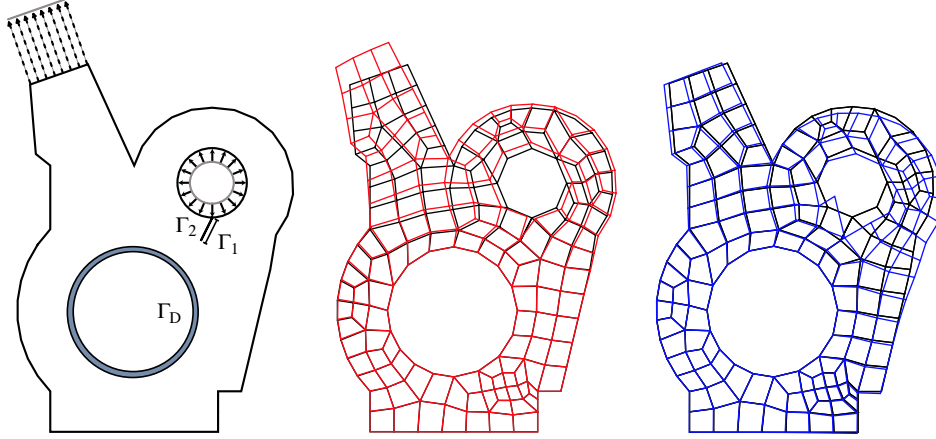


Fig. 13. Crack opening model problem (left), primal (center) and dual (right) solutions. Uniform mesh with 173 nodes (346 d.o.f).

Table 5  
Crack opening problem: energy norm estimates and effectivity indices

| d.o.f. | Primal                        |                      |                        | Dual                                    |                             |                               |
|--------|-------------------------------|----------------------|------------------------|---|-----------------------------|-------------------------------|
|        | $\frac{\ e_h\ }{\ u_h\ }$ (%) | $\rho(E_u[e_h])$ (%) | $\rho(E_1^G[e_h])$ (%) | $\frac{\ \epsilon_h\ }{\ \psi_h\ }$ (%) | $\rho(E_u[\epsilon_h])$ (%) | $\rho(E_1^G[\epsilon_h])$ (%) |
| 346    | 20.6                          | 18.1                 | −48.8                  | 61.6                                    | 16.4                        | −19.6                         |
| 1344   | 10.5                          | 14.1                 | −82.9                  | 25.5                                    | 19.1                        | −63.1                         |

Table 6  
Crack opening problem: estimates for  $z_h^\pm$

| d.o.f. | $\ z_h^\pm\ $ | $\rho(E_u[z_h^+])$ (%) | $\rho(E_1^G[z_h^+])$ (%) | $\ z_h^\pm\ $ (%) | $\rho(E_u[z_h^-])$ (%) | $\rho(E_1^G[z_h^-])$ (%) |
|--------|---------------|------------------------|--------------------------|-------------------|------------------------|--------------------------|
| 346    | 0.666         | 16.7                   | −30.9                    | 0.457             | 18.3                   | −49.7                    |
| 1344   | 0.313         | 17.2                   | −74.6                    | 0.232             | 15.2                   | −84.4                    |

Note that the opposite sides of the crack,  $\Gamma_1$  and  $\Gamma_2$ , have opposite normal unit outward vectors. Thus,  $l^\circ(\mathbf{u})$  is the average (integrated) crack opening, it is positive for opening and negative for penetration.

First, the analysis is performed with the coarse uniform mesh shown in Fig. 13. Energy norm error estimates for both the primal and dual problems are summarized in Table 5. Global effectivity indices are of the same order of magnitude as in the previous example for the upper bound estimates. Although the mesh is excessively coarse and the error is large (78% for the dual problem), the behavior of the upper bound is similar and the quality lower bound estimates is better.

Table 6 shows the energy norm estimates for the quantities  $z_h^+$  and  $z_h^-$ . Recall that these linear combinations of the primal and dual errors,  $z_h^\pm = \kappa e_h \pm \frac{1}{\kappa} \epsilon_h$ , are required to assess the error in the quantities of interest. Upper and lower bounds for the quantity of interest  $l^\circ(\mathbf{u}_h)$  are obtained by properly combining upper and lower bounds in the energy norm of  $z_h^+$  and  $z_h^-$ , see Fig. 1. Table 6 indicates that upper bound estimates of  $z_h^\pm$  present similar values of effectivity indices as in previous examples. Effectivity for lower bound estimates is again quite poor, specially for  $z_h^-$  with a value of 0.08 ( $\rho \approx -92\%$ ). Note however, that this poor lower bound effectivity does not drastically downgrade bounds of the desired functional output  $l^\circ(\mathbf{u}_h)$ , which are computed as indicated in Fig. 1 and Section 6. In fact, the obtained bounds are better than if the trivial lower bounds (equal to zero) are imposed.

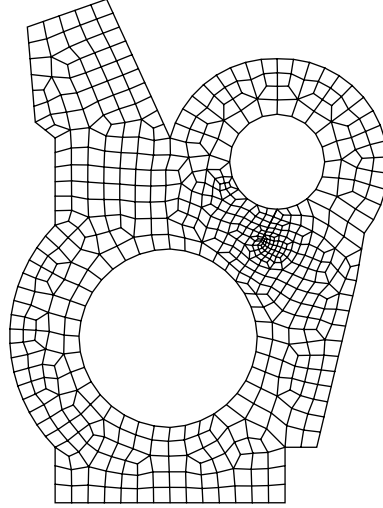


Fig. 14. Crack opening: heuristically adapted mesh with 672 nodes (1344 d.o.f.).

In fact, for the coarse case (346 d.o.f.), when upper and lower bounds for  $l^\theta(\mathbf{u}_h)$  are computed from data of Table 6, the following range is obtained:  $0.141 \leq l^\theta(\mathbf{u}_h) \leq 0.299$ . Note that the coarse mesh estimate is  $l^\theta(\mathbf{u}_H) = 0.161$  and the reference one is  $l^\theta(\mathbf{u}_h) = 0.220$ . In this case, the mesh is very inaccurate for the desired quantity of interest because the error of  $l^\theta(\mathbf{u}_H)$  is 26.8% with respect to the reference one. Nevertheless, the estimate  $0.220 \pm 0.079$  is a valuable information.

The same analysis is performed with the mesh (1344 d.o.f.) shown in Fig. 14. This mesh is uniformly densified and heuristically adapted concentrating elements in the neighborhood of the crack tip. As Table 6 shows, the energy errors are improved (10.5% and 26.4% for primal and dual problems) and the error of  $l^\theta(\mathbf{u}_H)$  is now 5.11% with respect to the reference one. Thus the output gap is reduced since  $0.200 \leq l^\theta(\mathbf{u}_h) \leq 0.249$  and the estimate is sharper  $0.225 \pm 0.025$ .

Fig. 15 shows the spatial distribution of the element-by-element contributions to  $l^\theta(\mathbf{e}_h)$ . That is, the local values for  $a_k(\mathbf{e}_h, \epsilon_h)$  in every element  $\Omega_k$  of the mesh are plotted. Note that this is not an estimate but the actual reference values. These local contributions may be either positive or negative. In order to better depict the areas that contribute to the error the distribution is represented by both the absolute value and the sign of the local contributions. The distribution of the absolute value shows that the main contributions to the error in the quantity of interest are related to elements in the neighborhood of the crack tip. Such a distribution of the error could guide an adaptive process.

The estimated spatial distribution of the error in the output of interest is also shown in Fig. 15. That is, the estimated values  $\tilde{e}$  and  $\tilde{\epsilon}$ , see for instance step 2 in Fig. 5, are used to evaluate  $a_k(\tilde{e}, \tilde{\epsilon})$  in every element  $\Omega_k$  of the mesh. The similarity of these distributions demonstrates the good agreement between the true and the estimated error distributions. Therefore, the introduced error estimators may be fairly used in a goal-oriented adaptive analysis.

#### 8.4. 3D mechanical problem

As previously observed, this approach easily accommodates a 3D analysis. Note that the modifications in code for the 3D analysis showed here were developed in four hours starting from the 2D implementation. Moreover, in 3D the conclusions drawn in [8] from a computational cost point of view are more critical.

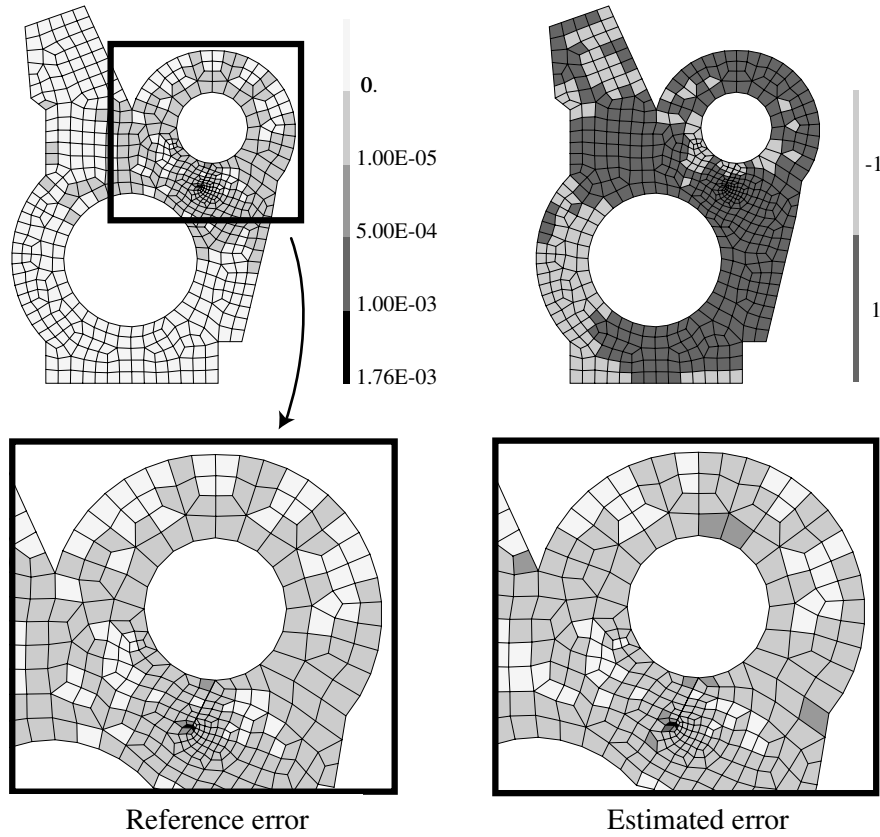


Fig. 15. Spatial distribution of the absolute value of the local contributions to  $l^0(e_h)$  (upper-left and zoom into the relevant zones). The zoom boxes compare the reference error distribution with the estimated error distribution. The upper-right plot describes the sign of these local contributions for the reference error distribution.

The geometry of the problem is inspired in an arch structure proposed in [20]. The structure is casted in the bottom bases and loaded with a uniform pressure on one lateral side, see Fig. 16. This figure also shows the mesh of 174 quadratic 10-noded tetrahedra used for the analysis. The error assessment is also per-

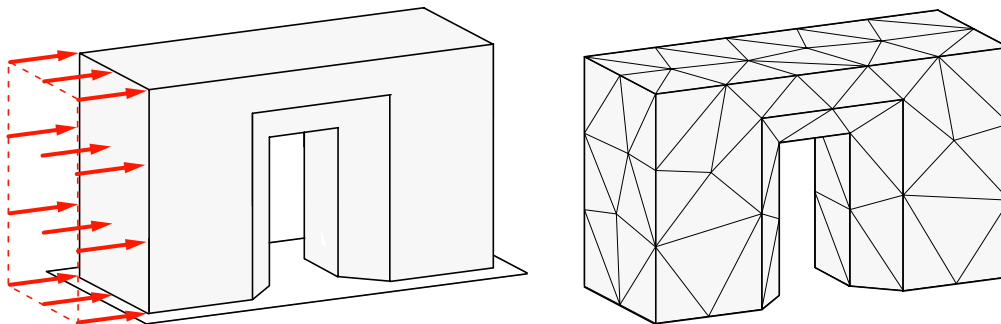


Fig. 16. 3D model problem (left) and uniform mesh (right) with 174 elements and 401 nodes.

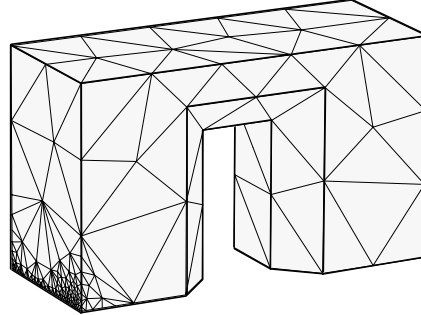


Fig. 17. 3D adapted mesh with 695 elements and 1495 nodes.

Table 7

3D upper and lower bounds for  $\|e_h\|$

| Mesh    | d.o.f. | $\ e_h\ $ | $\frac{\ e_h\ }{\ u_h\ }$ (%) | $\rho(E_u[e_h])$ (%) | $\rho(E_1^G[e_h])$ (%) | $\rho(\hat{E}_u[e_h])$ (%) |
|---------|--------|-----------|-------------------------------|----------------------|------------------------|----------------------------|
| Uniform | 1203   | 4.13      | 22.0                          | 8.18                 | -8.56                  | 53.7                       |
| Adapted | 4485   | 3.63      | 19.2                          | 9.67                 | -14.4                  | 58.6                       |

formed on the mesh shown in Fig. 17, which has been heuristically adapted by refining the elements along the loaded lateral side where stresses are larger. This adapted mesh contains 695 tetrahedra and 1495 nodes.

In this example, we compute the estimates  $E_u[e_h]$ , and  $E_1^G[e_h]$  introduced in this paper using the projection  $\pi^h$  for the r.h.s. term of the residual equation, as described in Section 7.1. This estimate is compared with  $\hat{E}_u[e_h]$  described in [5,7] (with  $w_i = \phi_i$ ).

The reference mesh is obtained dividing each tetrahedron in 8 tetrahedra. Table 7 shows that the estimate proposed here is one order of magnitude sharper than the one introduced in [5,7]. Note that lower bound estimates present similar effectivities.

Surprisingly, comparing results for the uniform mesh with 1203 d.o.f and the adapted mesh with 4485 d.o.f., the reference error is reduced only from 22.0% to 19.2%. This is due to the fact that the first mesh is too coarse and the corresponding reference mesh is not accurate enough. Using this reference mesh the norm of the reference error is 4.13, the exact error is however much larger. A finer reference mesh is built up by splitting each element into 64 tetrahedra (instead of 8). Due to the size of the problem, for this very fine reference mesh, the reference error can be estimated but it cannot be computed. The obtained values for the upper bound estimate  $E_u[e_h]$  and the lower bound estimate  $E_1^G[e_h]$  are 4.34 and 5.24 respectively. Thus the exact error is approximately 5 (probably larger). That means that for the uniform mesh the error is closer 27% than 22%.

## 9. Concluding remarks

This paper introduces a new approach to subdomain-based error estimates. The implementation is less cumbersome compared to hybrid-flux estimators where flux-equilibration algorithms must be implemented. Moreover, accuracy of the results (sharpness of the upper bound) are drastically improved compared to other flux-free estimates. In fact, it is at least comparable to hybrid-flux techniques.

The resulting estimates yield guaranteed (and sharp) upper bounds of the reference error. A simple and painless post-processing yields lower bounds of the error with a little extra computational cost. Upper and

lower bounds of the error are particularly interesting for assessing the error in quantities of interest, computing bounds for functional outputs.

The local problems that have to be solved in this context are flux-free, that is no flux equilibration is required. The flux-free property is specially significant when compared with the standard residual type error estimators (hybrid-flux approach). The local boundary conditions for the local problems in the standard estimators require flux equilibration and result in costly computations and complex programming, especially in 3D.

The distribution of the local contributions to the error are also accurately estimated, both for the energy norm of the error and for the error measured using some functional output. These estimates are therefore well suited to guide goal-oriented adaptive procedures.

### Acknowledgement

Sponsored by Ministerio de Ciencia y Tecnología (grants: DPI2004-3000 and CGL2004-06171-C03-01) and the Generalitat de Catalunya (grant: 2001SGR00257).

