# Finite Element Methods

Susanne C. Brenner[1] and Carsten Carstensen[2]

[1] *Louisiana State University, Baton Rouge, LA, USA* [2] *Humboldt-Universität zu Berlin, Berlin, Germany*

## ABSTRACT

This introductory chapter on the mathematical theory of finite element methods (FEMs) discusses its $h$-version for elliptic boundary value problems in the displacement formulation. Topics addressed range from a priori to a posteriori error estimates and also include weak forms of elliptic PDEs, Galerkin schemes, finite element spaces, and adaptive local mesh refinement. Nonconformities and variational crimes as well as algorithmic aspects conclude the chapter.

KEY WORDS: finite element, Ritz-Galerkin methods, a priori error estimate, a posteriori error estimate, adaptive local mesh refinement

## 1. Introduction

The finite element method is one of the most widely used techniques in computational mechanics. The mathematical origin of the method can be traced to a paper by Courant (1943). We refer the readers to the articles by Babuška (1994) and Oden (1991) for the history of the finite element method. In this chapter, we give a concise account of the $h$-version of the finite element method for elliptic boundary value problems in the displacement formulation, and refer the readers to The $p$-version of the Finite Element Method and Mixed Finite Element Methods for the theory of the $p$-version of the finite element method and the theory of mixed finite element methods.

This chapter is organized as follows. The finite element method for elliptic boundary value problems is based on the Ritz-Galerkin approach, which is discussed in Section 2. The construction of finite element spaces and the a priori error estimates for finite element methods are presented in Sections 3 and 4. The a posteriori error estimates for finite element methods and their applications to adaptive local mesh refinements are discussed in Sections 5 and 6. For

the ease of presentation, the contents of Sections 3 and 4 are restricted to symmetric problems on polyhedral domains using conforming finite elements. The extension of these results to more general situations is outlined in Section 7.

For the classical material in Sections 3, 4, and 7, we are content with highlighting the important results and pointing to the key literature. We also concentrate on basic theoretical results and refer the readers to other chapters in this encyclopedia for complications that may arise in applications. For the recent development of a posteriori error estimates and adaptive local mesh refinements in Sections 5 and 6, we try to provide a more comprehensive treatment. Owing to space limitations many significant topics and references are inevitably absent. For in-depth discussions of many of the topics covered in this chapter (and the ones that we do not touch upon), we refer the readers to the following survey articles and books (which are listed in alphabetical order) and the references therein (Ainsworth and Oden, 2000; Apel, 1999; Aziz, 1972; Babuška and Aziz, 1972; Babuška and Strouboulis, 2001; Bangerth and Rannacher, 2003; Bathe, 1996; Becker, Carey and Oden, 1981; Becker and Rannacher, 2001; Braess, 2001; Brenner and Scott, 2002; Ciarlet, 1978, 1991; Eriksson *et al*, 1995; Hughes, 2000; Oden and Reddy, 1976; Schatz, Thomée and Wendland, 1990; Strang and Fix, 1973; Szabó and Babuška, 1991; Verfürth, 1996; Wahlbin, 1991, 1995; Zienkiewicz and Taylor, 2000).

## 2. Ritz-Galerkin Methods for Linear Elliptic Boundary Value Problems

In this section, we set up the basic mathematical framework for the analysis of Ritz-Galerkin methods for linear elliptic boundary value problems. We will concentrate on symmetric problems. Nonsymmetric elliptic boundary value problems will be discussed in Section 7.1.

### 2.1. Weak problems

Let $\Omega$ be a bounded connected open subset of the Euclidean space $\mathbb{R}^d$ with a piecewise smooth boundary. For a positive integer $k$, the Sobolev space $H^k(\Omega)$ is the space of square integrable functions whose weak derivatives up to order $k$ are also square integrable, with the norm

$$\|v\|_{H^k(\Omega)} = \left( \sum_{|\alpha| \le k} \left\| \frac{\partial^\alpha v}{\partial x^\alpha} \right\|_{L_2(\Omega)}^2 \right)^{1/2}$$

The seminorm $\left( \sum_{|\alpha|=k} \|(\partial^\alpha v / \partial x^\alpha)\|_{L_2(\Omega)}^2 \right)^{1/2}$ will be denoted by $|v|_{H^k(\Omega)}$. We refer the readers to Nečas (1967), Adams (1995), Triebel (1978), Grisvard (1985), and Wloka (1987) for the properties of the Sobolev spaces. Here we just point out that $\|\cdot\|_{H^k(\Omega)}$ is a norm induced by an inner product and $H^k(\Omega)$ is complete under this norm, that is, $H^k(\Omega)$ is a Hilbert space. (We assume that the readers are familiar with normed and Hilbert spaces.)

Using the Sobolev spaces we can represent a large class of symmetric elliptic boundary value

problems of order $2m$ in the following abstract weak form:

Find $u \in V$, a closed subspace of a Sobolev space $H^m(\Omega)$, such that

$$a(u, v) = F(v) \qquad \forall\, v \in V \tag{1}$$

where $F\colon V \to \mathbb{R}$ is a bounded linear functional on $V$ and $a(\cdot, \cdot)$ is a symmetric bilinear form that is bounded and $V$-elliptic, that is,

$$
\begin{aligned}
\left| a(v_1, v_2) \right| &\leq C_1 \|v_1\|_{H^m(\Omega)} \|v_2\|_{H^m(\Omega)} &&\forall\, v_1, v_2 \in V \tag{2} \\
a(v, v) &\geq C_2 \|v\|_{H^m(\Omega)}^2 &&\forall\, v \in V \tag{3}
\end{aligned}
$$

**Remark** 1. We use $C$, with or without subscript, to represent a generic positive constant that can take different values at different occurrences.

**Remark** 2. Equation (1) is the Euler-Lagrange equation for the variational problem of finding the minimum of the functional $v \mapsto \frac{1}{2} a(v, v) - F(v)$ on the space $V$. In mechanics, this functional often represents an energy and its minimization follows from the Dirichlet principle. Furthermore, the corresponding Euler-Lagrange equations (also called first variation) (1) often represent the principle of virtual work.

It follows from conditions (2) and (3) that $a(\cdot, \cdot)$ defines an inner product on $V$ which is equivalent to the inner product of the Sobolev space $H^m(\Omega)$. Therefore the existence and uniqueness of the solution of (1) follow immediately from (2), (3), and the Riesz Representation Theorem (Yosida, 1995; Reddy, 1986; Oden and Demkowicz, 1996).

The following are typical examples from computational mechanics.

**Example 1.** *Let $a(\cdot, \cdot)$ be defined by*

$$a(v_1, v_2) = \int_\Omega \nabla v_1 \cdot \nabla v_2 \, \mathrm{d}x \tag{4}$$

*For $f \in L_2(\Omega)$, the weak form of the Poisson problem*

$$
\begin{aligned}
-\Delta u &= f &&on &&\Omega \\
u &= 0 &&on &&\Gamma \\
\frac{\partial u}{\partial n} &= 0 &&on &&\partial\Omega \setminus \Gamma
\end{aligned} \tag{5}
$$

*where $\Gamma$ is a subset of $\partial\Omega$ with a positive $(d-1)$-dimensional measure, is given by (1) with $V = \{v \in H^1(\Omega)\colon v|_\Gamma = 0\}$ and*

$$F(v) = \int_\Omega fv\,\mathrm{d}x = (f, v)_{L_2(\Omega)} \tag{6}$$

*For the pure Neumann problem where $\Gamma = \emptyset$, since the gradient vector vanishes for constant functions, an appropriate function space for the weak problem is $V = \{v \in H^1(\Omega)\colon (v, 1)_{L_2(\Omega)} = 0\}$.*

The boundedness of $F$ and $a(\cdot,\cdot)$ is obvious and the coercivity of $a(\cdot,\cdot)$ follows from the Poincaré-Friedrichs inequalities (Nečas, 1967) :

$$\|v\|_{L_2(\Omega)} \leq C \left( |v|_{H^1(\Omega)} + \left| \int_\Gamma v \mathrm{ds} \right| \right) \quad \forall\, v \in H^1(\Omega) \tag{7}$$

$$\|v\|_{L_2(\Omega)} \leq C \left( |v|_{H^1(\Omega)} + \left| \int_\Omega v \mathrm{dx} \right| \right) \quad \forall\, v \in H^1(\Omega) \tag{8}$$

**Example 2.** Let $\Omega \subset \mathbb{R}^d (d = 2,3)$ and $\boldsymbol{v} \in [H^1(\Omega)]^d$ be the displacement of an elastic body. The strain tensor $\boldsymbol{\epsilon}(\boldsymbol{v})$ is given by the $d \times d$ matrix with components

$$\epsilon_{ij}(\boldsymbol{v}) = \frac{1}{2} \left( \frac{\partial v_i}{\partial x_j} + \frac{\partial v_j}{\partial x_i} \right) \tag{9}$$

and the stress tensor $\boldsymbol{\sigma}(\boldsymbol{v})$ is the $d \times d$ matrix defined by

$$\boldsymbol{\sigma}(\boldsymbol{v}) = 2\mu\, \boldsymbol{\epsilon}(\boldsymbol{v}) + \lambda\, (\mathit{div}\ \boldsymbol{v})\, \boldsymbol{\delta} \tag{10}$$

where $\boldsymbol{\delta}$ is the $d \times d$ identity matrix and $\mu > 0$ and $\lambda > 0$ are the Lamé constants.

Let the bilinear form $a(\cdot,\cdot)$ be defined by

$$\begin{aligned} a(\boldsymbol{v}_1, \boldsymbol{v}_2) &= \int_\Omega \sum_{i,j=1}^d \sigma_{ij}(\boldsymbol{v}_1)\, \epsilon_{ij}(\boldsymbol{v}_2) \mathrm{dx} \\ &= \int_\Omega \boldsymbol{\sigma}(\boldsymbol{v}_1) \colon \boldsymbol{\epsilon}(\boldsymbol{v}_2) \mathrm{dx} \end{aligned} \tag{11}$$

For $\boldsymbol{f} \in [L_2(\Omega)]^d$, the weak form of the linear elasticity problem (Ciarlet, 1988)

$$\begin{aligned} \mathbf{div}\,[\boldsymbol{\sigma}(\boldsymbol{u})] &= \boldsymbol{f} & \text{on } \Omega \\ \boldsymbol{u} &= 0 & \text{on } \Gamma \\ [\boldsymbol{\sigma}(\boldsymbol{u})]\boldsymbol{n} &= 0 & \text{on } \partial\Omega \setminus \Gamma \end{aligned} \tag{12}$$

where $\Gamma$ is a subset of $\partial\Omega$ with a positive $(d-1)$-dimensional measure, is given by (1) with $V = \{\boldsymbol{v} \in [H^1(\Omega)]^d \colon \boldsymbol{v}|_\Gamma = 0\}$ and

$$F(\boldsymbol{v}) = \int_\Omega \boldsymbol{f} \cdot \boldsymbol{v} \mathrm{dx} = (\boldsymbol{f}, \boldsymbol{v})_{L_2(\Omega)} \tag{13}$$

For the pure traction problem where $\Gamma = \emptyset$, the strain tensor vanishes for all infinitesimal rigid motions, i.e., displacement fields of the form $\boldsymbol{m} = \boldsymbol{a} + \boldsymbol{\rho}\,\boldsymbol{x}$, where $\boldsymbol{a} \in \mathbb{R}^d$, $\boldsymbol{\rho}$ is a $d \times d$ antisymmetric matrix and $\boldsymbol{x} = (x_1, \ldots, x_d)^t$ is the position vector. In this case an appropriate function space for the weak problem is $V = \{\boldsymbol{v} \in [H^1(\Omega)]^d \colon \int_\Omega \nabla \times \boldsymbol{v} \mathrm{dx} = 0 = \int_\Omega \boldsymbol{v} \mathrm{dx}\}$.

The boundedness of $F$ and $a(\cdot,\cdot)$ is obvious and the coercivity of $a(\cdot,\cdot)$ follows from Korn's inequalities (Friedrichs, 1947; Duvaut and Lions, 1976; Nitsche, 1981) (see *Finite Element*

*Methods for Elasticity with Error-controlled Discretization and Model Adaptivity*) :

$$\|\boldsymbol{v}\|_{H^1(\Omega)} \ \leq \ C\left(\|\boldsymbol{\varepsilon}(\boldsymbol{v})\|_{L_2(\Omega)} + \left|\int_\Gamma \boldsymbol{v}\mathrm{ds}\right|\right)$$
$$\forall\, v \in [H^1(\Omega)]^d \tag{14}$$

$$\|\boldsymbol{v}\|_{H^1(\Omega)} \ \leq \ C\left(\|\boldsymbol{\varepsilon}(\boldsymbol{v})\|_{L_2(\Omega)} + \left|\int_\Omega \nabla \times \boldsymbol{v}\mathrm{ds}\right| + \left|\int_\Omega \boldsymbol{v}\mathrm{dx}\right|\right)$$
$$\forall\, v \in [H^1(\Omega)]^d \tag{15}$$

**Example 3.** *Let $\Omega$ be a domain in $\mathbb{R}^2$ and the bilinear form $a(\cdot,\cdot)$ be defined by*

$$a(v_1, v_2) \ = \ \int_\Omega \Bigg[\Delta v_1 \Delta v_2 + (1-\sigma) \\ \times\left(2\frac{\partial^2 v_1}{\partial x_1 \partial x_2}\frac{\partial^2 v_2}{\partial x_1 \partial x_2} - \frac{\partial^2 v_1}{\partial x_1^2}\frac{\partial^2 v_2}{\partial x_2^2} - \frac{\partial^2 v_1}{\partial x_2^2}\frac{\partial^2 v_2}{\partial x_1^2}\right)\Bigg]\mathrm{dx} \tag{16}$$

*where $\sigma \in (0, 1/2)$ is the Poisson ratio.*

*For $f \in L_2(\Omega)$, the weak form of the clamped plate bending problem (*Ciarlet, 1997*)*

$$\Delta^2 u = f \quad on\ \Omega, \quad u = \frac{\partial u}{\partial n} = 0 \quad on\ \partial\Omega \tag{17}$$

*is given by (1), where $V = \{v \in H^2(\Omega)\colon v = \partial v/\partial n = 0\ on\ \partial\Omega\} = H_0^2(\Omega)$ and $F$ is defined by (6). For the simply supported plate bending problem, the function space $V$ is $\{v \in H^2(\Omega)\colon v = 0\ on\ \partial\Omega\} = H^2(\Omega) \cap H_0^1(\Omega)$.*

*For these problems, the coercivity of $a(\cdot,\cdot)$ is a consequence of the following Poincaré-Friedrichs inequality (*Nečas, 1967*) :*

$$\|v\|_{H^1(\Omega)} \leq C|v|_{H^2(\Omega)} \qquad \forall\, v \in H^2(\Omega) \cap H_0^1(\Omega) \tag{18}$$

**Remark** 3. The weak formulation of boundary value problems for beams and shells can be found in Plates and Shells: Asymptotic Expansions and Hierarchic Models and Models and Finite Elements for Thin-walled Structures.

*2.2. Ritz-Galerkin methods*

In the Ritz-Galerkin approach for (1), a discrete problem is formulated as follows.

Find $\tilde{u} \in \widetilde{V}$ such that

$$a(\tilde{u}, \tilde{v}) = F(\tilde{v}) \qquad \forall\, \tilde{v} \in \widetilde{V} \tag{19}$$

where $\widetilde{V}$, the space of trial/test functions, is a finite-dimensional subspace of $V$.

The orthogonality relation

$$a(u - \tilde{u}, \tilde{v}) = 0 \qquad \forall\, \tilde{v} \in \widetilde{V} \tag{20}$$

follows by subtracting (19) from (1), and hence

$$\|u - \tilde{u}\|_a = \inf_{\tilde{v} \in V} \|u - \tilde{v}\|_a \tag{21}$$

where $\|\cdot\|_a = (a(\cdot, \cdot))^{1/2}$. Furthermore, (2), (3), and (21) imply that

$$\|u - \tilde{u}\|_{H^m(\Omega)} \le \left(\frac{C_1}{C_2}\right)^{1/2} \inf_{\tilde{v} \in \widetilde{V}} \|u - \tilde{v}\|_{H^m(\Omega)} \tag{22}$$

that is, the error for the approximate solution $\tilde{u}$ is quasi-optimal in the norm of the Sobolev space underlying the weak problem.

The abstract estimate (22), called Cea's lemma, reduces the error estimate for the Ritz-Galerkin method to a problem in approximation theory, namely, to the determination of the magnitude of the error of the best approximation of $u$ by a member of $\widetilde{V}$. The solution of this problem depends on the regularity (smoothness) of $u$ and the nature of the space $\widetilde{V}$.

One can also measure $u - \tilde{u}$ in other norms. For example, an estimate of $\|u - \tilde{u}\|_{L_2(\Omega)}$ can be obtained by the Aubin-Nitsche duality technique as follows. Let $w \in V$ be the solution of the weak problem

$$a(v, w) = \int_\Omega (u - \tilde{u}) v \mathrm{dx} \qquad \forall\, v \in V \tag{23}$$

Then we have, from (20), (23), and the Cauchy-Schwarz inequality,

$$\begin{aligned} \|u - \tilde{u}\|_{L_2(\Omega)}^2 &= a(u - \tilde{u}, w) = a(u - \tilde{u}, w - \tilde{v}) \\ &\le C_2 \|u - \tilde{u}\|_{H^m(\Omega)} \|w - \tilde{v}\|_{H^m(\Omega)} \quad \forall\, \tilde{v} \in \widetilde{V} \end{aligned}$$

which implies that

$$\|u - \tilde{u}\|_{L_2(\Omega)} \le C_2 \left( \inf_{\tilde{v} \in \widetilde{V}} \frac{\|w - \tilde{v}\|_{H^m(\Omega)}}{\|u - \tilde{u}\|_{L_2(\Omega)}} \right) \|u - \tilde{u}\|_{H^m(\Omega)} \tag{24}$$

In general, since $w$ can be approximated by members of $\widetilde{V}$ to high accuracy, the term inside the bracket on the right-hand side of (24) is small, which shows that the $L_2$ error is much smaller than the $H^m$ error.

The estimates (22) and (24) provide the basic a priori error estimates for the Ritz-Galerkin method in an abstract setting.

On the other hand, the error of the Ritz-Galerkin method can also be estimated in an a posteriori fashion. Let the computable linear functional (the residual of the approximate solution $\tilde{u}$) $R\colon V \to \mathbb{R}$ be defined by

$$R(v) = a(u - \tilde{u}, v) = F(v) - a(\tilde{u}, v) \tag{25}$$

The global a posteriori error estimate

$$\|u - \tilde{u}\|_{H^m(\Omega)} \le \frac{1}{C_2} \sup_{v \in V} \frac{|R(v)|}{\|v\|_{H^m(\Omega)}} \tag{26}$$

then follows from (3) and (25).

Let $D$ be a subdomain of $\Omega$ and $H_0^m(D)$ be the subspace of $V$ whose members vanish identically outside $D$. It follows from (25) and the local version of (2) that we also have a local a posteriori error estimate:

$$\|u - \tilde{u}\|_{H^m(D)} \geq \frac{1}{C_1} \sup_{v \in H_0^m(D)} \frac{|R(v)|}{\|v\|_{H^m(D)}} \tag{27}$$

The equivalence of the error norm with the dual norm of the residual will be the point of departure in Section 5.1.2 (cf. (70)).

### 2.3. Elliptic regularity

As mentioned above, the magnitude of the error of a Ritz-Galerkin method for an elliptic boundary value problem depends on the regularity of the solution. Here we give a brief description of elliptic regularity for the examples in Section 2.1.

If the boundary $\partial\Omega$ is smooth and the homogeneous boundary conditions are also smooth (i.e. the Dirichlet and Neumann boundary condition in (5) and the displacement and traction boundary conditions in (12) are defined on disjoint components of $\partial\Omega$), then the solution of the elliptic boundary value problems in Section 2.1 obey the classical *Shift Theorem* (Agmon, 1965; Nečas, 1967; Gilbarg and Trudinger, 1983; Wloka, 1987). In other words, if the right-hand side of the equation belongs to the Sobolev space $H^\ell(\Omega)$, then the solution of a $2m$-th order elliptic boundary problem belongs to the Sobolev space $H^{2m+\ell}(\Omega)$.

The Shift Theorem does not hold for domains with piecewise smooth boundary in general. For example, let $\Omega$ be the $L$-shaped domain depicted in Figure 1 and

$$u(x) = \phi(r)\, r^{2/3} \sin\left(\frac{2}{3}\left(\theta - \frac{\pi}{2}\right)\right) \tag{28}$$

where $r = (x_1^2 + x_2^2)^{1/2}$ and $\theta = \arctan(x_2/x_1)$ are the polar coordinates and $\phi$ is a smooth cut-off function that equals 1 for $0 \leq r < 1/2$ and 0 for $r > 3/4$. It is easy to check that $u \in H_0^1(\Omega)$ and $-\Delta u \in C^\infty(\overline{\Omega})$. Let $D$ be any open neighborhood of the origin in $\Omega$. Then $u \in H^2(\Omega \setminus \overline{D})$ but $u \notin H^2(D)$. In fact $u$ belongs to the Besov space $B_{2,\infty}^{5/3}(D)$ (Babuška and Osborn, 1991), which implies that $u \in H^{5/3-\epsilon}(D)$ for any $\epsilon > 0$, but $u \notin H^{5/3}(D)$ (see Triebel (1978) and Grisvard (1985) for a discussion of Besov spaces and fractional order Sobolev spaces). A similar situation occurs when the types of boundary condition change abruptly, such as the Poisson problem with mixed boundary conditions depicted on the circular domain in Figure 1, where the homogeneous Dirichlet boundary condition is assumed on the upper semicircle and the homogeneous Neumann boundary condition is assumed on the lower semicircle.

Therefore (Dauge, 1988), for the second (respectively fourth) order model problems in Section 2.1, the solution in general only belongs to $H^{1+\alpha}(\Omega)$ (respectively $H^{2+\alpha}(\Omega)$) for some $\alpha \in (0, 1]$ even if the right-hand side of the equation belongs to $C^\infty(\overline{\Omega})$.
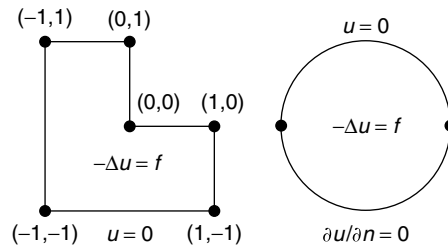
Figure 1: Singular points of two-dimensional elliptic boundary value problems.

For two-dimensional problems, the vertices of $\Omega$ and the points where the boundary condition changes type are the singular points (cf. Figure 1). Away from these singular points, the Shift Theorem is valid. The behavior of the solution near the singular points is also well understood. If the right-hand side function and its derivatives vanish to sufficiently high order at the singular points, then the Shift Theorem holds for certain weighted Sobolev spaces (Nazarov and Plamenevsky, 1994; Kozlov, Maz'ya and Rossman, 1997, 2001). Alternatively, one can represent the solution near a singular point as a sum of a regular part and a singular part (Grisvard, 1985; Dauge, 1988; Nicaise, 1993). For a $2m$-th order problem, the regular part of the solution belongs to the Sobolev space $H^{2m+k}(\Omega)$ if the right-hand side function belongs to $H^k(\Omega)$, and the singular part of the solution is a linear combination of special functions with less regularity, analogous to the function in (28).

The situation in three dimensions is more complicated due to the presence of edge singularities, vertex singularities, and edge-vertex singularities. The theory of three-dimensional singularities remains an active area of research.

## 3. Finite Element Spaces

Finite element methods are Ritz-Galerkin methods where the finite-dimensional trial/test function spaces are constructed by piecing together polynomial functions defined on (small) parts of the domain $\Omega$. In this section, we describe the construction and properties of finite element spaces. We will concentrate on conforming finite elements here and leave the discussion of nonconforming finite elements to Section 7.2.

### 3.1. The concept of a finite element

A $d$-dimensional finite element (Ciarlet, 1978; Brenner and Scott, 2002) is a triple $(K, \mathcal{P}_K, \mathcal{N}_K)$, where $K$ is a closed bounded subset of $\mathbb{R}^d$ with nonempty interior and a piecewise smooth boundary, $\mathcal{P}_K$ is a finite-dimensional vector space of functions defined on $K$ and $\mathcal{N}_K$ is a basis of the dual space $\mathcal{P}'_K$. The function space $\mathcal{P}_K$ is the space of the shape functions and the
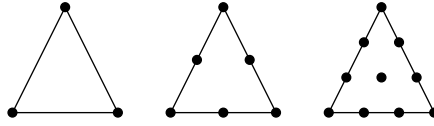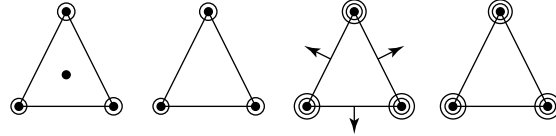
Figure 2: Lagrange elements.



Figure 3: Cubic Hermite element, Zienkiewicz element, fifth degree Argyris element and Bell element.

elements of $\mathcal{N}_K$ are the nodal variables (degrees of freedom).

The following are examples of two-dimensional finite elements.

**Example 4.** (*Triangular Lagrange Elements*) *Let $K$ be a triangle, $\mathcal{P}_K$ be the space $P_n$ of polynomials in two variables of degree $\leq n$, and let the set $\mathcal{N}_K$ consist of evaluations of shape functions at the nodes with barycentric coordinates $\lambda_1 = i/n$, $\lambda_2 = j/n$ and $\lambda_3 = k/n$, where $i, j, k$ are nonnegative integers and $i + j + k = n$. Then $(K, \mathcal{P}_K, \mathcal{N}_K)$ is the two-dimensional $P_n$ Lagrange finite element. The nodal variables for the $P_1$, $P_2$, and $P_3$ Lagrange elements are depicted in Figure 2, where $\bullet$ (here and in the following examples) represents pointwise evaluation of shape functions.*

**Example 5.** (*Triangular Hermite Elements*) *Let $K$ be a triangle. The cubic Hermite element is the triple $(K, P_3, \mathcal{N}_K)$ where $\mathcal{N}_K$ consists of evaluations of shape functions and their gradients at the vertices and evaluation of shape functions at the center of $K$. The nodal variables for the cubic Hermite element are depicted in the first figure in Figure 3, where $\circ$ (here and in the following examples) represents pointwise evaluation of gradients of shape functions.*

*By removing the nodal variable at the center (cf. the second figure in Figure 3) and reducing the space of shape functions to*

$$\left\{ v \in P_3 \colon 6v(c) - 2\sum_{i=1}^{3} v(p_i) + \sum_{i=1}^{3} (\nabla v)(p_i) \cdot (p_i - c) = 0 \right\} (\supset P_2)$$

*where $p_i$ $(i = 1, 2, 3)$ and $c$ are the vertices and center of $K$ respectively, we obtain the Zienkiewicz element.*

*The fifth degree Argyris element is the triple $(K, P_5, \mathcal{N}_K)$ where $\mathcal{N}_K$ consists of evaluations of the shape functions and their derivatives up to order two at the vertices and evaluations of the normal derivatives at the midpoints of the edges. The nodal variables for the Argyris element are depicted in the third figure in Figure 3, where $\bigcirc$ and $\uparrow$ (here and in the following*
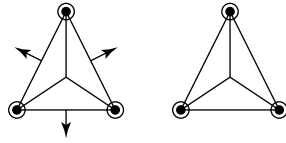
Figure 4: Hsieh-Clough-Tocher element and reduced Hsieh-Clough-Tocher element.
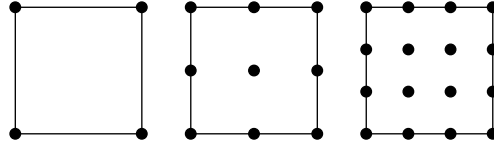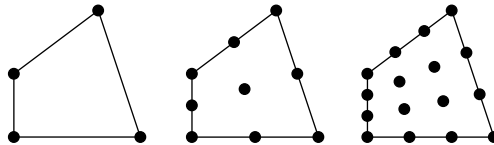


Figure 5: Tensor product elements.