### References

- [1] M. Paraschivoiu, J. Peraire, A.T. Patera, A posteriori finite element bounds for linear-functional outputs of elliptic partial differential equations, *Comput. Methods Appl. Mech. Engrg.* 150 (1–4) (1997) 289–312, Symposium on Advances in Computational Mechanics, Austin, TX, vol. 2, 1997.
- [2] M. Ainsworth, J.T. Oden, *A Posteriori Error Estimation in Finite Element Analysis*, John Wiley & Sons, Chichester, 2000.
- [3] J.T. Oden, S. Prudhomme, Goal-oriented error estimation and adaptivity for the finite element method, *Comput. Math. Appl.* 41 (5–6) (2001) 735–756.
- [4] P. Ladevèze, D. Leguillon, Error estimate procedure in the finite element method and applications, *SIAM J. Numer. Anal.* 20 (3) (1983) 485–509.
- [5] C. Carstensen, S.A. Funken, Fully reliable localized error control in the FEM, *SIAM J. Sci. Comput.* 21 (4) (1999/00) 1465–1484 (electronic).
- [6] L. Machiels, Y. Maday, A.T. Patera, A “flux-free” nodal Neumann subproblem approach to output bounds for partial differential equations, *C.R. Acad. Sci. Paris Sér. I Math.* 330 (3) (2000) 249–254.
- [7] P. Morin, R.H. Nochetto, K.G. Siebert, Local problems on stars: a posteriori error estimators, convergence, and performance, *Math. Comput.* 72 (243) (2003) 1067–1097.
- [8] H.-W. Choi, M. Paraschivoiu, Adaptive computations of a posteriori finite element output bounds: a comparison of the “hybrid-flux” approach and the “flux-free” approach, *Comput. Methods Appl. Mech. Engrg.* 193 (36–38) (2004) 4001–4033.
- [9] Y. Maday, A.T. Patera, J. Peraire, A general formulation for a posteriori bounds for output functionals of partial differential equations; application to the eigenvalue problem, *C.R. Acad. Sci. Paris Sér. I Math.* 328 (9) (1999) 823–828.
- [10] S. Prudhomme, J.T. Oden, On goal-oriented error estimation for elliptic problems: application to the control of pointwise errors, *Comput. Methods Appl. Mech. Engrg.* 176 (1–4) (1999) 313–331, *New advances in computational methods*, Cachan, 1997.
- [11] A.T. Patera, J. Peraire, A general Lagrangian formulation for the computation of a posteriori finite element bounds, in: *Error Estimation and Adaptive Discretization Methods in Computational Fluid Dynamics* Lect. Notes Comput. Sci. Engrg., vol. 25, Springer, Berlin, 2003, pp. 159–206.
- [12] I. Babuška, W.C. Rheinboldt, Error estimates for adaptive finite element computations, *SIAM J. Numer. Anal.* 15 (4) (1978) 736–754.
- [13] P. Díez, N. Parés, A. Huerta, Recovering lower bounds of the error by postprocessing implicit residual a posteriori error estimates, *Int. J. Numer. Methods Engrg.* 56 (10) (2003) 1465–1488.
- [14] A. Huerta, P. Díez, Error estimation including pollution assessment for nonlinear finite element analysis, *Comput. Methods Appl. Mech. Engrg.* 181 (1–3) (2000) 21–41.
- [15] S.C. Brenner, L.R. Scott, *The mathematical theory of finite element methods*, second ed. Texts in Applied Mathematics, vol. 15, Springer-Verlag, New York, 2002.
- [16] R.E. Bank, A. Weiser, Some a posteriori error estimators for elliptic partial differential equations, *Math. Comput.* 44 (170) (1985) 283–301.



- [17] M. Ainsworth, J.T. Oden, A unified approach to a posteriori error estimation using element residual methods, *Numer. Math.* 65 (1) (1993) 23–50.
- [18] I. Babuška, T. Strouboulis, C.S. Upadhyay, S.K. Gangaraj, K. Copps, Validation of a posteriori error estimators by numerical approach, *Int. J. Numer. Methods Engrg.* 37 (7) (1994) 1073–1123.
- [19] J. Peraire, A.T. Patera, Bounds for linear-functional outputs of coercive partial differential equations: local indicators and adaptive refinement, in: *Advances in Adaptive Computational Methods in Mechanics*, Cachan, 1997 *Stud. Appl. Mech.*, vol. 47, Elsevier, Amsterdam, 1998, pp. 199–216.
- [20] E. Florentin, L. Gallimard, J.P. Pelle, Evaluation of the local quality of stresses in 3D finite element analysis, *Comput. Methods Appl. Mech. Engrg.* 191 (39–40) (2002) 4441–4457.



# The computation of bounds for linear-functional outputs of weak solutions to the two-dimensional elasticity equations

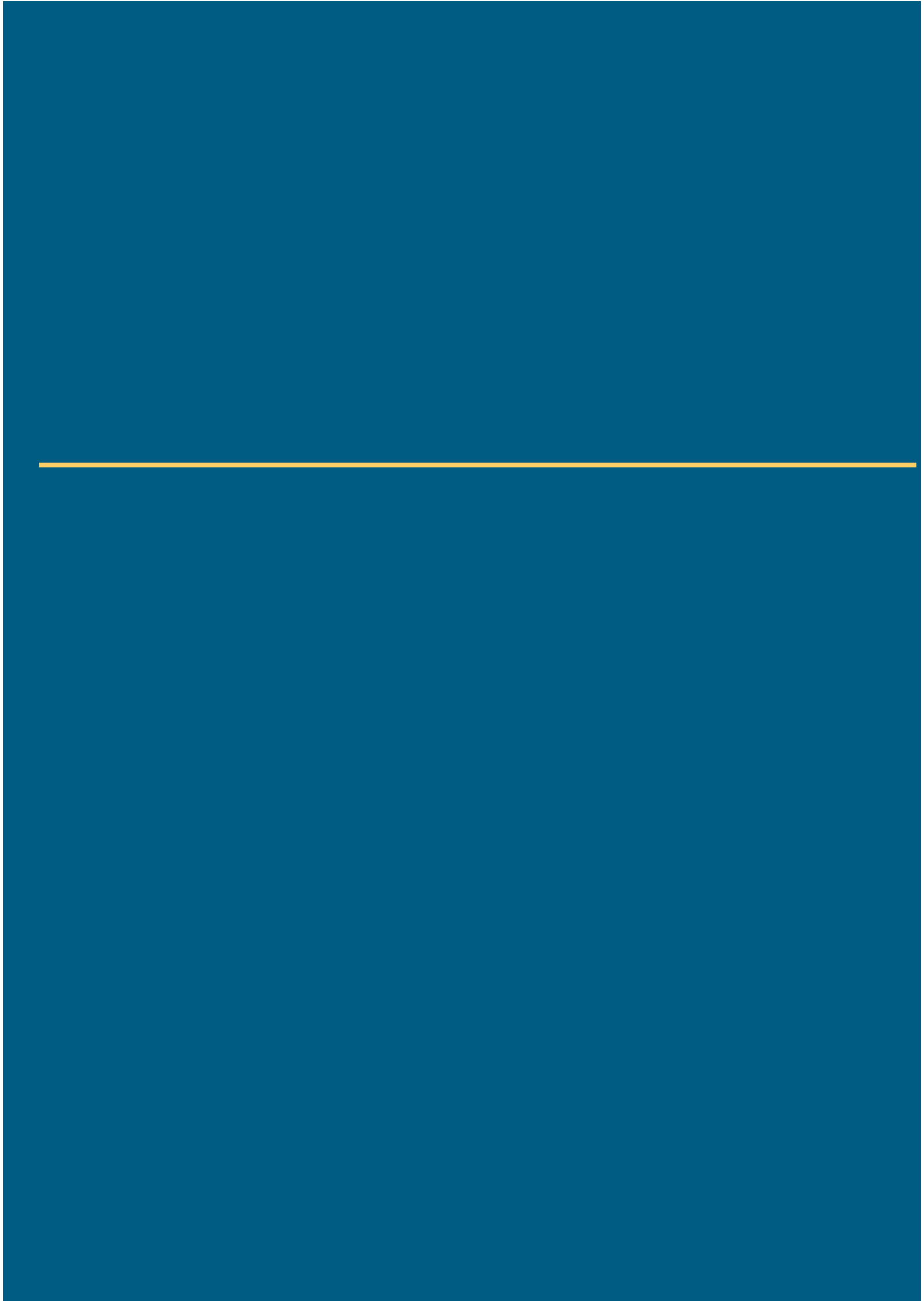
Parés N., Bonet J., Huerta A. and Peraire J.

---

*Computer Methods in Applied*

*Mechanics and Engineering*

**194.** In press.





Available online at [www.sciencedirect.com](http://www.sciencedirect.com)

SCIENCE @ DIRECT®

Comput. Methods Appl. Mech. Engrg. xxx (2005) xxx–xxx

**Computer methods  
in applied  
mechanics and  
engineering**

[www.elsevier.com/locate/cma](http://www.elsevier.com/locate/cma)

# The computation of bounds for linear-functional outputs of weak solutions to the two-dimensional elasticity equations

N. Parés<sup>a</sup>, J. Bonet<sup>b</sup>, A. Huerta<sup>a</sup>, J. Peraire<sup>c,\*</sup>

<sup>a</sup> *Laboratori de Càlcul Numèric, Universitat Politècnica de Catalunya, Barcelona E-08034, Spain*

<sup>b</sup> *Civil and Computational Engineering Centre, School of Engineering, University of Wales, Swansea SA28PP, United Kingdom*

<sup>c</sup> *Department of Aeronautics and Astronautics, Massachusetts Institute of Technology, Massachusetts Avenue 77, Building 37-451, Cambridge, MA 02139, United States*

Received 20 March 2004; received in revised form 15 September 2004; accepted 10 October 2004

---

## Abstract

We present a method for the computation of upper and lower bounds for linear-functional outputs of the exact solutions to the two dimensional elasticity equations. The method can be regarded as a generalization of the well known complementary energy principle. The desired output is cast as the supremum of a quadratic-linear convex functional over an infinite dimensional domain. Using duality the computation of an upper bound for the output of interest is reduced to a feasibility problem for the complementary, or dual, problem. In order to make the problem tractable from a computational perspective two additional relaxations that preserve the bounding properties are introduced. First, the domain is triangulated and a domain decomposition strategy is used to generate a sequence of independent problems to be solved over each triangle. The Lagrange multipliers enforcing continuity are approximated using piecewise linear functions over the edges of the triangulation. Second, the solution of the adjoint problem is approximated over the triangulation using a standard Galerkin finite element approach. A lower bound for the output of interest is computed by repeating the process for the negative of the output. Reversing the sign of the computed upper bound for the negative of the output yields a lower bound for the actual output. The method can be easily generalized to three dimensions. However, a constructive proof for the existence of feasible solutions for the outputs of interest is only given in two dimensions. The computed bound gaps are found to converge optimally, that is, at the same rate as the finite element approximation. An attractive feature of the proposed approach is that it allows for a data set to be generated that can be used to certify and document the computed bounds. Using this data set and a simple algorithm, the correctness of the computed bounds can be established without recourse to the original code used to compute them. In the present paper, only computational domains whose boundary is made up of straight sided segments and polynomially varying loads are considered. Two examples are given to illustrate the proposed methodology.

© 2005 Elsevier B.V. All rights reserved.

---

\* Corresponding author.

E-mail address: [peraire@mit.edu](mailto:peraire@mit.edu) (J. Peraire).

*Keywords:* Linear-functional outputs; Two-dimensional elasticity equations

---

## 1. Introduction

Linear elastic analysis is one of the most common tools used in practical computer aided engineering design. Many materials of practical interest can be adequately modelled as being linear elastic and the physical considerations underlying this assumption are usually well understood. From the theoretical point of view, the equations of linear elasticity have been studied in depth and a number of finite element algorithms exist that can be used to compute approximations to the solutions in an efficient and reliable manner. In particular, a priori error estimates can be established that guarantee the convergence of the computed solution to the exact solution when the mesh is suitably refined.

In practice, the computed solutions of the elasticity equations are used to determine approximations to quantities, or outputs, of practical interest such as displacements, forces or stresses. Once a solution has been computed on a given mesh, one is interested in determining the accuracy of the approximated outputs. In order to address this question a number of a posteriori error estimation methods have been proposed that attempt to quantify the error of the computed solution in either the energy norm [1,3,5,8], or more relevantly, in functional outputs of practical interest such as displacements or stresses [11,13,15]. These a posteriori approaches can be coupled with mesh adaptive strategies, e.g. [12], thus producing algorithms that, in principle, can be used to iterate from initial solution until a preset level of accuracy has been met.

Despite these numerous advances a fundamental issue still remains. Procedures that can be used to unambiguously certify the accuracy of the computed results have been elusive. The reasons for that are essentially twofold. First, existing a posteriori error estimation methods are considered to be quite reliable in practice but still involve uncertain ingredients. In some cases the expressions that bound the error involve continuity or interpolation constants which are non-computable and can be approximated accurately only when the solution is well resolved [13,15]. In other situations, the bounds are uniform for any level of mesh refinement, but in practice are only computable by assuming that the exact solution can be locally represented on a conservatively refined mesh [11,12]. Alternatively, numerical integration may be required to evaluate integrals involving analytic functions [19]. Second, the above solution algorithms are implemented in computer codes which can easily have hundredths of thousands of lines of code, the correctness of which is virtually impossible to verify in practice.

In this paper, we present a method to compute upper and lower bounds for linear outputs of interest of the exact weak solutions of the linear elasticity equations. The method is described in detail for two dimensions but the extension to the three dimensional case does not seem to present major difficulties. The approach presented can be interpreted as a generalization of the well-known complementary energy principle [7]. This principle, which in its original form only yields bounds for the energy norm of the solution, has been known for many decades. Here, an extension to linear outputs of interest is presented.

The starting point for our bounds procedure are finite element approximations to the displacement solution and to the output dependent adjoint solution. These approximations are then post-processed to yield the so called inter-element hybrid fluxes. The hybrid fluxes are then used as data for the computation of locally equilibrated stress fields. The final expression for the bounds is obtained by calculating appropriate norms of the stress fields. It is shown that the computed bounds are uniformly valid regardless of the size of the underlying coarse discretization, but as expected, their sharpness depends on the accuracy of the approximated solutions. A mesh adaptive procedure is also described which can be used to determine the bounds to a preset level of accuracy. Many of the components involved in this approach were presented in [17,18] for scalar equations. In this paper we emphasize those aspects of the method which are particular to the elasticity equations.

An attractive feature of the proposed approach is that the piecewise polynomial equilibrated stress-like fields, which are computed as part of the bound process, can be used as certificates to guarantee the correctness of the computed bounds. It turns out that given a stress field it is easy to check whether this field corresponds to a valid certificate, and in the affirmative case it is straightforward to determine the value of the output that it can certify. In particular, the stress fields need to satisfy continuity of normal tractions across elements, and membership of an appropriate space.

The idea of a certificate that is computed simultaneously with the solution has many attractive features. In particular, a certificate consisting of the data set necessary to describe the piecewise polynomial stress-like fields could be used to document the computed results. We note that exercising the certificate does not require access to the code used to compute it and can be done with a simple algorithm which does not require solving any system of equations. A very important point is that, if a certificate meets all the necessary conditions, which in turn are easy to verify, then there is no need to certify the code used to compute it. In practice, the size of these certificates depends on the required level of certainty. As expected, we shall find that high levels of certainty, i.e. small bound gaps, will often require longer certificates (larger data sets) than those required to certify less sharp claims.

## 2. Problem definition

The elasticity equations on a general two dimensional polygonal domain  $\Omega$  are considered. The boundary,  $\Gamma = \partial\Omega$ , is divided into two complementary disjoint parts  $\Gamma^D$  and  $\Gamma^N$ , where essential and Neumann boundary conditions are imposed, respectively. Furthermore, the boundary  $\Gamma^D$  is assumed to be a non empty set. The weak formulation of the problem consists of finding the displacements  $\bar{\mathbf{u}} \in \mathcal{U}$ , such that

$$a(\bar{\mathbf{u}}, \mathbf{v}) = l(\mathbf{v}) \quad \forall \mathbf{v} \in \mathcal{V},$$

where  $\mathcal{U} = \{\mathbf{v} \in [\mathcal{H}^1(\Omega)]^2, \mathbf{v}|_{\Gamma^D} = \mathbf{g}_D\}$  is the space of admissible displacement fields,  $\mathcal{V} = \{\mathbf{v} \in [\mathcal{H}^1(\Omega)]^2, \mathbf{v}|_{\Gamma^D} = \mathbf{0}\}$  is the space of test functions, and  $\mathcal{H}^1(\Omega)$  denotes the standard Sobolev space.

The linear forcing functional  $l: [\mathcal{H}^1(\Omega)]^2 \rightarrow \mathbb{R}$

$$l(\mathbf{v}) = \int_{\Omega} \mathbf{f} \cdot \mathbf{v} \, d\Omega + \int_{\Gamma^N} \mathbf{g} \cdot \mathbf{v} \, d\Gamma,$$

contains both the internal forces per unit volume  $\mathbf{f} \in [\mathcal{H}^{-1}(\Omega)]^2$  and the Neumann boundary tractions  $\mathbf{g} \in [\mathcal{H}^{-\frac{1}{2}}(\Gamma^N)]^2$ , and  $a: [\mathcal{H}^1(\Omega)]^2 \times [\mathcal{H}^1(\Omega)]^2 \rightarrow \mathbb{R}$  is the symmetric positive definite bilinear form given by

$$a(\mathbf{w}, \mathbf{v}) = \int_{\Omega} \boldsymbol{\sigma}(\mathbf{w}) : \boldsymbol{\varepsilon}(\mathbf{v}) \, d\Omega.$$

Here,  $\boldsymbol{\varepsilon}(\mathbf{v})$  is the second order deformation tensor, which is defined as the symmetric part of the gradient tensor  $\nabla \mathbf{v}$ , that is,  $\boldsymbol{\varepsilon}(\mathbf{v}) = \frac{1}{2}(\nabla \mathbf{v} + (\nabla \mathbf{v})^T)$ . The stress tensor  $\boldsymbol{\sigma}(\mathbf{v})$ , is related to the deformation tensor through a linear constitutive relation of the form,  $\boldsymbol{\sigma}(\mathbf{v}) = \mathbb{C} : \boldsymbol{\varepsilon}(\mathbf{v})$ , where  $\mathbb{C}$  is the fourth-order elasticity tensor. Throughout the paper the energy norm induced by the bilinear form  $a(\cdot, \cdot)$  is denoted by  $\|\cdot\|$ , that is,  $\|\mathbf{v}\|^2 = a(\mathbf{v}, \mathbf{v})$ .

Our goal is to compute bounds for output quantities of interest,  $\bar{s} = \ell^0(\bar{\mathbf{u}})$ , where  $\ell^0: [\mathcal{H}^1(\Omega)]^2 \rightarrow \mathbb{R}$  is a linear continuous functional defined as

$$\ell^0(\mathbf{v}) = \int_{\Omega} \mathbf{f}^0 \cdot \mathbf{v} \, d\Omega + \int_{\Gamma^N} \mathbf{g}^0 \cdot \mathbf{v} \, d\Gamma - a(\mathbf{u}^0, \mathbf{v}) \quad (1)$$

for given  $\mathbf{f}^0 \in [\mathcal{H}^{-1}(\Omega)]^2$ ,  $\mathbf{g}^0 \in [\mathcal{H}^{-\frac{1}{2}}(\Gamma^N)]^2$  and  $\mathbf{u}^0 \in [\mathcal{H}^{-1}(\Omega)]^2$ . Note that this form may easily incorporate, as particular cases, displacements or tractions integrated over arbitrary subdomains or boundary segments.

For any given  $\mathbf{u}^D \in \mathcal{U}$  we can write  $\bar{\mathbf{u}} = \mathbf{u}^D + \mathbf{u}$ , where  $\mathbf{u} \in \mathcal{V}$  is the solution of

$$a(\mathbf{u}, \mathbf{v}) = l(\mathbf{v}) - a(\mathbf{u}^D, \mathbf{v}) =: \ell(\mathbf{v}) \quad \forall \mathbf{v} \in \mathcal{V} \quad (2)$$

and therefore the output of interest can be rewritten as

$$\bar{s} = \ell^0(\mathbf{u}^D) + \ell^0(\mathbf{u}).$$

Working with  $\mathbf{u} \in \mathcal{V}$  rather than  $\bar{\mathbf{u}}$  has the advantage of avoiding the notational complexity introduced by the non-homogeneous Dirichlet boundary conditions. Thus, our goal will be to compute bounds for  $s = \ell^0(\mathbf{u})$ , from which we can easily evaluate the bounds for  $\bar{s} = \ell^0(\bar{\mathbf{u}})$ .

### 3. Output bounds

A key ingredient to our bound procedure, is the reformulation of our output of interest as a constrained minimization problem. We write the output of interest  $s = \ell^0(\mathbf{u})$  as

$$\begin{aligned} s &= \inf_{\mathbf{v} \in \mathcal{V}} \ell^0(\mathbf{v}) + \kappa^2(a(\mathbf{v}, \mathbf{v}) - \ell(\mathbf{v})) \\ \text{s.t. } a(\mathbf{v}, \boldsymbol{\varphi}) &= \ell(\boldsymbol{\varphi}) \quad \forall \boldsymbol{\varphi} \in \mathcal{V} \end{aligned} \quad (3)$$

for all  $\kappa \in \mathbb{R}$ . The above statement can be easily verified by noting that, from (2), the constraint  $a(\mathbf{v}, \boldsymbol{\varphi}) = \ell(\boldsymbol{\varphi}) \quad \forall \boldsymbol{\varphi} \in \mathcal{V}$  is only satisfied when  $\mathbf{v} = \mathbf{u}$  and clearly  $a(\mathbf{u}, \mathbf{u}) - \ell(\mathbf{u}) = 0$ . Now, the Lagrangian associated with the above constrained minimization problem is given by

$$L(\mathbf{v}, \boldsymbol{\varphi}) = \ell^0(\mathbf{v}) + \kappa^2(a(\mathbf{v}, \mathbf{v}) - \ell(\mathbf{v})) + \ell(\boldsymbol{\varphi}) - a(\mathbf{v}, \boldsymbol{\varphi}) \quad (4)$$

and problem (3) becomes

$$s = \inf_{\mathbf{v} \in \mathcal{V}} \sup_{\boldsymbol{\varphi} \in \mathcal{V}} L(\mathbf{v}, \boldsymbol{\varphi}). \quad (5)$$

A lower bound,  $s^-$ , for the output,  $s$ , can be easily found using the strong duality of convex minimization and the saddle point property of the Lagrange multipliers as

$$s = \inf_{\mathbf{v} \in \mathcal{V}} \sup_{\boldsymbol{\varphi} \in \mathcal{V}} L(\mathbf{v}, \boldsymbol{\varphi}) = \sup_{\boldsymbol{\varphi} \in \mathcal{V}} \inf_{\mathbf{v} \in \mathcal{V}} L(\mathbf{v}, \boldsymbol{\varphi}) \geq \inf_{\mathbf{v} \in \mathcal{V}} L(\mathbf{v}, \tilde{\boldsymbol{\varphi}}) \equiv s^- \quad \forall \tilde{\boldsymbol{\varphi}} \in \mathcal{V}, \quad (6)$$

where in order to obtain sharp bounds, it is important to use a good approximation  $\tilde{\boldsymbol{\varphi}}$  of the Lagrange multiplier. Thus, we see from (6), that the problem of computing a lower bound for the output of interest is cast as an unconstrained minimization problem.

The optimal value of for the Lagrange multiplier is obtained by solving the saddle problem (5) and is given by  $\boldsymbol{\varphi}^* = \kappa^2 \mathbf{u} + \boldsymbol{\psi}$  where  $\boldsymbol{\psi} \in \mathcal{V}$  is the solution of the problem,

$$a(\mathbf{v}, \boldsymbol{\psi}) = \ell^0(\mathbf{v}) \quad \forall \mathbf{v} \in \mathcal{V} \quad (7)$$

called dual or adjoint problem with respect to the selected output  $\ell^0(\cdot)$ . Note that due to the symmetry of  $a(\cdot, \cdot)$ , the only difference between the primal problem (2) for  $\mathbf{u}$ , and the dual problem (7) for  $\boldsymbol{\psi}$ , is only in the forcing data ( $\mathbf{f}^0$  instead of  $\mathbf{f}$ ,  $\mathbf{g}^0$  instead of  $\mathbf{g}$  and  $\mathbf{u}^0$  instead of  $\mathbf{u}^D$ ).

The solutions of the primal and dual problems are approximated by  $\mathbf{u}_h$  and  $\boldsymbol{\psi}_h$  respectively, which lie in a finite element interpolation space  $\mathcal{V}^h \subset \mathcal{V}$ , associated with a finite element mesh of characteristic size  $h$  and verify



$$a(\mathbf{u}_h, \mathbf{v}) = \ell(\mathbf{v}) \quad \mathbf{v} \in \mathcal{V}^h$$

and

$$a(\mathbf{v}, \boldsymbol{\psi}_h) = \ell^0(\mathbf{v}) \quad \mathbf{v} \in \mathcal{V}^h.$$

An approximation to the Lagrange multiplier  $\tilde{\boldsymbol{\varphi}}$ , is now obtained by setting  $\tilde{\boldsymbol{\varphi}} = \kappa^2 \mathbf{u}_h + \boldsymbol{\psi}_h$ . We note, however, that different options are also possible (see for instance [6,14]). With our choice for  $\tilde{\boldsymbol{\varphi}}$  the optimization over  $\mathbf{v}$  in (6), leads to

$$s^- = \frac{1}{4} \left\| \kappa \mathbf{u}_h + \frac{1}{\kappa} \boldsymbol{\psi}_h \right\|^2 - \frac{1}{4} \left\| \kappa \mathbf{u} - \frac{1}{\kappa} \boldsymbol{\psi} \right\|^2. \quad (8)$$

**Remark 1.** In the particular case when  $\ell^0(\cdot) = \ell(\cdot)$ , then  $s = \ell^0(\mathbf{u}) = \|\mathbf{u}\|^2$  and also  $\mathbf{u} = \boldsymbol{\psi}$  and  $\mathbf{u}_h = \boldsymbol{\psi}_h$ . In this case, the lower bound we obtain for  $\kappa = 1$ , is  $s^- = \|\mathbf{u}_h\|^2$ , which implies

$$\|\mathbf{u}_h\|^2 \leq \|\mathbf{u}\|^2.$$

This is the classical lower bound property of the energy norm of the finite element solution with respect to the exact solution norm.

An analogous expression for an upper bound,  $s^+$ , of  $s$ , is obtained by replacing  $\ell^0(\mathbf{u})$  by  $-\ell^0(\mathbf{u})$  in the original optimization problem (3) to obtain

$$\begin{aligned} -s &= \inf_{\mathbf{v} \in \mathcal{V}} -\ell^0(\mathbf{v}) + \kappa^2 (a(\mathbf{v}, \mathbf{v}) - \ell(\mathbf{v})) \\ \text{s.t. } &a(\mathbf{v}, \boldsymbol{\varphi}) = \ell(\boldsymbol{\varphi}) \quad \forall \boldsymbol{\varphi} \in \mathcal{V}. \end{aligned}$$

The optimal multiplier in this case is approximated by  $\tilde{\boldsymbol{\varphi}} = \kappa^2 \mathbf{u}_h - \boldsymbol{\psi}_h$ , and the optimization process yields

$$-s^+ = \frac{1}{4} \left\| \kappa \mathbf{u}_h - \frac{1}{\kappa} \boldsymbol{\psi}_h \right\|^2 - \frac{1}{4} \left\| \kappa \mathbf{u} + \frac{1}{\kappa} \boldsymbol{\psi} \right\|^2$$

which is equivalent to

$$s^+ = \frac{1}{4} \left\| \kappa \mathbf{u} + \frac{1}{\kappa} \boldsymbol{\psi} \right\|^2 - \frac{1}{4} \left\| \kappa \mathbf{u}_h - \frac{1}{\kappa} \boldsymbol{\psi}_h \right\|^2.$$

**Remark 2** (*Bounds for the output of the error*). If the particular lift of  $\bar{\mathbf{u}}$  is its finite element approximation, that is  $\mathbf{u}^D = \bar{\mathbf{u}}_h$ , then  $\mathbf{u}$  is the error in the finite element approximation,  $\mathbf{u} = \bar{\mathbf{u}} - \bar{\mathbf{u}}_h = \boldsymbol{\epsilon}$ , and  $\mathbf{u}_h = 0$ . In this case, the previous methodology would lead after some algebra to bounds for  $s = \ell^0(\boldsymbol{\epsilon})$

$$\begin{aligned} s^- &= -\frac{1}{4} \left\| \kappa \boldsymbol{\epsilon} - \frac{1}{\kappa} \boldsymbol{\epsilon} \right\|^2, \\ s^+ &= \frac{1}{4} \left\| \kappa \boldsymbol{\epsilon} + \frac{1}{\kappa} \boldsymbol{\epsilon} \right\|^2 \end{aligned}$$

where  $\boldsymbol{\epsilon} = \boldsymbol{\psi} - \boldsymbol{\psi}_h$  is the error in the finite element approximation of the dual problem.

Writing together the expressions for the upper and lower bounds we have

$$\frac{1}{4} \left\| \kappa \mathbf{u}_h + \frac{1}{\kappa} \boldsymbol{\psi}_h \right\|^2 - \frac{1}{4} \left\| \kappa \mathbf{u} - \frac{1}{\kappa} \boldsymbol{\psi} \right\|^2 \leq s \leq \frac{1}{4} \left\| \kappa \mathbf{u} + \frac{1}{\kappa} \boldsymbol{\psi} \right\|^2 - \frac{1}{4} \left\| \kappa \mathbf{u}_h - \frac{1}{\kappa} \boldsymbol{\psi}_h \right\|^2.$$

It is clear that these expressions are non-computable, since they depend on the exact solution of both the primal and dual problems. However, they illustrate our basic approach to obtaining bounds for outputs of interest: if we can compute upper bounds  $\|\kappa \mathbf{u} \pm (1/\kappa) \boldsymbol{\psi}\|_{\text{UB}}^2$  for  $\|\kappa \mathbf{u} \pm (1/\kappa) \boldsymbol{\psi}\|^2$ , then, we can write computable expressions for the output bounds as

$$\frac{1}{4} \left\| \kappa \mathbf{u}_h + \frac{1}{\kappa} \boldsymbol{\psi}_h \right\|^2 - \frac{1}{4} \left\| \kappa \mathbf{u} - \frac{1}{\kappa} \boldsymbol{\psi} \right\|_{\text{UB}}^2 \leq s \leq \frac{1}{4} \left\| \kappa \mathbf{u} - \frac{1}{\kappa} \boldsymbol{\psi} \right\|_{\text{UB}}^2 - \frac{1}{4} \left\| \kappa \mathbf{u}_h + \frac{1}{\kappa} \boldsymbol{\psi}_h \right\|^2. \quad (9)$$

In the next section, we present an approach for computing upper bounds for the energy norm of the solution of the elasticity equations. This result is then generalized in Section 5 to compute the upper bounds for the linear combination of the primal and dual functions,  $\kappa \mathbf{u} \pm (1/\kappa) \boldsymbol{\psi}$ .

#### 4. Upper bounds for the energy norm

Consider the generalized elasticity problem with Neumann and homogeneous Dirichlet boundary conditions written in weak form as: find  $\mathbf{z} \in \mathcal{V}$  such that

$$a(\mathbf{z}, \mathbf{v}) = \ell^*(\mathbf{u}) \quad \forall \mathbf{v} \in \mathcal{V}, \quad (10)$$

where

$$\ell^*(\mathbf{v}) = \int_{\Omega} \mathbf{f}^* \cdot \mathbf{v} \, d\Omega + \int_{\Gamma^N} \mathbf{g}^* \cdot \mathbf{v} \, d\Gamma - a(\mathbf{u}^*, \mathbf{v}).$$

It is clear that any linear combination of the primal and dual solutions,  $\alpha \mathbf{u} + \beta \boldsymbol{\psi}$ ,  $\alpha, \beta \in \mathbb{R}$ , is the solution of problem (10) with  $\mathbf{f}^* = \alpha \mathbf{f} + \beta \mathbf{f}^0$ ,  $\mathbf{g}^* = \alpha \mathbf{g} + \beta \mathbf{g}^0$  and  $\mathbf{u}^* = \alpha \mathbf{u}^D + \beta \mathbf{u}^0$ . In particular, the choice  $\alpha = \kappa$ ,  $\beta = \pm 1/\kappa$  will be used later to obtain the required upper bounds for  $\|\kappa \mathbf{u} \pm (1/\kappa) \boldsymbol{\psi}\|^2$ .

In this section we consider the problem of computing an upper bound for  $\|\mathbf{z}\|^2$ . We recall that  $\|\mathbf{z}\|^2$  can be obtained as the solution of the optimization procedure

$$\|\mathbf{z}\|^2 = \sup_{\mathbf{v} \in \mathcal{V}} 2\ell^*(\mathbf{v}) - a(\mathbf{v}, \mathbf{v}). \quad (11)$$

The above problem is to be considered over an infinite dimensional space of functions which are defined over the whole domain  $\Omega$ . In order to come up with a computable expression for an upper bound of  $\|\mathbf{z}\|^2$ , two relaxations are introduced. First, a domain decomposition strategy is used to transform the maximization problem over functions in  $\Omega$ , into a number of independent problems which are defined over subdomains (triangles in our case). Second, duality is exploited to transform each of the convex maximization problems into a feasibility problem for the dual functional which is shown to yield an upper bound for the optimal solution.