Figure 6: $Q_n$ quadrilateral elements.

*examples) represent pointwise evaluation of second order derivatives and the normal derivative of the shape functions, respectively.*

*By removing the nodal variables at the midpoints of the edges (cf. the fourth figure in Figure 3) and reducing the space of shape functions to $\{v \in P_5 \colon (\partial v/\partial n)\big|_e \in P_3(e)$ for each edge $e\}$, we obtain the Bell element.*

**Example 6.** (*Triangular Macro Elements*) *Let $K$ be a triangle that is subdivided into three subtriangles by the center of $K$, $\mathcal{P}_K$ be the space of piecewise cubic polynomials with respect to this subdivision that belong to $C^1(K)$, and let the set $\mathcal{N}_K$ consist of evaluations of the shape functions and their first-order derivatives at the vertices of $K$ and evaluations of the normal derivatives of the shape functions at the midpoints of the edges of $K$. Then $(K, \mathcal{P}_K, \mathcal{N}_K)$ is the Hsieh-Clough-Tocher macro element. The nodal variables for this element are depicted in the first figure in Figure 4.*

*By removing the nodal variables at the midpoints of the edges (cf. the second figure in Figure 4) and reducing the space of shape functions to $\{v \in C^1(K) \colon v$ is piecewise cubic and $(\partial v/\partial n)\big|_e \in P_1(e)$ for each edge $e\}$, we obtain the reduced Hsieh-Clough-Tocher macro element.*

**Example 7.** (*Rectangular Tensor Product Elements*) *Let $K$ be the rectangle $[a_1, b_1] \times [a_2, b_2]$, $\mathcal{P}_K$ be the space spanned by the monomials $x_1^i x_2^j$ for $0 \le i, j \le n$, and the set $\mathcal{N}_K$ consist of evaluations of shape functions at the nodes with coordinates $\big(a_1 + i(b_1 - a_1)/n, a_2 + j(b_2 - a_2)/n\big)$ for $0 \le i, j \le n$. Then $(K, \mathcal{P}_K, \mathcal{N}_K)$ is the two-dimensional $Q_n$ tensor product element. The*
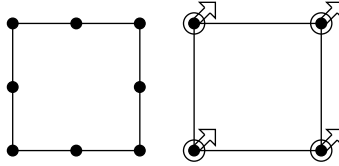
Figure 7: Serendipity and Bogner-Fox-Schmit elements.

*nodal variables of the $Q_1$, $Q_2$ and $Q_3$ elements are depicted in Figure 5.*

**Example 8.** (*Quadrilateral $Q_n$ Elements*) *Let $K$ be a convex quadrilateral; then there exists a bilinear map $(x_1, x_2) \mapsto B(x_1, x_2) = (a_1 + b_1 x_1 + c_1 x_2 + d_1 x_1 x_2,\ a_2 + b_2 x_1 + c_2 x_2 + d_2 x_1 x_2)$ f rom the biunit square $S$ with vertices $(\pm 1, \pm 1)$ onto $K$. The space of shape functions is defined by $v \in \mathcal{P}_K$ if and only if $v \circ B \in Q_n$ and $\mathcal{N}_K$ consists of pointwise evaluations of the shape functions at the nodes of $K$ corresponding under the map $B$ to the nodes of the $Q_n$ tensor product element on $S$. The nodal variables of the $Q_1$, $Q_2$ and $Q_3$ quadrilateral elements are depicted in Figure 6.*

**Example 9.** (*Other Rectangular Elements*) *Let $K$ be the rectangle $[a_1, b_1] \times [c_1, d_1]$;*

$$\mathcal{P}_K = \left\{ v \in Q_2 : 4v(c) + \sum_{i=1}^{4} v(p_i) - 2\sum_{i=1}^{4} v(m_i) = 0 \right\} (\supset P_2)$$

*where the $p_i$'s are the vertices of $K$, the $m_i$'s are the midpoints of the edges of $K$ and $c$ is the center of $K$; and $\mathcal{N}_K$ consist of evaluations of the shape functions at the vertices and the midpoints (cf. the first figure in Figure 7). Then $(K, \mathcal{P}_K, \mathcal{N}_K)$ is the 8-node serendipity element.*

*If we take $\mathcal{P}_K$ to be the space of bicubic polynomials spanned by $x_1^i x_2^j$ for $0 \le i, j \le 3$ and $\mathcal{N}_K$ to be the set consisting of evaluations at the vertices of $K$ of the shape functions, their first-order derivatives and their second-order mixed derivatives, then we have the Bogner-Fox-Schmit element. The nodal variables for this element are depicted in the second figure in Figure 7, where the tilted arrows represent pointwise evaluations of the second-order mixed derivatives of the shape functions.*

**Remark** 4. The triangular $P_n$ elements and the quadrilateral $Q_n$ elements, which are suitable for second order elliptic boundary value problems, can be generalized to any dimension in a straightforward manner. The Argyris element, the Bell element, the macro elements, and the Bogner-Fox-Schmit element are suitable for fourth-order problems in two space dimensions.

*3.2. Triangulations and finite element spaces*

We restrict $\Omega \subset \mathbb{R}^d$ ($d = 1, 2, 3$) to be a polyhedral domain in this and the following sections. The case of curved domains will be discussed in Section 7.4.
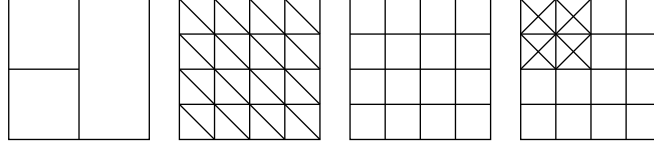
Figure 8: Partitions and triangulations.

A *partition* of $\Omega$ is a collection $\mathcal{P}$ of polyhedral subdomains of $\Omega$ such that

$$\overline{\Omega} = \bigcup_{D \in \mathcal{P}} \overline{D} \quad \text{and} \quad D \cap D' = \emptyset \quad \text{if } D, D' \in \mathcal{P}, D \neq D'$$

where we use $\overline{\Omega}$ and $\overline{D}$ to represent the closures of $\Omega$ and $D$.

A *triangulation* of $\Omega$ is a partition where the intersection of the closures of two distinct subdomains is either empty, a common vertex, a common edge or a common face. For $d = 1$, every partition is a triangulation. But the two concepts are different when $d \geq 2$. A partition that is not a triangulation is depicted in the first figure in Figure 8, where the other three figures represent triangulations. Below we will concentrate on triangulations consisting of triangles or convex quadrilaterals in two dimensions and tetrahedrons or convex hexahedrons in three dimensions.

The shape regularity of a triangle (or tetrahedron) $D$ can be measured by the parameter

$$\gamma(D) = \frac{\operatorname{diam} D}{\text{diameter of the largest ball in } \overline{D}} \tag{29}$$

which will be referred to as the aspect ratio of the triangle (tetrahedron). We say that a family of triangulations of triangles (or tetrahedrons) $\{\mathcal{T}_i : i \in I\}$ is *regular* (or *nondegenerate*) if the aspect ratios of all the triangles (tetrahedrons) in the triangulations are bounded, that is, there exists a positive constant $C$ such that

$$\gamma(D) \leq C \qquad \text{for all} \quad D \in \mathcal{T}_i \quad \text{and} \quad i \in I$$

The shape regularity of a convex quadrilateral (or hexahedron) $D$ can be measured by the parameter $\gamma(D)$ defined in (29) and the parameter

$$\sigma(D) = \max\left\{\frac{|e_1|}{|e_2|} : e_1 \text{ and } e_2 \text{ are any two edges of } D\right\} \tag{30}$$

We will refer to the number $\max(\gamma(D), \sigma(D))$ as the aspect ratio of the convex quadrilateral (hexahedron). We say that a family of triangulations of convex quadrilaterals (or hexahedrons) $\{\mathcal{T}_i : i \in I\}$ is *regular* if the aspect ratios of all the quadrilaterals in the triangulations are bounded, that is, there exists a positive constant $C$ such that

$$\gamma(D), \sigma(D) \leq C \qquad \text{for all} \quad D \in \mathcal{T}_i \quad \text{and} \quad i \in I$$

A family of triangulations is *quasi-uniform* if it is regular and there exists a positive constant $C$ such that

$$h_i \leq C \operatorname{diam} D \qquad \forall\, D \in \mathcal{T}_i, \quad i \in I \tag{31}$$

where $h_i$ is the maximum of the diameters of the subdomains in $\mathcal{T}_i$.

**Remark** 5. For a triangle or a tetrahedron $D$, a lower bound for the angles of $D$ can lead to an upper bound for $\gamma(D)$ (and vice versa). Therefore, the *regularity* of a family of simplicial triangulations (i.e. triangulations consisting of triangles or tetrahedrons) is equivalent to the following *minimum angle condition*: There exists $\theta_* > 0$ such that the angles of the simplexes in all the triangulations $\mathcal{T}_i$ are bounded below by $\theta_*$.

**Remark** 6. A family of triangulations obtained by successive uniform subdivisions of an initial triangulation is quasi-uniform. A family of triangulations generated by a *local refinement* strategy is usually regular but not quasi-uniform.

Let $\mathcal{T}$ be a triangulation of $\Omega$, and a finite element $(\overline{D}, \mathcal{P}_{\overline{D}}, \mathcal{N}_{\overline{D}})$ be associated with each subdomain $D \in \mathcal{T}$. We define the corresponding finite element space to be

$$
\begin{aligned}
FE_{\mathcal{T}} \;=\; & \{v \in L_2(\Omega) \colon v_{\overline{D}} = v|_{\overline{D}} \in \mathcal{P}_{\overline{D}} \quad \forall\, D \in \mathcal{T}, \text{ and} \\
& v_{\overline{D}}, v_{\overline{D}'} \text{ share the same nodal values on } \overline{D} \cap \overline{D}'\}
\end{aligned} \tag{32}
$$

We say that $FE_{\mathcal{T}}$ is a $C^r$ finite element space if $FE_{\mathcal{T}} \subset C^r(\overline{\Omega})$. For example, the finite element spaces constructed from the Lagrange finite elements (Example 4), the tensor product elements (Example 7), the cubic Hermite element (Example 5), the Zienkiewicz element (Example 5) and the serendipity element (Example 9) are $C^0$ finite element spaces, and those constructed from the quintic Argyris element (Example 5), the Bell element (Example 5), the macro elements (Example 6) and the Bogner-Fox-Schmit element (Example 9) are $C^1$ finite element spaces.

Note that a $C^r$ finite element space is automatically a subspace of the Sobolev space $H^{r+1}(\Omega)$ and therefore appropriate for elliptic boundary value problems of order $2(r+1)$.

### 3.3. Element nodal interpolation operators and interpolation error estimates

Let $(K, \mathcal{P}_K, \mathcal{N}_K)$ be a finite element. Denote the nodal variables in $\mathcal{N}_K$ by $N_1, \ldots, N_n$ ($n = \dim \mathcal{P}_K$) and the dual basis of $\mathcal{P}_K$ by $\phi_1, \ldots, \phi_n$, that is,

$$N_i(\phi_j) = \delta_{ij} = \begin{cases} 1 & \text{if} \quad i = j \\ 0 & \text{if} \quad i \neq j \end{cases}$$

Assume that $\zeta \mapsto N_i(\zeta)$ is well-defined for $\zeta \in H^s(K)$ (where $s$ is a sufficiently large positive number), then we can define the *element nodal interpolation operator* $\Pi_K \colon H^s(K) \to \mathcal{P}_K$ by

$$\Pi_K \zeta = \sum_{j=1}^{n} N_j(\zeta) \phi_j \tag{33}$$

Note that (33) implies

$$\Pi_K v = v \qquad \forall\, v \in \mathcal{P}_K \tag{34}$$

For example, by the Sobolev embedding theorem (Adams, 1995; Nečas, 1967; Wloka, 1987; Gilbarg and Trudinger, 1983), the nodal interpolation operators associated with the Lagrange finite elements (Example 4), the tensor product finite elements (Example 7), and the serendipity element (Example 9) are well-defined on $H^s(K)$ for $s > 1$ if $K \subset \mathbb{R}^2$ and for $s > 3/2$ if $K \subset \mathbb{R}^3$. On the other hand the nodal interpolation operators associated with the Zienkiewicz element (Example 5) and the macro elements (Example 6) are well-defined on $H^s(K)$ for $s > 2$, while the interpolation operators for the quintic Argyris element (Example 5), the Bell element (Example 5) or the Bogner-Fox-Schmit (Example 9) are well-defined on $H^s(K)$ for $s > 3$.

The error of the element nodal interpolation operator for a triangular (tetrahedral) or convex quadrilateral (hexagonal) element $(K, \mathcal{P}_K, \mathcal{N}_K)$ can be controlled in terms of the shape regularity of $K$. Let $\hat{K}$ be the image of $K$ under the scaling map

$$x \mapsto \mathcal{H}(x) = (\operatorname{diam} K)^{-1} x \tag{35}$$

Then $\widehat{K}$ is a domain of unit diameter and we can define a finite element $(\widehat{K}, \mathcal{P}_{\widehat{K}}, \mathcal{N}_{\widehat{K}})$ as follows: (i) $\hat{v} \in \mathcal{P}_{\widehat{K}}$ if and only if $\hat{v} \circ \mathcal{H} \in \mathcal{P}_K$, and (ii) $N \in \mathcal{N}_{\widehat{K}}$ if and only if the linear functional $v \mapsto N(v \circ \mathcal{H}^{-1})$ on $\mathcal{P}_K$ belongs to $\mathcal{N}_K$. It follows that the dual basis $\hat{\phi}_1, \ldots, \hat{\phi}_n$ of $\mathcal{P}_{\widehat{K}}$ is related to the dual basis $\phi_1, \ldots, \phi_n$ of $\mathcal{P}_K$ through the relation $\hat{\phi}_i \circ \mathcal{H} = \phi_i$, and (33) implies that

$$(\Pi_K \zeta) \circ \mathcal{H}^{-1} = \Pi_{\widehat{K}}(\zeta \circ \mathcal{H}^{-1}) \tag{36}$$

for all sufficiently smooth functions $\zeta$ defined on $K$. Moreover, for the functions $\hat{\zeta}$ and $\zeta$ related by $\zeta(x) = \hat{\zeta}(\mathcal{H}(x))$, we have

$$|\hat{\zeta}|^2_{H^s(\widehat{K})} = (\operatorname{diam} K)^{2s-d} |\zeta|^2_{H^s(K)} \tag{37}$$

where $d$ is the spatial dimension.

Assuming that $\mathcal{P}_{\widehat{K}} \supseteq P_m$ (equivalently $\mathcal{P}_K \supseteq P_m$), we have, by (34),

$$\begin{aligned}
\|\hat{\zeta} - \Pi_{\widehat{K}}\hat{\zeta}\|_{H^m(\widehat{K})} &= \|(\hat{\zeta} - p) - \Pi_{\widehat{K}}(\hat{\zeta} - p)\|_{H^m(\widehat{K})} \\
&\leq 2\|\Pi_{\widehat{K}}\|_{m,s} \|\hat{\zeta} - p\|_{H^s(\widehat{K})} \quad \forall\, p \in P_m
\end{aligned}$$

where $\|\Pi_{\widehat{K}}\|_{m,s}$ is the norm of the operator $\Pi_{\widehat{K}} : H^s(\widehat{K}) \to H^m(\widehat{K})$, and hence

$$\|\hat{\zeta} - \Pi_{\widehat{K}}\hat{\zeta}\|_{H^m(\widehat{K})} \leq 2\|\Pi_{\widehat{K}}\|_{m,s} \inf_{p \in P_m} \|\hat{\zeta} - p\|_{H^s(\widehat{K})} \tag{38}$$

Since $K$ is convex, the following estimate (Verfürth, 1999) holds provided $m$ is the largest integer strictly less than $s$:

$$\inf_{p \in P_m} \|\hat{\zeta} - p\|_{H^s(\widehat{K})} \leq C_{s,d} |\hat{\zeta}|_{H^s(\widehat{K})} \qquad \forall\, \hat{\zeta} \in H^s(\widehat{K}) \tag{39}$$

where the positive constant $C_{s,d}$ depends only on $s$ and $d$.

Combining (38) and (39) we find

$$\|\hat{\zeta} - \Pi_{\widehat{K}}\hat{\zeta}\|_{H^m(\widehat{K})} \leq 2C_{s,d}\|\Pi_{\widehat{K}}\|_{m,s}|\hat{\zeta}|_{H^s(\widehat{K})} \quad \forall \hat{\zeta} \in H^s(\widehat{K}) \tag{40}$$

We have therefore reduced the error estimate for the element nodal interpolation operator to an estimate of $\|\Pi_{\widehat{K}}\|_{m,s}$. Since diam $\widehat{K} = 1$, the norm $\|\Pi_{\widehat{K}}\|_{m,s}$ is a constant depending only on the shape of $\widehat{K}$ (equivalently of $K$), if we considered $s$ and $m$ to be fixed for a given type of element.

For triangular elements, we can use the concept of *affine-interpolation-equivalent* elements to obtain a more concrete description of the dependence of $\|\Pi_{\widehat{K}}\|_{m,s}$ on the shape of $\widehat{K}$. A $d$-dimensional nondegenerate affine map is a map of the form $x \mapsto Ax+b$ where $A$ is a nonsingular $d \times d$ matrix and $b \in \mathbb{R}^d$. We say that two finite elements $(K_1, \mathcal{P}_{K_1}, \mathcal{N}_{K_1})$ and $(K_2, \mathcal{P}_{K_2}, \mathcal{N}_{K_2})$ are affine-interpolation-equivalent if (i) there exists a nondegenerate affine map $\Phi$ that maps $K_1$ onto $K_2$, (ii) $v \in \mathcal{P}_{K_2}$ if and only if $v \circ \Phi \in \mathcal{P}_{K_1}$ and (iii)

$$(\Pi_{K_2}\zeta) \circ \Phi = \Pi_{K_1}(\zeta \circ \Phi) \tag{41}$$

for all sufficiently smooth functions $\zeta$ defined on $K_2$. For example, any triangular elements in one of the families (except the Bell element and the reduced Hsieh-Clough-Tocher element) described in Section 3.1 are affine- interpolation-equivalent to the corresponding element on the standard simplex $S$ with vertices $(0,0)$, $(1,0)$ and $(0,1)$.

Assuming $(\widehat{K}, \mathcal{P}_{\widehat{K}}, \mathcal{N}_{\widehat{K}})$ (or equivalently $(K, \mathcal{P}_K, \mathcal{N}_K)$) is affine-interpolation-equivalent to the element $(S, \mathcal{P}_S, \mathcal{N}_S)$ on the standard simplex, it follows from (41) and the chain rule that

$$\|\Pi_{\widehat{K}}\|_{m,s} \leq C\|\Pi_S\|_{m,s} \tag{42}$$

where the positive constant depends only on the Jacobian matrix of the affine map $\hat{\Phi} \colon S \to \widehat{K}$ and thus depends only on an upper bound of the parameter $\gamma(\widehat{K})$ (cf. (29)) which is identical with $\gamma(K)$.

Combining (36), (37), (40) and (42), we find

$$\sum_{k=0}^{m}(\operatorname{diam} K)^k|\zeta - \Pi_K\zeta|_{H^m(\widehat{K})} \leq C(\operatorname{diam} K)^s|\zeta|_{H^s(K)}$$
$$\forall \zeta \in H^s(K) \tag{43}$$

where the positive constant $C$ depends only on $s$ and an upper bound of the parameter $\gamma(K)$ (the aspect ratio of $K$), provided that (i) the element nodal interpolation operator is well-defined on $H^s(K)$, (ii) the triangular element $(K, \mathcal{P}_K, \mathcal{N}_K)$ is affine-interpolation-equivalent to a reference element $(S, \mathcal{P}_S, \mathcal{N}_S)$ on the standard simplex, (iii) $\mathcal{P} \supseteq P_m$, and (iv) $m$ is the largest integer $< s$.

For convex quadrilateral elements, we can similarly obtain a concrete description of the dependence of $\|\Pi_{\widehat{K}}\|_{m,s}$ on the shape of $\widehat{K}$ by assuming that there is a reference

element $(S, \mathcal{P}_S, \mathcal{N}_S)$ defined on the biunit square $S$ with vertices $(\pm 1, \pm 1)$ and a bilinear homeomorphism $\hat{\Phi}$ from $S$ onto $\widehat{K}$ with the following properties: $\hat{v} \in \mathcal{P}_{\widehat{K}}$ if and only if $v \circ \Phi \in \mathcal{P}_S$ and $(\Pi_{\widehat{K}} \hat{\zeta}) \circ \hat{\Phi} = \Pi_S(\hat{\zeta} \circ \hat{\Phi})$ for all sufficiently smooth functions $\hat{\zeta}$ defined on $\widehat{K}$. Note that because of (36) this is equivalent to the existence of a bilinear homeomorphism from $S$ onto $K$ such that

$$v \in \mathcal{P}_K \Longleftrightarrow v \circ \Phi \in \mathcal{P}_S$$
$$\text{and} \quad (\Pi_K \zeta) \circ \Phi = \Pi_S(\zeta \circ \Phi) \tag{44}$$

for all sufficiently smooth functions $\zeta$ defined on $K$. The estimate (42) holds again by the chain rule, where the positive constant $C$ depends only on the Jacobian matrix of $\hat{\Phi}$ and thus depends only on upper bounds for the parameters $\gamma(\widehat{K})$ and $\sigma(\widehat{K})$ (cf. (30)), which are identical with $\gamma(K)$ and $\sigma(K)$. We conclude that the estimate (43) also holds for convex quadrilateral elements where the positive constant $C$ depends on upper bounds of $\gamma(K)$ and $\sigma(K)$ (equivalently an upper bound of the aspect ratio of $K$) provided condition (ii) is replaced by (44). For example, the estimate (43) is valid for the quadrilateral $Q_n$ element in Example 8.

**Remark** 7. The general estimate (40) can be refined to yield *anisotropic* error estimates for certain reference elements. For example, in two dimensions, the following estimates (Apel and Dobrowolski, 1992; Apel, 1999) hold for the $P_n$ Lagrange elements on the reference simplex $S$ and the $Q_n$ tensor product elements on the reference square $S$:

$$\left\| \frac{\partial}{\partial x_j} (\zeta - \Pi_S \zeta) \right\|_{L_2(S)}$$
$$\leq C \left( \left\| \frac{\partial^2 \zeta}{\partial x_1 \partial x_j} \right\|_{L_2(S)} + \left\| \frac{\partial^2 \zeta}{\partial x_2 \partial x_j} \right\|_{L_2(S)} \right) \tag{45}$$

for $j = 1, 2$ and for all $\zeta \in H^2(S)$. We refer the readers to Interpolation in *h*-version Finite Element Spaces, for more details.

**Remark** 8. The analysis of the quadrilateral serendipity elements is more subtle. A detailed discussion can be found in Arnold, Boffi and Falk (2002).

**Remark** 9. The estimate (43) can be generalized naturally to 3-D tetrahedral $P_n$ elements and hexahedral $Q_n$ elements.

**Remark** 10. Let $n$ be a nonnegative integer and $n < s \leq n + 1$. The estimate

$$\inf_{p \in P_n(\Omega)} \|\zeta - p\|_{H^s(\Omega)} \leq C_{\Omega,s} |\zeta|_{H^s(\Omega)} \qquad \forall \zeta \in H^s(\Omega) \tag{46}$$

for general $\Omega$ follows from generalized Poincaré-Friedrichs inequalities (Nečas, 1967). In the case where $\Omega$ is convex, the constant $C_{\Omega,s}$ depends only on $s$ and the dimension of $\Omega$, but not on the shape of $\Omega$, as indicated by the estimate (39). For nonconvex domains, the constant $C_{\Omega,s}$ does depend on the shape of $\Omega$ (Dupont and Scott, 1980, Verfürth, 1999).

Let $F$ be a bounded linear functional on $H^s(\Omega)$ with norm $\|F\|$ such that $F(p) = 0$ for all

$p \in P_n(\Omega)$. It follows from (46) that

$$
\begin{aligned}
|F(\zeta)| &\leq \inf_{p \in P_n(\Omega)} |F(\zeta - p)| \leq \|F\| \inf_{p \in P_n(\Omega)} \|\zeta - p\|_{H^s(\Omega)} \\
&\leq (C_{\Omega,s}\|F\|)|\zeta|_{H^s(\Omega)}
\end{aligned}
\tag{47}
$$

for all $\zeta \in H^s(\Omega)$. The estimate (47), known as the Bramble-Hilbert lemma (Bramble and Hilbert, 1970), is useful for deriving various error estimates.

### 3.4. Some discrete estimates

The finite element spaces in Section 3.2 are designed to be subspaces of Sobolev spaces so that they can serve as the trial/test spaces for Ritz-Galerkin methods. On the other hand, since finite element spaces are constructed by piecing together finite-dimensional function spaces, there are *discrete* estimates valid on the finite element spaces but not the Sobolev spaces.

Let $(K, \mathcal{P}_K, \mathcal{N}_K)$ be a finite element such that $\mathcal{P}_K \subset H^k(K)$ for a nonnegative integer $k$. Since any seminorm on a finite-dimensional space is continuous with respect to a norm, we have, by scaling, the following *inverse* estimate:

$$
|v|_{H^k(K)} \leq C(\operatorname{diam} K)^{\ell-k}\|v\|_{H^\ell(K)} \quad \forall\, v \in \mathcal{P}_K,\ 0 \leq \ell \leq k
\tag{48}
$$

where the positive constant $C$ depends on the domain $\widehat{K}$ (the image of $K$ under the scaling map $\mathcal{H}$ defined by (35)) and the space $\mathcal{P}_K$.

For finite elements whose shape functions can be pulled back to a fixed finite-dimensional function space on a reference element, the constant $C$ depends only on the shape regularity of the element domain $K$ and global versions of (48) can be easily derived. For example, for a quasi-uniform family $\{\mathcal{T}_i : i \in I\}$ of simplicial or quadrilateral triangulations of a polygonal domain $\Omega$, we have

$$
|v|_{H^1(\Omega)} \leq Ch_i^{-1}\|v\|_{L_2(\Omega)} \qquad \forall\, v \in V_i \quad \text{and} \quad i \in I
\tag{49}
$$

where $V_i \subset H^1(\Omega)$ is either the $P_n$ triangular finite element space or the $Q_n$ quadrilateral finite element space associated with $\mathcal{T}_i$. Note that $V_i \subset H^s(\Omega)$ for any $s < 3/2$ and a bit more work shows that the following inverse estimate (Ben Belgacem and Brenner, 2001) also holds:

$$
|v|_{H^s(\Omega)} \leq C_s h_i^{1-s}\|v\|_{H^1(\Omega)} \quad \forall\, v \in V_i,\ i \in I
\tag{50}
$$

and $1 \leq s < 3/2$, where the positive constant $C_s$ can be uniformly bounded for $s$ in a compact subset of $[1, 3/2)$.