##### 4.1. Domain decomposition

We consider a triangulation of the computational domain  $\Omega$  into  $n_{e1}$  triangles and denote by  $\Omega_k$  a generic triangle,  $k = 1, \dots, n_{e1}$ . Let  $\Gamma_h$  be the set of all the edges in the mesh, and  $\boldsymbol{\Lambda} = \prod_{k=1}^{n_{e1}} [\mathcal{H}^{-\frac{1}{2}}(\partial\Omega_k)]^2$  the space of integrable tractions in  $\Gamma_h$ . The set of all the interior edges of the mesh is denoted by  $\Gamma^I$ , that is  $\Gamma_h = \Gamma \cup \Gamma^I$ . For each edge  $\gamma \in \Gamma_h$  a unit normal direction,  $\mathbf{n}^\gamma$ , is assigned such that, if  $\gamma$  is an exterior edge,

$\mathbf{n}^\gamma$  coincides with the outward unit normal to  $\Gamma$ . Similarly, given an element  $\Omega_k$  and an edge of this element  $\gamma \in \partial\Omega_k$ , the outward normal to the element associated to  $\gamma$  is denoted by  $\mathbf{n}_k^\gamma$ . Then,  $\tau_k$  is defined as  $\tau_k|_\gamma = \mathbf{n}_k^\gamma \cdot \mathbf{n}^\gamma$ , that is

$$\tau_k|_\gamma = \mathbf{n}_k^\gamma \cdot \mathbf{n}^\gamma = \begin{cases} 1 & \text{if } \mathbf{n}_k^\gamma = \mathbf{n}^\gamma, \\ -1 & \text{if } \mathbf{n}_k^\gamma = -\mathbf{n}^\gamma. \end{cases}$$

Note that if  $\gamma = \partial\Omega_k \cap \partial\Omega_l$ , then  $\tau_k|_\gamma + \tau_l|_\gamma = 0$ .

The broken space  $\widehat{\mathcal{V}}$  is introduced by relaxing in  $\mathcal{V}$  both the Dirichlet homogeneous boundary conditions and the continuity of the functions across the edges of  $\Gamma_h$ , that is,

$$\widehat{\mathcal{V}} = \{ \hat{\mathbf{v}} \in [\mathcal{L}^2(\Omega)]^2, \hat{\mathbf{v}}|_{\Omega_k} \in [\mathcal{H}^1(\Omega_k)]^2 \quad \forall \Omega_k \in \Omega \}.$$

Given a function in the broken space  $\hat{\mathbf{v}} \in \widehat{\mathcal{V}}$ , the jump of  $\hat{\mathbf{v}}$  across the mesh edges is defined as

$$[\hat{\mathbf{v}}]_\gamma = \begin{cases} \hat{\mathbf{v}}|_{\Omega_k} \tau_k|_\gamma + \hat{\mathbf{v}}|_{\Omega_l} \tau_l|_\gamma, & \text{if } \gamma = \partial\Omega_k \cap \partial\Omega_l \in \Gamma^I, \\ \hat{\mathbf{v}}, & \text{if } \gamma \in \Gamma, \end{cases}$$

where the definition of the jump depends on the arbitrary choice of the edge normals. Note that if  $\hat{\mathbf{v}}$  is a continuous function verifying the Dirichlet homogeneous boundary conditions,  $\hat{\mathbf{v}} \in \mathcal{V}$ , then  $[\hat{\mathbf{v}}] = 0$  in  $\Gamma^I \cup \Gamma^D$ .

Then, given a broken function  $\hat{\mathbf{v}} \in \widehat{\mathcal{V}}$ , the continuity at inter-elemental edges and Dirichlet homogeneous boundary conditions in  $\Gamma^D$  can be enforced weakly through the bilinear form  $b: \widehat{\mathcal{V}} \times \Lambda \rightarrow \mathbb{R}$

$$b(\hat{\mathbf{v}}, \boldsymbol{\lambda}) = \sum_{\gamma \in \Gamma^I \cup \Gamma^D} \int_\gamma \boldsymbol{\lambda} \cdot [\hat{\mathbf{v}}] d\Gamma = \sum_{k=1}^{n_{e1}} \int_{\partial\Omega_k \setminus \Gamma^N} \tau_k \boldsymbol{\lambda} \cdot \hat{\mathbf{v}}|_{\Omega_k} d\Gamma$$

by imposing  $b(\hat{\mathbf{v}}, \boldsymbol{\lambda}) = 0$  for all  $\boldsymbol{\lambda} \in \Lambda$ . Therefore, the space of test functions  $\mathcal{V}$  can be recovered as

$$\mathcal{V} = \{ \hat{\mathbf{v}} \in \widehat{\mathcal{V}}, b(\hat{\mathbf{v}}, \boldsymbol{\lambda}) = 0 \quad \forall \boldsymbol{\lambda} \in \Lambda \}.$$

Let us denote by  $\mathcal{V}_k$  the restriction of the broken space  $\widehat{\mathcal{V}}$  to the element  $\Omega_k$ , that is,  $\mathcal{V}_k = \widehat{\mathcal{V}}|_{\Omega_k} = [\mathcal{H}^1(\Omega_k)]^2$ . Formally, any function  $\mathbf{v}_k \in \mathcal{V}_k$  is not defined in the whole domain  $\Omega$  but only in the element  $\Omega_k$ . In the following, any function  $\mathbf{v}_k \in \mathcal{V}_k$  is naturally extended to  $\Omega$  by setting the values outside  $\Omega_k$  to zero. Then, a function  $\hat{\mathbf{v}} \in \widehat{\mathcal{V}}$  can be decomposed as the sum of its restrictions to each element  $\mathbf{v}_k = \hat{\mathbf{v}}|_{\Omega_k} \in \mathcal{V}_k$ , that is,  $\hat{\mathbf{v}} = \sum_{k=1}^{n_{e1}} \mathbf{v}_k$ , and  $\widehat{\mathcal{V}} = \bigoplus_{k=1}^{n_{e1}} \mathcal{V}_k$ .

We can now rewrite the maximization problem of Eq. (11) as a constrained saddle problem defined over functions in  $\widehat{\mathcal{V}}$  and Lagrange multipliers in  $\Lambda$  as,

$$\|\mathbf{z}\|^2 = \sup_{\hat{\mathbf{v}} \in \widehat{\mathcal{V}}} \inf_{\boldsymbol{\lambda} \in \Lambda} J(\hat{\mathbf{v}}, \boldsymbol{\lambda}),$$

where  $J(\hat{\mathbf{v}}, \boldsymbol{\lambda})$  is the quadratic-linear Lagrangian which can be expressed using the local restrictions  $\mathbf{v}_k$  of  $\hat{\mathbf{v}}$  as

$$J(\hat{\mathbf{v}}, \boldsymbol{\lambda}) = \sum_{k=1}^{n_{e1}} J_k(\mathbf{v}_k, \boldsymbol{\lambda}),$$

where

$$J_k(\mathbf{v}_k, \boldsymbol{\lambda}) = 2\ell_k^*(\mathbf{v}_k) - a_k(\mathbf{v}_k, \mathbf{v}_k) + 2b_k(\mathbf{v}_k, \boldsymbol{\lambda}). \tag{12}$$

Here, the subscript  $k$  denotes the restriction of the linear and bilinear forms to the element  $\Omega_k$ , that is,

$$a_k(\mathbf{w}, \mathbf{v}) = \int_{\Omega_k} \boldsymbol{\sigma}(\mathbf{w}): \boldsymbol{\varepsilon}(\mathbf{v}) d\Omega,$$

$$\ell_k^*(\mathbf{v}) = \int_{\Omega_k} \mathbf{f}^* \cdot \mathbf{v} d\Omega + \int_{\Gamma^N \cap \partial\Omega_k} \mathbf{g}^* \cdot \mathbf{v} d\Gamma - a_k(\mathbf{u}^*, \mathbf{v}),$$

and

$$b_k(\mathbf{v}, \boldsymbol{\lambda}) = \int_{\partial\Omega_k \setminus \Gamma^N} \tau_k \boldsymbol{\lambda} \cdot \mathbf{v}|_{\Omega_k} d\Gamma.$$

Then, using standard duality arguments an upper bound for  $\|\mathbf{z}\|^2$  is obtained as

$$\|\mathbf{z}\|^2 = \sup_{\hat{\mathbf{v}} \in \hat{\mathcal{V}}} \inf_{\boldsymbol{\lambda} \in \Lambda} J(\hat{\mathbf{v}}, \boldsymbol{\lambda}) = \inf_{\boldsymbol{\lambda} \in \Lambda} \sup_{\hat{\mathbf{v}} \in \hat{\mathcal{V}}} J(\hat{\mathbf{v}}, \boldsymbol{\lambda}) \leq \sup_{\hat{\mathbf{v}} \in \hat{\mathcal{V}}} J(\hat{\mathbf{v}}, \tilde{\boldsymbol{\lambda}}) \quad \forall \tilde{\boldsymbol{\lambda}} \in \Lambda. \quad (13)$$

Clearly, the Lagrange multiplier  $\tilde{\boldsymbol{\lambda}}$  has to be properly chosen in order to ensure that the resulting maximization is bounded from above and that the resulting optimum is accurate. The importance of the weak imposition of the continuity requirement and the approximation of the Lagrange multiplier, is that once the Lagrange multiplier is fixed, the lagrangian  $J(\hat{\mathbf{v}}, \tilde{\boldsymbol{\lambda}})$  decomposes into local elementary contributions, and the maximization in (13) decomposes into local maximization problems in the elements of the mesh.

In order to simplify the notation, we will rewrite  $J_k(\cdot, \tilde{\boldsymbol{\lambda}})$  in a simpler way. We note that given a Lagrange multiplier  $\tilde{\boldsymbol{\lambda}} \in \Gamma_h$ , the values of  $\tilde{\boldsymbol{\lambda}}$  in  $\Gamma^N$  do not contribute to  $J_k(\cdot, \tilde{\boldsymbol{\lambda}})$ . Therefore we can define

$$\tilde{\boldsymbol{\lambda}}|_{\Gamma^N} = \mathbf{g}^*, \quad (14)$$

so that,

$$\int_{\Gamma^N \cap \partial\Omega_k} \mathbf{g}^* \cdot \mathbf{v}_k d\Gamma + \int_{\partial\Omega_k \setminus \Gamma^N} \tau_k \tilde{\boldsymbol{\lambda}} \cdot \mathbf{v}_k d\Gamma = \int_{\partial\Omega_k} \tau_k \tilde{\boldsymbol{\lambda}} \cdot \mathbf{v}_k d\Gamma$$

and therefore,

$$J_k(\mathbf{v}_k, \tilde{\boldsymbol{\lambda}}) = 2 \int_{\Omega_k} \mathbf{f}^* \cdot \mathbf{v}_k d\Omega + 2 \int_{\partial\Omega_k} \tau_k \tilde{\boldsymbol{\lambda}} \cdot \mathbf{v}_k d\Gamma - 2a_k(\mathbf{u}^*, \mathbf{v}_k) - a_k(\mathbf{v}_k, \mathbf{v}_k). \quad (15)$$

Thus the global maximization of Eq. (13) can be decomposed as,

$$\sup_{\hat{\mathbf{v}} \in \hat{\mathcal{V}}} J(\hat{\mathbf{v}}, \tilde{\boldsymbol{\lambda}}) = \sum_{k=1}^{n_{e1}} \sup_{\mathbf{v}_k \in \mathcal{V}_k} J_k(\mathbf{v}_k, \tilde{\boldsymbol{\lambda}}), \quad (16)$$

allowing to obtain an upper bound for  $\|\mathbf{z}\|^2$ , maximizing the the local functionals  $J_k(\cdot, \tilde{\boldsymbol{\lambda}})$  in each element of the mesh independently.

This local maximization problems, although local, can not be solved exactly because  $\mathcal{V}_k$  is an infinite dimensional space. Moreover, if we replace  $\mathcal{V}_k$  with a finite dimensional subspace, the upper bound property is lost.

#### 4.2. Complementary energy relaxation

We consider now the problem of finding computable upper bounds for the local maximization problems (16), that is, find  $v_k \in \mathbb{R}$  such that,

$$\sup_{\mathbf{v}_k \in \mathcal{V}_k} J_k(\mathbf{v}_k, \tilde{\boldsymbol{\lambda}}) \leq v_k, \quad (17)$$

so that a global upper bound for  $\|\mathbf{z}\|^2$  will be recovered as

$$\|\mathbf{z}\|^2 \leq \sum_{k=1}^{n_{e1}} \sup_{\mathbf{v}_k \in \mathcal{V}_k} J_k(\mathbf{v}_k, \tilde{\boldsymbol{\lambda}}) \leq \sum_{k=1}^{n_{e1}} v_k.$$

The upper bounds  $v_k$  are computed using a standard duality argument which transforms the problem of finding the maximum over the infinite dimensional space  $\mathcal{V}_k$  to a problem of finding a feasible solution in an appropriate finite dimensional space.

Let  $\mathcal{S}_k$  denote the space of componentwise square-integrable stress fields in  $\Omega_k$ , that is,  $\mathcal{S}_k$  contains all the second-order tensors with  $\sigma_{ij} \in \mathcal{L}^2(\Omega_k)$ ,  $\forall i, j$ . Then,  $\mathcal{S}_k^{\text{eq}}$  denotes the subset of  $\mathcal{S}_k$  which contains all the equilibrated stress fields with respect to  $\mathbf{f}^*$ ,  $\tilde{\boldsymbol{\lambda}}$  and  $\mathbf{u}^*$ , that is,  $\boldsymbol{\sigma}_k \in \mathcal{S}_k^{\text{eq}}$  verifies

$$\int_{\Omega_k} \boldsymbol{\sigma}_k : \boldsymbol{\varepsilon}(\mathbf{v}_k) d\Omega = \int_{\Omega_k} \mathbf{f}^* \cdot \mathbf{v}_k d\Omega + \int_{\partial\Omega_k} \tau_k \tilde{\boldsymbol{\lambda}} \cdot \mathbf{v}_k d\Gamma - a_k(\mathbf{u}^*, \mathbf{v}_k) \quad \forall \mathbf{v}_k \in \mathcal{V}_k. \quad (18)$$

The stress fields in  $\mathcal{S}_k^{\text{eq}}$  are usually referred to as being statically admissible. In addition, we define the complementary energy of a stress field  $\boldsymbol{\sigma}_k \in \mathcal{S}_k$ , as the value given by the the functional  $J_k^c : \mathcal{S}_k \rightarrow \mathbb{R}$ ,

$$J_k^c(\boldsymbol{\sigma}_k) = \int_{\Omega_k} \boldsymbol{\sigma}_k : \mathbb{C}^{-1} : \boldsymbol{\sigma}_k d\Omega.$$

**Lemma 1.** *If  $\tilde{\boldsymbol{\lambda}}$  is such that  $J_k(\mathbf{v}_k, \tilde{\boldsymbol{\lambda}})$  is bounded from above for all  $\mathbf{v}_k \in \mathcal{V}_k$ , then the following duality relation holds:*

$$\sup_{\mathbf{v}_k \in \mathcal{V}_k} J_k(\mathbf{v}_k, \tilde{\boldsymbol{\lambda}}) = \inf_{\boldsymbol{\sigma}_k \in \mathcal{S}_k^{\text{eq}}} J_k^c(\boldsymbol{\sigma}_k).$$

**Proof.** Let  $\boldsymbol{\sigma}_k \in \mathcal{S}_k^{\text{eq}}$  and  $\mathbf{v}_k \in \mathcal{V}_k$ , then

$$\begin{aligned} 0 &\leq \int_{\Omega_k} (\boldsymbol{\sigma}_k - \boldsymbol{\sigma}(\mathbf{v}_k)) : \mathbb{C}^{-1} : (\boldsymbol{\sigma}_k - \boldsymbol{\sigma}(\mathbf{v}_k)) d\Omega \\ &= \int_{\Omega_k} \boldsymbol{\sigma}_k : \mathbb{C}^{-1} : \boldsymbol{\sigma}_k d\Omega + \int_{\Omega_k} \boldsymbol{\sigma}(\mathbf{v}_k) : \boldsymbol{\varepsilon}(\mathbf{v}_k) d\Omega - 2 \int_{\Omega_k} \boldsymbol{\sigma}_k : \boldsymbol{\varepsilon}(\mathbf{v}_k) d\Omega \\ &= J_k^c(\boldsymbol{\sigma}_k) + a_k(\mathbf{v}_k, \mathbf{v}_k) - 2 \int_{\Omega_k} \boldsymbol{\sigma}_k : \boldsymbol{\varepsilon}(\mathbf{v}_k) d\Omega. \end{aligned}$$

Now, since  $\boldsymbol{\sigma}_k \in \mathcal{S}_k^{\text{eq}}$  Eq. (18) holds true and

$$2 \int_{\Omega_k} \boldsymbol{\sigma}_k : \boldsymbol{\varepsilon}(\mathbf{v}_k) d\Omega - a_k(\mathbf{v}_k, \mathbf{v}_k) = J_k(\mathbf{v}_k, \tilde{\boldsymbol{\lambda}}),$$

leading to

$$0 \leq J_k^c(\boldsymbol{\sigma}_k) - J_k(\mathbf{v}_k, \tilde{\boldsymbol{\lambda}}),$$

which implies  $J_k(\mathbf{v}_k, \tilde{\boldsymbol{\lambda}}) \leq J_k^c(\boldsymbol{\sigma}_k)$ .