It is well-known that in two dimensions the Sobolev space $H^1(\Omega)$ is not a subspace of $C(\overline{\Omega})$. However, the $P_n$ triangular finite element space and the $Q_n$ quadrilateral finite element space do belong to $C(\overline{\Omega})$ and it is possible to bound the $L_\infty$ norm of the finite element function by its $H^1$ norm. Indeed, it follows from Fourier transform and extension theorems (Adams, 1995, Wloka, 1987) that, for $\epsilon > 0$,

$$
\|v\|_{L_\infty(\Omega)} \leq C\epsilon^{-1/2}\|v\|_{H^{1+\epsilon}(\Omega)} \qquad \forall\, v \in H^{1+\epsilon}(\Omega)
\tag{51}
$$

By taking $\epsilon = (1 + |\ln h_i|)^{-1}$ in (51) and applying (50), we arrive at the following *discrete Sobolev inequality*:

$$\|v\|_{L_\infty(\Omega)} \le C(1 + |\ln h_i|)^{1/2}\|v\|_{H^1(\Omega)} \quad \forall\, v \in V_i \tag{52}$$

where the positive constant $C$ is independent of $i \in I$.

The discrete Sobolev inequality and the Poincaré-Friedrichs inequality (8) imply immediately the following *discrete Poincaré inequality*:

$$
\begin{aligned}
\|v\|_{L_\infty(\Omega)} &\le \|v - \bar{v}\|_{L_\infty(\Omega)} + \|\bar{v}\|_{L_\infty(\Omega)} \\
&\le 2\|v - \bar{v}\|_{L_\infty(\Omega)} \\
&\le C(1 + |\ln h_i|)^{1/2}\|v - \bar{v}\|_{H^1(\Omega)} \\
&\le C(1 + |\ln h_i|)^{1/2}|v|_{H^1(\Omega)}
\end{aligned}
\tag{53}
$$

for all $v \in V_i$ that vanishes at a given point in $\overline{\Omega}$ and with mean $\bar{v} = \int_\Omega v\,\mathrm{d}x/|\Omega|$.

**Remark** 11. The discrete Sobolev inequality can also be established directly using calculus and inverse estimates (Bramble, Pasciak and Schatz, 1986; Brenner and Scott, 2002), and both (52) and (53) are sharp (Brenner and Sung, 2000).

## 4. A Priori Error Estimates for Finite Element Methods

Let $\mathcal{T}$ be a triangulation of $\Omega$ and a finite element $(\overline{D}, \mathcal{P}_{\overline{D}}, \mathcal{N}_{\overline{D}})$ be associated with each subdomain $D \in \mathcal{T}$ so that the resulting finite element space $FE_\mathcal{T}$ (cf. (32)) is a subspace of $C^{m-1}(\overline{\Omega}) \subset H^m(\Omega)$. By imposing appropriate boundary conditions, we can obtain a subspace $V_\mathcal{T}$ of $FE_\mathcal{T}$ such that $V_\mathcal{T} \subset V$, the subspace of $H^m(\Omega)$, where the weak problem (1) is formulated. The corresponding finite element method for (1) is:

Find $u_\mathcal{T} \in V_\mathcal{T}$ such that

$$a(u_\mathcal{T}, v) = F(v) \qquad \forall\, v \in V_\mathcal{T} \tag{54}$$

In this section, we consider a priori estimates for the discretization error $u - u_\mathcal{T}$. We will discuss the second-order and fourth-order cases separately. We use the letter $C$ to denote a generic positive constant that can take different values at different appearances.

Let us also point out that the asymptotic error analysis carried out in this section is not sufficient for parameter-dependent problems (e.g. thin structures and nearly incompressible elasticity) that can experience *locking* (Babuška and Suri, 1992). We refer the readers to other chapters in this encyclopedia that are devoted to such problems for the discussion of the techniques that can overcome locking.

### 4.1. Second-order problems

We will devote most of our discussion to the case where $\Omega \subset \mathbb{R}^2$ and only comment briefly on the 3-D case. For preciseness, we also assume the right-hand side of the elliptic boundary value problem to be square integrable. We first consider the case where $V \subset H^1(\Omega)$ is defined by homogeneous Dirichlet boundary conditions (cf. Section 2.1) on $\Gamma \subset \partial\Omega$. Such problems can be discretized by triangular $P_n$ elements (Example 4) or quadrilateral $Q_n$ elements (Example 8).

Let $\mathcal{T}$ be a triangulation of $\Omega$ by triangles (convex quadrilaterals) and each triangle (quadrilateral) in $\mathcal{T}$ be equipped with the $P_n$ $(n \geq 1)$ Lagrange element ($Q_n$ quadrilateral element). The resulting finite element space $FE_{\mathcal{T}}$ is a subspace of $C^0(\overline{\Omega}) \subset H^1(\Omega)$. We assume that $\Gamma$ is the union of the edges of the triangles (quadrilaterals) in $\mathcal{T}$ and take $V_{\mathcal{T}} = V \cap FE_{\mathcal{T}}$, the subspace defined by the homogeneous Dirichlet boundary condition on $\Gamma$.

We know from the discussion in Section 2.3 that $u \in H^{1+\alpha(D)}(D)$ for each $D \in \mathcal{T}$, where the number $\alpha(D) \in (0,1]$ and $\alpha(D) = 1$ for $D$ away from the singular points. Hence, the element nodal interpolation operator $\Pi_{\overline{D}}$ is well-defined on $u$ for all $D \in \mathcal{T}$. We can therefore piece together a global nodal interpolant $\Pi_{\mathcal{T}}^N u \in V_{\mathcal{T}}$ by the formula

$$\left(\Pi_{\mathcal{T}}^N u\right)\big|_{\overline{D}} = \Pi_{\overline{D}}\left(u\big|_D\right) \tag{55}$$

From the discussion in Section 3.3, we know that (43) is valid for both the triangular $P_n$ element and the quadrilateral $Q_n$ element. We deduce from (43) and (55) that

$$\|u - \Pi_{\mathcal{T}}^N u\|_{H^1(\Omega)}^2 \leq C \sum_{D \in \mathcal{T}} (\operatorname{diam} D)^{2\alpha(D)} |u|_{H^{1+\alpha(D)}(D)}^2 \tag{56}$$

where $C$ depends only on the maximum of the aspect ratios of the element domains in $\mathcal{T}$. Combining (22) and (56) we have the a priori discretization error estimate

$$\|u - u_{\mathcal{T}}\|_{H^1(\Omega)} \leq C \left( \sum_{D \in \mathcal{T}} (\operatorname{diam} D)^{2\alpha(D)} |u|_{H^{1+\alpha(D)}(D)}^2 \right)^{1/2} \tag{57}$$

where $C$ depends only on the constants in (2) and (3) and the maximum of the aspect ratios of the element domains in $\mathcal{T}$.

Hence, if $\{\mathcal{T}_i : i \in I\}$ is a regular family of triangulations, and the solution $u$ of (1) belongs to the Sobolev space $H^{1+\alpha}(\Omega)$ for some $\alpha \in (0,1]$, then we can deduce from (57) that

$$\|u - u_{\mathcal{T}_i}\|_{H^1(\Omega)} \leq C h_i^\alpha |u|_{H^{1+\alpha}(\Omega)} \tag{58}$$

where $h_i = \max_{D \in \mathcal{T}_i} \operatorname{diam} D$ is the mesh size of $\mathcal{T}_i$ and $C$ is independent of $i \in I$. Note that the solution $w$ of (23) with $\tilde{u}$ replaced by $u_{\mathcal{T}}$ also belongs to $H^{1+\alpha}(\Omega)$ and satisfies the elliptic regularity estimate

$$\|w\|_{H^{1+\alpha}(\Omega)} \leq C \|u - u_{\mathcal{T}}\|_{L_2(\Omega)}$$

Therefore, we have

$$
\begin{aligned}
\inf_{v \in V_{\mathcal{T}}} \|w - v\|_{H^1(\Omega)} &\leq \|w - \Pi_{\mathcal{T}}^N w\|_{H^1(\Omega)} \\
&\leq Ch_i^{\alpha} |w|_{H^{1+\alpha}(\Omega)} \\
&\leq Ch_i^{\alpha} \|u - u_{\mathcal{T}}\|_{L_2(\Omega)}
\end{aligned} \tag{59}
$$

The abstract estimate (24) with $\tilde{u}$ replaced by $u_{\mathcal{T}}$ and (59) yield the following $L_2$ estimate:

$$
\|u - u_{\mathcal{T}_i}\|_{L_2(\Omega)} \leq Ch_i^{2\alpha} |u|_{H^{1+\alpha}(\Omega)} \tag{60}
$$

where $C$ is also independent of $i \in I$.

**Remark** 12. In the case where $\alpha = 1$ (for example, when $\Gamma = \partial\Omega$ in Example 1 and $\Omega$ is convex), the estimate (58) is optimal and it is appropriate to use a quasi-uniform family of triangulations. In the case where $\alpha(D) < 1$ for $D$ next to singular points, the estimate (57) allows the possibility of improvement by graded meshes (cf. Section 6).

**Remark** 13. In the derivations of (58) and (60) above for the triangular $P_n$ elements, we have used the minimum angle condition (cf. Remark 5.). In view of the anisotropic estimates (45), these estimates also hold for triangular $P_n$ elements under the *maximum angle condition* (Babuška and Aziz, 1976; Jamet, 1976; Ženišek, 1995; Apel, 1999): there exists $\theta_* < \pi$ such that all the angles in the family of triangulations are $\leq \theta_*$. The estimates (58) and (60) are also valid for $Q_n$ elements on parallelograms satisfying the maximum angle condition. They can also be established for certain *thin* quadrilateral elements (Apel, 1999).

The 2-D results above also hold for 3-D tetrahedral $P_n$ elements and 3-D hexagonal $Q_n$ elements if the solution $u$ of (1) belongs to $H^{1+\alpha}(\Omega)$ where $1/2 < \alpha \leq 1$, since the nodal interpolation operator are then well-defined by the Sobolev embedding theorem. This is the case, for example, if $\Gamma = \partial\Omega$ in Example 1. However, new interpolation operators that require less regularity are needed if $0 < \alpha \leq 3/2$. Below, we construct an interpolation operator $\Pi_{\mathcal{T}}^A \colon H^1(\Omega) \to V_{\mathcal{T}}$ using the local averaging technique of Scott and Zhang (1990).
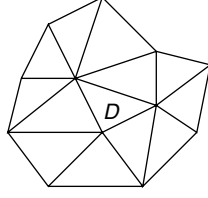
For simplicity, we take $V_{\mathcal{T}}$ to be a tetrahedral $P_1$ finite element space. Therefore, we only need to specify the value of $\Pi_{\mathcal{T}}^A \zeta$ at the vertices of $\mathcal{T}$ for a given function $\zeta \in H^1(\Omega)$. Let $p$ be a vertex. We choose a face (or edge in 2-D) $\mathcal{F}$ of a subdomain in $\mathcal{T}$ such that $p \in \overline{\mathcal{F}}$. The choice of $\mathcal{F}$ is of course not unique. But we always choose $\mathcal{F} \subset \partial\Omega$ if $p \in \partial\Omega$ so that the resulting interpolant will satisfy the appropriate Dirichlet boundary condition. Let $\{\psi_j\}_{j=1}^d \subset P_1(\mathcal{F})$ be biorthogonal to the nodal basis $\{\phi_j\}_{j=1}^d \subset P_1(\mathcal{F})$ with respect to the $L_2(\mathcal{F})$ inner product. In other words $\phi_j$ equals 1 at the $j$th vertex of $\mathcal{F}$ and vanishes at the other vertices, and

$$
\int_{\mathcal{F}} \psi_i \phi_j \, \mathrm{d}s = \delta_{ij} \tag{61}
$$

Suppose $p$ corresponds to the $j$th vertex of $\mathcal{F}$. We then define

$$
(\Pi_{\mathcal{T}}^A \zeta)(p) = \int_{\mathcal{F}} \psi_j \zeta \, \mathrm{d}s \tag{62}
$$

where the integral is well-defined because of the trace theorem.

Figure 9: A two-dimensional example of $S(D)$.

It is clear in view of (61) and (62) that $\Pi_\mathcal{T}^A v = v$ for all $v \in FE_\mathcal{T}$ and $\Pi_\mathcal{T}^A \zeta = 0$ on $\Gamma$ if $\zeta = 0$ on $\Gamma$. Note also that $\Pi_\mathcal{T}^A$ is not a local operator, i.e., $(\Pi_\mathcal{T}^A \zeta)\big|_D$ is in general determined by $\zeta\big|_{S(D)}$, where $S(D)$ is the polyhedral domain formed by the subdomains in $\mathcal{T}$ sharing (at least) a vertex with $D$ (cf. Figure 9 for a 2-D example). It follows that the interpolation error estimate for $\widetilde{\Pi}_\mathcal{T}$ takes the following form:

$$
\begin{aligned}
\|\zeta - \Pi_\mathcal{T}^A \zeta\|_{L_2(D)}^2 &+ (\operatorname{diam} D)^2 |\zeta - \Pi_\mathcal{T}^A \zeta|_{H^1(D)}^2 \\
&\leq C \big(\operatorname{diam} D\big)^{2(1+\alpha(S(D)))} |\zeta|_{H^{1+\alpha(S(D))}(S(D))}^2
\end{aligned}
\tag{63}
$$

where $C$ depends on the shape regularity of $\mathcal{T}$, provided that $\zeta \in H^{1+\alpha(S(D))}(S(D))$ for some $\alpha(S(D)) \in (0,1]$. The estimates (58) and (60) for tetrahedral $P_1$ elements can be derived for general $\alpha \in (0,1]$ and regular triangulations using the estimate (63).

**Remark** 14. The interpolation operator $\Pi_\mathcal{T}^A$ can be defined for general finite elements (Scott and Zhang, 1990; Girault and Scott, 2002) and anisotropic estimates can be obtained for $\Pi_\mathcal{T}^A$ f or certain triangulations (Apel, 1999). There also exist interpolation operators for less regular functions (Clément, 1975; Bernardi and Girault, 1998).

Next, we consider the case where $V$ is a closed subspace of $H^1(\Omega)$ with finite codimension $n < \infty$, as in the case of the Poisson problem with pure homogeneous Neumann boundary condition (where $n = 1$) or the elasticity problem with pure homogeneous traction boundary condition (where $n = 1$ when $d = 1$, and $n = 3(d-1)$ when $d = 2$ or 3). The key assumption here is that there exists a bounded linear projection $Q$ from $H^1(\Omega)$ onto an $n$ dimensional subspace of $FE_\mathcal{T}$ such that $\zeta \in H^1(\Omega)$ belongs to $V$ if and only if $Q\zeta = 0$. We can then define an interpolation operator $\widetilde{\Pi}_\mathcal{T}$ from appropriate Sobolev spaces onto $V_\mathcal{T}$ by

$$
\widetilde{\Pi}_\mathcal{T} = (I - Q)\Pi_\mathcal{T}
$$

where $\Pi_\mathcal{T}$ is either the nodal interpolation operator $\Pi_\mathcal{T}^N$ or the Scott-Zhang averaging interpolation operator $\Pi_\mathcal{T}^A$ introduced earlier. Observe that, since the weak solution $u$ belongs to $V$,

$$
u - \widetilde{\Pi}_\mathcal{T} u = u - Qu - (I - Q)\Pi_\mathcal{T} u = (I - Q)(u - \Pi_\mathcal{T} u)
$$

and the interpolation error of $\widetilde{\Pi}_\mathcal{T}$ can be estimated in terms of the norm of $Q \colon H^1(\Omega) \to H^1(\Omega)$ and the interpolation error of $\Pi_\mathcal{T}$. Therefore, the a priori discretization error estimates for Dirichlet or Dirichlet/Neumann boundary value problems also hold for this second type of (pure Neumann) boundary value problems.

For the Poisson problem with homogeneous Neumann boundary condition, we can take

$$Q\zeta = \frac{1}{|\Omega|} \int_{\Omega} \zeta \mathrm{dx}$$

the mean of $\zeta$ over $\Omega$. For the elasticity problem with pure homogeneous traction boundary condition, the operator $Q$ from $[H^1(\Omega)]^d$ onto the space of infinitesimal rigid motions is defined by

$$\int_{\Omega} Q\boldsymbol{\zeta}\mathrm{dx} = \int_{\Omega} \boldsymbol{\zeta}\mathrm{dx} \quad \text{and}$$

$$\int_{\Omega} \nabla \times Q\boldsymbol{\zeta}\mathrm{dx} = \int_{\Omega} \nabla \times \boldsymbol{\zeta}\mathrm{dx} \qquad \forall \boldsymbol{\zeta} \in [H^1(\Omega)]^d$$

In both cases, the norm of $Q$ is bounded by a constant $C_{\Omega}$.

**Remark** 15. In the case where $f \in H^k(\Omega)$ for $k > 0$, the solution $u$ belongs to $H^{2+k}(\Omega)$ away from the geometric or boundary data singularities and, in particular, away from $\partial\Omega$. Therefore, it is advantageous to use higher-order elements in certain parts of $\Omega$, or even globally (with curved elements near $\partial\Omega$) if singularities are not present. In the case where $f \in L_2(\Omega)$, the error estimate (57) indicates that the order of the discretization error for the triangular $P_n$ element or the quadrilateral $Q_n$ element is independent of $n \geq 1$. However, the convergence of the finite element solutions to a particular solution as $h \downarrow 0$ can be improved by using higher-order elements because of the existence of *nonuniform error estimates* (Babuška and Kellogg, 1975).

### 4.2. Fourth-order problems

We restrict the discussion of fourth-order problems to the two-dimensional plate bending problem of Example 3.

Let $\mathcal{T}$ be a triangulation of $\Omega$ by triangles and each triangle in $\mathcal{T}$ be equipped with the Hsieh-Clough-Tocher macro element (cf. Example 6). The finite element space $FE_{\mathcal{T}}$ defined by (32) is a subspace of $C^1(\overline{\Omega}) \subset H^2(\Omega)$. We take $V_{\mathcal{T}}$ to be $V \cap FE_{\mathcal{T}}$, where $V = H_0^2(\Omega)$ for the clamped plate and $V = H_0^1(\Omega) \cap H^2(\Omega)$ for the simply supported plate.

The solution $u$ of the plate-bending problem belongs to $H^{2+\alpha(D)}(D)$ for each $D \in \mathcal{T}$, where $\alpha(D) \in (0, 2]$ and $\alpha(D) = 2$ for $D$ away from the corners of $\Omega$. The elemental nodal interpolation operator $\Pi_{\overline{D}}$ is well-defined on $u$ for all $D \in \mathcal{T}$. We can therefore define a global nodal interpolation operator $\Pi_{\mathcal{T}}^N$ by the formula (55). Since the Hsieh-Clough-Tocher macro element is affine-interpolation-equivalent to the reference element on the standard simplex, we deduce from (55) and (43) that

$$\|u - \Pi_{\mathcal{T}}^N u\|_{H^2(\Omega)}^2 \leq C \sum_{D \in \mathcal{T}} (\mathrm{diam}\, D)^{2\alpha(D)} |u|_{H^{2+\alpha(D)}(D)}^2 \tag{64}$$

where $C$ depends only on the maximum of the aspect ratios of the triangles in $\mathcal{T}$ (or equivalently

the minimum angle of $\mathcal{T}$). From (22) and (64), we have

$$\|u - u_{\mathcal{T}}\|_{H^2(\Omega)} \leq C \left( \sum_{D \in \mathcal{T}} (\operatorname{diam} D)^{2\alpha(D)} |u|^2_{H^{2+\alpha(D)}(D)} \right)^{1/2} \tag{65}$$

where $C$ depends only on the constants in (2) and (3) and the minimum angle of $\mathcal{T}$.

Hence, if $\{\mathcal{T}_i : i \in I\}$ is a regular family of triangulations, and the solution $u$ of the plate bending problem belongs to the Sobolev space $H^{2+\alpha}(\Omega)$ for some $\alpha \in (0, 2]$, we can deduce from (65) that

$$\|u - u_{\mathcal{T}_i}\|_{H^2(\Omega)} \leq C h_i^\alpha |u|_{H^{2+\alpha}(\Omega)} \tag{66}$$

where $h_i = \max_{D \in \mathcal{T}_i} \operatorname{diam} D$ is the mesh size of $\mathcal{T}_i$ and $C$ is independent of $i \in I$. Since the solution $w$ of (23) also belongs to $H^{2+\alpha}(\Omega)$, the abstract estimate (24) combined with an error estimate for $w$ in the $H^2$-norm analogous to (66) yields the following $L_2$ estimate:

$$\|u - u_{\mathcal{T}_i}\|_{L_2(\Omega)} \leq C h_i^{2\alpha} |u|_{H^{2+\alpha}(\Omega)} \tag{67}$$

where $C$ is also independent of $i \in I$.

**Remark** 16. The analysis of general triangular and quadrilateral $C^1$ macro elements can be found in Douglas *et al* (1979).

The plate-bending problem can also be discretized by the Argyris element (cf. Example 5). If $\alpha(D) > 1$ for all $D \in \mathcal{D}$, then the nodal interpolation operator $\Pi_{\mathcal{T}}^N$ is well-defined for the Argyris finite element space. If $\alpha(D) \leq 1$ for some $D \in \mathcal{T}$, then the nodal interpolation operator $\Pi_{\mathcal{T}}^N$ must be replaced by an interpolation operator constructed by the technique of local averaging. In either case, the estimates (66) and (67) remain valid for the Argyris finite element solution.

## 5. A Posteriori Error Estimates and Analysis

In this section, we review explicit and implicit estimators as well as averaging and multilevel estimators for a posteriori finite element error control.

Throughout this section, we adopt the notation of Sections 2.1 and 2.2 and recall that $u$ denotes the (unknown) exact solution of (1) while $\tilde{u} \in \widetilde{V}$ denotes the discrete and given solution of (19). It is the aim of Section 5.1-5.6 to estimate the error $e := u - \tilde{u} \in V$ in the energy norm $\| \cdot \|_a = (a(\cdot, \cdot))^{1/2}$ in terms of computable quantities while Section 5.7 concerns other error norms or goal functionals.

Throughout this section, we assume $0 < \|e\|_a$ to exclude the exceptional situation $u = \tilde{u}$.

*5.1. Aims and concepts in a posteriori finite element error control*

The following five sections introduce the notation, the concepts of efficiency and reliability, the definitions of residual and error, a posteriori error control and adaptive algorithms, and comment on some relevant literature.

*5.1.1. Error estimators, efficiency, reliability, asymptotic exactness.* Regarded as an approximation to the (unknown) error norm $\|e\|_a$, a (computable) quantity $\eta$ is called *a posteriori error estimator*, or *estimator* for brevity, if it is a function of the known domain $\Omega$ and its boundary $\Gamma$, the quantities of the right-hand side $F$, cf. (6) and (13), as well as of the (given) discrete solution $\tilde{u}$, or the underlying triangulation.

An estimator $\eta$ is called *reliable* if

$$\|e\|_a \le \mathrm{C_{rel}}\, \eta + \mathrm{h.o.t._{rel}} \tag{68}$$

An estimator $\eta$ is called *efficient* if

$$\eta \le \mathrm{C_{eff}}\, \|e\|_a + \mathrm{h.o.t._{eff}} \tag{69}$$

An estimator is called *asymptotically exact* if it is reliable and efficient in the sense of (68)-(69) with $\mathrm{C_{rel}} = \mathrm{C_{eff}^{-1}}$.

Here, $\mathrm{C_{rel}}$ and $\mathrm{C_{eff}}$ are multiplicative constants that do not depend on the mesh size of an underlying finite element mesh $\mathcal{T}$ for the computation of $\tilde{u}$ and h.o.t. denotes higher-order terms. The latter are generically much smaller than $\eta$ or $\|e\|_a$, but usually, this depends on the (unknown) smoothness of the exact solution or the (known) smoothness of the given data. The readers are warned that, in general, h.o.t. may not be neglected; in case of high oscillations they may even dominate (68) or (69).

*5.1.2. Error and residual.*   Abstract examples for estimators are (26) and (27), which involve dual norms of the residual (25). Notice carefully that $R := F - a(\tilde{u}, \cdot)$ is a bounded linear functional in $V$, written $R \in V^*$, and hence the dual norm

$$\|R\|_{V^*} := \sup_{v \in V \setminus \{0\}} \frac{R(v)}{\|v\|_a} = \sup_{v \in V \setminus \{0\}} \frac{a(e, v)}{\|v\|_a} = \|e\|_a < \infty \tag{70}$$

The second equality immediately follows from (25). A Cauchy inequality in (70) with respect to the scalar product $a$ results in $\|R\|_{V^*} \le \|e\|_a$ while $v = e$ in (70) yields finally the equality $\|R\|_{V^*} = \|e\|_a$.

That is, the error (estimation) in the energy norm *is equivalent* to the (computation of the) dual norm of the given residual. Furthermore, it is even of *comparable computational effort* to

compute an optimal $v = e$ in (70) or to compute $e$. The proof of (70) yields even a stability estimate: The relative error of $R(v)$ as an approximation to $\|e\|_a$ equals

$$\frac{(\|e\|_a - R(v))}{\|e\|_a} \quad = \quad \frac{1}{2} \left\| v - \frac{e}{\|e\|_a} \right\|_a^2$$
$$\text{for all} \quad v \in V \quad \text{with} \quad \|v\|_a = 1 \tag{71}$$

In fact, given any $v \in V$ with $\|v\|_a = 1$, the identity (71) follows from

$$1 - a\left(\frac{e}{\|e\|_a}, v\right) \quad = \quad \frac{1}{2}\, a\left(\frac{e}{\|e\|_a}, \frac{e}{\|e\|_a}\right) - a\left(\frac{e}{\|e\|_a}, v\right)$$
$$+ \frac{1}{2}\, a(v, v) = \frac{1}{2} \left\| v - \frac{e}{\|e\|_a} \right\|_a^2$$

The error estimate (71) implies that the maximizing $v$ in (70) (i.e. $v \in V$ with maximal $R(v)$ subject to $\|v\|_a \leq 1$) is unique and equals $e/\|e\|_a$. As a consequence, the computation of the maximizing $v$ in (70) is equivalent to and indeed equally expensive as the computation of the unknown $e/\|e\|_a$ and so (since $\tilde{u}$ is known) of the exact solution $u$. Therefore, a posteriori error analysis aims to compute lower and upper bounds of $\|R\|_{V^*}$ rather than its exact value.

*5.1.3. Error estimators and error control.* For an idealized termination procedure, one is given a tolerance Tol > 0 and interested in a stopping criterion (of successively adapted mesh refinements)

$$\|e\|_a \leq \text{Tol}$$

Since the error $\|e\|_a$ is unknown, it is replaced by its upper bound (68) and then leads to

$$C_{\text{rel}}\eta + \text{h.o.t.}_{\text{rel}} \leq \text{Tol} \tag{72}$$

For a verification of (72), in practice, one requires not only $\eta$ but also $C_{\text{rel}}$ and h.o.t.$_{\text{rel}}$. The later quantity cannot be dropped; it is not sufficient to know that h.o.t.$_{\text{rel}}$ is (possibly) negligible for sufficient small mesh-sizes.

Section 5.6 presents numerical examples and further discussions of this aspect.

*5.1.4. Adaptive mesh-refining algorithms.* Error estimators are used in adaptive mesh-refining algorithms to motivate a *refinement rule*, which determines whether an element or edge and so on shall be refined or coarsened. This will be discussed in Section 6 below.

At this stage two remarks are in order. First, one should be precise in the language and distinguish between error estimators, which are usually global and fully involve constants and higher-order terms and (local) refinement indicators used in refinement rules. Second, constants and higher-order terms might be seen as less important and are often omitted in the usage as refinement indicators for the step MARK in Section 6.2.

*5.1.5. Literature.*   Amongst the most influential pioneering publications on a posteriori error control are Babuška and Rheinboldt (1978), Ladeveze and Leguillon (1983), Bank and Weiser (1985), Babuška and Miller (1987), Eriksson and Johnson (1991), followed by many others. The readers may find it rewarding to study the survey articles of Eriksson *et al* (1995), Becker and Rannacher (2001) and the books of Verfürth (1996), Ainsworth and Oden (2000), Babuška and Strouboulis (2001), Bangerth and Rannacher (2003), Repin (2008), and Verfürth (2013) for a first insight and further references.

## 5.2. Explicit residual-based error estimators

The most frequently considered and possibly easiest class of error estimators consists of local norms of explicitly given volume and jump residuals multiplied by mesh-depending weights.

To derive them for a general class of abstract problems from Section 2.1, let $u \in V$ be an exact solution of the problem (1) and let $\tilde{u} \in \widetilde{V}$ be its Galerkin approximation from (19) with residual $R(v)$ from (25). Moreover, as in Example 1 or 2, it is supposed throughout this chapter that the strong form of the equilibration associated with the weak form (19) is of the form

$$-\mathrm{div}p = f \quad \text{for some flux or stress} \quad p \in L^2(\Omega; \mathbb{R}^{m \times n})$$

The discrete analog $\tilde{p}$ is piecewise smooth but, in general, discontinuous; at several places below, it is a $\mathcal{T}$ piecewise constant $m \times n$ matrix as it is proportional to the gradient of some (piecewise) $P_1$ FE function $\tilde{u}$. The description of the residuals is based on the weak form of $f + \mathrm{div}p = 0$.

*5.2.1. Residual representation formula.*   It is the aim of this section to recast the residual in the form

$$R(v) = \sum_{T \in \mathcal{T}} \int_T r_T \cdot v \mathrm{dx} - \sum_{E \in \mathcal{E}} \int_E r_E \cdot v \mathrm{ds} \tag{73}$$

of a sum of integrals over all element domains $T \in \mathcal{T}$ plus a sum of integrals over all edges or faces $E \in \mathcal{E}$ and to identify the explicit volume residual $r_T$ and the jump residual $r_E$.

The boundary $\partial T$ of each finite element domain $T \in \mathcal{T}$ is a union of edges or faces, which form the set $\mathcal{E}(T)$, written $\partial T = \cup \mathcal{E}(T)$. Each edge or face $E \in \mathcal{E}$ in the set of all possible edges or faces $\mathcal{E} = \cup \{\mathcal{E}(T) : T \in \mathcal{T}\}$ is associated with a unit normal vector $\nu_E$, which is unique up to an orientation $\pm\nu_E$, which is globally f ixed. By convention, the unit normal $\nu$ on the domain $\Omega$ or on an element $T$ points outwards.

For the ease of exploration, suppose that the underlying boundary value problem allows the bilinear form $a(\tilde{u}, v)$ to equal the sum over all $\int_T \tilde{p}_{jk} D_j v_k \mathrm{dx}$ with given fluxes or stresses $\tilde{p}_{jk}$. Moreover, Neumann data are excluded from the description in this section and hence only interior edges contribute with a jump residual. An integration by parts on $T$ with outer unit

normal $\nu_T$ yields

$$\int_T \tilde{p}_{jk}\, D_j v_k \mathrm{dx} = \int_{\partial T} \tilde{p}_{jk}\, v_k\, \nu_{T,j}\mathrm{ds} - \int_T v_k\, D_j \tilde{p}_{jk}\mathrm{dx}$$

which, with the divergence operator div and proper evaluation of $\tilde{p}\nu$, reads

$$a(\tilde{u}, v) + \sum_{T\in\mathcal{T}} \int_T v \cdot \operatorname{div}\tilde{p}\mathrm{dx} = \sum_{T\in\mathcal{T}} \int_{\partial T} (\tilde{p}\nu)\cdot v \mathrm{ds}$$

Each boundary $\partial T$ is rewritten as a sum of edges or faces. Each such edge or face $E$ belongs either to the boundary $\partial\Omega$, written $E\in\mathcal{E}_{\partial\Omega}$, or is an interior edge, written $E\in\mathcal{E}_\Omega$. For $E\in\mathcal{E}_{\partial\Omega}$ there exists exactly one element $T$ with $E\in\mathcal{E}(T)$ and one defines $T_+ = T$, $T_- = E \subset \partial\Omega$, $\omega_E = \operatorname{int}(T)$ and $\nu_E := \nu_T = \nu_\Omega$. Any $E\in\mathcal{E}(\Omega)$ is the intersection of exactly two elements, which we name $T_+$ and $T_-$ and which essentially determine the patch $\omega_E := \operatorname{int}(T_+\cup T_-)$ of $E$. This description of $T_\pm$ is unique up to the order that is fixed in the sequel by the convention that $\nu_E = \nu_{T_+}$ is exterior to $T_+$. Then,

$$\sum_{T\in\mathcal{T}} \int_{\partial T} (\tilde{p}\nu)\cdot v\mathrm{ds} = \sum_{E\in\mathcal{E}(\Omega)} \int_E [\tilde{p}\nu_E]\cdot v\mathrm{ds}$$

where $[\tilde{p}\nu_E] := (\tilde{p}|_{T_+} - \tilde{p}|_{T_-})\nu_E$ for $E = \partial T_+ \cap \partial T_- \in \mathcal{E}(\Omega)$ and $[\tilde{p}\nu_E] := 0$ for $E\in\mathcal{E}(T)\cap\mathcal{E}_{\partial\Omega}$. Altogether, one obtains the error residual error representation formula (73) with the

$$
\begin{aligned}
\textit{volume residuals} \quad r_T &:= \quad f + \operatorname{div}\tilde{p} \quad &&\text{in } T\in\mathcal{T} \\
\textit{jump residuals} \quad r_E &:= \quad [\tilde{p}\nu_E] \quad &&\text{along } E\in\mathcal{E}(\Omega)
\end{aligned}
$$

*5.2.2. Weak approximation operators.* In terms of the residual $R$, the orthogonality condition (20) is rewritten as $R(\tilde{v}) = 0$ for all $\tilde{v}\in\widetilde{V}$. Hence, given any $v\in V$ with norm $\|v\|_a = 1$, there holds $R(v) = R(v - \tilde{v})$.

Explicit error estimators rely on the design of $\tilde{v} := \Pi^A_\mathcal{T}(v)$ as a function of $v$, $\Pi^A_\mathcal{T}$ is called *approximation operator* as in (61)-(63) and discussed further in Section 4. See also (Carstensen, 1999; Carstensen and Funken, 2000; Nochetto and Wahlbin, 2002). For the understanding of this section, it suffices to know that there are several choices of $\tilde{v}\in\widetilde{V}$ that satisfy first-order approximation and stability properties in the sense of

$$
\begin{aligned}
\sum_{T\in\mathcal{T}} \|h_T^{-1}(v - \tilde{v})\|^2_{L_2(T)} \quad &+ \quad \sum_{E\in\mathcal{E}(\Omega)} \|h_E^{-1/2}(v - \tilde{v})\|^2_{L_2(E)} \\
&+ \quad |v - \tilde{v}|^2_{H^1(\Omega)} \le C\, |v|^2_{H^1(\Omega)}
\end{aligned}
\tag{74}
$$

Here, $h_T$ and $h_E$ denotes the diameter of an element $T\in\mathcal{T}$ and an edge $E\in\mathcal{E}$, respectively. The multiplicative constant $C$ is independent of the mesh-sizes $h_T$ or $h_E$, but depends on the shape of the element domains through their minimal angle condition (for simplices) or aspect ratio (for tensor product elements).

*5.2.3. Reliability.*    Given the explicit volume and jump residuals $r_T$ and $r_E$ in (73), one defines the *explicit residual-based estimator* $\eta_{R,R}$,

$$\eta_{R,R}^2 := \sum_{T \in \mathcal{T}} h_T^2 \left\| r_T \right\|_{L_2(T)}^2 + \sum_{E \in \mathcal{E}(\Omega)} h_E \left\| r_E \right\|_{L_2(E)}^2 \tag{75}$$

which is reliable, that is

$$\left\| e \right\|_a \le C \, \eta_{R,R} \tag{76}$$

The proof of (76) follows from (73)-(75) and Cauchy-Schwarz inequalities:

$$
\begin{aligned}
R(v) \;\; = \;\; & R(v - \tilde{v}) = \sum_{T \in \mathcal{T}} \int_T r_T \cdot (v - \tilde{v}) \mathrm{d}x \\
& - \sum_{E \in \mathcal{E}(\Omega)} \int_E r_E \cdot (v - \tilde{v}) \mathrm{d}s \\
\le \;\; & \sum_{T \in \mathcal{T}} (h_T \left\| r_T \right\|_{L_2(T)})(h_T^{-1} \left\| v - \tilde{v} \right\|_{L_2(T)}) \\
& + \sum_{E \in \mathcal{E}(\Omega)} (h_E^{1/2} \left\| r_E \right\|_{L_2(E)})(h_E^{-1/2} \left\| v - \tilde{v} \right\|_{L_2(E)}) \\
\le \;\; & \left( \sum_{T \in \mathcal{T}} h_T^2 \left\| r_T \right\|_{L_2(T)}^2 \right)^{1/2} \left( \sum_{T \in \mathcal{T}} h_T^{-2} \left\| v - \tilde{v} \right\|_{L_2(T)}^2 \right)^{1/2} \\
& + \left( \sum_{E \in \mathcal{E}(\Omega)} h_E \left\| r_E \right\|_{L_2(E)}^2 \right)^{1/2} \left( \sum_{E \in \mathcal{E}(\Omega)} h_E^{-1} \left\| v - \tilde{v} \right\|_{L_2(E)}^2 \right)^{1/2} \\
\le \;\; & C \, \eta_{R,R} \, |v|_{H^1(\Omega)}
\end{aligned}
$$

For first-order finite element methods in the situation of Example 1 or 2, the volume term $r_T = f$ can be substituted by the higher-order term of oscillations, that is

$$\left\| e \right\|_a^2 \le C \left( \mathrm{osc}(f)^2 + \sum_{E \in \mathcal{E}(\Omega)} h_E \left\| r_E \right\|_{L_2(E)}^2 \right) \tag{77}$$

For each node $z \in \mathcal{N}$ with nodal basis function $\varphi_z$ and patch $\omega_z := \{ x \in \Omega : \varphi_z(x) \ne 0 \}$ of diameter $h_z$ and the source term $f \in L_2(\Omega)^m$ with integral mean $f_z := |\omega_z|^{-1} \int_{\omega_z} f(x) \mathrm{d}x \in \mathbb{R}^m$, the oscillations of $f$ are defined by

$$\mathrm{osc}(f) := \left( \sum_{z \in \mathcal{N}} h_z^2 \left\| f - f_z \right\|_{L_2(\omega_z)}^2 \right)^{1/2}$$

Notice for $f \in H^1(\Omega)^m$ and the mesh size $h_{\mathcal{T}} \in P_0(\mathcal{T})$ there holds

$$\mathrm{osc}(f) \le C \left\| h_{\mathcal{T}}^2 \, Df \right\|_{L_2(\Omega)}$$

and so $\mathrm{osc}(f)$ is of quadratic and hence of higher-order. We refer to Carstensen and Verfürth (1999), Nochetto (1993), Becker and Rannacher (1996), and Rodriguez (1994b) for further details on and proofs of (77).

*5.2.4. Efficiency.* Following a technique with inverse estimates due to Verfürth (1996), this section investigates the proof of efficiency of $\eta_{R,R}$ in a local form, namely,

$$h_T\,\|r_T\|_{L_2(T)} \quad \leq \quad C\left(\|e\|_{H^1(T)} + \operatorname{osc}(f,T)\right) \tag{78}$$

$$h_E^{1/2}\,\|r_E\|_{L_2(E)} \quad \leq \quad C\left(\|e\|_{H^1(\omega_E)} + \operatorname{osc}(f,\omega_E)\right) \tag{79}$$

where $\tilde{f}$ denotes an elementwise polynomial (best-) approximation of $f$ and

$$\operatorname{osc}(f,T) \quad := \quad h_T\,\|f - \tilde{f}\|_{L_2(T)} \tag{80}$$

and

$$\operatorname{osc}(f,\omega_E) \quad := \quad h_E\,\|f - \tilde{f}\|_{L_2(\omega_E)} \tag{81}$$

The main tools in the proof of (79) and (78) are bubble functions $b_E$ and $b_T$ based on an edge or face $E \in \mathcal{E}$ and an element $T \in \mathcal{T}$ with nodes $\mathcal{N}(E)$ and $\mathcal{N}(T)$, respectively. Given a nodal basis $(\varphi_z : z \in \mathcal{N})$ of a first-order finite element method with respect to $\mathcal{T}$ define, for any $T \in \mathcal{T}$ and $E \in \mathcal{E}(\Omega)$, the element- and edge-bubble functions

$$b_T := \prod_{z \in \mathcal{N}(T)} \varphi_z \in H_0^1(T) \quad \text{and} \quad b_E := \prod_{z \in \mathcal{N}(E)} \varphi_z \in H_0^1(\omega_E) \tag{82}$$

$b_E$ and $b_T$ are nonnegative and continuous piecewise polynomials $\leq 1$ with support $\operatorname{supp} b_E = \overline{\omega_E} = T_+ \cup T_-$ (for $T_\pm \in \mathcal{T}$ with $E = T_+ \cap T_-$) and $\operatorname{supp} b_T = T$.

Utilizing the bubble functions (82), the proof of (78)-(79) essentially consists in the design of test functions $w_T \in H_0^1(T)$, $T \in \mathcal{T}$, and $w_E \in H_0^1(\omega_E)$, $E \in \mathcal{E}(\Omega)$, with the properties

$$|w_T|_{H^1(T)} \quad \leq \quad Ch_T\,\|r_T\|_{L_2(T)}$$

and

$$|w_E|_{H^1(\omega_E)} \quad \leq \quad Ch_E^{1/2}\,\|r_E\|_{L_2(E)} \tag{83}$$

$$h_T^2\,\|r_T\|_{L_2(T)}^2 \quad \leq \quad C_1\,R(w_T) + C_2\operatorname{osc}(f,T)^2 \tag{84}$$

$$h_E\,\|r_E\|_{L_2(E)}^2 \quad \leq \quad C_1\,R(w_E) + C_2\operatorname{osc}(f,\omega_E)^2 \tag{85}$$

In fact, (83)-(85), the definition of the residual $R = a(e,\cdot)$, and Cauchy-Schwarz inequalities with respect to the scalar product $a$ prove (78)-(79).

To construct the test function $w_T$, $T \in \mathcal{T}$, recall $\operatorname{div} p + f = 0$ and $r_T = f + \operatorname{div} \tilde{p}$ and set $\tilde{r}_T := \tilde{f} + \operatorname{div} \tilde{p}$ for some polynomial $\tilde{f}$ on $T$ such that $\tilde{r}_T$ is a best approximation of $r_T$ in some finite-dimensional (polynomial) space with respect to $L_2(T)$. Since

$$h_T\,\|\tilde{r}_T\|_{L_2(T)} \quad \leq \quad h_T\,\|r_T\|_{L_2(T)} \leq h_T\,\|\tilde{r}_T\|_{L_2(T)} + h_T\,\|f - \tilde{f}\|_{L_2(T)}$$

it remains to bound $\tilde{r}_T$, which belongs to a finite-dimensional space and hence satisfies an inverse inequality

$$h_T\,\|\tilde{r}_T\|_{L_2(T)} \leq Ch_T\,\|b_T^{1/2}\tilde{r}_T\|_{L_2(T)}$$

This motivates the estimation of

$$
\begin{aligned}
\|b_T^{1/2}\tilde{r}_T\|_{L_2(T)}^2 &= \int_T b_T\tilde{r}_T \cdot (\tilde{r}_T - r_T)\mathrm{dx} + \int_T \mathrm{b_T}\tilde{r}_T \cdot \mathrm{r_T}\mathrm{dx} \\
&\leq \|b_T^{1/2}\tilde{r}_T\|_{L_2(T)}\|b_T^{1/2}(f - \tilde{f})\|_{L_2(T)} \\
&\quad + \int_T b_T\tilde{r}_T \cdot \mathrm{div}\,(\tilde{p} - p)\mathrm{dx}
\end{aligned}
$$

The combination of the preceding estimates results in

$$
\begin{aligned}
h_T^2\,\|r_T\|_{L_2(T)}^2 &\leq C_1 \int_T (h_T^2 b_T\tilde{r}_T) \cdot \mathrm{div}\,(\tilde{p} - p)\mathrm{dx} \\
&\quad + C_2\,\mathrm{osc}(f,T)^2
\end{aligned}
$$

An integration by parts concludes the proof of (84) for

$$
w_T := h_T^2 b_T\tilde{r}_T \tag{86}
$$

the proof of (83) for this $w_T$ is immediate.

Given an interior edge $E = T_+ \cap T_- \in \mathcal{E}(\Omega)$ with its neighboring elements $T_+$ and $T_-$, simultaneously addressed as $T_\pm \in \mathcal{T}$, extend the edge residual $r_E$ from the edge $E$ to its patch $\omega_E = \mathrm{int}(T_+ \cup T_-)$ such that

$$
\begin{aligned}
\|b_E r_E\|_{L_2(\omega_E)} + h_E|b_E r_E|_{H^1(\omega_E)} &\leq C_1 h_E^{1/2}\,\|r_E\|_{L_2(E)} \\
&\leq C_2 h_E^{1/2}\,\|b_E^{1/2}r_E\|_{L_2(E)}
\end{aligned} \tag{87}
$$

(with an inverse inequality at the end). The choice of the two real constants

$$
\alpha_\pm = \frac{\displaystyle\int_{T_\pm} h_E b_E\,\tilde{r}_{T_\pm} \cdot r_E\mathrm{dx}}{\displaystyle\int_{T_\pm} w_{T_\pm} \cdot \tilde{r}_{T_\pm}\mathrm{dx}}
$$

in the definition

$$
w_E := \alpha_+ w_{T_+} + \alpha_- w_{T_-} - h_E b_E r_E \tag{88}
$$

yields $\int_{T_\pm} w_E \cdot \tilde{r}_{T_\pm}\mathrm{dx} = 0$. Since $\int_{T_\pm} w_{T_\pm} \cdot \tilde{r}_{T_\pm}\mathrm{dx} = \mathrm{h}_{T_\pm}^2\,\|\mathrm{b}_{T_\pm}^{1/2}\tilde{r}_{T_\pm}\|_{L_2(T_\pm)}^2$, one eventually deduces $|\alpha_\pm|\,|w_{T_\pm}|_{H^1(T_\pm)} \leq Ch_E^{1/2}\,\|r_E\|_{L_2(E)}$ and then concludes (83). An integration by parts shows

$$
\begin{aligned}
C^{-2}h_E\,\|r_E\|_{L_2(E)}^2 &\leq h_E\|b_E^{1/2}r_E\|_{L_2(E)}^2 \\
&= -\int_E w_E \cdot r_E\mathrm{ds} = \int_E w_E \cdot [(p - \tilde{p}) \cdot \nu_E]\mathrm{ds} \\
&= \int_{\omega_E} (p - \tilde{p}) : Dw_E\mathrm{dx} + \int_{\omega_E} w_E \cdot \mathrm{div}_{\mathcal{T}}\,(p - \tilde{p})\mathrm{dx} \\
&= R(w_E) - \int_{\omega_E} w_E \cdot (f + \mathrm{div}_{\mathcal{T}}\tilde{p})\mathrm{dx} \\
&= R(w_E) - \int_{\omega_E} w_E \cdot (f - \tilde{f})\mathrm{dx}
\end{aligned}
$$

(with $\int_{T_\pm} w_E \cdot \tilde{r}_{T_\pm} \mathrm{dx} = 0$ in the last step). A Friedrichs inequality $\|w_E\|_{L_2(\omega_E)} \leq Ch_E \, |w_E|_{H^1(\omega_E)}$ and (83) then conclude the proof of (85).

### 5.3. Implicit error estimators

Implicit error estimators are based on a local norm of a solution of a localized problem of a similar type with the residual terms on the right-hand side. This section introduces two different versions based on a partition of unity and based on an equilibration technique.

*5.3.1. Localization by partition of unity.* Given a nodal basis $(\varphi_z : z \in \mathcal{N})$ of a first-order finite element method with respect to $\mathcal{T}$, there holds the partition of unity property

$$\sum_{z \in \mathcal{N}} \varphi_z = 1 \quad \text{in} \quad \Omega$$

Given the residual $R = F - a(\tilde{u}, \cdot) \in V^*$, we observe that $R_z(v) := R(\varphi_z v)$ defines a bounded linear functional $R_z$ on a localized space called $V_z$ and specified below.

The bilinear form $a$ is an integral over $\omega$ on some integrand. The latter may be weighted with $\varphi_z$ to define some (localized) bilinear form $a_z : V_z \times V_z \to \mathbb{R}$. Supposing that $a_z$ is $V_z$-elliptic one defines the norm $\| \cdot \|_{a_z}$ on $V_z$ and considers

$$\eta_z := \sup_{v \in V_z \setminus \{0\}} \frac{R_z(v)}{\|v\|_{a_z}} \tag{89}$$

The dual norm is as in (70)-(71) and hence equivalent to the computation of the norm $\|e_z\|_{a_z}$ of a local solution

$$e_z \in V_z \quad \text{with} \quad a_z(e_z, \cdot) = R_z \in V_z^* \tag{90}$$

(The proof of $\|e_z\|_{a_z} = \eta_z$ follows the arguments of Section 5.1.2 and hence is omitted.)

**Example 10.** *Adopt notation from Example 1 and let $(\varphi_z : z \in \mathcal{N})$ be the first-order finite element nodal basis functions. Then define $R_z$ and $a_z$ by*

$$R_z(v) \quad := \quad \int_\Omega \varphi_z \, fv \mathrm{dx} - \int_\Omega \nabla \tilde{u} \cdot \nabla(\varphi_z v)\mathrm{dx} \quad \forall v \in V$$

$$a_z(v_1, v_2) \quad := \quad \int_\Omega \varphi_z \nabla v_1 \cdot \nabla v_2 \mathrm{dx} \quad \forall v_1, v_2 \in V$$

*Let $V_z$ denote the completion of $V$ under the norm given by the scalar product $a_z$ when $\varphi_z \not\equiv 0$ on $\Gamma$ or otherwise its quotient space with $\mathbb{R}$, i.e.*

$$V_z = \begin{cases} \{v \in H^1_{\mathrm{loc}}(\omega_z) : \ a_z(v, v) < \infty, \ \varphi_z v = 0 \ \text{on} \ \Gamma \cap \partial\omega_z\} \\ \qquad \text{if} \quad \varphi_z \not\equiv 0 \ \text{on} \ \Gamma \\ \{v \in H^1_{\mathrm{loc}}(\omega_z) : \ a_z(v, v) < \infty, \ \int_\Omega \varphi_z v \mathrm{dx} = 0\} \\ \qquad \text{if} \quad \varphi_z \equiv 0 \ \text{on} \ \Gamma \end{cases}$$

*Notice that $R_z(1) = 0$ for a free node $z$ such that (90) has a unique solution and hence $\eta_z < \infty$.*

In practical applications, the solution $e_z$ of (90) has to be approximated by some finite element approximation $\tilde{e}_z$ on a discrete space $\widetilde{V}_z$ based on a finer mesh or of higher order. (Arguing as in the stability estimate (71), leads to an error estimate for an approximation $\|\tilde{e}_z\|_{a_z}$ of $\eta_z$.)

Suppose that $\eta_z$ is known exactly (or computed with high and controlled accuracy) and that the bilinear form $a$ is localized through the partition of unity such that (e.g. in Example 10)

$$a(u,v) = \sum_{z \in \mathcal{N}} a_z(u,v) \qquad \forall u, v \in V \tag{91}$$

Then the implicit error estimator $\eta_L$ is reliable with $C_{\mathrm{rel}} = 1$ and $\mathrm{h.o.t.}_{\mathrm{rel}} = 0$,

$$\|e\|_a \leq \eta_L := \left( \sum_{z \in \mathcal{N}} \eta_z^2 \right)^{1/2} \tag{92}$$

The proof of (92) follows from the definition of $R_z$, $\eta_z$, and $e_z$ and Cauchy-Schwarz inequalities:

$$
\begin{aligned}
\|e\|_a^2 &= R(e) = \sum_{z \in \mathcal{N}} R_z(e) \leq \sum_{z \in \mathcal{N}} \eta_z \|e\|_{a_z} \\
&\leq \left( \sum_{z \in \mathcal{N}} \eta_z^2 \right)^{1/2} \left( \sum_{z \in \mathcal{N}} \|e\|_{a_z}^2 \right)^{1/2} = \eta_L \|e\|_a
\end{aligned}
$$

Notice that $\|\tilde{e}_z\|_{a_z} := \tilde{\eta}_z \leq \eta_z$ for any approximated local solution

$$\tilde{e}_z \in \widetilde{V}_z \quad \text{with} \quad a_z(\tilde{e}_z, \cdot) = R_z \in \widetilde{V}_z^* \tag{93}$$

and all of them are efficient estimators. The proof of efficiency is based on a weighted Poincaré or Friedrichs inequality which reads

$$\|\varphi_z v\|_a \leq C \|v\|_{a_z} \quad \forall v \in V_z \tag{94}$$

In fact, in Example 1, 2, and 3, one obtains efficiency in a more local form than indicated in

$$\eta_z \leq C \|e\|_a \quad \text{with} \ \mathrm{h.o.t.}_{\mathrm{eff}} = 0 \tag{95}$$

(This follows immediately from (94):

$$
\begin{aligned}
R_z(v) &= R(\varphi_z v) = a(e, \varphi_z v) \leq \|e\|_a \|\varphi_z v\|_a \\
&\leq C \|e\|_a \|v\|_{a_z})
\end{aligned}
$$

In the situation of Example 10, the estimator $\eta_L$ dates back to Babuška and Miller (1987); the use of weights was established in Carstensen and Funken (1999/00). A reliable computable estimator $\tilde{\eta}_L$ is introduced in Morin, Nochetto and Siebert (2003a) based on a proper finite-dimensional space $\widetilde{V}_z$ of some piecewise quadratic polynomials on $\omega_z$.

*5.3.2. Equilibration estimators.* The nonoverlapping domain decomposition schemes employ artificial unknowns $g_T \in L_2(\partial T)^m$ for each $T \in \mathcal{T}$ at the interfaces, which allow a representation of the form

$$
\begin{aligned}
R(v) \quad &= \quad \sum_{T \in \mathcal{T}} R_T(v) \quad \text{where} \\
R_T(v) \quad &:= \quad \int_T f \cdot v \mathrm{dx} - \int_T \tilde{p} : Dv \mathrm{dx} + \int_{\partial T} g_T \cdot v \mathrm{ds}
\end{aligned}
\tag{96}
$$

Adopting the notation from Section 5.2.1, the new quantities $g_T$ satisfy

$$
g_{T_+} + g_{T_-} = 0 \quad \text{along } E = \partial T_+ \cap \partial T_- \in \mathcal{E}(\Omega)
$$

(where $T_\pm$ and $T$ denote neighboring element domains) to guarantee (96). (There are nondisplayed modifications on any Neumann boundary edge $E \subset \partial\Omega$.) Moreover, the bilinear form $a$ is expanded in an elementwise form

$$
a(u,v) = \sum_{T \in \mathcal{T}} a_T(u,v) \qquad \forall u, v \in V
\tag{97}
$$

Under the *equilibration condition* $R_T(c) = 0$ for all kernel functions $c$ (namely, the constant functions for the Laplace model problem), the resulting *local problem* reads

$$
e_T \in V_T \quad \text{with } a_T(e_T, \cdot) = R_T \in V_T^*
\tag{98}
$$

and is equivalent to the computation of

$$
\eta_T := \sup_{v \in V_T \backslash \{0\}} \frac{R_T(v)}{\|v\|_{a_T}} = \|e_T\|_{a_T}
\tag{99}
$$

The sum of all local contributions defines the reliable *equilibration error estimator* $\eta_{EQ}$,

$$
\|e\|_a \leq \eta_{EQ} := \Big( \sum_{T \in \mathcal{T}} \eta_T^2 \Big)^{1/2}
\tag{100}
$$

(The immediate proof of (100) is analogous to that of (92) and hence is omitted.)