Now, let  $\bar{\mathbf{v}}_k$  be the point at which  $J_k(\mathbf{v}_k, \tilde{\boldsymbol{\lambda}})$  is maximum, that is,

$$\|\bar{\mathbf{v}}_k\|_k^2 = \sup_{\mathbf{v}_k \in \mathcal{V}_k} J_k(\mathbf{v}_k, \tilde{\boldsymbol{\lambda}})$$

where  $\|\mathbf{v}_k\|_k^2 = a_k(\mathbf{v}_k, \mathbf{v}_k)$ . Moreover, the gradient condition for  $\bar{\mathbf{v}}_k$  leads to

$$a_k(\bar{\mathbf{v}}_k, \mathbf{v}_k) = \int_{\Omega_k} \mathbf{f}^* \cdot \mathbf{v}_k d\Omega + \int_{\partial\Omega_k} \tau_k \tilde{\boldsymbol{\lambda}} \cdot \mathbf{v}_k d\Gamma - a_k(\mathbf{u}^*, \mathbf{v}_k) \quad \forall \mathbf{v}_k \in \mathcal{V}_k$$

from where it follows that  $\bar{\boldsymbol{\sigma}}_k = \boldsymbol{\sigma}(\bar{\mathbf{v}}_k) \in \mathcal{S}_k^{\text{eq}}$ .

Thus,  $\forall \boldsymbol{\sigma}_k \in \mathcal{S}_k^{\text{eq}}, \forall \mathbf{v}_k \in \mathcal{V}_k$

$$J_k(\mathbf{v}_k, \tilde{\boldsymbol{\lambda}}) \leq J_k(\bar{\mathbf{v}}_k, \tilde{\boldsymbol{\lambda}}) = \|\bar{\mathbf{v}}_k\|_k^2 = J_k^c(\bar{\boldsymbol{\sigma}}_k) \leq J_k^c(\boldsymbol{\sigma}_k)$$

and the lemma is proved.  $\square$

Lemma 1 provides the key to the obtention of the local upper bounds  $v_k$ . It is sufficient to compute a statically admissible stress field in  $\boldsymbol{\sigma}_k^* \in \mathcal{S}_k^{\text{eq}}$ , and then evaluate its complementary energy

$$\sup_{\mathbf{v}_k \in \mathcal{V}_k} J_k(\mathbf{v}_k, \tilde{\boldsymbol{\lambda}}) = \inf_{\boldsymbol{\sigma}_k \in \mathcal{S}_k^{\text{eq}}} J_k^c(\boldsymbol{\sigma}_k) \leq J_k^c(\boldsymbol{\sigma}_k^*) = v_k. \quad (19)$$

The remainder of the section is devoted to show that one can chose the statically admissible stress field to be piecewise polynomial provided that the forcing data  $\mathbf{f}^*, \mathbf{g}^*$  and the displacement fields  $\mathbf{u}^*$  are piecewise polynomial functions.

#### 4.3. Upper bound computation

In order to construct the statically admissible stress field required in Eq. (19), it is first necessary to evaluate the Lagrange multiplier  $\tilde{\boldsymbol{\lambda}}$  satisfying the necessary constraints. We will then construct the stress fields  $\boldsymbol{\sigma}_k^*$  inside the elements.

##### 4.3.1. Lagrange multiplier approximation

In order to obtain a sharp upper bound, the choice of the Lagrange multiplier is critical. In particular, the maximization in Eq. (17) must be bounded from above.

From Eq. (15), we note that the Lagrange multiplier  $\tilde{\boldsymbol{\lambda}}$  is precisely the Neumann boundary condition for the local problems. That is, the traction distribution applied on the boundary of each element. When integrated over each element these tractions must therefore be equilibrated so that the local problems are solvable. This is equivalent to saying that

$$\ell_k^*(\mathbf{v}) + b_k(\mathbf{v}, \tilde{\boldsymbol{\lambda}}) = 0 \quad \forall \mathbf{v} \in \mathbf{P}_{\text{sm}}, \quad (20)$$

where  $\mathbf{P}_{\text{sm}}$ , is the space of solid motions which includes any combination of translations and rotation. Moreover, since the optimal traction distribution is given by the tractions of the exact solution  $\mathbf{z}$  over the edges of the elements, the Lagrange multipliers have to be both equilibrated and a good approximation to the tractions of the exact solution.

There are several known choices for the Lagrange multipliers which are approximations to the continuous tractions of the exact solution  $\mathbf{z}$  at the inter element boundaries. Here we follow the the strategy proposed by Ladeveze et al. [8].

The approximated Lagrange multiplier is denoted by  $\boldsymbol{\lambda}_h$  and it is a linear function in each edge of the mesh verifying

$$b(\hat{\mathbf{v}}, \boldsymbol{\lambda}_h) = a(\mathbf{z}_h, \hat{\mathbf{v}}) - \ell^*(\hat{\mathbf{v}}) \quad \forall \hat{\mathbf{v}} \in \widehat{\mathcal{V}}^h, \quad (21)$$

where  $\mathbf{z}_h$  is the standard Galerkin finite element approximation of  $\mathbf{z}$ . We note that for any continuous  $\hat{\mathbf{v}}$ ,  $b(\hat{\mathbf{v}}, \boldsymbol{\lambda}_h) = 0$ , and therefore  $a(\mathbf{z}_h, \hat{\mathbf{v}}) - \ell^*(\hat{\mathbf{v}}) = 0$  thus highlighting that  $\mathbf{z}_h$  is indeed the finite element approximation to  $\mathbf{z}$ .

The above equations do not determine the Lagrange multiplier  $\boldsymbol{\lambda}_h$  on  $\Gamma^{\text{N}}$ . Therefore,  $\boldsymbol{\lambda}_h$  is extended into  $\Gamma^{\text{N}}$  using Eq. (14).

**Lemma 2.** *If  $\boldsymbol{\lambda}_h$  verifies the equilibration conditions given in Eq. (21), then the local problems*

$$\sup_{\mathbf{v}_k \in \mathcal{V}_k} J_k(\mathbf{v}_k, \boldsymbol{\lambda}_h)$$

*are bounded from above.*

**Proof.** The null space of the bilinear form  $a_k(\cdot, \cdot)$  is the three dimensional space of the rigid solid motions in the element  $\Omega_k$  (translations and rotations), that is,  $\mathbf{w}_k$  is a rigid solid body motion if and only if  $a(\mathbf{v}_k, \mathbf{w}_k) = 0, \forall \mathbf{v}_k \in \mathcal{V}_k$ . Then,  $J_k(\mathbf{w}_k, \boldsymbol{\lambda}_h)$  must vanish for any rigid motion  $\mathbf{w}_k$ , otherwise, given a rigid solid motion  $\mathbf{w}_k$ , for any  $\alpha \in \mathbb{R}, \alpha \mathbf{w}_k \in \mathcal{V}_k$ ,

$$J_k(\alpha \mathbf{w}_k, \boldsymbol{\lambda}_h) = 2\alpha(\ell_k^*(\mathbf{w}_k) + b_k(\mathbf{w}_k, \tilde{\boldsymbol{\lambda}})) = \alpha J_k(\mathbf{w}_k, \boldsymbol{\lambda}_h)$$

and this will lead to an unbounded maximization problem. Let us verify that  $J_k(\mathbf{w}_k, \boldsymbol{\lambda}_h) = 0$  for any  $\mathbf{w}_k$  in the null space of  $a_k(\cdot, \cdot)$ .

Since the Lagrange multiplier  $\boldsymbol{\lambda}_h$  is equilibrated, Eq. (21) is satisfied; thus, for any rigid solid motion  $\mathbf{w}_k$  in the element  $\Omega_k$

$$b_k(\mathbf{w}_k, \boldsymbol{\lambda}_h) = a_k(\mathbf{z}_h, \mathbf{w}_k) - \ell_k^*(\mathbf{w}_k) = -\ell_k^*(\mathbf{w}_k),$$

since  $\mathbf{w}_k \in \widehat{\mathcal{V}}^h$ . Therefore,

$$J_k(\mathbf{w}_k, \boldsymbol{\lambda}_h) = 2\ell_k^*(\mathbf{w}_k) + 2b_k(\mathbf{w}_k, \boldsymbol{\lambda}_h) - a_k(\mathbf{w}_k, \mathbf{w}_k) = 0. \quad \square$$

#### 4.3.2. Construction of an equilibrated stress field $\boldsymbol{\sigma}_k \in \mathcal{S}_k^{\text{eq}}$

Once the Lagrange multipliers  $\boldsymbol{\lambda}_h$  have been determined, an equilibrated stress field  $\boldsymbol{\sigma}_k \in \mathcal{S}_k^{\text{eq}}$  has to be evaluated in order to obtain an upper bound for  $\|\mathbf{z}\|^2$ .

The existence of a piecewise polynomial stress field is established in the following theorem.

**Theorem 1.** For any given forcing function  $\mathbf{f}^*|_{\Omega_k} \in [P_r(\Omega_k)]^2$  and any equilibrated Lagrange multiplier  $\boldsymbol{\lambda}_h|_{\partial\Omega_k} \in [P_p(\partial\Omega_k)]^2$ , that is

$$\int_{\Omega_k} \mathbf{f}^* \cdot \mathbf{v} d\Omega + \int_{\partial\Omega_k} \tau_k \boldsymbol{\lambda}_h \cdot \mathbf{v} d\Gamma = 0 \quad \forall \mathbf{v} \in \mathbf{P}_{\text{sm}}, \quad (22)$$

there exists at least one dual feasible solution  $\boldsymbol{\sigma}_k \in \mathcal{S}_k = \{\boldsymbol{\sigma}, \sigma_{ij} \in \mathcal{L}^2(\Omega_k) \forall i, j\}$ , verifying

$$\int_{\Omega_k} \boldsymbol{\sigma}_k : \boldsymbol{\varepsilon}(\mathbf{v}) d\Omega = \int_{\Omega_k} \mathbf{f}^* \cdot \mathbf{v} d\Omega + \int_{\partial\Omega_k} \tau_k \boldsymbol{\lambda}_h \cdot \mathbf{v} d\Gamma \quad \forall \mathbf{v} \in \mathcal{V}_k, \quad (23)$$

which is piecewise polynomial of degree  $q$  in each component, with  $q \geq p$  and  $q > r$ .

A constructive proof of this theorem, based on some results presented in [4,9], is given in Appendix A.

**Remark 3.** For a linear equilibrated stress field (useful for linear applied tractions and constant force term), the local equilibrated stress fields are uniquely determined. Otherwise,  $q > 1$ , there are extra degrees of freedom associated both to the internal boundaries and to the space of divergence free and trace free stress fields. This degrees of freedom are used to sharpen the bounds, that is, to minimize the value of  $J_k^c(\boldsymbol{\sigma}_k)$ .

The computation of the upper bound is summarized in the box shown in Fig. 1.

## 5. Bounds for the output of interest $s$

We note that  $\mathbf{z}^\pm = \kappa \mathbf{u} \pm (1/\kappa)\boldsymbol{\psi} \in \mathcal{V}$  is the solution of the boundary value problem

$$a(\mathbf{z}^\pm, \mathbf{v}) = \kappa \ell(\mathbf{v}) \pm \frac{1}{\kappa} \ell^0(\mathbf{v}) \quad \forall \mathbf{v} \in \mathcal{V}, \quad (24)$$

which is a particular case of (10) with  $\ell^*(\mathbf{v}) = \kappa \ell(\mathbf{v}) \pm (1/\kappa)\ell^0(\mathbf{v})$ .

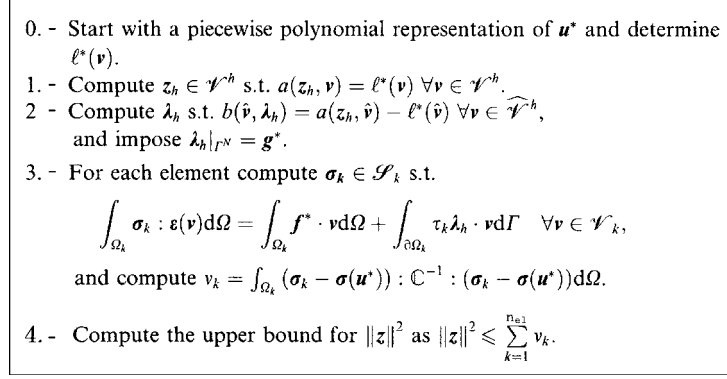


Fig. 1. Upper bounds for the squared energy norm of the error.

Therefore, the approach described in the previous section we can used compute upper bounds for the energy norm of  $\mathbf{z}^\pm = \kappa \mathbf{u} \pm (1/\kappa) \boldsymbol{\psi} \in \mathcal{V}$ . These can then be used in expression (9) to yield computable expressions for the upper and lower bounds for  $s$ .

First, we compute approximations  $\mathbf{u}_h, \boldsymbol{\psi}_h \in \mathcal{V}$  by solving,

$$\begin{aligned} a(\mathbf{u}_h, \mathbf{v}) &= \ell(\mathbf{v}) \quad \forall \mathbf{v} \in \mathcal{V}^h, \\ a(\mathbf{v}, \boldsymbol{\psi}_h) &= \ell^0(\mathbf{v}) \quad \forall \mathbf{v} \in \mathcal{V}^h, \end{aligned}$$

respectively, and set  $\mathbf{z}_h^\pm = \kappa \mathbf{u}_h \pm (1/\kappa) \boldsymbol{\psi}_h$ . Here, we assume that  $\mathbf{u}^D$  and  $\mathbf{u}^0$  are piecewise polynomial over the elements of the working triangulation. For the particular case in which  $\mathbf{u}^D$  and  $\mathbf{u}^0$  are the finite element approximations to  $\mathbf{u}$  and  $\boldsymbol{\psi}$ , respectively, we will have  $\mathbf{u}_h = \boldsymbol{\psi}_h = 0$ .

Second, using the strategy described in [8], compute Lagrange multipliers by equilibrating the primal and dual problems, namely, find  $\lambda_h^u$  and  $\lambda_h^\psi$ , such that

$$\begin{aligned} b(\hat{\mathbf{v}}, \lambda_h^u) &= a(\mathbf{u}_h, \hat{\mathbf{v}}) - \ell(\hat{\mathbf{v}}) \quad \forall \mathbf{v} \in \widehat{\mathcal{V}}^h, \\ b(\hat{\mathbf{v}}, \lambda_h^\psi) &= a(\hat{\mathbf{v}}, \boldsymbol{\psi}_h) - \ell^0(\hat{\mathbf{v}}) \quad \forall \mathbf{v} \in \widehat{\mathcal{V}}^h. \end{aligned}$$

Extend  $\lambda_h^u$  and  $\lambda_h^\psi$  at the Neumann boundaries according to  $\lambda_h^u|_{\Gamma^N} = \mathbf{g}$  and  $\lambda_h^\psi|_{\Gamma^N} = \mathbf{g}^0$ , respectively. Finally, set  $\lambda_h^\pm = \kappa \lambda_h^u \pm (1/\kappa) \lambda_h^\psi$ .

Third, for each element in the mesh, we determine an equilibrated stress field  $\sigma_k^\pm$  verifying the equivalent of Eq. (23). That is, we compute  $\sigma_k^u$  and  $\sigma_k^\psi$  such that

$$\begin{aligned} \int_{\Omega_k} \sigma_k^u : \varepsilon(\mathbf{v}) d\Omega &= \int_{\Omega_k} \mathbf{f} \cdot \mathbf{v} d\Omega + \int_{\partial\Omega_k} \tau_k \lambda_h^u \cdot \mathbf{v} d\Gamma \quad \forall \mathbf{v} \in \mathcal{V}_k, \\ \int_{\Omega_k} \sigma_k^\psi : \varepsilon(\mathbf{v}) d\Omega &= \int_{\Omega_k} \mathbf{f}^0 \cdot \mathbf{v} d\Omega + \int_{\partial\Omega_k} \tau_k \lambda_h^\psi \cdot \mathbf{v} d\Gamma \quad \forall \mathbf{v} \in \mathcal{V}_k \end{aligned}$$

and set  $\sigma_k^\pm = \kappa(\sigma_k^u - \sigma_k^D) \pm (1/\kappa)(\sigma_k^\psi - \sigma_k^0)$ , where  $\sigma_k^D = \sigma(\mathbf{u}^D)|_{\Omega_k}$  and  $\sigma_k^0 = \sigma(\mathbf{u}^0)|_{\Omega_k}$ .