**Example 11.** *In the situation of Example 10 there holds $\eta_T < \infty$ if and only if either $\Gamma \cap \partial T$ has positive surface measure (with $V_T = \{v \in H^1(T) : v = 0 \text{ on } \Gamma \cap \partial T\}$) or otherwise $R_T(1) = 0$ (with $V_T = \{v \in H^1(T) : \int_T v \mathrm{dx} = 0\}$). Ladeveze and Leguillon (1983) suggested a certain choice of the interface corrections to guarantee this and even higher-order equilibrations are established. Details on the implementation are given in Ainsworth and Oden (2000); a detailed error analysis with higher-order equilibrations and the perturbation by a finite element simulation of the local problems with corrections can be found in Ainsworth and Oden (2000) and Babuška and Strouboulis (2001).*

The error estimator $\eta = \eta_{EQ}$ is efficient in the sense of (69) with higher-order terms h.o.t.$(T)$ on $T$ that depend on the given data provided

$$
\begin{aligned}
h_E^{1/2} \|g_T - \tilde{p}\nu_E\|_{L_2(E)} \quad &\leq \quad C \|e\|_{a_T} + \text{h.o.t.}(T) \\
&\quad \text{for all } E \in \mathcal{E}(T)
\end{aligned}
\tag{101}
$$

(Recall that $\mathcal{E}(T)$ denotes the set of edges or faces of $T$.) This stability property depends on the design of $g_T$; a positive example is given in Theorem 6.2 of Ainsworth and Oden (2000) for Example 1. Given Inequality (101), the efficiency of $\eta_T$ follows with standard arguments, for example, an integration by parts, a trace and Poincaré or Friedrichs inequality $h_T^{-1}\|v\|_{L_2(T)} + h_T^{-1/2}\|v\|_{L_2(\partial T)} \leq C\,\|v\|_{a_T}$ for $v \in V_T$:

$$
\begin{aligned}
R_T(v) &= \int_T r_T \cdot v \mathrm{dx} + \int_{\partial \mathrm{T}} (\mathrm{g_T} - \tilde{\mathrm{p}}\nu) \cdot \mathrm{vds} \\
&\leq C\left(h_T\|r_T\|_{L_2(T)} + h_T^{1/2}\|g_T - \tilde{p}\nu\|_{L_2(\partial T)}\right)\|v\|_{a_T}
\end{aligned}
$$

followed by (79) and (101).

Further examples and the relation to the hypercircle identity can be found in Braess (2001). Some refinement with a postprocessing is suggested in Carstensen and Merdon (2014).

### 5.4. Multilevel error estimators

While the preceding estimators evaluate or estimate the residual of one finite element solution $u_H$, multilevel estimators concern at least two meshes $\mathcal{T}_H$ and $\mathcal{T}_h$ with associated discrete spaces $V_H \subset V_h \subset V$ and two discrete solutions $u_H = \tilde{u}$ and $u_h$. The interpretation is that $p_h$ is computed on a finer mesh (e.g. $\mathcal{T}_h$ is a refinement of $\mathcal{T}_H$) or that $p_h$ is computed with higher polynomial order than $p_H = \tilde{p}$.

*5.4.1. Error-reduction property and multilevel error estimator.* Let $V_H \subset V_h \subset V$ denote two nested finite element spaces in $V$ with coarse and fine finite element solution $u_H = \tilde{u} \in V_H = \widetilde{V}$ and $u_h \in V_h$ of the discrete problem (19), respectively, and with the exact solution $u$. Let $p_H = \tilde{p}$, $p_h$, and $p$ denote the respective fluxes and let $\|\cdot\|$ be a norm associated to the energy norm, for example, a norm with $\|p - \tilde{p}\| = \|u - \tilde{u}\|_a$ and $\|p - p_h\| = \|u - u_h\|_a$. Then, the *multilevel error estimator*

$$\eta_{ML} := \|p_h - p_H\| = \|u_h - u_H\|_a \tag{102}$$

is simply the norm of the difference of the two discrete solutions. The interpretation is that the error $\|p - p_h\|$ of the finer discrete solution is systematically smaller than the error $\|e\|_a = \|p - p_H\|$ of the coarser discrete solution in the sense of an *error-reduction property*: For some constant $\varrho < 1$, there holds

$$\|p - p_h\| \leq \varrho\,\|p - p_H\| \tag{103}$$

Notice the bound $\varrho \leq 1$ for Galerkin errors in the energy norm (because of the best-approximation property). The point is that $\varrho < 1$ in (103) is bounded away from one. Then, the error-reduction property (103) immediately implies reliability and efficiency of $\eta_{ML}$:

$$(1 - \varrho)\|p - p_H\| \leq \eta_{ML} = \|p_h - p_H\| \leq (1 + \varrho)\|p - p_H\| \tag{104}$$

(The immediate proof of (104) utilizes (103) and a simple triangle inequality.)
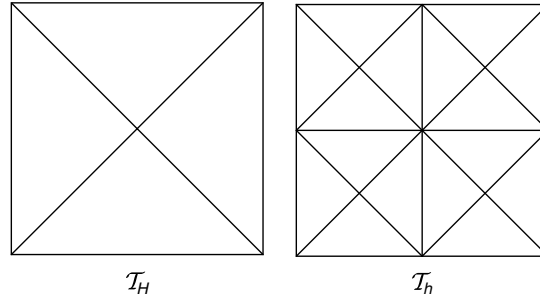
Figure 10: Two triangulations $\mathcal{T}_H$ and $\mathcal{T}_h$ with equal discrete $P_1$ finite element solutions $u_H = u_h$ for a Poisson problem with right-hand side $f = 1$. The refinement $\mathcal{T}_h$ of $\mathcal{T}_H$ is generated by two (newest-vertex) bisections per element (initially with the interior node in $\mathcal{T}_H$ as newest vertex).

Four remarks on the error-reduction conclude this section: Efficiency of $\eta_{ML}$ in the sense of (69) is robust in $\varrho \to 1$, but reliability is not: The reliability constant $C_{rel} = (1 - \varrho)^{-1}$ in (68) tends to infinity as $\varrho$ approaches 1.

Higher-order terms are invisible in (104): h.o.t.$_{rel} = 0 = $ h.o.t.$_{eff}$. This is unexpected when compared to all the other error estimators and hence indicates that (103) should fail to hold for heavily oscillating right-hand sides.

The error-reduction property (103) is often observed in practice for fine meshes and can be monitored during the calculation. For coarse meshes and in the preasymptotic range, (103) may fail to hold.

The error-reduction property (103) is often called *saturation assumption* in the literature and frequently has a status of an unproven hypothesis.

*5.4.2. Counterexample for error-reduction.*   The error-reduction property (103) may fail to hold even if $f$ shows *no* oscillations: Figure 10 displays two triangulations, $\mathcal{T}_H$ and its refinement $\mathcal{T}_h$, with one and five free nodes, respectively. If the right-hand side is constant and if the problem has homogeneous Dirichlet conditions for the Poisson problem

$$1 + \Delta u = 0 \quad \text{in } \Omega := (0,1)^2 \quad \text{and} \quad u = 0 \quad \text{on } \partial\Omega$$

then the corresponding $P_1$ finite element solutions coincide: $u_H = u_h$. A direct proof is based on the nodal basis function $\Phi_1$ of the free node in $V_H := FE_{\mathcal{T}_H}$ (the first-order finite element space with respect to the coarse mesh $\mathcal{T}_H$) and the nodal basis functions $\varphi_2, \ldots, \varphi_5 \in V_h := \mathcal{S}_0^1(\mathcal{T}_h)$ of the new free nodes in $\mathcal{T}_h$. Then,

$$\Phi_2 := \Phi_1 - (\varphi_2 + \cdots + \varphi_5) \in V_h := FE_{\mathcal{T}_h}$$

satisfies (since $\int_E \Phi_2 ds = 0$ for all edges $E$ in $\mathcal{T}_H$ and $\int_\Omega \Phi_2 dx = 0$)

$$R(\Phi_2) = 0$$

Thus $u_H$ is the finite element solution in $W_h := \text{span}\{\Phi_1, \Phi_2\} \subset V_h$. Since, by symmetry, the finite element solution $u_h$ in $V_h$ belongs to $W_h$, there holds $u_H = u_h$.

The characterization of the saturation property in Carstensen, Gallistl and Gedicke (2016) illustrates that only a few very particular triangulations lead to counterexamples of this type.

*5.4.3. Affirmative example for error-reduction.* Adopt notation from Section 5.2.4 with a coarse discrete space $V_H = \tilde{V}$ and consider the fine space $V_h := V_H \oplus W_h$ for

$$W_h := \text{span}\{\tilde{r}_T \, b_T : T \in \mathcal{T}\} \oplus \text{span}\{r_E \, b_E : E \in \mathcal{E}(\Omega)\} \subset V \qquad (105)$$

Then there holds the error-reduction property up to higher-order terms

$$\text{osc}(f) := \left( \sum_{T \in \mathcal{T}} h_T^2 \|f - \tilde{f}\|_{L_2(T)}^2 \right)^{1/2}$$

namely

$$\|p - p_h\|^2 \le \varrho \, \|p - p_H\|^2 + \text{osc}(f)^2 \qquad (106)$$

The constant $\varrho$ in (106) is uniformly smaller than one, independent of the mesh size, and depends on the shape of the elements and the type of ansatz functions through constants in (83)-(85).

The proof of (106) is based on the test functions $w_T$ and $w_E$ in (86) and (88) of Section 5.2.4 and

$$w_h := \sum_{T \in \mathcal{T}} w_T + \sum_{E \in \mathcal{E}(\Omega)} w_E \in W_h \subset V$$

Utilizing (83)-(85) one can prove

$$
\begin{aligned}
\|w_h\|_a^2 &\le C \left( \sum_{T \in \mathcal{T}} h_T^2 \|\tilde{r}_T\|_{L_2(T)}^2 + \sum_{E \in \mathcal{E}(\Omega)} h_E \|r_E\|_{L_2(E)}^2 \right) \\
&\le C \left( \sum_{T \in \mathcal{T}} R(w_T) + \sum_{E \in \mathcal{E}(\Omega)} R(w_E) + \text{osc}(f)^2 \right) \\
&= C \left( R(w_h) + \text{osc}(f)^2 \right)
\end{aligned}
$$

Since $w_h$ belongs to $V_h$ and $u_h$ is the finite element solution with respect to $V_h$ there holds

$$R(w_h) = a(u_h - u_H, w_h) \le \|u_h - u_H\|_a \|w_h\|_a$$

The combination of the preceding inequalities yields the key inequality

$$
\begin{aligned}
\eta_{R,R}^2 &= \sum_{T \in \mathcal{T}} h_T^2 \|\tilde{r}_T\|_{L_2(T)}^2 + \sum_{E \in \mathcal{E}(\Omega)} h_E \|r_E\|_{L_2(E)}^2 \\
&\le C(\|u_h - u_H\|_a^2 + \text{osc}(f)^2)
\end{aligned}
$$

This, the Galerkin orthogonality $\|u - u_H\|_a^2 = \|u - u_h\|_a^2 + \|u_h - u_H\|_a^2$, and the reliability of $\eta_R$ show

$$\|u - u_H\|_a^2 \le CC_{\mathrm{rel}}^2 \big( \|u - u_H\|_a^2 - \|u - u_h\|_a^2 + \mathrm{osc}(f)^2 \big)$$

and so imply (106) with $\varrho = 1 - C^{-1}C_{\mathrm{rel}}^{-2} < 1$.

**Example 12.** *In the Poisson problem with $P_1$ finite elements and discrete space $V_H$, let $W_h$ consist of the quadratic and cubic bubble functions (82). Then (106) holds; cf. also Example 14 below.*

**Example 13.** *Other affirmative examples for the Poisson problem consist of the $P_1$ and $P_2$ finite element spaces $V_H$ and $V_h$ over one regular triangulation $\mathcal{T}$ or of the $P_1$ finite element spaces with respect to a regular triangulation $\mathcal{T}_H$ and its red-refinement $\mathcal{T}_h$. The observation that the element-bubble functions are in fact redundant is due to Dörfler and Nochetto (2002).*

*5.4.4. Hierarchical error estimator.* Given a smooth right-hand side $f$ and based on the example of the previous section, the multilevel estimator (102) is reliable and efficient up to higher-order terms. The costly calculation of $u_h$, however, exclusively allows for an accurate error control of $u_H$ (and no reasonable error control for $u_h$). Instead of (102), cheaper versions are favored where $u_h$ is replaced by some quantity computed by a localized problem. One reliable and efficient version is the *hierarchical error estimator*

$$\eta_H := \left( \sum_{T \in \mathcal{T}} \eta_T^2 + \sum_{E \in \mathcal{E}} \eta_E^2 \right)^{1/2} \tag{107}$$

where, for each $T \in \mathcal{T}$ and $E \in \mathcal{E}$ and their test functions (86) and (88),

$$\eta_T := \frac{R(w_T)}{\|w_T\|_a} \quad \text{and} \quad \eta_E := \frac{R(w_E)}{\|w_E\|_a} \tag{108}$$

(The proof of reliability and efficiency follows from (83)-(85) by the arguments from Section 5.2.4.)

**Example 14.** *In the Poisson problem with $P_1$ finite elements and discrete space $V_H$, let $W_h$ consist of the quadratic and cubic bubble functions (82). Then,*

$$\eta_H := \left( \sum_{T \in \mathcal{T}} \frac{R(b_T)^2}{\|b_T\|_a^2} + \sum_{E \in \mathcal{E}(\Omega)} \frac{R(b_E)^2}{\|b_E\|_a^2} \right)^{1/2} \tag{109}$$

*is a reliable and efficient hierarchical error estimator. With the error-reduction property of $P_1$ and $P_2$ finite elements due to Dörfler and Nochetto (2002),*

$$\eta_H := \left( \sum_{E \in \mathcal{E}(\Omega)} \frac{R(b_E)^2}{\|b_E\|_a^2} \right)^{1/2} \tag{110}$$

*is reliable and efficient as well. The same is true if, for each edge $E \in \mathcal{E}$, $b_E$ defines a hat function of the midpoint of $E$ with respect to a red-refinement $\mathcal{T}_h$ of $\mathcal{T}_H$ (that is, each edge is halved and each triangle is divided into four congruent subtriangles; cf. Figure 15, left).*

*5.5. Averaging error estimators*

Averaging techniques, also called (*gradient*) *recovery estimators*, focus on one mesh and one known low-order flux approximation $\tilde{p}$ and the difference to a piecewise polynomial $\tilde{q}$ in a finite-dimensional subspace $\widetilde{Q} \subset L_2(\Omega; \mathbb{R}^{m \times n})$ of higher polynomial degrees and more restrictive continuity conditions than those generally satisfied by $\tilde{p}$. Averaging techniques are universal in the sense that there is no need for any residual or partial differential equation in order to apply them.

*5.5.1. Definition of averaging error estimators.* The procedure consists of taking a piecewise smooth $\tilde{p}$ and approximating it by some globally continuous piecewise polynomials (denoted by $\widetilde{Q}$) of higher degree $A\tilde{p}$. A simple example, frequently named after Zienkiewicz and Zhu and sometimes even called the ZZ estimator, reads as follows: For each node $z \in \mathcal{N}$ and its patch $\omega_z$ let

$$(A\tilde{p})(z) = \frac{\displaystyle\int_{\omega_z} \tilde{p}\mathrm{dx}}{\displaystyle\int_{\omega_z} 1\mathrm{dx}} \in \mathbb{R}^{m \times n} \tag{111}$$

be the integral mean of $\tilde{p}$ over $\omega_z$. Then, define $A\tilde{p}$ by interpolation with (conforming, i.e. globally continuous) hat functions $\varphi_z$, for $z \in \mathcal{N}$,

$$A\tilde{p} = \sum_{z \in \mathcal{N}} (A\tilde{p})(z)\, \varphi_z \in \tilde{Q}$$

Let $\widetilde{Q} = \mathrm{span}\{\varphi_z \colon z \in \mathcal{N}\}$ denote the (conforming) first-order finite element space and let $\|\cdot\|$ be the norm for the fluxes. Then the averaging estimator is defined by

$$\eta_A := \|\tilde{p} - A\tilde{p}\| \tag{112}$$

Notice that there is a minimal version

$$\eta_M := \min_{\tilde{q} \in \tilde{Q}} \|\tilde{p} - \tilde{q}\| \leq \eta_A \tag{113}$$

The efficiency of $\eta_M$ follows from a triangle inequality, namely

$$\eta_M \leq \|p - \tilde{p}\| + \|p - \tilde{q}\| \quad \text{for all } \tilde{q} \in \widetilde{Q} \tag{114}$$

and the fact that $\|p - \tilde{p}\| = \mathcal{O}(h)$ while (in all the examples of this chapter)

$$\min_{\tilde{q} \in \tilde{Q}} \|p - \tilde{q}\| = \mathrm{h.o.t.}(p) =: \mathrm{h.o.t._{eff}}$$

This is of higher order for smooth $p$ and efficiency follows for $\eta = \eta_M$ and $\mathrm{C_{eff}} = 1$.

It turns out that $\eta_A$ and $\eta_M$ are very close and accurate estimators in many numerical examples; cf. Section 5.6.4 below. This and the fact that the calculation of $\eta_A$ is an easy postprocessing made $\eta_A$ extremely popular.

For proper treatment of Neumann boundary conditions, we refer to Carstensen and Bartels (2002) and for applications in solid and fluid mechanics to Alberty and Carstensen (2003) and Carstensen and Funken (2001a,b) and for higher-order FEM to Bartels and Carstensen (2002).

Multigrid smoothing steps may be successfully employed as averaging procedures as proposed in Bank and Xu (2003).

*5.5.2. All averaging error estimators are reliable.*   The first proof of reliability dates back to Rodriguez (1994a,b) and we refer to Carstensen (2004) for an overview. A simplified reliability proof for $\eta_M$ and hence for *all* averaging techniques (Carstensen, Bartels and Klose, 2001) is outlined in the sequel.

First let $\Pi$ be the $L_2$ projection onto the first-order finite element space $\widetilde{V}$ and let $\tilde{q}$ be arbitrary in $\widetilde{Q}$, that is, each of the $m \times n$ components of $\tilde{q}$ is a first-order finite element function in $FE_{\mathcal{T}}$. The Galerkin orthogonality shows for the error $e := u - \tilde{u}$ and $\tilde{p} := \nabla \tilde{u}$ in the situation of Example 1 that

$$
\begin{aligned}
\|e\|_a^2 \;=\; & \int_\Omega (\nabla u - \tilde{q}) \cdot \nabla(e - \Pi e)\mathrm{d}x \\
& + \int_\Omega (\tilde{q} - \tilde{p}) \cdot \nabla(e - \Pi e)\mathrm{d}x
\end{aligned}
$$

A Cauchy-Schwarz inequality in the latter term is combined with the $H^1$-stability of $\Pi$, namely,

$$
\|\nabla(e - \Pi e)\|_{L_2(\Omega)} \le \mathrm{C}_{\mathrm{stab}}\|\nabla e\|_{L_2(\Omega)}
$$

(For sufficient conditions for this we refer to Crouzeix and Thomée (1987), Bramble, Pasciak and Steinbach (2002), Carstensen (2002, 2003b), and Carstensen and Verfürth (1999).) The $H^1$-stability in the second term and an integration by parts in the first term on the right-hand side show

$$
\begin{aligned}
\|e\|_a^2 \;\le\; & \int_\Omega f \cdot (e - \Pi e)\mathrm{d}x + \int_\Omega (e - \Pi e) \cdot \mathrm{div}\,\tilde{q}\mathrm{d}x \\
& + \mathrm{C}_{\mathrm{stab}}\|\nabla e\|_{L_2(\Omega)}\|\tilde{p} - \tilde{q}\|_{L_2(\Omega)}
\end{aligned}
$$

Since $e - \Pi e$ is $L_2$-orthogonal onto $f_h := \Pi f \in \widetilde{V}$,

$$
\begin{aligned}
& \int_\Omega f \cdot (e - \Pi e)\mathrm{d}x \\
& = \int_\Omega (f - f_h) \cdot (e - \Pi e)\mathrm{d}x \\
& \le \|h_{\mathcal{T}}^{-1}(e - \Pi e)\|_{L_2(\Omega)}\|h_{\mathcal{T}}(f - \Pi f)\|_{L_2(\Omega)}
\end{aligned}
$$

Notice that, despite possible boundary layers, $\|h_{\mathcal{T}}(f - \Pi f)\|_{L_2(\Omega)} = $ h.o.t. is of higher order. The first-order approximation property of the $L_2$ projection,

$$
\|h_{\mathcal{T}}^{-1}(e - \Pi e)\|_{L_2(\Omega)} \le \mathrm{C}_{\mathrm{apx}}\|\nabla e\|_{L_2(\Omega)}
$$

follows from the $H^1$-stability (cf. e.g. Carstensen and Verfürth (1999) for a proof). Similar arguments for the remaining term $\operatorname{div}\tilde{q}$ and the $\mathcal{T}$-piecewise divergence operator $\operatorname{div}_\mathcal{T}$ with $\operatorname{div}\tilde{q} = \operatorname{div}_\mathcal{T}\tilde{q} = \operatorname{div}_\mathcal{T}(\tilde{q}-\tilde{p})$ (recall that $\tilde{u}$ is of first-order and hence $\Delta\tilde{u}=0$ on each element) lead to

$$\int_\Omega (e - \Pi e)\cdot\operatorname{div}\tilde{q}\mathrm{dx}$$

$$\leq \mathrm{C}_{\mathrm{apx}}\|\nabla e\|_{L_2(\Omega)}\|h_\mathcal{T}\operatorname{div}_\mathcal{T}(\tilde{q}-\tilde{p})\|_{L_2(\Omega)}$$

An inverse inequality $h_T\|\operatorname{div}_\mathcal{T}(\tilde{q}-\tilde{p})\|_{L_2(T)} \leq \mathrm{C}_{\mathrm{inv}}\|\tilde{q}-\tilde{p}\|_{L_2(T)}$ (cf. Section 3.4) shows

$$\int_\Omega (e-\Pi e)\cdot\operatorname{div}\tilde{q}\mathrm{dx} \leq \mathrm{C}_{\mathrm{apx}}\mathrm{C}_{\mathrm{inv}}\|\nabla e\|_{L_2(\Omega)}\|\tilde{q}-\tilde{p}\|_{L_2(\Omega)}$$

The combination of all established estimates plus a division by $\|e\|_a = \|\nabla e\|_{L_2(\Omega)}$ yield the announced reliability result

$$\|e\|_a \leq (\mathrm{C}_{\mathrm{stab}} + \mathrm{C}_{\mathrm{apx}}\mathrm{C}_{\mathrm{inv}})\|\tilde{p}-\tilde{q}\|_{L_2(\Omega)} + \text{h.o.t.}$$

In the second step, one designs a more local approximation operator $J$ to substitute $\Pi$ as in Carstensen and Bartels (2002); the essential properties are the $H^1$-stability, the first-order approximation property, and a local form of the orthogonality condition; we omit the details.

*5.5.3. Averaging error estimators and edge contributions.* There is a local equivalence of the estimators $\eta_A$ from a local averaging process (111) and the edge estimator

$$\eta_E := \left(\sum_{E\in\mathcal{E}(\Omega)} h_E\|[\tilde{p}\nu_E]\|_{L_2(E)}^2\right)^{1/2}$$

The observation that, with some mesh-size-independent constant $C$,

$$C^{-1}\,\eta_E \leq \eta_A \leq C\,\eta_E \tag{115}$$

dates back to Rodriguez (1994a) and can be found in Verfürth (1996). The proof of (115) for piecewise linears in $\widetilde{V}$ is based on the equivalence of the two seminorms

$$\varrho_1(\tilde{q}) \quad := \quad \min_{r\in\mathbb{R}^{m\times n}}\|\tilde{q}-r\|_{L_2(\omega_z)}$$

and

$$\varrho_2(\tilde{q}) \quad := \quad \left(\sum_{E\in\mathcal{E}(z)} h_E\|[\tilde{q}\nu_E]\|_{L_2(E)}^2\right)^{1/2}$$

for piecewise constant vector-valued functions $\tilde{q}$ in $P_z$, the set of possible fluxes $\tilde{q}$ restricted on the patch $\omega_z$ of a node $z$, and with the set of edges $\mathcal{E}(z) := \{E\in\mathcal{E} : z\in E\}$. The main observation is that $\varrho_1$ and $\varrho_2$ vanish exactly for constants functions $\mathbb{R}^{m\times n}$ and hence they are norms on the quotient space $P_z/\mathbb{R}^{m\times n}$. By the equivalence of norms on any finite-dimensional space $P_z/\mathbb{R}^{m\times n}$, there holds

$$C^{-1}\,\varrho_1(\tilde{q}) \leq \varrho_2(\tilde{q}) \leq C\,\varrho_1(\tilde{q}) \quad \forall\tilde{q}\in P_z$$

This is a local version of (115), which eventually implies (115) by localization and composition; we omit the details.

For triangles and tetrahedra and piecewise linear finite element functions, it is proved in Carstensen (2004) that

$$\eta_M \leq \eta_A \leq C_d\,\eta_M \tag{116}$$

with universal constants $C_2 = \sqrt{10}$ and $C_3 = \sqrt{15}$ for 2-D and 3-D, respectively. This equivalence holds for a larger class of elements and (first-order) averaging operators and then proves efficiency for $\eta_A$ whereas efficiency of $\eta_M$ follows from a triangle inequality in (114).

### 5.6. Comparison of error bounds in benchmark example

This section is devoted to a numerical comparison of energy errors and its a posteriori error estimators for an elliptic model problem.

*5.6.1. Benchmark example.* The numerical comparisons are computed for the Poisson problem

$$1 + \Delta u = 0 \text{ in } \Omega \quad \text{and} \quad u = 0 \text{ on } \partial\Omega \tag{117}$$

on the $L$-shaped domain $\Omega = (-1, +1)^2 \backslash ([0,1] \times [-1,0])$ and its boundary $\partial\Omega$. The first mesh $\mathcal{T}_1$ consists of 17 free nodes and 48 elements and is obtained by, first, a decomposition of $\Omega$ in 12 congruent squares of size $1/2$ and, second, a decomposition of each of the boxes along its two diagonals into 4 congruent triangles. The subsequent meshes are successively red-refined (i.e. each triangle is partitioned into 4 congruent subtriangles of Figure 15, left). This defines the (conforming) $P_1$ finite element spaces $V_1 \subset V_2 \subset V_3 \subset \cdots \subset V := H_0^1(\Omega)$. The error $e := u - u_j$ of the finite element solution $u_j = \tilde{u}$ in $V_j = \widetilde{V}$ of dimension $N = \dim(V_j)$ is measured in the energy norm (the Sobolev seminorm in $H^1(\Omega)$)

$$
\begin{aligned}
|e|_{1,2} \quad &:= \quad |e|_{H^1(\Omega)} := \left( \int_\Omega |De|^2 \mathrm{dx} \right)^{1/2} = a(u - \tilde{u}, u + \tilde{u})^{1/2} \\
&= \quad \left( |u|^2_{H^1(\Omega)} - |\tilde{u}|^2_{H^1(\Omega)} \right)^{1/2}
\end{aligned}
$$

(by the Galerkin orthogonality) computed with the approximation $|u|^2_{H^1(\Omega)} \approx 0.21407315683398$.

Figure 11 displays the computed values of $|e|_{1,2}$ for a sequence of uniformly refined meshes $\mathcal{T}_1, \mathcal{T}_2, \ldots, \mathcal{T}_6$ as a function of the number of degrees of freedom $N = 17, 81, 353, 1473, 6017, 24321, 97793$. It is this curve that will be estimated by computable upper and lower bounds explained in the subsequent sections where we use the fact that each element is a right isosceles triangle.

Notice that the two axes in Figure 11 scale logarithmically such that any algebraic curve of growth $\alpha$ is mapped into a straight line of slope $-\alpha$. The experimental convergence rate is $2/3$
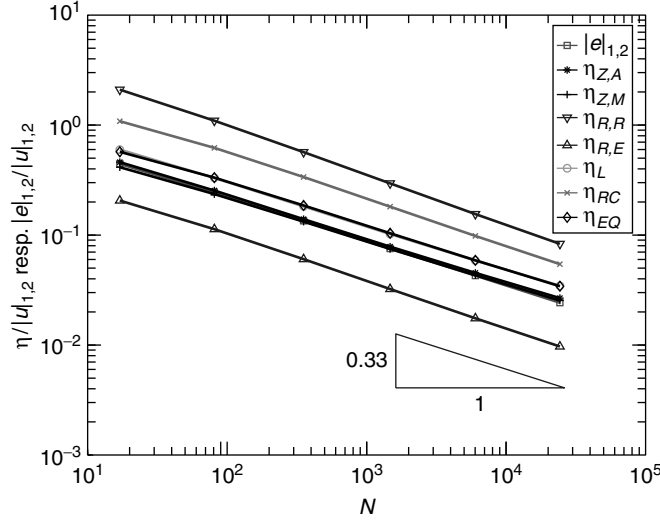
Figure 11: Experimental results for the benchmark problem (117) with meshes $\mathcal{T}_1, \ldots, \mathcal{T}_6$. The relative error $|e|_{1,2}/|u|_{1,2}$ and various estimators $\eta/|u|_{1,2}$ are plotted as functions of the number of degrees of freedom $N$. Both axes are in a logarithmic scaling such that an algebraic curve $N^{-\alpha}$ is visible as a straight line with slope $-\alpha$. A triangle with slope $-0.33$ is displayed for comparison.

in agreement with the (generic) singularity of the domain and resulting theoretical predictions. More details can be found in Carstensen, Bartels and Klose (2001).

*5.6.2. Explicit error estimators.* For the benchmark problem of Section 5.6.1, the error estimator (75) can be written in the form

$$\eta_{R,R} \quad := \quad \left( \sum_{T \in \mathcal{T}} h_T^2 \|1\|_{L_2(T)}^2 \right)^{1/2} + \left( \sum_{E \in \mathcal{E}_\Omega} h_E \int_E \left[ \frac{\partial \tilde{u}}{\partial \nu_E} \right]^2 \mathrm{ds} \right)^{1/2} \tag{118}$$

and is reliable with $\mathrm{C}_{\mathrm{rel}} = 1$ and h.o.t.$_{\mathrm{rel}} = 0$ (Carstensen and Funken, 2001a). Figure 11 displays $\eta_{R,R}$ as a function of the number of unknowns $N$ and illustrates $|e|_{1,2} \leq \eta_{R,R}$.

Guaranteed lower error bounds, i.e. with the efficiency constant $\mathrm{C}_{\mathrm{eff}}$, are less established and the higher-order terms h.o.t.$_{\mathrm{eff}}$ usually involve $\|f - f_T\|_{L_2(T)}$ for the elementwise integral mean $f_T$ of the right-hand side. Here, $f \equiv 1$ and so h.o.t.$_{\mathrm{eff}} = 0$. Following a derivation in Carstensen, Bartels and Klose (2001), Figure 11 also displays an efficient variant $\eta_{R,E}$ as a lower error bound: $\eta_{R,E} \leq |e|_{1,2}$.

The guaranteed lower and upper bounds with explicit error estimators leave a very large region for the true error. Our interpretation is that various geometries (the shape of the patches) lead to different constants and $\mathrm{C}_{\mathrm{rel}} = 1$ reflects the worst possible situation for every

patch in the current mesh.

A more efficient reliable explicit error estimator $\eta_{R,C}$ from Carstensen and Funken (2001a) displayed in Figure 11 requires the computation of local (patchwise) analytical eigenvalues and hence is very expensive. However, the explicit estimators $\eta_{R,C}$ and $\eta_{R,E}$ still overestimate and underestimate the true error by a huge factor (up to 10 and even more) in the simple situation of the benchmark. One conclusion is that the involved constants estimate a worst-case scenario with respect to every right-hand side or every exact solution.

This experimental evidence supports the design of more elaborate estimators: The stopping criterion (72) with the reliable explicit estimators may appear very cheap and easy. But the decision (72) may have too costly consequences.

*5.6.3. Implicit estimators.* For comparison, the two implicit estimators $\eta_L$ and $\eta_{EQ}$ are displayed in Figure 11 as functions of $N$. It is stressed that both estimators are efficient and reliable (Carstensen and Funken, 1999/00)

$$|e|_{1,2} \le \eta_L \le 2.37\,|e|_{1,2}$$

The practical performance of $\eta_L$ and $\eta_{EQ}$ in Figure 11 is comparable and in fact is much sharper than that of $\eta_{R,E}$ and $\eta_{R,R}$.

*5.6.4. Averaging estimator.* The averaging estimators $\eta_A$ and $\eta_M$ are as well displayed in Figure 11 as a function of $N$. Here, $\eta_M$ is efficient up to higher-order terms (since the exact solution $u \in H^{5/3-\varepsilon}(\Omega)$ is singular, this is not really guaranteed) while its reliability is open, i.e. the corresponding constants have not been computed. Nevertheless, the behavior of $\eta_A$ and $\eta_M$ is exclusively seen here from an experimental point of view. The striking numerical result is an amazing high accuracy of $\eta_M \approx \eta_A$ as an empirical guess of $|e|_{1,2}$. If we took $C_{\mathrm{rel}}$ into account, this effect would be destroyed: The high accuracy is an empirical observation in this example (and possibly many others) but does not yield an accurate guaranteed error bound.

*5.6.5. Adapted meshes.* The benchmark in Figure 11 is based on a sequence of uniform meshes and hence results in an experimental convergence rate 2/3 according to the corner singularity of this example. Adaptive mesh-refining algorithms, described below in more detail, are empirically studied also in Carstensen, Bartels and Klose (2001). The observations can be summarized as follows: The quality of the estimators and their relative accuracy is similar to what is displayed in Figure 11 even though the convergence rates are optimally improved to one.

*5.7. Goal-oriented error estimators*

This section provides a brief introduction to goal-oriented error control.

*5.7.1. Goal functionals.* Given the Sobolev space $V = H_0^1(\Omega)$ with a finite-dimensional subspace $\widetilde{V} \subset V$, a bounded and $V$-elliptic bilinear form $a : V \times V \to \mathbb{R}$, a bounded linear form $F : V \to \mathbb{R}$, there exists an exact solution $u \in V$ and a discrete solution $\tilde{u} \in \widetilde{V}$ of

$$a(u, v) = F(v) \quad \forall v \in V \quad \text{and} \quad a(\tilde{u}, \tilde{v}) = F(\tilde{v}) \quad \forall \tilde{v} \in \widetilde{V} \tag{119}$$

The previous sections concern estimations of the error $e := u - \tilde{u}$ in the energy norm, equivalent to the Sobolev norm in $V$. Other norms are certainly of some interest as well as the error with respect to a certain goal functional. The latter is some given bounded and linear functional $J : V \to \mathbb{R}$ with respect to which one aims to monitor the error, that is, one wants to find computable lower and upper bounds for the (unknown) quantity

$$|J(u) - J(\tilde{u})| = |J(e)|$$

Typical examples of goal functionals are described by $L_2$ functions, for example,

$$J(v) = \int_\Omega \varrho \, v \mathrm{dx} \quad \forall v \in V$$

for a given $\varrho \in L_2(\Omega)$ or as contour integrals.

In many cases, the main interest is on a point value and then $J(v)$ is given by a mollification $\varrho$ of a singular measure in order to guarantee the boundedness of $J : V \to \mathbb{R}$.

*5.7.2. Duality technique.* To bound or approximate $J(e)$ one considers the *dual problem*

$$a(v, z) = J(v) \quad \forall v \in V \tag{120}$$

with exact solution $z \in V$ (guaranteed by the Lax-Milgram lemma) and the discrete solution $\tilde{z} \in \widetilde{V}$ of

$$a(\tilde{v}, \tilde{z}) = J(\tilde{v}) \quad \forall \tilde{v} \in \widetilde{V}$$

Set $f := z - \tilde{z}$. On the basis of the Galerkin orthogonality $a(e, \tilde{z}) = 0$ one infers

$$J(e) = a(e, z) = a(e, z - \tilde{z}) = a(e, f) \tag{121}$$

As a result of (121) and the boundedness of $a$ one obtains the a posteriori estimate

$$|J(e)| \leq \|a\| \, \|e\|_V \, \|f\|_V \leq \|a\| \, \eta_u \eta_z$$

Indeed, utilizing the primal and dual residual $R_u$ and $R_z$ in $V^*$, defined by

$$R_u := F - a(\tilde{u}, \cdot) \quad \text{and} \quad R_z := J - a(\cdot, \tilde{z})$$

computable upper error bounds for $\|e\|_V \leq \eta_u$ and $\|f\|_V \leq \eta_z$ can be found by the arguments of the energy error estimators of the previous sections. This yields a computable upper error bound $\|a\| \, \eta_u \eta_z$ for $|J(e)|$ which is global, that is, the interaction of $e$ and $f$ is *not* reflected. One might therefore speculate that the upper bound is often too coarse and inappropriate for goal-oriented adaptive mesh refinement.

*5.7.3. Upper and lower bounds of $J(e)$.* Throughout the rest of this section, let the bilinear form $a$ be symmetric and positive definite; hence a scalar product with induced norm $\|\cdot\|_a$. Then, the parallelogram rule shows

$$2\,J(e) = 2\,a(e,f) = \|e+f\|_a^2 - \|e\|_a^2 - \|f\|_a^2$$

This right-hand side can be written in terms of residuals, in the spirit of (70), namely, $\|e\|_a = \|\mathrm{Res}_u\|_{V^*}$, $\|f\|_a = \|\mathrm{Res}_z\|_{V^*}$, and

$$
\begin{aligned}
\|e+f\|_a &= \|\mathrm{Res}_{u+z}\|_{V^*} \text{ for } \mathrm{Res}_{u+z} := F + J - a(\tilde{u}+\tilde{z},\cdot)\\
&= \mathrm{Res}_u + \mathrm{Res}_z \in V^*
\end{aligned}
$$

Therefore, the estimation of $J(e)$ is reduced to the computation of lower and upper error bounds for the three residuals $\mathrm{Res}_u$, $\mathrm{Res}_z$, and $\mathrm{Res}_{u+z}$ with respect to the energy norm. This illustrates that the energy error estimation techniques of the previous sections may be employed for goal-oriented error control.

For more details and examples of a refined estimation see Ainsworth and Oden (2000) and Babuška and Strouboulis (2001).

*5.7.4. Computing an approximation to $J(e)$.* An immediate consequence of (121) is

$$J(e) = R(z)$$

and hence $J(e)$ is easily computed once the dual solution $z$ of (120) is known or at least approximated to sufficient accuracy. An upper error bound for $|J(e)| = |R(z)|$ is obtained following the methodology of Becker and Rannacher (1996, 2001) and Bangerth and Rannacher (2003).

To outline this methodology, consider the residual representation formula (73) following the notation of Section 5.2.1. Suppose that $z \in H^2(\Omega)$ (e.g. for a $H^2$ regular dual problem) and let $Iz$ denote the nodal interpolant in the lowest-order finite element space $\widetilde{V}$. With some interpolation constant $C_I > 0$, there holds, for any element $T \in \mathcal{T}$,

$$h_T^{-2}\|z - Iz\|_{L_2(T)} + h_T^{-3/2}\|z - Iz\|_{L_2(\partial T)} \le C_I |z|_{H^2(T)}$$

The combination of this with (73) shows

$$
\begin{aligned}
J(e) &= \mathrm{Res}(z - Iz)\\
&= \sum_{T \in \mathcal{T}} \int_T r_T \cdot (z - Iz)\mathrm{dx} - \sum_{\mathrm{E} \in \mathcal{E}} \int_{\mathrm{E}} \mathrm{r}_{\mathrm{E}} \cdot (\mathrm{z} - \mathrm{Iz})\mathrm{ds}\\
&\le \sum_{T \in \mathcal{T}} \Big( \|r_T\|_{L_2(T)} \|z - Iz\|_{L_2(T)}\\
&\qquad + \|r_E\|_{L_2(\partial T)} \|z - Iz\|_{L_2(\partial T)} \Big)\\
&\le \sum_{T \in \mathcal{T}} C_I \left( h_T^2 \|r_T\|_{L^2(T)} + h_T^{3/2} \|r_E\|_{L_2(\partial T)} \right) |z|_{H^2(T)}
\end{aligned}
$$

The influence of the goal functional in this upper bound is through the unknown $H^2$ seminorm $|z|_{H^2(T)}$, which is to be replaced by some discrete analog based on a computed approximation $z_h$. The justification of some substitute $|D_h^2 z_h|_{L_2(T)}$ (postprocessed with some averaging technique) for $|z|_{H^2(T)}$ is through striking numerical evidence; we refer to Becker and Rannacher (2001) and Bangerth and Rannacher (2003) for details and numerical experiments.

## 6. Local Mesh Refinement

This section is devoted to the mesh-design task in the finite element method based on a priori and a posteriori information. Examples of the former type are graded meshes or geometric meshes with an a priori choice of refinement toward corner singularities briefly mentioned in Section 6.1. Examples of the latter type are adaptive algorithms for automatic mesh refinement (or mesh coarsening) strategies with a successive call of the steps

$$\text{SOLVE} \Rightarrow \text{ESTIMATE} \Rightarrow \text{MARK} \Rightarrow \text{REFINE}$$

Given the current triangulation, one has to compute the finite element solution in step SOLVE; cf. Section 7.8 for a MATLAB realization of that. The accuracy of this finite element approximation is checked in the step ESTIMATE. On the basis of the refinement indicators of Section 6.2 the step MARK identifies the elements, edges or patches in the current mesh in need of refinement (or coarsening). The new data structure is generated in step REFINE where a partition is given and a *closure algorithm* computes a triangulation described in Section 6.3. The convergence and optimality of the adaptive algorithm is discussed in Section 6.5.

### 6.1. A priori mesh design

The singularities of the exact solution of the Laplace equation on domains with corners (cf. Figure 1) are reasonably well understood and motivate the (possibly anisotropic) mesh refinement toward vertices or edges. This section aims a short introduction for two-dimensional $P_1$ finite elements–Interpolation in $h$-version Finite Element Spaces, will report on three-dimensional examples.

Given a polygonal domain with a coarse triangulation into triangles (which specify the geometry), macro elements can be used to fill the domain with graded meshes. Figure 12(a) displays a macro element described in the sequel while Figure 12(b) illustrates the resulting fine mesh for an $L$-shaped domain.

The description is restricted to the geometry on the reference element $T_{\text{ref}}$ with vertices $(0,0)$, $(1,0)$, and $(0,1)$ of Figure 12(a). The general situation is then obtained by an affine transformation illustrated in Figure 12(b). The macro element $T_{\text{ref}}$ is generated as follows: Given a grading parameter $\beta > 0$ for a grading function $g(t) = t^\beta$, and given a natural number $N$, set $\xi_j := g(j/N)$ and draw line segments aligned to the antidiagonal through $(0, \xi_j)$ and
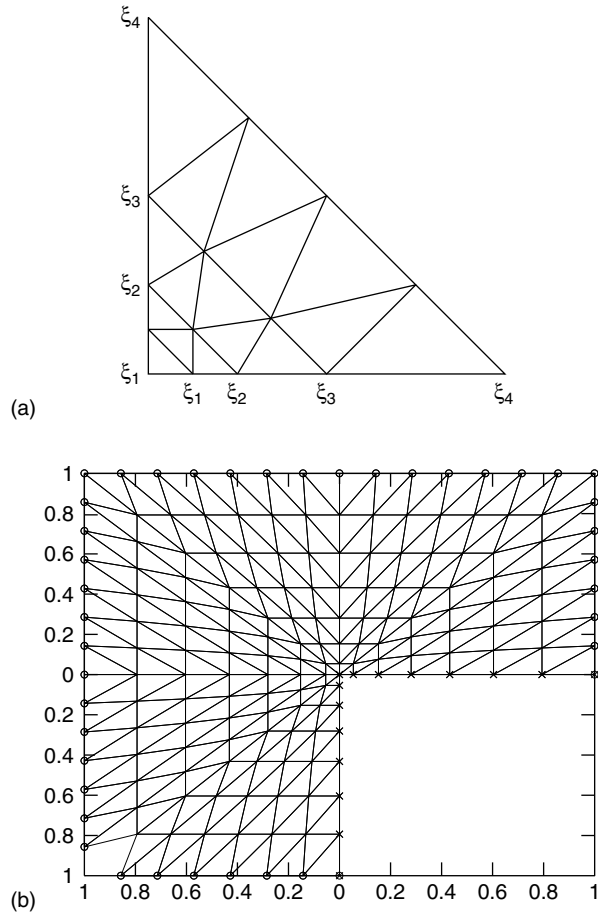
Figure 12: (a) Reference domain $T_{\mathrm{ref}}$ with graded mesh for $\beta = 3/2$ and $N = 4$. (b) Graded mesh on $L$-shaped domain with refinement toward origin and uniform refinement far away from the origin. Notice that the outer boundaries of the macro elements show a uniform distribution and so match each other in one global regular triangulation.

$(\xi_j, 0)$ for $j = 0, 1, \ldots, N$. Each of these segments is divided into $j$ uniform edges and so define the set of nodes $(0,0)$ and $\xi_j/j \, (j - k, k)$ for $k = 0, \ldots, j$ and $j = 1, \ldots, N$. The elements are then given by the vertices $\xi_j/j \, (j - k, k)$ and $\xi_j/j \, (j - k - 1, k + 1)$ aligned with the antidiagonal and the vertex $\xi_{j-1}/(j - 1) \, (j - k - 1, k)$ on the finer and $\xi_{j+1}/(j + 1) \, (j - k, k + 1)$ on the coarser neighboring segment, respectively. The finest element is $\mathrm{conv}\{(0,0), (0, \xi_1), (\xi_1, 0)\}$ of diameter $\sqrt{(2)} \, g(1/N) \approx N^{-\beta}$.

Figure 12(b) displays a triangulation of the $L$-shaped domain with a refinement toward the origin designed by a union of transformed macro elements with $\beta = 3/2$ and $N = 7$. The other vertices of the $L$-shaped domain yield higher singularities, which are not important for the first-order Courant finite element.
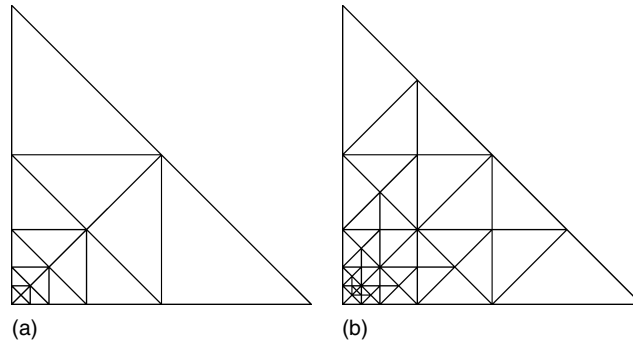
(a)                                  (b)

Figure 13: (a) Reference domain $T_{\mathrm{ref}}$ with geometric mesh for parameter $\beta = 1/2$ and $N = 4$. This mesh can also be generated by an adaptive red-green-blue refinement of Section 6.3. (b) Illustration of the closure algorithm. The refinement triangulation with 50 element domains is obtained from the mesh (a) with 18 element domains by marking one edge (namely the second along the antidiagonal) in the mesh (a).

The geometric mesh depicted in Figure 13 yields a finer refinement toward some corners of the polygonal domain. Given a parameter $\beta > 0$ in this type of triangulation, the nodes $\xi_0 := 0$ and $\xi_j := \beta^{N-j}$ for $j = 1, \ldots, N$ define antidiagonals through $(\xi_j, 0)$ and $(0, \xi_j)$, which are in turn bisected. For such a triangulation, the polynomial degrees $p_T$ on each triangle $T$ are distributed as follows: $p_T = 1$ for the two triangles $T$ with vertex $(0,0)$ and $p_T = j + 2$ for the four elements in the convex quadrilateral with the vertices $(\xi_j, 0)$, $(0, \xi_j)$, $(\xi_{j+1}, 0)$, and $(0, \xi_{j+1})$ for $j = 0, \ldots, N - 1$. Figure ?? compares experimental convergence rates of the error in $H^{-1}$-seminorm $|e|_{H^1}$ for various graded meshes for the $P_1$ finite element method, the $p$-and $hp$-finite element method, and for the adaptive algorithm of Section 6.4. The $P_1$ finite element method on graded meshes with $\beta = 3/2$, $\beta = 2$ and $h$-adaptivity recover optimality in the convergence rate as opposite to the uniform refinement ($\beta = 1$), leading only to a sub-optimal convergence due to the corner singularity. The $hp$-finite element method performs better convergence rate compared to the $p$-finite element method.

Tensor product meshes are more appropriate for smaller values of $\beta$; the one-dimensional model analysis of Babuška and Guo (1986) suggests $\beta = (2)^{1/2} - 1 \approx 0.171573$.

### 6.2. Adaptive mesh-refining algorithms

The automatic mesh refinement for regular triangulations called MARK and RE-FINE/COARSEN frequently consists of three stages:

  (i)  the marking of elements or edges for refinement;
 (ii)  the closure algorithm to ensure that the resulting triangulation is (or remains) regular;
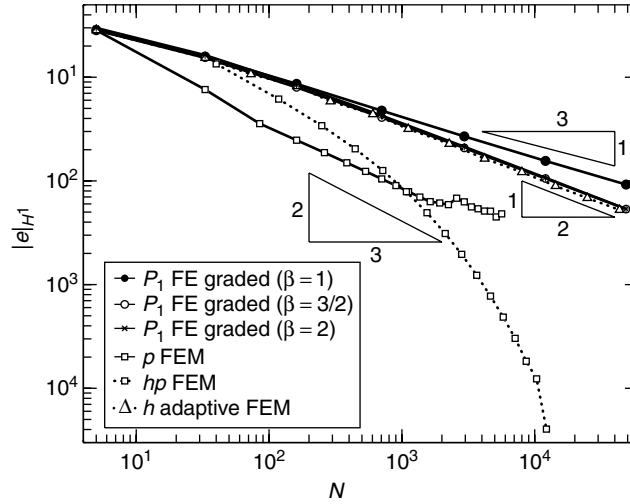(iii)  the refinement itself, i.e. the change of the underlying data structures.

Figure 14: Experimental convergence rates for various graded meshes for the $P_1$ finite element method, the $p$- and $hp$-finite element method, and for the adaptive algorithm of Section 6.5 for the Poisson problem on the $L$-shaped domain.

This section will focus on the marking strategies (i) and the subsequent one on (ii)-(iii).

In a model situation with a sum over all elements $T \in \mathcal{T}$ (or over all edges, faces, or nodes), the a posteriori error estimators of the previous section give rise to a lower or upper error bound $\eta = \sqrt{\sum_{T \in \mathcal{T}} \eta_T^2}$. The marking strategy is an algorithm selects a subset $\mathcal{M}$ of $\mathcal{T}$ called the *marked* elements; these are marked with the intention of being refined during the later refinement algorithm.

A typical algorithm computes a threshold $L$, a positive real number, and then utilizes the *refinement rule* or *marking criterion*

$$\text{mark } T \in \mathcal{T} \quad \text{if} \quad L \leq \eta_T$$

Therein, $\eta_T$ is referred to as the *refinement indicator* whereas $L$ is the *threshold*; that is,

$$\mathcal{M} := \{T \in \mathcal{T} \colon L \leq \eta_T\}$$

Typical examples for the computation of a threshold $L$ are the *maximum criterion*

$$L := \Theta \max\{\eta_T \colon T \in \mathcal{T}\}$$

or the *bulk criterion*, where $L$ is the largest value such that

$$(1 - \Theta)^2 \sum_{T \in \mathcal{T}} \eta_T^2 \leq \sum_{T \in \mathcal{M}} \eta_T^2$$

The parameter $\Theta$ is chosen with $0 \leq \Theta \leq 1$; $\Theta = 0$ corresponds to an almost uniform refinement and $\Theta = 1$ to a raw refinement of just a few elements (no refinements in the bulk criterion).
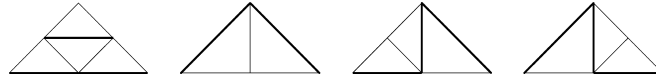
Figure 15: Red-, green-, blue-(left)- and blue-(right) refinement with reference edge on bottom of a triangle (from left to right) into four, two, and three subtriangles. The bold lines opposite the newest vertex indicate the next reference edge for further refinements.

A different strategy is possible if the error estimator gives rise to a quantitative bound of a new meshsize. For instance, the explicit error estimator can be rewritten as $\eta_R =: \|h_{\mathcal{T}} R\|_{L_2(\Omega)}$ with a given function $R \in L_2(\Omega)$ and the local mesh size $h_{\mathcal{T}}$ (when edge contributions are recast as volume contributions). Then, given a tolerance Tol and using the heuristic that $R$ would not change (at least not dramatically increase) during a refinement, the new local mesh size $h^{\mathrm{new}}$ can be calculated from the condition

$$\|h^{\mathrm{new}} R\|_{L_2(\Omega)} = \mathrm{Tol}$$

upon the equi-distribution hypothesis $h^{\mathrm{new}} \propto \mathrm{Tol}/R$. Another approach that leads to a requested mesh-size distribution is based on sharp a priori error bounds, such as $\|h_{\mathcal{T}} D^2 u\|_{L_2(\Omega)}$ where $D^2 u$ denotes the matrix of all second derivatives of the exact solution $u$. Since $D^2 u$ is unknown, it has to be approximated by postprocessing a finite element solution with some averaging technique.

The aforementioned refinement rules for the step MARK (intended for conforming FEM) ignore further decisions such as the particular type of anisotropic refinement or the increase of polynomial degrees versus the mesh refinements for $hp$-FEM.

### 6.3. Mesh refining of regular triangulations

Given a marked set of objects such as nodes, edges, faces, or elements, the refinement of the element (triangles or tetrahedra) plus further refinements (closure algorithm) for the design of regular triangulations are considered in this section.

*6.3.1. Refinement of a triangle.*   Triangular elements in two dimensions are refined into two, three, or four subtriangles as indicated in Figure 15. All these divisions are based on hidden information on some reference edge. Rivara (1984) assumed the longest edge in the triangle as base of the refinement strategies while the one below is based on the explicit marking of a reference edge. In Figure 15, the bottom edge of the original triangle acts as the *reference edge* and, in any refinement, is halved. The four divisions displayed correspond to the bisection of one (called green-refinement), of two (the two versions of blue-refinement), or of all three edges (for red-refinement) as long as the bottom edge is refined. In Figure 15, the reference edges in the generated subtriangles are drawn with a bold line.

*6.3.2. Closure algorithms.* The bisection of some set of marked elements does not always lead to a regular triangulation–the new vertices may be hanging nodes for the neighboring elements. Further refinements are necessary to make those nodes regular. The default procedure within the class of bisection algorithms is to work on a set of marked elements $\mathcal{M}(k)$.