Finally, the upper bound is found as

$$\|\mathbf{z}^\pm\|^2 \leq \sum_{k=1}^{n_{el}} J_k^c(\sigma_k^\pm) = \kappa^2 v^u + \frac{1}{\kappa^2} v^\psi \pm 2v^{u\psi} \equiv \|\mathbf{z}^\pm\|_{UB}^2, \quad (25)$$



where

$$v^u = \sum_{k=1}^{n_{e1}} J_k^c(\sigma_k^u - \sigma_k^D),$$

$$v^\psi = \sum_{k=1}^{n_{e1}} J_k^c(\sigma_k^\psi - \sigma_k^\theta),$$

$$v^{u\psi} = \sum_{k=1}^{n_{e1}} \int_{\Omega_k} (\sigma_k^u - \sigma(u^D)) : \mathbb{C}^{-1}(\sigma_k^\psi - \sigma(u^\theta)) d\Omega.$$

Introducing expression (25) into (9) leads, after some algebra, to the following expressions for the upper and lower bounds of  $s$ :

$$s^+ = \frac{1}{2}s_h + \frac{1}{2}v^{u\psi} + \frac{\kappa^2}{4}(v^u - \|\mathbf{u}_h\|^2) + \frac{1}{4\kappa^2}(v^\psi - \|\boldsymbol{\psi}_h\|^2),$$

$$s^- = \frac{1}{2}s_h + \frac{1}{2}v^{u\psi} - \frac{\kappa^2}{4}(v^u - \|\mathbf{u}_h\|^2) - \frac{1}{4\kappa^2}(v^\psi - \|\boldsymbol{\psi}_h\|^2),$$

where  $s_h = \ell^\theta(\mathbf{u}_h)$ .

Following [10,11], the bounds are optimized with respect to the arbitrary parameter  $\kappa$ . The optimal value is given by  $\bar{\kappa}^2 = \left(\sqrt{v^\psi - \|\boldsymbol{\psi}_h\|^2}\right) / \left(\sqrt{v^u - \|\mathbf{u}_h\|^2}\right)$ . The resulting procedure to determine the bounds for  $s$  is summarized in the box of Fig. 2.

### 6. Adaptive mesh refinement

Once upper and lower bounds for the output quantity  $s$  are computed, one can compute the bound average

$$\underline{s} = \frac{1}{2}(s^+ + s^-) = s_h + \frac{1}{2}v^{u\psi}$$

and the bound gap

$$\Delta = s^+ - s^- = \sqrt{v^u - \|\mathbf{u}_h\|^2} \sqrt{v^\psi - \|\boldsymbol{\psi}_h\|^2}.$$

The bound average  $\underline{s}$  is a new estimate of the output  $s$ , where its error with respect to  $s$  can be easily bounded since

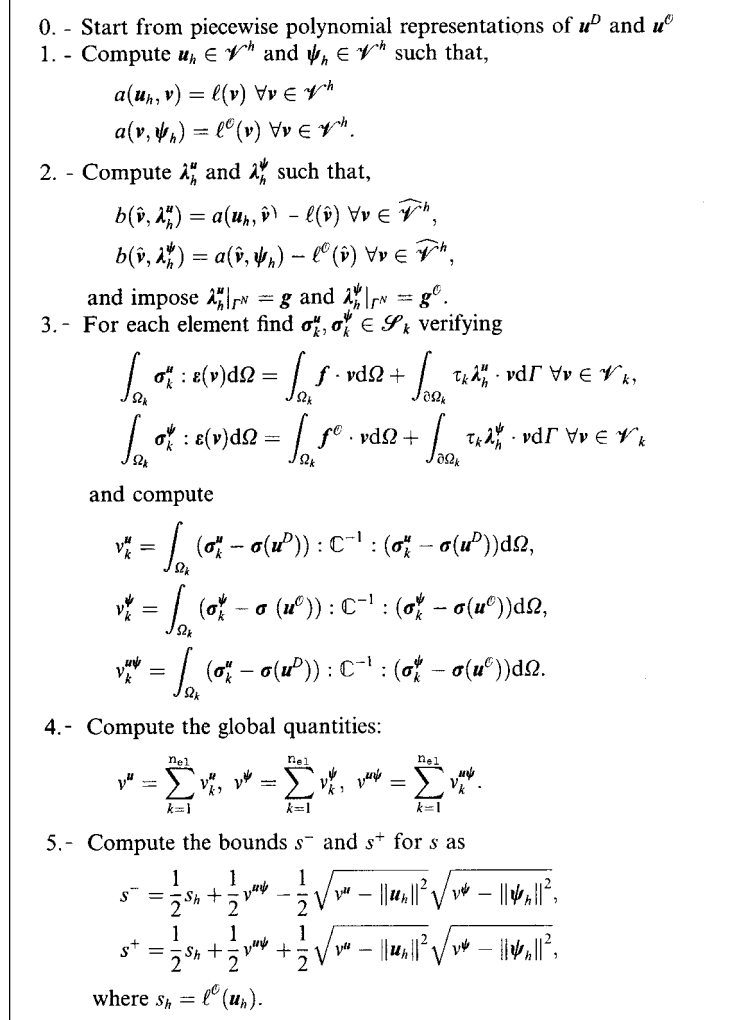
$$|s - \underline{s}| \leq \frac{1}{2}\Delta.$$

If this error meets the desired requirements of accuracy the computation is concluded. On the other hand, if the level of precision does not meet the requirements a mesh adaptive procedure can be easily devised [12].

The bound gap can be written as

$$\Delta = \sum_{k=1}^{n_{e1}} \frac{\bar{\kappa}^2}{2} \left(v_k^u - \|\mathbf{u}_k\|^2\right) + \frac{1}{2\bar{\kappa}^2} \left(v_k^\psi - \|\boldsymbol{\psi}_k\|^2\right) = \sum_{k=1}^{n_{e1}} \Delta_k,$$

where  $\bar{\kappa}$  is the optimal value of the parameter  $\kappa$  which optimizes the bounds. The above expression identifies the elemental contributions  $\Delta_k$ . These contributions can be shown to be always positive (since the comple-

Fig. 2. Bounds for the output of interest  $s$ .

mentary energy of an equilibrated stress field and the energy of an arbitrary displacement field, are upper and lower bounds to the energy of the exact solution, respectively) and can therefore be used as a local refinement indicator.

Then, given a target bound gap  $\Delta_{\text{tol}}$ , at each level of refinement, the elements with  $\Delta_k \geq (\Delta_{\text{tol}}/n_{e1})$  are refined. Numerical experimentation indicates that this strategy leads to a robust and reliable procedure to achieve the desired accuracy. The refined meshes are obtained using the mesh generator presented in [16].

## 7. Numerical examples

The presented method is illustrated with two numerical examples: a linearly forced square which has a regular solution for which an analytical expression exists, and a square plate with two interior rectangular cut-outs, the solution of which, has corner singularities. The outputs of interest are in both cases displace-

ments and reaction forces integrated over parts of the boundary. Linear finite elements approximations have been used for the adjoint and the hybrid fluxes have also been interpolated linearly over each edge. The equilibrated stress fields in the dual problem are also taken to be linearly varying in space.

The coarse mesh problems are solved using triangular linear finite elements, and the local equilibrated stress fields are taken to be piecewise linear in each triangle of the mesh. Four estimates of  $s$  are considered: the upper and lower bounds ( $s^+$  and  $s^-$ , respectively), their average,  $\underline{s} = (s^+ + s^-)/2$ , and also the output given by the finite element approximation, denoted by  $s_h = \ell^c(\mathbf{u}_h)$ .

In the first example the analytical solution of the problem is known and the quality of the different estimates is measured with the following effectivity indices  $\rho^\pm = (s^\pm/s) - 1$ ,  $\underline{\rho} = (\underline{s}/s) - 1$ , and  $\rho_h = (s_h/s) - 1$ . Another measure of the accuracy of the bounds is given by the relative half bound gap

$$\rho_G = \frac{1}{2} \frac{s^+ - s^-}{|s|} \geq 0.$$

Since  $s^+$  and  $s^-$  are upper and lower bounds of  $s$ , the index  $\rho_G$  is an upper bound of the relative error between the approximation  $\underline{s}$  and the exact output  $s$ , that is

$$\frac{|\underline{s} - s|}{|s|} \leq \rho_G.$$

In the second example, where the analytical solution is not known, the bound accuracy is measured in terms of the relative half bound gap,  $\rho_G$ , which is re-defined as

$$\rho_G = \frac{1}{2} \frac{s^+ - s^-}{|\underline{s}|},$$

where the exact output is replaced by the average estimate.

### 7.1. Linearly forced square

The plane stress elasticity equations are considered in the unity square  $[0, 1]^2$ . On the left edge of the square,  $x_1 = 0$ , Dirichlet homogeneous boundary conditions are imposed in the  $x_2$  direction, and in the left-lower corner,  $(0, 0)$ , both the  $x_1$  and  $x_2$  displacements are prescribed to zero. Also, a linear normal traction,  $\mathbf{g} = (x_2, 0)^T$ , is applied at the right edge,  $x_1 = 1$ .

The analytical solution of the problem  $\mathbf{u} = (u_1, u_2)$ , is given by

$$u_1(x_1, x_2) = \frac{1}{E} x_1 x_2, \quad u_2(x_1, x_2) = -\frac{1}{2E} (v x_2^2 + x_1^2),$$

where  $E$  and  $\nu$  are the Young's modulus and the Poisson's ratio.

The output considered is the weighted average normal displacement at the right edge,

$$s = \int_0^1 x_2 u_1(1, x_2) dx_2 = \frac{1}{3E}.$$

It turns out that for this particular forcing and output, the primal and adjoint problems are the same. For this case, called compliance, the output is proportional to the energy norm of the solution and the finite element approximation directly provides a lower bound. The numerical results demonstrate that our method, while more expensive, leads to the same lower bound, doing no worse than the inherent bound of the finite element approximation.

Four uniform triangular meshes have been considered, the initial one with 18 elements ( $h = 1/3$ ). The other meshes are obtained by uniformly subdividing each element of the previous mesh into four new elements. The results are summarized in the Table 1.

Table 1  
Bounds and effectivity indices in a series of uniformly refined meshes

| $h$  | $s_h$ | $s^-$ | $s^+$ | $\underline{s}$ | $\rho^-$ | $\rho^+$ | $\underline{\rho}$ | $\rho_G$ |
|------|-------|-------|-------|-----------------|----------|----------|--------------------|----------|
| 1/3  | .3124 | .3124 | .5621 | .4372           | -.062740 | .686210  | .311720            | .3745    |
| 1/6  | .3264 | .3264 | .4370 | .3817           | -.020710 | .310880  | .145070            | .1658    |
| 1/12 | .3314 | .3314 | .3653 | .3484           | -.005710 | .095990  | .045140            | .0508    |
| 1/24 | .3328 | .3328 | .3419 | .3374           | -.001480 | .025640  | .012080            | .0136    |
| 1/48 | .3332 | .3332 | .3355 | .3344           | -.000370 | .006530  | .003080            | .0034    |

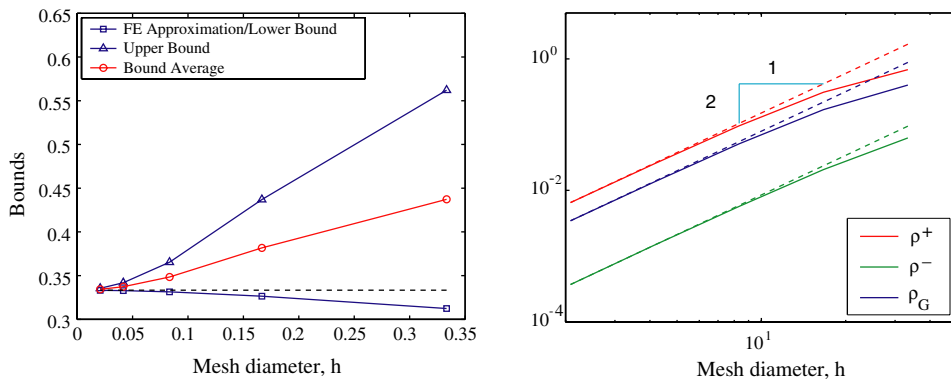


Fig. 3. Computed bounds for a uniform  $h$ -refinement process (left) and its convergence (right).

Fig. 3 displays the results graphically and also illustrates the convergence rate of the bounds. The results for both the upper and lower bounds, average, and relative half bound gap,  $\rho_G$ , asymptotically approach the finite element convergence rate of  $\mathcal{O}(h^2)$ .

## 7.2. Square plate

A square thin plate with two rectangular holes is considered. Normal tractions are applied on the left and right sides of the plate [12]. Since the problem is symmetric, only one fourth of the plate is considered, as shown in Fig. 4.

Two outputs of interest are considered: the average normal displacement over the boundary  $\Gamma_0$ , and the integrated normal component of the traction in  $\Gamma_1$ , that is,

$$\ell_0^\emptyset(\mathbf{v}) = \int_{\Gamma_0} \mathbf{v} \cdot \mathbf{n} d\Gamma, \quad \ell_1^\emptyset(\mathbf{v}) = \int_{\Gamma_1} \mathbf{n} \cdot \boldsymbol{\sigma}(\mathbf{v}) \mathbf{n} d\Gamma. \quad (26)$$

**Remark 4.** The first output is already in the form of Eq. (1) with  $\mathbf{g}^\emptyset = \mathbf{n}|_{\Gamma_0}$  and  $\mathbf{g}^\emptyset = \mathbf{0}$  elsewhere. The second output, on the other hand, does not have the same form. In order to transform this output into the form (1) considered here, we introduce a continuous function  $\chi$  such that  $\chi|_{\Gamma_1} = 1$  and is equal to zero at all the other vertical boundaries. Then, if  $\mathbf{n}^1 = \mathbf{n}|_{\Gamma_1}$ , we have

$$s = \ell_1^\emptyset(\mathbf{u}) = \int_{\Gamma_1} \mathbf{n} \cdot \boldsymbol{\sigma}(\mathbf{u}) \mathbf{n} d\Gamma = a(\mathbf{u}, \chi \mathbf{n}^1) =: \tilde{\ell}_1^\emptyset(\mathbf{u})$$

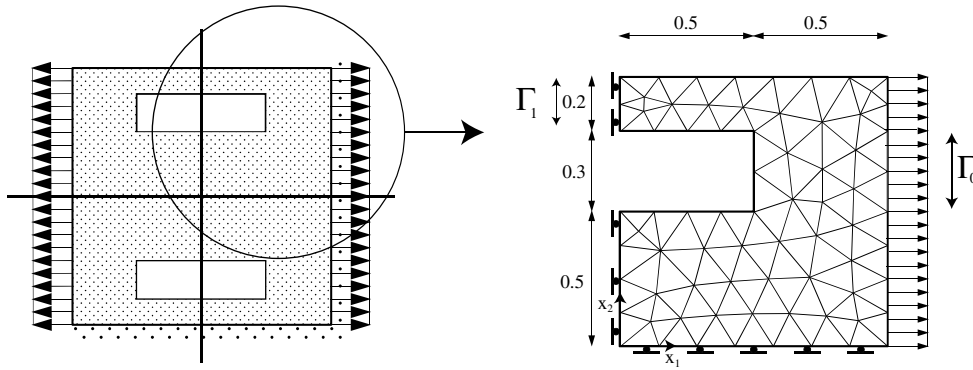


Fig. 4. Example 2: model problem (left) and initial mesh (right).

and instead of working with the functional  $\ell_1^0(\cdot)$ , we work with  $\tilde{\ell}_1^0(\cdot)$ . This is much easier since this corresponds to  $\mathbf{u}^0 = -\chi \mathbf{n}^1$  in Eq. (1).

Fig. 5 and Table 2 show the bounds obtained in this example. A nested sequence of meshes is considered. The initial mesh ( $h_{ini}$ ) is shown in Fig. 4, and the refined meshes are obtained, as in the first example, dividing each element into 4 new ones. The function  $\chi$  required in  $\tilde{\ell}_1^0(\cdot)$ , is defined on the initial mesh by setting all the nodal values equal to zero except for those nodes on  $\Gamma_1$  which are given a value of unity.

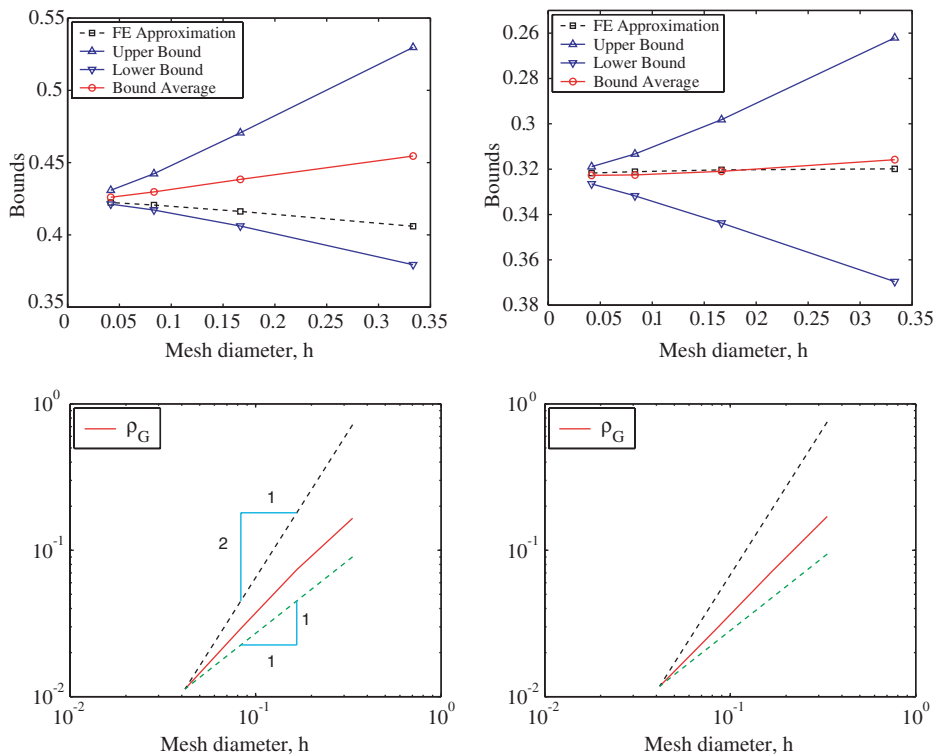


Fig. 5. Bounds convergence for a uniform  $h$ -refinement (up) and for the displacement output  $\ell_0^c(\mathbf{u})$  (left) and for the reaction output  $\ell_0^c(\mathbf{u})$  (right).