**Closure Algorithm for Bisection.** Input a regular triangulation $\mathcal{T}^{(0)} := \mathcal{T}$ and an initial subset $\mathcal{M}^{(0)} := \mathcal{M} \subset \mathcal{T}^{(0)}$ of marked elements, set $\mathcal{Z} = \emptyset$ and $k := 0$.
While $\mathcal{M}(k) \neq \emptyset$ repeat (i)-(iv):

(i) choose some element $K$ in $\mathcal{M}(k)$ with reference edge $E$ and initiate its midpoint $z_K$ as new vertex, set $\mathcal{Z} := \mathcal{Z} \cup \{z_K\}$;

(ii) bisect $E$ and divide $K$ into $K_+$ and $K_-$ and set $\mathcal{T}^{(k+1)} := \{K_+, K_-\} \cup (\mathcal{T}^{(k)} \setminus \{K\})$;

(iii) find all elements $T_1, \ldots, T_{m_K} \in \mathcal{T}^{(k+1)}$ with han- ging node $z \in \mathcal{Z}$ (if any) and set $\mathcal{M}^{(k+1)} := \{T_1, \ldots, T_{m_K}\} \cup (\mathcal{M}(k) \setminus \{K\})$;

(iv) update $k := k + 1$ and go to (i).

Output a refined triangulation $\mathcal{T}^{(k)}$.

According to step (iii) of the closure algorithm, any termination leads to a regular triangulation $\mathcal{T}^{(k)}$. The remaining essential detail is to guarantee that there will always occur a termination via $\mathcal{M}(k) = \emptyset$ for some $k$. In the newest-vertex bisection, for instance, the reference edges are inherited in such a way that any element $K \in \mathcal{T}$, the initial regular triangulation, is refined only by subdivisions which, at most, halve each edge of $K$. Since the closure algorithm only halves some edge in $\mathcal{T}$ and prohibits any further refinements, any intermediate (irregular) $\mathcal{T}^{(k)}$ remains coarser than or equal to some regular uniform refinement $\widehat{\mathcal{T}}$ of $\mathcal{T}$. This is the main argument to prove that the closure algorithm cannot refine forever and stops after a finite number of steps.

Figure 13(b) shows an example where, given the initial mesh of Figure 13(a), only one edge, namely the second on the diagonal, is marked for refinement and the remaining refinement is induced by the closure algorithm. Nevertheless, the number of new elements can be bounded in terms of the initial triangulation and the number of marked elements (Binev, Dahmen and DeVore, 2004).

The closure algorithm for the red-green-blue refinement in 2-D is simpler when the focus is on marking of edges. One main ingredient is that each triangle $K$ is assigned a reference edge $E(K)$. If we are given a set of marked elements, let $\mathcal{M}$ denote the set of corresponding assigned reference edges.

**Closure Algorithm for Red-Green-Blue Refinement.** Input a regular triangulation $\mathcal{T}$ with a set of edges $\mathcal{E}$ and an initial subset $\mathcal{M} \subset \mathcal{E}$ of marked edges, set $k := 0$ and

$\mathcal{M}^{(0)} := \mathcal{N}^{(0)} := \mathcal{M}$.
While $\mathcal{N}^{(k)} \neq \emptyset$ repeat (i)-(iv):

(i) choose some edge $E$ in $\mathcal{N}^{(k)}$ and let $T_{\pm} \in \mathcal{T}$ denote the (at most) two triangles that share the edge $E \subset \partial T_{\pm}$;

(ii) set $\mathcal{M}^{(k+1)} := \mathcal{M}(k) \cup \{E_+, E_-\}$ with the reference edge $E_{\pm} := E(T_{\pm})$ of $T_{\pm}$;

(iii) if $\mathcal{M}^{(k+1)} = \mathcal{M}(k)$ set $\mathcal{N}^{(k+1)} := \mathcal{N}^{(k)} \setminus \{E\}$ else set $\mathcal{N}^{(k+1)} := (\mathcal{N}^{(k)} \cup \{E_+, E_-\}) \setminus \{E\}$;

(iv) update $k := k + 1$ and go to (i).

Bisect the marked edges $\{E \in \mathcal{M}(k) : E \subset \partial T\}$ of each element $T \in \mathcal{T}$ and refine $T$ by one of the red-green-blue refinement rules to generate elementwise a partition $\widehat{\mathcal{T}}$ as output.

The closure algorithm for red-green-blue refinement terminates as $\mathcal{N}^{(k)}$ is decreasing and $\mathcal{M}(k)$ is increasing and outputs a set $\widehat{\mathcal{M}} := \mathcal{M}(k)$ of marked edges with the following closure property: Any element $T \in \mathcal{T}$ with an edge in $\mathcal{M}$ satisfies $E(T) \in \mathcal{M}$, i.e. if $T$ is marked, then at least its reference edge will be halved. This property allows the application of one properly chosen refinement of Figure 15 and leads to a regular triangulation.

The reference edge $E(K)$ in the closure algorithm is assigned to each element $K$ of the initial triangulation and then is inherited according to the rules of Figure 15. For newest-vertex bisection, each triangle with vertices of global numbers $j$, $k$, and $\ell$ has the reference edge opposite to the vertex number $\max\{j, k, \ell\}$.

On the basis of refinement rules that inherit a reference edge to the generated elements, one can prove that a finite number of affine-equivalent elements domains occur.

### 6.4. Newest vertex bisection (NVB)

The newest vertex bisection for simplicial finite element domains in $\mathbb{R}^n$ for any space dimension $n = 1, 2, 3, \ldots$ dates back to Maubach (1995) and Traxler (1997); this subsection follows the description of Stevenson (2008) to which we refer for proofs.

A *tagged simplex* $(z_0, \ldots, z_n; \gamma)$ is an $(n+2)$-tuple of vertices $z_0, \ldots, z_n \in \mathbb{R}^n$ and of a type $\gamma \in \{0, \ldots, n-1\}$. The simplex $\mathrm{dom}(z_0, \ldots, z_n; \gamma) := \mathrm{conv}\{z_0, \ldots, z_n\}$ is a compact subset of $\mathbb{R}^n$ supposed to have a positive $n$-dimensional volume; in other words, the vertices $z_0, \ldots, z_n \in \mathbb{R}^n$ do not belong to an $(n-1)$-dimensional hyperplane. The *nodes* $\mathcal{N}(T)$ and *sides* $\mathcal{E}(T)$ of a tagged simplex $(z_0, \ldots, z_n; \gamma)$ are

$$\mathcal{N}(T) := \{z_0, \ldots, z_n\} \quad \text{and} \quad \mathcal{E}(T) := \{\mathrm{conv}\{z_0, \ldots, z_{j-1}, z_{j+1}, \ldots, z_n\} : j = 0, \ldots, n\}.$$

It sometimes appears convenient to identify the tagged simplex $(z_0, \ldots, z_n; \gamma)$ with its domain $T := \mathrm{dom}(z_0, \ldots, z_n; \gamma) := \mathrm{conv}\{z_0, \ldots, z_n\}$; but this suppresses further information on the $n$ different types of simplices with the same domain and with a different order of the vertices, which conduct the mesh-refinement.

The *bisection* of a tagged simplex $(z_0, \ldots, z_n; \gamma)$ bisects the *refinement edge* $\mathrm{conv}\{z_0, z_n\}$ and so generates the two tagged simplices

$$
\begin{aligned}
&\left(z_0, \frac{z_0 + z_n}{2}, z_1, \ldots, z_\gamma, z_{\gamma+1}, \ldots, z_{n-1}; \gamma'\right), \\
&\left(z_n, \frac{z_0 + z_n}{2}, z_1, \ldots, z_\gamma, z_{n-1}, \ldots, z_{\gamma+1}; \gamma'\right).
\end{aligned}
\tag{122}
$$

of the type $\gamma' := (\gamma + 1) \pmod{n}$ one higher than $\gamma$ with the convention of $n$-periodicity. (The list $z_{n-1}, \ldots, z_{\gamma+1}$ represents $z_{n-1}, z_{n-2}, z_{n-3}, \ldots, z_{\gamma+2}, z_{\gamma+1}$ and this list is empty, whence neglected, for $n - 1 < \gamma + 1$.) The two new tagged simplices (122) are called the *children* of the tagged simplex $(z_0, \ldots, z_n; \gamma)$. A tagged simplex generated from $T$ by a finite number of applications of (122) is called a descendant of (the tagged simplex) $T$. In particular, any child of some child is called *grandchild* (of a tagged simplex). Notice that, given a tagged simplex $T = (z_0, \ldots, z_n; \gamma)$, the tagged simplex

$$
T_R := (z_n, z_1, \ldots, z_\gamma, z_{n-1}, \ldots, z_{\gamma+1}, z_0; \gamma)
$$

has the same children as $T$. Notice that the type $\gamma$ does not play any role for $n \le 2$.

A *regular triangulation* $\mathcal{T}$ of a polyhedral bounded Lipschitz domain $\Omega \subset \mathbb{R}^n$ into tagged simplices is a finite set of tagged simplices with the following properties (a)–(c). The simplices cover the domain $\Omega$ in the sense that (a) their union is equal to the closure $\overline{\Omega}$ of the domain

$$
\bigcup_{T \in \mathcal{T}} \mathrm{dom}(T) = \overline{\Omega}
$$

the simplices are non-overlapping in the sense that (b) the interiors of two distinct tagged simplices $T$ and $T'$ are disjoint

$$
\mathrm{int}(\mathrm{dom}(T)) \cap \mathrm{int}(\mathrm{dom}(T')) = \emptyset
$$

(c) two non-disjoint tagged simplices $T = (y_0, \ldots, y_n; \gamma)$ and $T' := (z_0, \ldots, z_n; \gamma')$ share some hyper-surface in that there exists $0 \le j_1 < \cdots < j_N \le n$ and $0 \le k_1 < \cdots < k_N \le n$ for some $N \in \{1, \ldots, n\}$ such that

$$
\mathrm{dom}(T) \cap \mathrm{dom}(T') = \mathrm{conv}\{y_{j_1}, \ldots, y_{j_N}\} = \mathrm{conv}\{z_{k_1}, \ldots, z_{k_N}\}
\tag{123}
$$

The initial condition of the input triangulation of Stevenson, 2008 involves the concept of a reflected neighbor. Two distinct tagged simplices $T$ and $T'$ in $\mathcal{T}$ are *neighbors* if they share a side in that $\mathcal{F}(T) \cap \mathcal{F}(T') \ne \emptyset$. The regularity of $\mathcal{T}$ shows that $T$ and $T'$ share exactly the $n$ vertices $\mathcal{N}(T) \cap \mathcal{N}(T')$. If, in addition to that, the positions of the common vertices in the respective lists of vertices in $T$ and $T'$ or in $T_R$ and $T'$ coincide in all but exactly one position, then $T$ and $T'$ are called *reflected neighbors*. In other words, $N = n$ holds in (123) and, moreover, there is exactly one index $j \in \{0, \ldots, n\}$ with $y_j \ne z_j$ such that $(y_0, \ldots, y_{j-1}, y_j, y_{j+1}, \ldots, y_n; \gamma) \in \{T, T_R\}$ and $T' := (y_0, \ldots, y_{j-1}, z_j, y_{j+1}, \ldots, y_n; \gamma')$.

The following condition on the initial triangulation $\mathcal{T}_0$ is assumed throughout the newest-vertex bisection (NVB). The regular triangulation $\mathcal{T}_0$ satisfies the *matching condition* or *initial*

*condition* if the following conditions (a)–(c) hold. (a) $\mathcal{T}_0$ is a regular triangulation of $\Omega$ into tagged simplices. (b) All simplices in $\mathcal{T}_0$ are of the same type $\gamma$. (c) Any two neighbouring tagged simplices $T = (y_0, \ldots, y_n; \gamma)$ and $T' = (z_0, \ldots, z_n; \gamma)$ in $\mathcal{T}_0$ satisfy (c1)–(c2).
(c1) If $\mathrm{conv}\{y_0, y_n\} \subseteq T \cap T'$ or $\mathrm{conv}\{z_0, z_n\} \subseteq T \cap T'$, then $T$ and $T'$ are reflected neighbours.
(c2) If $\mathrm{conv}\{y_0, y_n\} \nsubseteq T \cap T' \neq \emptyset$ and $\mathrm{conv}\{z_0, z_n\} \nsubseteq T \cap T'$, then any child $S$ of $T$ and any child $S'$ of $T'$ (each understood as in (122)) are either reflected neighbors or no neighbors.

This initial condition guarantees that uniform refinements of a triangulation are regular; in particular there exist refinements of $\mathcal{T}_0$ which are regular, cf. Stevenson, 2008, Theorem 4.3 for further details. In other words, the set $\mathbb{T}$ of all such *admissible refinements* is non-trivial. Given the initial triangulation $\mathcal{T}_0$ which satisfies the aforementioned initial condition then $\mathcal{T}$ is an *admissible triangulation*, written $\mathcal{T} \in \mathbb{T}$, if $\mathcal{T}$ is a regular triangulation and if there is a finite sequence of successive bisections of the type (122) that applies to $\mathcal{T}_0$ and generates $\mathcal{T}$ in the following sense. There exists a nonnegative integer $L$ and sets $\mathcal{T}_0, \ldots, \mathcal{T}_L$ of tagged simplices ($\mathcal{T}_1, \ldots, \mathcal{T}_{L-1}$ may not be regular triangulations) such that $\mathcal{T}_L = \mathcal{T}$ and for each $\ell = 0, \ldots, L-1$ it holds $\mathcal{T}_{\ell+1} = \{T_1^{(\ell)}, T_2^{(\ell)}\} \cup \mathcal{T}_\ell \setminus \{T^{(\ell)}\}$ for exactly one $T^{(\ell)} \in \mathcal{T}_\ell$ bisected into its children $T_1^{(\ell)}$ and $T_2^{(\ell)}$ by (122).

This implied concept of successive bisections equips $\mathbb{T}$ with a partial ordering $\leq$ and defines a lattice $(\mathbb{T}, \leq)$. Hence there exists the smallest common refinement $\vee$ and the greatest common coarsening $\wedge$ of a finite number of admissible triangulation. The smallest common refinement $\mathcal{T} \otimes \mathcal{T}' := \mathcal{T} \bigvee \mathcal{T}'$ of two admissible triangulations $\mathcal{T}, \mathcal{T}' \in \mathbb{T}$, also called their *overlay*, satisfies

$$\mathrm{card}(\mathcal{T} \otimes \mathcal{T}') + \mathrm{card}(\mathcal{T}_0) \leq \mathrm{card}(\mathcal{T}) + \mathrm{card}(\mathcal{T}')$$

The mesh-refinement in adaptive finite element algorithms is usually driven by some marking followed by refinement. Given an admissible triangulation $\mathcal{T} \in \mathbb{T}$ and a subset $\mathcal{M} \subset \mathcal{T}$, let $\hat{\mathcal{T}} =: \mathrm{refine}(\mathcal{T}, \mathcal{M})$ be the smallest admissible refinement of $\mathcal{T}$ with $\mathcal{M} \cap \hat{\mathcal{T}} = \emptyset$. There exists effective algorithms to compute $\mathrm{refine}(\mathcal{T}, \mathcal{M})$ for a singleton $\mathcal{M} = \{T\}$ for $T \in \mathcal{T}$ and by successive calls of those routines also for the computation of

$$\mathrm{refine}(\mathcal{T}, \mathcal{M}) = \bigvee_{M \in \mathcal{M}} \mathrm{refine}(\mathcal{T}, \{M\})$$

The reader is referred to p.235 in Stevenson (2008) for algorithms and proofs. Two further properties of $\mathbb{T}$ will be employed in the optimality analysis outlined in Subsection 6.5 below.

The fundamental overhead control due to Binev, Dahmen and DeVore (2004) for $n = 2$ and Stevenson (2008) for $n \geq 3$ reads, in the notation of an adaptive algorithm below with $\mathcal{T}_{\ell+1} := \mathrm{refine}(\mathcal{T}_\ell, \mathcal{M}_\ell)$ for any $\mathcal{M}_\ell \subset \mathcal{T}_\ell$ and $\ell = 0, 1, 2 \ldots$, as

$$\mathrm{card}(\mathcal{T}_k) - \mathrm{card}(\mathcal{T}_0) \leq C_{BDV} \sum_{j=0}^{k-1} \mathrm{card}(\mathcal{M}_j)$$

for all $k = 0, 1, 2 \ldots$ with some constant $C_{BDV}$, which solely depends on $\mathcal{T}_0$. The example of Figure 13 illustrates that this result cannot be proven by mathematical induction over the

levels; it requires a deeper insight in the refinement of simplices and their size as well as their distances. It is a nontrivial observation of Gallistl, Schedensack and Stevenson (2014) that the number of children in each of the above calls refine($\mathcal{T}_\ell, \mathcal{M}_\ell$) is indeed controlled by card($\mathcal{T}_\ell$) times some universal constant.

### 6.5. Convergence of adaptive algorithms

The overwhelming practical success of adaptive mesh-refining algorithms has been justified in one-dimensional examples by Babuska and his collaborators in the eighties while the first convergence proof dates back to Dörfler (1996). The first optimal convergence rates are obtained in Binev, Dahmen and DeVore (2004) with an embedded coarsening step, before it became clear that the standard adaptive algorithm leads to asymptotically optimal convergence rates in Stevenson (2007). For more discussions on the history and a larger overview about the literature, the reader is referred to Carstensen, Feischl, Page and Praetorius (2014). The remaining parts of this subsection focus on an outline of abstract conditions called axioms of adaptivity that are sufficient for optimal asymptotic convergence rates of the adaptive algorithm (CAFM).

The adaptive algorithm based on collective marking is abstractly described in terms of an estimator $\eta$ which is defined for any admissible triangulation $\mathcal{T} \in \mathbb{T}$ in the notation of the previous subsection for newest vertex bisection as follows. Given any $\mathcal{T} \in \mathbb{T}$, there is a map $\eta(\mathcal{T}, \bullet) : \mathcal{T} \to [0, \infty)$ (written $\eta(\mathcal{T}, \bullet) \in [0, \infty)^{\mathcal{T}}$) which fulfills certain properties, called axioms below. In other words, $\eta(\mathcal{T}, K)$ is a computable nonnegative real number with square $\eta^2(\mathcal{T}, K) \equiv \eta(\mathcal{T}, K)^2$ for any $K \in \mathcal{T} \in \mathbb{T}$.

Throughout this subsection, an *estimator* is such a family $\eta := (\eta(\mathcal{T}, \bullet) : \mathcal{T} \in \mathbb{T}) \in \prod_{\mathcal{T} \in \mathbb{T}} [0, \infty)^{\mathcal{T}}$. The adaptive algorithm (CAFEM) driven the an estimator $\eta$ and the collective Dörfler marking reads as follows.

**CAFEM.** INPUT  initial coarse triangulation $\mathcal{T}_0$ and bulk parameter $0 < \theta < 1$

FOR $\ell = 0, 1, 2, \ldots$ DO

COMPUTE $\eta(\mathcal{T}_\ell, K)$ for all $K \in \mathcal{T}_\ell$ and their norm $\eta_\ell := \sqrt{\sum_{K \in \mathcal{T}_\ell} \eta^2(\mathcal{T}_\ell, K)}$

SELECT  a subset $\mathcal{M}_\ell \subset \mathcal{T}_\ell$ of (almost) minimal cardinality with

$$\theta \eta_\ell^2 \le \sum_{K \in \mathcal{M}_\ell} \eta^2(\mathcal{T}_\ell, K) =: \eta^2(\mathcal{T}_\ell, \mathcal{M}_\ell) \tag{124}$$

CALL  $\mathcal{T}_{\ell+1} := \text{refine}(\mathcal{T}_\ell, \mathcal{M}_\ell)$ from NVB  OD

OUTPUT  sequence of triangulations $(\mathcal{T}_\ell)$ and estimated values $(\eta_\ell)$

The convergence rates in CAFEM concern the decay of the estimated values $\eta_\ell$ towards zero and the growth of the number $N_\ell := \text{card}(\mathcal{T}_\ell) - \text{card}(\mathcal{T}_0)$ of extra element domains on the level $\ell$ towards infinity as $\ell \to \infty$. The *asymptotic convergence rate is at least $s > 0$* provided that

$$\sup_{\ell \in \mathbb{N}_0} (N_\ell + 1)^s \eta_\ell < \infty$$

the supremum of all $s$ with this property is *the asymptotic convergence rate*. The comparison of the asymptotic convergence rate of CAFEM is with respect to an optimal choice of a triangulation with respect to the values of the estimators. Given any $N \in \mathbb{N}_0$, this amounts in the minimization of the quantity

$$\eta(\mathcal{T}) := \eta(\mathcal{T}, \mathcal{T}) := \sqrt{\sum_{K \in \mathcal{T}} \eta^2(\mathcal{T}, K)}$$

for all admissible triangulations $\mathcal{T} \in \mathbb{T}(N)$, where

$$\mathbb{T}(N) := \{ \mathcal{T} \in \mathbb{T} : \mathrm{card}(\mathcal{T}) - \mathrm{card}(\mathcal{T}_0) \le N \}$$

The *optimal convergence rate of an estimator* $\eta \in \prod_{\mathcal{T} \in \mathbb{T}} [0, \infty)^{\mathcal{T}}$ is the supremum of all $s > 0$ such that

$$\sup_{N \in \mathbb{N}} (N + 1)^s \min_{\mathcal{T} \in \mathbb{T}(N)} \eta(\mathcal{T}) < \infty$$

Based on an observation due to Gallistl, Schedensack and Stevenson (2014), it is obvious that the asymptotic convergence rate of the CAFEM is always smaller than or equal to the optimal convergence rate of an estimator. Optimality of the adaptive algorithm means that equality holds.

   Given (1) the set of admissible triangulations $\mathbb{T}$ with the partial ordering $\le$ for refinement from the previous subsection, given (2) an estimator $\eta \in \prod_{\mathcal{T} \in \mathbb{T}} [0, \infty)^{\mathcal{T}}$, and given (3) a function $\delta : \{ (\mathcal{T}, \widehat{\mathcal{T}}) \in \mathbb{T}^2 : \mathcal{T} \le \widehat{\mathcal{T}} \} \to [0, \infty)$ suppose the following four conditions with universal positive constants $\Lambda_1, \ldots, \Lambda_4, \Lambda_3' < \infty$ and $\rho_2 < 1$. In (A1)–(A3), $\widehat{\mathcal{T}}$ is an arbitrary refinement of $\mathcal{T}$ and the subset $\mathcal{R} \subset \mathcal{T}$ of $\mathcal{T}$ includes the refined triangles $\mathcal{T} \setminus \widehat{\mathcal{T}} \subset \mathcal{R}$ but not too many further triangles such that $\mathrm{card}(\mathcal{R}) \le \Lambda_3' \mathrm{card}(\mathcal{T} \setminus \widehat{\mathcal{T}})$. The CAFEM algorithm outputs the triangulations $\mathcal{T}_\ell$ and estimated values $\eta_\ell := \eta(\mathcal{T}_\ell)$ for any $\ell = 0, 1, 2, \ldots$ which arise in (A4).

$$|\eta(\widehat{\mathcal{T}}, \mathcal{T} \cap \widehat{\mathcal{T}}) - \eta(\mathcal{T}, \mathcal{T} \cap \widehat{\mathcal{T}})| \le \Lambda_1 \delta(\mathcal{T}, \widehat{\mathcal{T}}) \tag{A1}$$

$$\eta(\widehat{\mathcal{T}}, \widehat{\mathcal{T}} \setminus \mathcal{T}) \le \rho_2 \eta(\mathcal{T}, \mathcal{T} \setminus \widehat{\mathcal{T}}) + \Lambda_2 \delta(\mathcal{T}, \widehat{\mathcal{T}}) \tag{A2}$$

$$\delta^2(\mathcal{T}, \widehat{\mathcal{T}}) \le \Lambda_3 \eta^2(\mathcal{T}, \mathcal{R}) \tag{A3}$$

$$\sum_{k=\ell}^{\infty} \delta^2(\mathcal{T}_k, \mathcal{T}_{k+1}) \le \Lambda_4 \eta_\ell^2 \quad \text{for all } \ell \in \mathbb{N}_0 \tag{A4}$$

The axioms (A1)-(A4) are sufficient for convergence $\eta_\ell \to 0$ as $\ell \to \infty$ and allow an optimal asymptotic convergence rate for bulk parameters $\theta < \theta_0 := (1 + \Lambda_1^2 \Lambda_3)^{-1}$. For any $s > 0$ and $\theta < \theta_0$, the equivalence

$$\sup_{\ell \in \mathbb{N}_0} (N_\ell + 1)^s \eta_\ell \approx \sup_{N \in \mathbb{N}} (N + 1)^s \min_{\mathcal{T} \in \mathbb{T}(N)} \eta(\mathcal{T})$$

is contained in Carstensen, Feischl, Page and Praetorius (2014) and Carstensen and Rabus (2016) to which we refer to an overview of the literature as well as for many examples. Unlike the former contributions in the literature, the critical bound $\theta_0$ does not depend on any notion of efficiency. In fact, the efficiency is solely required in a global version to lead to optimal convergence rates of the error as well. Easy examples cover lowest-order conforming,

nonconforming or mixed adaptive finite element methods with $\mathcal{R} \subset \mathcal{T} \setminus \widehat{\mathcal{T}}$ in (A3) and the error $\delta(\mathcal{T}, \widehat{\mathcal{T}}) := ||\widehat{p} - p||_{L^2(\Omega)}$ of the flux approximation $p$ (resp. $\widehat{p}$) computed for the triangulation $\mathcal{T}$ (resp. its refinement $\widehat{\mathcal{T}}$). For nonconforming or mixed schemes, the quasi-orthogonality (A4) requires a little modification called (A4$_\varepsilon$).

A generalization with separate marking appears necessary for least-squares or mixed finite element methods when the data approximation term does not allow any weight by the mesh-size. Details on the algorithm SAFEM and the generalization of the axioms for optimal asymptotic convergence rates appear in Carstensen and Rabus (2016).

## 7. Other Aspects

In this section, we discuss briefly several topics in finite element methodology. Some of the discussions involve the Sobolev space $W_p^k(\Omega)$ ($1 \le p \le \infty$), which is the space of functions in $L_p(\Omega)$ whose weak derivatives up to order $k$ also belong to $L_p(\Omega)$, with the norm

$$\|v\|_{W_p^k(\Omega)} = \left( \sum_{|\alpha| \le k} \left\| \frac{\partial^\alpha v}{\partial x^\alpha} \right\|_{L_p(\Omega)} \right)^{1/p}$$

for $1 \le p < \infty$ and

$$\|v\|_{W_\infty^k(\Omega)} = \max_{|\alpha| \le k} \left\| \left( \frac{\partial^\alpha v}{\partial x^\alpha} \right) \right\|_{L_\infty(\Omega)}$$

For $1 \le p < \infty$, the seminorm $\left( \sum_{|\alpha|=k} \|(\partial^\alpha v/\partial x^\alpha)\|_{L_p(\Omega)} \right)^{1/p}$ will be denoted by $|v|_{W_p^k(\Omega)}$, and the seminorm $\max_{|\alpha|=k} \|(\partial^\alpha v/\partial x^\alpha)\|_{L_\infty(\Omega)}$ will be denoted by $|v|_{W_\infty^k(\Omega)}$.