Table 2

Example 2: Bounds and relative bound gap in a series of uniformly refined  $h$ -meshes both for  $\ell_0^c(\mathbf{u})$  and  $\ell_1^c(\mathbf{u})$ 

| $h$                 | Displacement average |       |       |                 |          | Reaction average |        |        |                 |          |
|---------------------|----------------------|-------|-------|-----------------|----------|------------------|--------|--------|-----------------|----------|
|                     | $s_h$                | $s^-$ | $s^+$ | $\underline{s}$ | $\rho_G$ | $s_h$            | $s^-$  | $s^+$  | $\underline{s}$ | $\rho_G$ |
| $h_{\text{ini}}$    | .4060                | .3794 | .5297 | .4546           | .1654    | -.3199           | -.3696 | -.2621 | -.3158          | .1702    |
| $1/2h_{\text{ini}}$ | .4163                | .4061 | .4706 | .4384           | .0736    | -.3203           | -.3438 | -.2982 | -.3210          | .0710    |
| $1/4h_{\text{ini}}$ | .4207                | .4172 | .4423 | .4298           | .0292    | -.3211           | -.3318 | -.3133 | -.3225          | .0286    |
| $1/8h_{\text{ini}}$ | .4224                | .4213 | .4309 | .4261           | .0113    | -.3217           | -.3265 | -.3189 | -.3227          | .0118    |

Unlike the first example, the outputs in (26) are general and the finite element approximation can no longer be guaranteed to provide a lower bound. This example shows that the bounds behave well even for problems with singularities. However, it is also observed that the convergence rate for the bounds, the finite element approximation  $s_h$  and the bound average, is no longer  $\mathcal{O}(h^2)$ , although it is still faster than linear.

For the reaction output,  $\ell_1^c(\mathbf{u})$ , an adaptive procedure has been employed starting with the mesh shown in Fig. 4) where the bound gap  $\Delta_{\text{ini}}$  is 0.1075, and two target bound gaps have been considered  $\Delta_{\text{tol}} = \frac{1}{2}\Delta_{\text{ini}}$  and  $\Delta_{\text{tol}} = \frac{1}{10}\Delta_{\text{ini}}$ .

In order to achieve  $\Delta_{\text{tol}} = \frac{1}{2}\Delta_{\text{ini}}$  four new meshes are generated, where the bound gap for the last mesh is  $\Delta_{\text{r}} = 0.0471$ . The resulting sequence of meshes can be seen in Fig. 6, where the local elementary contributions to the global bound gap are plotted in each element of the mesh. As can be seen not only the zone where the output is measured ( $\Gamma_1$ ) is refined, but also the corners where the solution is singular.

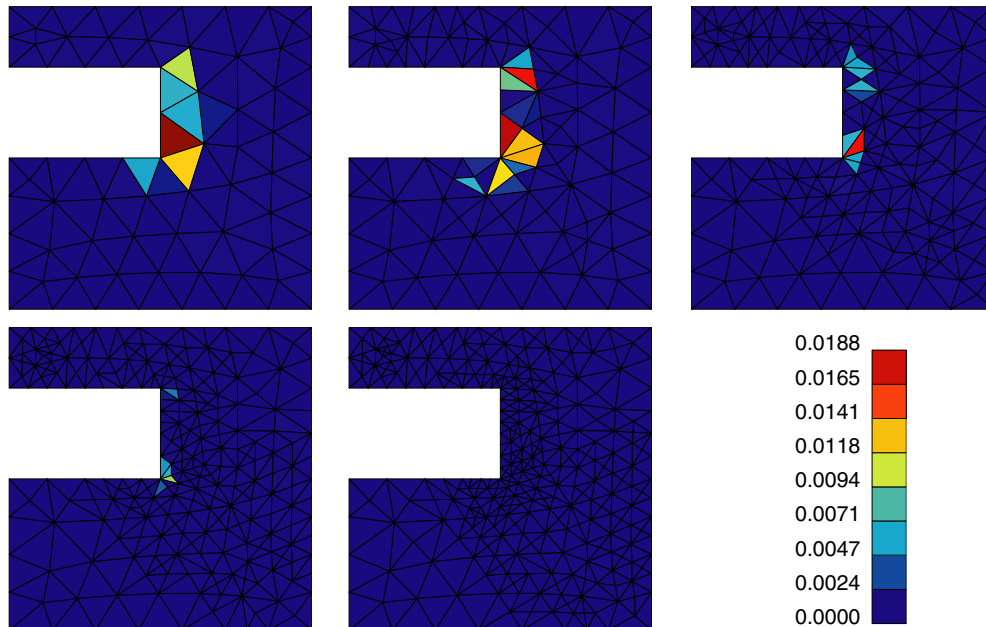


Fig. 6. Example 2: sequence of adapted meshes for the output  $\ell_1^c(\mathbf{u})$  with desired final gap  $\Delta_{\text{tol}} = \frac{1}{2}\Delta_{\text{ini}}$  with  $n_{e1} = 108, 165, 280, 405$  and 538.

Table 3

Example 2: bounds in a series of adaptively  $h$ -refined meshes both for  $\ell_1^e(\mathbf{u})$  with desired final gap  $A_{\text{tol}} = \frac{1}{10} A_{\text{ini}}$ 

| $n_{e1}$ | $\Delta$ | $s_1$   | $s_u$   |
|----------|----------|---------|---------|
| 108      | .10749   | –.36957 | –.26208 |
| 222      | .18215   | –.38940 | –.20725 |
| 433      | .12171   | –.36880 | –.24709 |
| 811      | .07199   | –.35089 | –.27891 |
| 1387     | .03755   | –.33750 | –.29995 |
| 1966     | .02428   | –.33392 | –.30964 |
| 2532     | .01574   | –.32922 | –.31348 |
| 3069     | .01172   | –.32826 | –.31654 |
| 3564     | .00834   | –.32627 | –.31793 |

The values of the bounds for the adaptive procedure with the desired final gap  $A_{\text{tol}} = \frac{1}{10} A_{\text{ini}}$  are shown in Table 3.

## 8. Conclusions

We have presented a method for the computation of bounds for linear-functional outputs of weak solutions to the linear elasticity equations. A distinctive feature of this method is that the computed bounds are strict with respect to the output of the exact solution. The numerical experiments presented show that the computed bounds are sharp and converge at the same rate as the finite element solution that would be obtained with a comparable amount of work. To our knowledge, this is the only published approach that can certify the certainty of the the computed bounds. We believe this feature is of clear interest in real engineering practice. The method has been presented for the two dimensional elasticity equations, but we expect that the extension to three dimensions will not present any additional difficulties. The major computational cost, in addition to a standard finite element solution, is the computation of an adjoint for each output considered. All other operations are local and result in a low computational overhead. Two limitations in the presented approach are the need for the forcing function to be of piecewise polynomial form, and the requirement for the computational domain to have piecewise straight boundaries. Future work will focus on relaxing these constraints.

## Acknowledgements

Nuria Pares would like to acknowledge the Departament d'Universitats, Recerca i Societat de la Informació of the Generalitat de Catalunya for the support provided during her visit to MIT. Jaume Peraire would like to acknowledge the generous support provided by the Singapore-MIT Alliance.

## Appendix A. Proof of Theorem 1

In this appendix we present a constructive proof of Theorem 1 which shows the existence of piecewise polynomial equilibrated stress fields. Towards this end some preliminary notation and results are required (see [4] for details).

**Lemma 1.** Given a triangle  $T$ , consider the following interpolation spaces:

$$P_q(T) = \{\text{polynomial functions of degree less or equal to } q \text{ in } T\},$$

$$\mathbf{SP}_q(T) = \{\text{stress fields with } \sigma_{xx}, \sigma_{xy}, \sigma_{yy} \in P_q(T)\},$$

$$R_q(\partial T) = \{\text{polynomial functions of degree less or equal to } q \text{ on each } \gamma_i \in \partial T\},$$

$$\mathbf{R}_q^c(\partial T) = \{\mathbf{g} \in [R_q(\partial T)]^2, \text{ s.t. } \exists \boldsymbol{\sigma} \in \mathbf{SP}_q(T), \boldsymbol{\sigma} \cdot \mathbf{n} = \mathbf{g} \text{ on } \partial T\},$$

$$\boldsymbol{\Phi}_q(T) = \{\boldsymbol{\zeta}_q \in \mathbf{SP}_q(T), \nabla \cdot \boldsymbol{\zeta}_q = 0, \boldsymbol{\zeta}_q \cdot \mathbf{n}|_{\partial T} = 0\},$$

$$\mathbf{P}_{\text{sm}} = \{\text{rigid solid body motions in } T\}$$

and  $\mathbf{P}_{\text{sm}}^\perp$  be the orthogonal complement of  $\mathbf{P}_{\text{sm}}$  with respect to the standard scalar product in  $[P_{q-1}(T)]^2$ , that is, every  $\mathbf{p} \in [P_{q-1}(T)]^2$  can be written uniquely as  $\mathbf{p} = \mathbf{p}_{\text{sm}} + \mathbf{p}^\perp$ , with  $\mathbf{p}_{\text{sm}} \in \mathbf{P}_{\text{sm}}$  and  $\mathbf{p}^\perp \in \mathbf{P}_{\text{sm}}^\perp$ . We note that for the case  $q = 1$ , the only member of  $\mathbf{P}_{\text{sm}}^\perp$  is the null function.

Then, for  $q \geq 1$ , and for any  $\boldsymbol{\sigma} \in \mathbf{SP}_q(T)$  the following relations imply  $\boldsymbol{\sigma} = 0$ :

$$\int_{\partial T} (\boldsymbol{\sigma} \cdot \mathbf{n}) \cdot \mathbf{p}_q \, d\Gamma = 0 \quad \forall \mathbf{p}_q \in \mathbf{R}_q^c(\partial T), \quad (\text{A.1})$$

$$\int_T \boldsymbol{\sigma} : \boldsymbol{\varepsilon}(\mathbf{p}_{q-1}) \, d\Omega = 0 \quad \forall \mathbf{p}_{q-1} \in \mathbf{P}_{\text{sm}}^\perp, \quad (\text{A.2})$$

$$\int_T \boldsymbol{\sigma} : \mathbb{C}^{-1} : \boldsymbol{\zeta}_q \, d\Omega = 0 \quad \forall \boldsymbol{\zeta}_q \in \boldsymbol{\Phi}_q(T). \quad (\text{A.3})$$

**Proof.** First let us check that Eqs. (A.1) and (A.2) imply that  $\boldsymbol{\sigma} \in \boldsymbol{\Phi}_q(T)$ . Indeed, on one hand, since  $\boldsymbol{\sigma} \cdot \mathbf{n}|_{\partial T} \in \mathbf{R}_q^c(\partial T)$ , from Eq. (A.1),  $\int_{\partial T} (\boldsymbol{\sigma} \cdot \mathbf{n})^2 \, d\Gamma = 0$ , which implies that  $\boldsymbol{\sigma} \cdot \mathbf{n} = 0$  in  $\partial T$ . On the other hand, the following integration by parts:

$$\int_T (\nabla \cdot \boldsymbol{\sigma}) \cdot (\nabla \cdot \boldsymbol{\sigma}) \, d\Omega = \int_{\partial T} (\boldsymbol{\sigma} \cdot \mathbf{n}) \cdot (\nabla \cdot \boldsymbol{\sigma}) \, d\Gamma - \int_T \boldsymbol{\sigma} : \boldsymbol{\varepsilon}(\nabla \cdot \boldsymbol{\sigma}) \, d\Omega,$$

plus the fact that  $\boldsymbol{\varepsilon}(\mathbf{p}_{\text{sm}}) = 0$  for  $\mathbf{p} \in \mathbf{P}_{\text{sm}}$ , leads to

$$\int_T (\nabla \cdot \boldsymbol{\sigma}) \cdot (\nabla \cdot \boldsymbol{\sigma}) \, d\Omega = \int_T \boldsymbol{\sigma} : \boldsymbol{\varepsilon}(\nabla \cdot \boldsymbol{\sigma} - \pi_{\text{sm}}(\nabla \cdot \boldsymbol{\sigma})) \, d\Omega,$$

where  $\pi_{\text{sm}}(\cdot)$  is the projection operator from  $[P_{q-1}(T)]^2$  onto the space  $\mathbf{P}_{\text{sm}}$ . Then, since  $\nabla \cdot \boldsymbol{\sigma} - \pi_{\text{sm}}(\nabla \cdot \boldsymbol{\sigma}) \in \mathbf{P}_{\text{sm}}^\perp$ , Eq. (A.2) implies

$$\int_T (\nabla \cdot \boldsymbol{\sigma}) \cdot (\nabla \cdot \boldsymbol{\sigma}) \, d\Omega = 0 \Rightarrow \nabla \cdot \boldsymbol{\sigma} = 0 \text{ in } \Omega,$$

which shows that (A.1) and (A.2) imply that  $\boldsymbol{\sigma} \in \boldsymbol{\Phi}_q(T)$ . Finally, using Eq. (A.3),

$$\int_T \boldsymbol{\sigma} : \mathbb{C}^{-1} : \boldsymbol{\sigma} \, d\Omega = 0 \Rightarrow \boldsymbol{\sigma} = 0,$$

which ends the proof.  $\square$



**Lemma 2.** Let  $\{\mathbf{p}_q^i\}_{i=1\dots I}$ ,  $\{\mathbf{p}_{q-1}^j\}_{j=1\dots J}$  and  $\{\zeta_q^l\}_{l=1\dots L}$  denote the elements of a basis of  $\mathbf{R}_q^c(\partial T)$ ,  $\mathbf{P}_{\text{sm}}^\perp$  and  $\Phi_q(T)$  respectively, where  $I, J$  and  $L$  simply denote the dimensions of each space. Then, any  $\boldsymbol{\sigma} \in \mathbf{SP}_q(T)$  is uniquely determined by the following degrees of freedom:

$$\begin{aligned} \int_{\partial T} (\boldsymbol{\sigma} \cdot \mathbf{n}) \cdot \mathbf{p}_q^i \, d\Gamma, \quad i \in I, \\ \int_T \boldsymbol{\sigma} : \boldsymbol{\varepsilon}(\mathbf{p}_{q-1}^j) \, d\Omega, \quad j \in J, \\ \int_T \boldsymbol{\sigma} : \mathbb{C}^{-1} : \zeta_q^l \, d\Omega, \quad l \in L. \end{aligned}$$

**Proof.** Lemma 1 states that any stress field  $\boldsymbol{\sigma} \in \mathbf{SP}_q(T)$  can be described giving the values of the previous degrees of freedom. However, this description it is not necessary unique, that is, different values of the previous degrees of freedom can yield the same stress field. In order to see that this description is unique, it is sufficient to see that the number of degrees of freedom coincides with the dimension of  $\mathbf{SP}_q(T)$ , where  $\dim(\mathbf{SP}_q(T)) = \frac{3}{2}(q+1)(q+2)$ .

Let us consider first the case  $q > 1$ . It is clear that  $\dim(\mathbf{R}_q^c(\partial T)) = 6(q+1) - 3$ . Now, a basis of  $[P_{q-1}(T)]^2$ , determined by  $q(q+1)$  elements, defines only  $q(q+1) - 3$  degrees of freedom of the form  $\int_T \boldsymbol{\sigma} : \boldsymbol{\varepsilon}(\mathbf{p}_{q-1}) \, d\Omega$ , since  $\boldsymbol{\varepsilon}(\mathbf{t}_x) = \boldsymbol{\varepsilon}(\mathbf{t}_y) = \boldsymbol{\varepsilon}(\mathbf{r}) = 0$ , for  $\mathbf{t}_x, \mathbf{t}_y$  and  $\mathbf{r}$  the three rigid solid body motions, that is,  $\dim(\mathbf{P}_{\text{sm}}^\perp) = q(q+1) - 3$ .

Then, the only remaining part is to determine the dimension of  $\Phi_q(T)$ . Any  $\zeta \in \Phi_q(T)$  can be rewritten as  $\zeta_{xx} = \partial^2 b / \partial^2 y$ ,  $\zeta_{yy} = \partial^2 b / \partial x^2$ ,  $\zeta_{xy} = \zeta_{yx} = -\partial^2 b / \partial x \partial y$ , where  $b \in b_T^2 P_{q-4}(T)$  for  $b_T$  the cubic bubble function on  $T$  vanishing on  $\partial T$  and achieving a maximum value of unity on  $T$ , see [2]. Therefore,  $\dim(\Phi_q(T)) = \dim(P_{q-4}(T)) = \frac{1}{2}(q-2)(q-3)$ .

Finally, it is trivial to check that

$$\underbrace{\frac{3}{2}(q+1)(q+2)}_{\dim(\mathbf{SP}_q(T))} = \underbrace{6(q+1) - 3}_{\dim(\mathbf{R}_q^c(\partial T))} + \underbrace{q(q+1) - 3}_{\dim(\mathbf{P}_{\text{sm}}^\perp)} + \underbrace{\frac{1}{2}(q-2)(q-3)}_{\dim(\Phi_q(T))}.$$

For the particular case  $q = 1$ ,  $[P_0(T)]^2 = \langle \mathbf{t}_x, \mathbf{t}_y \rangle$ , and equations  $\int_T \boldsymbol{\sigma} : \boldsymbol{\varepsilon}(\mathbf{p}_0) \, d\Omega = 0$  do not characterize any degree of freedom. Moreover, in this case  $\dim(\Phi_1(T)) = 0$ . Then the nine boundary degrees of freedom  $\mathbf{R}_1^c(\partial T)$  determine uniquely the linear stress field in the triangle.  $\square$

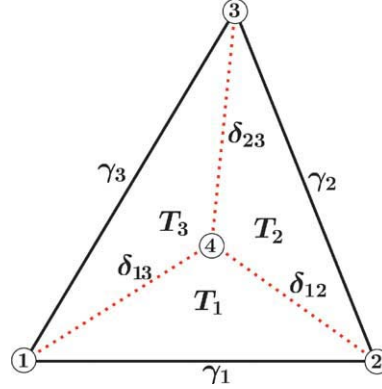
We can now proceed directly with the proof of Theorem 1. We will assume that the initial element  $\Omega_k$  is a triangle, but the strategy can be easily extended to quadrilateral elements. In order to find  $\boldsymbol{\sigma}_k$ , we follow [9] and divide the initial triangle into three new triangles,  $\Omega_k = T_1 \cup T_2 \cup T_3$  by adding a point in the triangle centroid, as indicated in Fig. 1.

Let,  $\boldsymbol{\sigma}_\lambda$  be a stress field in  $\Omega_k$ , where  $\boldsymbol{\sigma}_\lambda|_{T_i} \in \mathbf{SP}_q(T_i)$ ,  $i = 1, 2, 3$ , that is,  $\boldsymbol{\sigma}_\lambda$  is a polynomial stress field of degree  $q$  in each subtriangle, such that,  $\boldsymbol{\sigma}_\lambda|_{\partial T} = \tau_k \boldsymbol{\lambda}_h$ , and with continuous normal tractions at the internal edges of the partition ( $\delta_{ij} = T_i \cap T_j$ ,  $i, j = 1 \dots 3$ ).

Then, setting  $\boldsymbol{\sigma}_k = \boldsymbol{\sigma}_\lambda + \boldsymbol{\sigma}_0$ , the initial problem reduces to finding a piecewise polynomial stress field  $\boldsymbol{\sigma}_0$  verifying

$$\int_{\Omega_k} \boldsymbol{\sigma}_0 : \boldsymbol{\varepsilon}(\mathbf{v}) \, d\Omega = \int_{\Omega_k} (\mathbf{f}^* + \nabla \cdot \boldsymbol{\sigma}_\lambda) \cdot \mathbf{v} \, d\Omega \quad \forall \mathbf{v} \in \mathcal{V}_k. \tag{A.4}$$

Let  $\boldsymbol{\sigma}_0$  be a piecewise polynomial stress field, where in each triangle  $T_i$ , the stress field  $\boldsymbol{\sigma}_0^i = \boldsymbol{\sigma}_0|_{T_i}$  is assumed to be polynomial of degree  $q$  in each component, that is,  $\boldsymbol{\sigma}_0^i \in \mathbf{SP}_q(T_i)$ .

Fig. 1. Local subdivision of an element  $\Omega_k$  into the triangles  $T_1$ ,  $T_2$  and  $T_3$ .

Then,  $\sigma_0$  is uniquely determined by the degrees of freedom characterized in Lemma 2. First, there are the  $18(q+1) - 9$  degrees of freedom determining the value of  $\sigma_0$  at the edges of  $\Omega_k$  and at the internal edges, namely

$$\int_{\partial T_i} (\sigma_0^i \cdot \mathbf{n}) \cdot \mathbf{p} \, d\Gamma \quad \forall \mathbf{p} \in \mathbf{R}_q^c(\partial T_i), \quad i = 1 \dots 3.$$

Second, we have the degrees of freedom related to the divergence of  $\sigma_0$ , that is

$$\int_{T_i} \sigma_0^i : \varepsilon(\mathbf{p}) \, d\Omega \quad \forall \mathbf{p} \in \mathbf{P}_{\text{sm}}(T_i)^\perp, \quad i = 1 \dots 3.$$

And finally, the degrees of freedom associated to  $\Phi_q(T_i)$ ,  $i = 1, 2, 3$ , which can be set arbitrarily.

The proof ends with the construction of  $\sigma_0$  that verifies (A.4). This construction also follows the classification of Lemma 2: (i) the boundary degrees of freedom, (ii) the divergence ones and (iii) the related to  $\Phi_q(T_i)$ , which are not detailed because they are arbitrary.

(i) The  $18(q+1) - 9$  boundary degrees of freedom are determined in two steps. First, the  $12(q+1)$  constraints to enforce compatibility are imposed

$$\int_{\gamma_i} (\sigma_0^i \cdot \mathbf{n}) \cdot \mathbf{p} \, d\Gamma = 0 \quad \forall \mathbf{p} \in [R_q(\gamma_i)]^2, \quad i = 1 \dots 3, \quad (\text{A.5})$$

where  $\gamma_i = \partial T_i \cap \partial T$ , and

$$\int_{\delta_{ij}} (\sigma_0^i \cdot \mathbf{n}) \cdot \mathbf{p} \, d\Gamma = \int_{\delta_{ij}} (\sigma_0^j \cdot \mathbf{n}) \cdot \mathbf{p} \, d\Gamma \quad \forall \mathbf{p} \in [R_q(\delta_{ij})]^2, \quad i, j = 1 \dots 3, \quad i < j. \quad (\text{A.6})$$

Then, the  $6q - 3$  remainder degrees of freedom are used to impose the following additional constraints,

$$\int_{\partial T_i} (\sigma_0^i \cdot \mathbf{n}) \cdot \mathbf{v} \, d\Omega = - \int_{T_i} (\mathbf{f}^* + \nabla \cdot \sigma_i) \cdot \mathbf{v} \, d\Omega \quad \forall \mathbf{v} \in \mathbf{P}_{\text{sm}}(T_i), \quad i = 1 \dots 3. \quad (\text{A.7})$$

It is important to note that since  $\lambda_n$  and  $\mathbf{f}^*$  verify Eq. (22), some of the previous equations are redundant. For  $q > 1$ , Eq. (A.7) represents nine constraints but only 6 degrees of freedom are required to impose them because Eq. (22) is scalar and  $\dim(\mathbf{P}_{\text{sm}}) = 3$ . In the case  $q = 1$ , Eq. (A.7) represents six constraints but only three of them are independent for the same reason. Note that for  $q = 1$ , the boundary degrees of freedom are uniquely determined, while for  $q > 1$ , there are  $6q - 9$  degrees of freedom left associated to the internal boundaries which can be set arbitrarily.

(ii) Once the boundary degrees of freedom are fixed, we impose those related to the divergence of  $\sigma_0$ , namely

$$\int_{T_i} \sigma_0^i : \varepsilon(\mathbf{p}) d\Omega = \int_{T_i} (\mathbf{f}^* + \nabla \cdot \sigma_\lambda) \cdot \mathbf{p} d\Omega + \int_{\partial T_i} (\sigma_0^i \cdot \mathbf{n}) \cdot \mathbf{p} d\Gamma \tag{A.8}$$

for all  $\mathbf{p} \in \mathbf{P}_{\text{sm}}(T_i)^\perp$ ,  $i = 1 \dots 3$ .

Once (A.7) and in (A.8) have been imposed,  $\sigma_0$  verifies

$$\int_{T_i} \sigma_0^i : \varepsilon(\mathbf{p}) d\Omega = \int_{T_i} (\mathbf{f}^* + \nabla \cdot \sigma_\lambda) \cdot \mathbf{p} d\Omega + \int_{\partial T_i} (\sigma_0^i \cdot \mathbf{n}) \cdot \mathbf{p} d\Gamma \tag{A.9}$$

for all  $\mathbf{p} \in [P_{q-1}(T_i)]^2$ ,  $i = 1 \dots 3$ .

To conclude the proof it only remains to show that the stress field  $\sigma_0$  indeed verifies (A.4). On one hand, a simple integration by parts shows that Eq. (A.9) is equivalent to

$$\int_{T_i} (\nabla \cdot \sigma_0^i - (\mathbf{f}^* + \nabla \cdot \sigma_\lambda)) \cdot \mathbf{p} d\Omega = 0 \quad \forall \mathbf{p} \in [P_{q-1}(T_i)]^2, \quad i = 1 \dots 3. \tag{A.10}$$

Since  $\nabla \cdot \sigma_0^i, \mathbf{f}^*$  and  $\nabla \cdot \sigma_\lambda \in [P_{q-1}(T_i)]^2$ , we have  $\nabla \cdot \sigma_0^i = \mathbf{f}^* + \nabla \cdot \sigma_\lambda$ , and thus Eqs. (A.9) and (A.10) hold not only for  $\mathbf{p} \in [P_{q-1}(T_i)]^2$  but for  $\mathbf{p} \in [\mathcal{H}^1(T_i)]^2$ . On the other hand, using a similar reasoning, we have that  $\sigma_0^i \in \mathbf{SP}_q(T_i)$  and Eqs. (A.5), (A.6) give  $\sigma_0^i \cdot \mathbf{n}|_{\gamma_i} = 0$ ,  $i = 1 \dots 3$  and  $(\sigma_0^i - \sigma_0^j) \cdot \mathbf{n}|_{\delta_{ij}} = 0$ ,  $i, j = 1 \dots 3$ ,  $i < j$ . Thus,

$$\int_{\gamma_i} (\sigma_0^i \cdot \mathbf{n}) \cdot \mathbf{p} d\Gamma = 0 \quad \forall \mathbf{p} \in [R_q(\gamma_i)]^2, \quad i = 1 \dots 3, \tag{A.11}$$

where  $\gamma_i = \partial T_i \cap \partial T$ , and

$$\int_{\delta_{ij}} (\sigma_0^i \cdot \mathbf{n}) \cdot \mathbf{p} d\Gamma = \int_{\delta_{ij}} (\sigma_0^j \cdot \mathbf{n}) \cdot \mathbf{p} d\Gamma \quad \forall \mathbf{p} \in [R_q(\delta_{ij})]^2, \quad i, j = 1 \dots 3, \quad i < j \tag{A.12}$$

hold, not only for  $\mathbf{p} \in [R_q^c(\gamma_i)]^2$  and  $\mathbf{p} \in [R_q^c(\delta_{ij})]^2$ , but for  $\mathbf{p} \in [\mathcal{H}^{\frac{1}{2}}(\gamma_i)]^2$  and  $\mathbf{p} \in [\mathcal{H}^{\frac{1}{2}}(\delta_{ij})]^2$ , respectively. Finally, for any  $\mathbf{v} \in \mathcal{V}_k = [\mathcal{H}^1(T)]^2$  using the infinite dimensional versions of Eqs. (A.5), (A.6) and (A.9), and the fact that  $\mathbf{v}$  is continuous on  $\delta_{ij}$ ,

$$\begin{aligned} \int_T \sigma_0 : \varepsilon(\mathbf{v}) d\Omega &= \sum_{i=1}^3 \int_{T_i} \sigma_0^i : \varepsilon(\mathbf{v}) d\Omega = \int_T (\mathbf{f}^* + \nabla \cdot \sigma_\lambda) \cdot \mathbf{v} d\Omega + \underbrace{\sum_{i=1}^3 \int_{\partial T_i} (\sigma_0^i \cdot \mathbf{n}) \cdot \mathbf{v} d\Gamma}_0 \\ &= \int_T (\mathbf{f}^* + \nabla \cdot \sigma_\lambda) \cdot \mathbf{v} d\Omega. \quad \square \end{aligned}$$

## References

- [1] M. Ainsworth, J.T. Oden, A posteriori error estimation in finite element analysis, *Comput. Methods Appl. Mech. Engrg.* 142 (1997) 1–88.
- [2] D.N. Arnold, R. Winther, Mixed finite element for elasticity, *Numer. Math.* 92 (2002) 401–419.
- [3] R.E. Bank, A. Weiser, Some a posteriori error indicators for elliptic partial differential equations, *Math. Comput.* 44 (170) (1895) 283–301.
- [4] F. Brezzi, M. Fortin, *Mixed and Hybrid Finite Element Methods* Springer series in computational mathematics, Springer-Verlag, Berlin, 1991.
- [5] P. Destuynder, B. Métivet, Explicit error bounds in a conforming finite element method, *Math. Comput.* 68 (1999) 1379–1396.

- [6] P. Díez, N. Parés, A. Huerta, Recovering lower bounds of the error postprocessing implicit residual a posteriori error estimates, *Int. J. Numer. Methods Engrg.* 56 (2003) 1465–1488.
- [7] J. Fraeijs de Veubeke, Displacement and equilibrium models in the finite element method. B.M. Fraeijs de Veubeke Memorial Volume of Selected Papers 1980, *Int. J. Numer. Methods Engrg. Classical Reprint Ser.* 52 (2001) 287–342.
- [8] P. Ladevèze, D. Leguillon, Error estimate procedure in the finite element method and applications, *SIAM J. Numer. Anal.* 20 (1983) 485–509.
- [9] P. Ladevèze, J.P. Pelle, P. Rougeot, Error estimation and mesh optimization for classical finite elements, *Engrg. Comput.* 8 (1991) 69–80.
- [10] M. Paraschivoiu, A.T. Patera, A hierarchical duality approach to bounds for the outputs of partial differential equations, *Comput. Methods Appl. Mech. Engrg.* 158 (3-4) (1998) 389–407.
- [11] M. Paraschivoiu, J. Peraire, A.T. Patera, A posteriori finite element bounds for linear functional outputs of elliptic partial differential equations, *Comput. Methods Appl. Mech. Engrg.* 1580 (1997) 289–312.
- [12] J. Peraire, A.T. Patera, Bounds for linear-functional outputs of coercive partial differential equations: local indicators and adaptive refinement, in: P. Ladeveze, J. Oden (Eds.), *Proceedings of the Workshop on New Advances in Adaptive Computational Methods in Mechanics*, Cachan, September 17–19, Elsevier, Amsterdam, 1997.
- [13] S. Prudhomme, J.T. Oden, On goal-oriented error estimation for elliptic problems: application to the control of pointwise errors, *Comput. Methods Appl. Mech. Engrg.* 176 (1999) 313–331.
- [14] S. Prudhomme, J.T. Oden, T. Westermann, J. Bass, M.E. Botkin, Practical methods for a posteriori error estimation in engineering applications, *Int. J. Numer. Methods Engrg.* 56 (2003) 1193–1224.
- [15] R. Rannacher, F.-T. Suttmeier, A feed-back approach to error control in finite element methods: application to linear elasticity, *Comput. Mech.* 19 (1997) 434–446.
- [16] X. Roca, J. Sarrate, A. Huerta, Una librería orientada al objeto para el refinamiento de triangulos y tetraedros, *Aplicaciones al cálculo adaptado*. In: *Proceedings of Métodos Computacionais em Engenharia*, Lisbon 2004.
- [17] A.M. Sauer-Budge, J. Bonet, A. Huerta, J. Peraire, Computing bounds for linear functionals of exact weak solutions to Poisson's equation, *SIAM J. Numer. Anal.* 42 (4) (2004) 1610–1630.
- [18] A.M. Sauer-Budge, J. Peraire, Computing bounds for linear functionals of exact weak solutions to the advection–diffusion–reaction equation, *SIAM J. Sci. Comput.* 26 (2) (2003) 636–652.
- [19] T. Strouboulis, I. Babuska, Guaranteed computable bounds for the exact error in the finite element solution—Part II: bounds for the energy norm of the error in two dimensions, *Int. J. Numer. Meth. Engrn.* 47 (2000) 427–475.