### 7.1. Nonsymmetric/indefinite problems

The results in Section 4 can be extended to the case where the bilinear form $a(\cdot, \cdot)$ in the weak problem (1) is nonsymmetric and/or indefinite due to lower order terms in the partial differential equation. We assume that $a(\cdot, \cdot)$ is bounded (cf. (2)) on the closed subspace $V$ of the Sobolev space $H^m(\Omega)$ and replace (3) by the condition that

$$a(v, v) + L\|v\|_{L_2(\Omega)}^2 \ge C_3 \|v\|_{H^m(\Omega)}^2 \qquad \forall\, v \in V \tag{125}$$

where $L$ is a positive constant.

**Example 15.** *Let $a(\cdot, \cdot)$ be defined by*

$$a(v_1, v_2) = \int_\Omega \nabla v_1 \cdot \nabla v_2 \mathrm{dx} + \sum_{j=1}^d \int_\Omega b_j(x) \frac{\partial v_1}{\partial x_j} v_2 \mathrm{dx}$$

$$+ \int_\Omega c(x) v_1\, v_2 \mathrm{dx} \tag{126}$$

*for all $v_1, v_2 \in H^1(\Omega)$, where $b_j(x)$ $(1 \le j \le d)$, $c(x) \in L_\infty(\Omega)$. If we take $V = \{v \in H^1(\Omega) : v\big|_\Gamma = 0\}$ and $F$ is defined by (6), then (1) is the weak form of the nonsymmetric boundary value problem*

$$-\Delta u + \sum_{j=1}^{d} b_j \frac{\partial u}{\partial x_j} + cu = f, \quad u \ = \ 0 \quad on \quad \Gamma,$$

$$\frac{\partial u}{\partial n} \ = \ 0 \quad on \quad \partial\Omega \setminus \Gamma \tag{127}$$

*and the coercivity condition (125) follows from the well-known Gårding's inequality (Agmon, 1965).*

Unlike the symmetric positive definite case, we need to assume that the weak problem (1) has a unique solution, and that the adjoint problem is also uniquely solvable, that is, given any $G \in V^*$ there is a unique $w \in V$ such that

$$a(v, w) = G(v) \qquad \forall\, v \in V \tag{128}$$

Furthermore, we assume that the solution $w$ of (128) enjoys some elliptic regularity when $G(v) = (g, v)_{L_2(\Omega)}$ for $g \in L_2(\Omega)$, i.e., $w \in H^{m+\alpha}(\Omega)$ for some $\alpha > 0$ and

$$\|w\|_{H^{m+\alpha}(\Omega)} \le C\|g\|_{L_2(\Omega)} \tag{129}$$

Let $\mathcal{T}$ be a triangulation of $\Omega$ with mesh size $h_\mathcal{T} = \max_{T \in \mathcal{T}} \operatorname{diam} T$ and $V_\mathcal{T} \subset V$ be a finite element space associated with $\mathcal{T}$ such that the following approximation property is satisfied:

$$\inf_{v \in V_\mathcal{T}} \|w - v\|_{H^m(\Omega)} \le \epsilon_\mathcal{T} \|w\|_{H^{m+\alpha}(\Omega)} \qquad \forall\, w \in H^{m+\alpha}(\Omega) \tag{130}$$

where

$$\epsilon_\mathcal{T} \downarrow 0 \quad as \quad h_\mathcal{T} \downarrow 0 \tag{131}$$

The discrete problem is then given by (54).

Following Schatz (1974) the well-posedness of the discrete problem and the error estimate for the finite element approximate solution can be addressed simultaneously. Assume for the moment that $u_\mathcal{T} \in V_\mathcal{T}$ is a solution of (54). Then we have

$$a(u - u_\mathcal{T}, v) = 0 \qquad \forall\, v \in V_\mathcal{T} \tag{132}$$

We use (132) and a duality argument to estimate $\|u - u_\mathcal{T}\|_{L_2(\Omega)}$ in terms of $\|u - u_\mathcal{T}\|_{H^m(\Omega)}$. Let $w \in V$ satisfy

$$a(v, w) = (u - u_\mathcal{T}, v)_{L_2(\Omega)} \qquad \forall\, v \in V_\mathcal{T} \tag{133}$$

We obtain, from (2), (132), and (133), the following analog of (24):

$$\|u - u_\mathcal{T}\|_{L_2(\Omega)}^2 = a(u - u_\mathcal{T}, w)$$

$$\le C \left( \inf_{v \in V_\mathcal{T}} \|w - v\|_{H^m(\Omega)} \right) \|u - u_\mathcal{T}\|_{H^m(\Omega)} \tag{134}$$

and hence, by (129) and (130),

$$\|u - u_{\mathcal{T}}\|_{L_2(\Omega)} \le \epsilon_{\mathcal{T}} \|u - u_{\mathcal{T}}\|_{H^m(\Omega)} \tag{135}$$

It follows from (125) and (135) that

$$\|u - u_{\mathcal{T}}\|^2_{H^m(\Omega)} \le a(u - u_{\mathcal{T}}, u - u_{\mathcal{T}}) + C\epsilon^2_{\mathcal{T}}\|u - u_{\mathcal{T}}\|^2_{H^m(\Omega)}$$

which together with (131) implies, for $h_{\mathcal{T}}$ sufficiently small,

$$\|u - u_{\mathcal{T}}\|^2_{H^m(\Omega)} \le a(u - u_{\mathcal{T}}, u - u_{\mathcal{T}}) \tag{136}$$

For the special case where $F = 0$ and $u = 0$, any solution $u_{\mathcal{T}}$ of the homogeneous discrete problem

$$a(u_{\mathcal{T}}, v) = 0 \qquad \forall\, v \in V_{\mathcal{T}}$$

must satisfy, by (136),

$$\|u_{\mathcal{T}}\|^2_{H^m(\Omega)} \le 0$$

We conclude that any solution of the homogeneous discrete problem must be trivial and hence the discrete problem (19) is uniquely solvable provided $h_{\mathcal{T}}$ is sufficiently small. Under this condition, we also obtain immediately from (2), (132), and (136), the following analog of (22):

$$\|u - u_{\mathcal{T}}\|_{H^m(\Omega)} \le C \inf_{v \in V_{\mathcal{T}}} \|u - v\|_{H^m(\Omega)} \tag{137}$$

Concrete error estimates now follow from (137), (134), and the results in Section 3.3.

### 7.2. Nonconforming finite elements

When the finite element space $FE_{\mathcal{T}}$ defined by (32) does not belong to the Sobolev space $H^m(\Omega)$ where the weak problem (1) is posed, it is referred to as a *nonconforming* finite element space. Nonconforming finite element spaces are more flexible and are useful for problems with constrains where conforming finite element spaces are more difficult to construct.

**Example 16.** *(Triangular Nonconforming Elements) Let $K$ be a triangle. If the set $\mathcal{N}_K$ consists of evaluations of the shape functions at the midpoints of the edges of $K$ (Figure 16a), then $(K, P_1, \mathcal{N}_K)$ is the nonconforming $P_1$ element of Crouzeix and Raviart (Crouzeix and Raviart, 1973). It is the simplest triangular element that can be used to solve the incompressible Stokes equation. We refer the readers to Brenner (2015) for many applications of this interesting element.*

*If the set $\mathcal{N}_K$ consists of evaluations of the shape functions at the vertices of $K$ and the evaluations of the normal derivatives of the shape functions at the midpoints of the edges of $K$ (Figure 16b), then $(K, P_2, \mathcal{N}_K)$ is the Morley element (Morley, 1968; Shi, 1990). It is the simplest triangular element that can be used to solve the plate bending problem. Higher dimensional analogs of Morley can be found in Ruas (1988) and Wang and Xu (2006).*
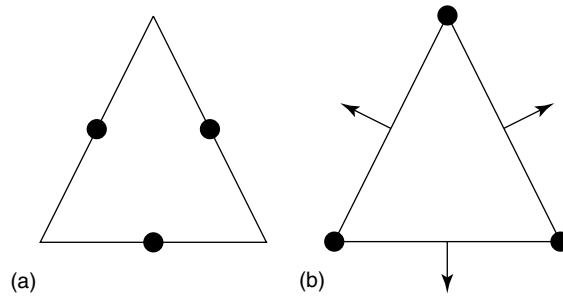
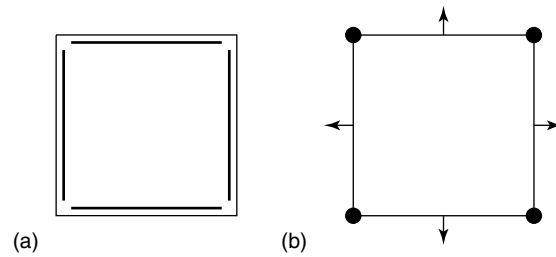Figure 16: Triangular nonconforming finite elements.



Figure 17: Rectangular nonconforming finite elements.

**Example 17.** *(Rectangular Nonconforming Elements) Let $K$ be a rectangle. If $\mathcal{P}_K$ is the space spanned by the functions $1, x_1, x_2$ and $x_1^2 - x_2^2$ and the set $\mathcal{N}_K$ consists of the mean values of the shape functions on the edges of $K$, then $(K, \mathcal{P}_K, \mathcal{N}_K)$ is the rotated $Q_1$ element of Rannacher and Turek (Rannacher and Turek, 1992) (Figure 17(a), where the thick lines represent mean values over the edges). It is the simplest rectangular element that can be used to solve the incompressible Stokes equation.*

*If $\mathcal{P}_K$ is the space spanned by the functions $1, x_1, x_2, x_1^2, x_1 x_2, x_2^2, x_1^2 x_2$ and $x_1 x_2^2$ and the set $\mathcal{N}_K$ consists of evaluations of the shape functions at the vertices of $K$ and evaluations of the normal derivatives at the midpoints of the edges (Figure 17b), then $(K, \mathcal{P}_K, \mathcal{N}_K)$ is the incomplete $Q_2$ element (Shi, 1986). It is the simplest rectangular element that can be used to solve the plate bending problem.*

**Remark** 17. Poincaré-Friedrichs and Korn's inequalities are responsible for coercivity in Examples 1-3. They have been extended to piecewise $H^1$ functions (Brenner, 2003; Brenner and Sung, 2015), piecewise $H^2$ functions (Brenner, Wang and Zhao, 2004), and piecewise $H^1$ vector fields (Brenner, 2004; Mardal and Winther, 2006; Brenner and Sung, 2015). These results can be applied in particular to nonconforming finite element spaces.

Consider the weak problem (1) for a symmetric positive definite boundary value problem, where $F$ is defined by (6) for a function $f \in L_2(\Omega)$. Let $V_\mathcal{T}$ be a nonconforming finite element space associated with the triangulation $\mathcal{T}$. We assume that there is a (mesh-dependent)

symmetric bilinear form $a_{\mathcal{T}}(\cdot, \cdot)$ defined on $V + V_{\mathcal{T}}$ such that (i) $a_{\mathcal{T}}(v, v) = a(v, v)$ for $v \in V$, (ii) $a_{\mathcal{T}}(\cdot, \cdot)$ is positive definite on $V_{\mathcal{T}}$. The discrete problem for the nonconforming finite element method reads: Find $u_{\mathcal{T}} \in V_{\mathcal{T}}$ such that

$$a_{\mathcal{T}}(u_{\mathcal{T}}, v) = \int_{\Omega} f v \mathrm{dx} \qquad \forall\, v \in V_{\mathcal{T}} \tag{138}$$

**Example 18.** *The Poisson problem in Example 1 can be solved by the nonconforming $P_1$ finite element method in which the finite element space is $V_{\mathcal{T}} = \{v \in L_2(\Omega) : v|_T \in P_1(T)$ for every triangle $T \in \mathcal{T}$, $v$ is continuous at the midpoints of the edges of $\mathcal{T}$ and $v$ vanishes at the midpoints of the edges of $\mathcal{T}$ along $\Gamma\}$ and the bilinear form $a_{\mathcal{T}}$ is defined by*

$$a_{\mathcal{T}}(v_1, v_2) = \sum_{T \in \mathcal{T}} \int_T \nabla v_1 \cdot \nabla v_2 \mathrm{dx} \tag{139}$$

Below we will discuss the *a priori* error analysis of nonconforming finite element methods and refer the readers to Dari, Duran, Padra and Vampa (1996), Carstensen and Hoppe (2006), Becker, Mao and Shi (2010) and Hu, Shi and Xu (2012) for the *a posteriori* analysis.

*7.2.1. Convergence analysis of nonconforming finite element methods: standard approach.* The nonconforming Ritz-Galerkin method (138) can be analyzed as follows. Let $\tilde{u}_{\mathcal{T}} \in V_{\mathcal{T}}$ be defined by

$$a_{\mathcal{T}}(\tilde{u}_{\mathcal{T}}, v) = a_{\mathcal{T}}(u, v) \qquad \forall\, v \in V_{\mathcal{T}}$$

Then we have

$$\|u - \tilde{u}_{\mathcal{T}}\|_{a_{\mathcal{T}}} = \inf_{v \in V_{\mathcal{T}}} \|u - v\|_{a_{\mathcal{T}}}$$

where $\|w\|_{a_{\mathcal{T}}} = (a_{\mathcal{T}}(w, w))^{1/2}$ is the nonconforming energy norm defined on $V + V_{\mathcal{T}}$, and we arrive at the following generalization (Berger, Scott and Strang, 1972) of (21):

$$
\begin{aligned}
\|u - u_{\mathcal{T}}\|_{a_{\mathcal{T}}} &\leq \|u - \tilde{u}_{\mathcal{T}}\|_{a_{\mathcal{T}}} + \|\tilde{u}_{\mathcal{T}} - u_{\mathcal{T}}\|_{a_{\mathcal{T}}} \\
&= \inf_{v \in V_{\mathcal{T}}} \|u - v\|_{a_{\mathcal{T}}} + \sup_{v \in V_{\mathcal{T}} \setminus \{0\}} \frac{a_{\mathcal{T}}(\tilde{u}_{\mathcal{T}} - u_{\mathcal{T}}, v)}{\|v\|_{a_{\mathcal{T}}}} \\
&= \inf_{v \in V_{\mathcal{T}}} \|u - v\|_{a_{\mathcal{T}}} + \sup_{v \in V_{\mathcal{T}} \setminus \{0\}} \frac{a_{\mathcal{T}}(u - u_{\mathcal{T}}, v)}{\|v\|_{a_{\mathcal{T}}}}
\end{aligned}
\tag{140}
$$

**Remark** 18. The second term on the right-hand side of (140), which vanishes in the case of conforming Ritz-Galerkin methods, measures the *consistency errors* of nonconforming methods.

As an example, we analyze the nonconforming $P_1$ finite element method for the Poisson problem in Example 18. For simplicity we assume $\Gamma = \partial\Omega$. For each $T \in \mathcal{T}$, we define an interpolation operator $\Pi_T : H^1(T) \to P_1(T)$ by

$$(\Pi_T \zeta)(m_E) = \frac{1}{h_E} \int_E \zeta \mathrm{ds}$$

where $m_E$ is the midpoint for the edge $E$ of $T$ and $h_E$ is the diameter of $E$. The interpolation operator $\Pi_T$ satisfies the estimate (43) for $m = 1$, and they can be pieced together to form an interpolation operator $\Pi : H^1(\Omega) \to V_{\mathcal{T}}$. Since the solution $u$ of (5) belongs to $H^{1+\alpha(T)}(T)$, where $\frac{1}{2} < \alpha(T) \le 1$ (Grisvard, 1985; Dauge, 1988), the first term on the right-hand side of (140) satisfies the estimate

$$
\begin{aligned}
\inf_{v \in V_{\mathcal{T}}} \|u - v\|_{a_{\mathcal{T}}} &\le \|u - \Pi u\|_{a_{\mathcal{T}}} \\
&\le C \left( \sum_{T \in \mathcal{T}} (\operatorname{diam} T)^{2\alpha(T)} |u|^2_{H^{1+\alpha(T)}(T)} \right)^{1/2}
\end{aligned}
\tag{141}
$$

where the constant $C$ depends only on the minimum angle in $\mathcal{T}$.

To analyze the second term on the right-hand side of (140), we write, using (5), (138), and (139),

$$
a_{\mathcal{T}}(u - u_{\mathcal{T}}, v) = - \sum_{E \in \mathcal{E}(\mathcal{T})} \int_E \frac{\partial u}{\partial n_E}[v]_E \mathrm{d}s
\tag{142}
$$

where $\mathcal{E}(\mathcal{T})$ is the set of edges in $\mathcal{T}$ that are not on $\Gamma$, $n_E$ is a unit vector normal to $E$, and $[v]_E = v_+ - v_-$ is the jump of $v$ across $E$ ($n_E$ is pointing from the minus side to the plus side and $v = 0$ outside $\Omega$). Note that (142) is well-defined because $u$ belongs to $H^{1+\alpha}(\Omega)$ for some $\alpha \in (\frac{1}{2}, 1]$ .

Since $[v]_E$ vanishes at the midpoint of $E \in \mathcal{E}(\mathcal{T})$, we have

$$
\int_E \frac{\partial u}{\partial n_E}[v]_E \mathrm{d}s = \int_E \frac{\partial(u - p)}{\partial n_E}[v]_E \mathrm{d}s \qquad \forall\, p \in P_1
\tag{143}
$$

Let $\mathcal{T}_E = \{T \in \mathcal{T} : E \subset \partial T\}$. It follows from (143), the trace theorem and the Bramble-Hilbert lemma (cf. Remark 10.) that

$$
\begin{aligned}
\left| \int_e \frac{\partial u}{\partial n_E}[v]_E \mathrm{d}s \right| &\le C \inf_{p \in P_1} \left[ \sum_{T \in \mathcal{T}_E} \left( |u - p|^2_{H^1(T)} \right. \right. \\
&\quad \left. \left. + (\operatorname{diam} T)^{2\alpha(T)} |u|^2_{H^{1+\alpha(T)}(T)} \right) \right]^{1/2} \left( \sum_{T \in \mathcal{T}_E} |v|^2_{H^1(T)} \right)^{1/2} \\
&\le C \left( \sum_{T \in \mathcal{T}_E} (\operatorname{diam} T)^{2\alpha(T)} |u|^2_{H^{1+\alpha(T)}(T)} \right)^{1/2} \left( \sum_{T \in \mathcal{T}_E} |v|^2_{H^1(T)} \right)^{1/2}
\end{aligned}
\tag{144}
$$

We conclude from (142) and (144) that

$$
\sup_{v \in V_{\mathcal{T}}} \frac{a_{\mathcal{T}}(u - u_{\mathcal{T}}, v)}{\|v\|_{a_{\mathcal{T}}}} \le C \left( \sum_{T \in \mathcal{T}} (\operatorname{diam} T)^{2\alpha(T)} |u|^2_{H^{1+\alpha(T)}(T)} \right)^{1/2}
\tag{145}
$$

where the constant $C$ depends only on the minimum angle in $\mathcal{T}$.

Combining (140), (141), and (145) we have the following analog of (57)

$$\sum_{T \in \mathcal{T}} |u - u_{\mathcal{T}}|_{H^1(T)}^2 \le C \sum_{T \in \mathcal{T}} (\operatorname{diam} T)^{2\alpha(T)} |u|_{H^{1+\alpha(T)}(T)}^2 \tag{146}$$

**Remark** 19. Estimates of $\|u - u_{\mathcal{T}}\|_{L_2(\Omega)}$ can also be obtained for nonconforming finite element methods (Crouzeix and Raviart, 1973). There is also a close connection between certain nonconforming methods and mixed methods (Arnold and Brezzi, 1982).

*7.2.2. Convergence analysis of nonconforming finite element methods: nonstandard approach.* The analysis of the $P_1$ nonconforming finite element method in Section 7.2.1 relies on the expression (142), which is not well-defined for $u \in H^1(\Omega)$. Therefore we had to use the elliptic regularity $u \in H^{1+\alpha}(\Omega)$ for some $\alpha \in (\frac{1}{2}, 1]$ in order to proceed. This is different from the convergence analysis of the conforming $P_1$ finite element method where the derivation of the error estimate (21) does not use any information on the regularity of $u$ other than the fact that it belongs to $H^1(\Omega)$.

The following estimate was established in Gudi (2010) without using any regularity of $u$ beyond $H^1(\Omega)$:

$$\|u - u_{\mathcal{T}}\|_{a_{\mathcal{T}}} \le C \Big( \inf_{v \in V_{\mathcal{T}}} \|u - v\|_{a_{\mathcal{T}}} + \sum_{T \in \mathcal{T}} \operatorname{Osc}(f, \mathcal{T}) \Big) \tag{147}$$

where

$$\operatorname{Osc}(f, \mathcal{T}) = \Big( \sum_{T \in \mathcal{T}} \operatorname{osc}(f, T)^2 \Big)^{\frac{1}{2}}$$

and $\operatorname{osc}(f, T)$ is defined in (80). The estimate (147) shows that Cea's lemma is valid for nonconforming finite element methods up to data oscillations. It puts the convergence analysis of nonconforming finite element methods on the same footing as that for conforming finite element methods.

The proof of (147) uses local efficiency estimates in *a posteriori* error analysis and an enriching $E_{\mathcal{T}} : V_{\mathcal{T}} \longrightarrow H_0^1(\Omega)$ with the following property:

$$\sum_{T \in \mathcal{T}} h_T^{-2} \|v - E_{\mathcal{T}} v\|_{L_2(T)}^2 + |E_{\mathcal{T}} v|_{H^1(\Omega)}^2 \le C \|v\|_{a_{\mathcal{T}}}^2 \tag{148}$$

for all $v \in V_{\mathcal{T}}$. The operator $E_{\mathcal{T}}$ can be constructed by using the $P_2$ Lagrange finite element space associated with $\mathcal{T}$ and averaging (Brenner, 1994).

We begin with the following analog of (140) :

$$\|u - u_{\mathcal{T}}\|_{a_{\mathcal{T}}} \le \|u - v\|_{a_{\mathcal{T}}} + \sup_{w \in V_{\mathcal{T}} \setminus \{0\}} \frac{a_{\mathcal{T}}(v - u_{\mathcal{T}}, w)}{\|w\|_{a_{\mathcal{T}}}} \tag{149}$$

for any $v \in V_{\mathcal{T}}$, and rewrite the numerator on the right-hand side of (149) as

$$a_{\mathcal{T}}(v - u_{\mathcal{T}}, w) = a_{\mathcal{T}}(v, w - E_{\mathcal{T}} w) + a_{\mathcal{T}}(v - u, E_{\mathcal{T}} w)$$

$$+ \int_{\Omega} f(E_{\mathcal{T}} w - w) \mathrm{d}x \tag{150}$$

Using integration by parts and the continuity of $w$ at the midpoints of the edges, the first term on the right-hand side of (150) can be written as

$$a_{\mathcal{T}}(v, w - E_{\mathcal{T}}w) = \sum_{E \in \mathcal{E}(\mathcal{T})} \int_E [\partial v/\partial n]_E \{w - E_{\mathcal{T}}w\}_E \mathrm{ds} \tag{151}$$

where $[\partial v/\partial n]_E$ is the jump of the normal derivative of $v$ across the edge $E$ and $\{w - E_{\mathcal{T}}w\}_E$ is the average of $w - E_{\mathcal{T}}w$ across $E$. Note that the integration by parts that led to (151) is legitimate because $v$, $w$ and $E_{\mathcal{T}}w$ are piecewise polynomial functions.

It follows from (151) that

$$a_{\mathcal{T}}(v, w - E_{\mathcal{T}}w) \leq \Big( \sum_{E \in \mathcal{E}(\mathcal{T})} h_E \|[\partial v/\partial n]_E\|^2_{L_2(E)} \Big)^{\frac{1}{2}}$$

$$\times \Big( \sum_{E \in \mathcal{E}(\mathcal{T})} h_E^{-1} \|\{w - E_{\mathcal{T}}w\}_E\|^2_{L_2(E)} \Big)^{\frac{1}{2}}$$

$$\leq \Big( \sum_{E \in \mathcal{E}(\mathcal{T})} h_E \|[\partial v/\partial n]_E|^2_{L_2(E)} \Big)^{\frac{1}{2}} \tag{152}$$

$$\times \Big( \sum_{T \in \mathcal{T}} h_T^{-2} \|w - E_{\mathcal{T}}w\|^2_{L_2(T)} \Big)^{\frac{1}{2}}$$

The second and third terms on the right-hand side of (150) are bounded by

$$a_{\mathcal{T}}(v - u, E_{\mathcal{T}}w) \leq \|v - u\|_{a_{\mathcal{T}}} |E_{\mathcal{T}}w|_{H^1(\Omega)} \tag{153}$$

$$\int_\Omega f(E_{\mathcal{T}}w - w)\mathrm{dx} \leq \Big( \sum_{T \in \mathcal{T}} h_T^2 \|f\|^2_{L_2(T)} \Big)^{\frac{1}{2}}$$

$$\times \Big( \sum_{T \in \mathcal{T}} h_T^{-2} \|w - E_{\mathcal{T}}w\|^2_{L_2(T)} \Big)^{\frac{1}{2}} \tag{154}$$

Finally we observe that the estimates

$$\Big( \sum_{T \in \mathcal{T}} h_T^2 \|f\|^2_{L_2(T)} \Big)^{\frac{1}{2}} \leq C \big( \|u - v\|_{a_{\mathcal{T}}} + \mathrm{Osc}(f, \mathcal{T}) \big) \tag{155}$$

$$\Big( \sum_{E \in \mathcal{E}(\mathcal{T})} h_E \|[\partial v/\partial n]_E\|^2_{L_2(E)} \Big)^{\frac{1}{2}} \leq C \big( \|u - v\|_{a_{\mathcal{T}}} + \mathrm{Osc}(f, \mathcal{T}) \big) \tag{156}$$

can be established by the same bubble function techniques that led to (78) and (79).

The estimate (147) follows from (148)-(156).

**Remark** 20. Since the derivation of (147) uses both techniques from *a priori* analysis and *a posteriori* analysis, convergence analysis based on (147) is also known as the *medius* analysis.

**Remark** 21. It follows from (147) and the density of $C_c^\infty(\Omega)$ in $H_0^1(\Omega)$ that the nonconforming $P_1$ finite element method converges to $u$ in $\|\cdot\|_{a_\mathcal{T}}$ as the mesh size of $\mathcal{T}$ decreases to 0. Therefore the *medius* analysis establishes the convergence of the $P_1$ nonconforming finite element method without using any elliptic regularity theory. The elliptic regularity of $u$ is only needed if we want to obtain a convergence rate. Indeed $O(h^\alpha)$ convergence follows immediately from (147) and $u \in H^{1+\alpha}(\Omega)$.

### 7.3. Effects of numerical integration

The explicit form of the finite element equation (54) involves the evaluations of integrals which, in general, cannot be computed exactly. Thus, the effects of numerical integration must be taken into account in the error analysis. We will illustrate the ideas in terms of finite element methods for simplicial triangulations. The readers are referred to Davis and Rabinowitz (1984) for a comprehensive treatment of numerical integration.

Consider the second order elliptic boundary value problem (5) in Example 1 on a bounded polyhedral domain $\Omega \subset \mathbb{R}^d$ ($d = 2$ or 3). Let $\mathcal{T}$ be a simplicial triangulation of $\Omega$ such that $\Gamma$ is a union of the edges (faces) of $\mathcal{T}$ and $V_\mathcal{T} \subset V$ be the corresponding $P_n$ Lagrange finite element space.

In Section 4, the finite element approximate solution $u_\mathcal{T} \in V_\mathcal{T}$ is defined by (54). But in practice, the integral $F(v) = \int_\Omega fv\mathrm{d}x$ is evaluated by a quadrature scheme and the approximate solution $u_\mathcal{T} \in V_\mathcal{T}$ is actually defined by

$$a(u_\mathcal{T}, v) = F_\mathcal{T}(v) \qquad \forall\, v \in V_\mathcal{T} \tag{157}$$

where $F_\mathcal{T}(v)$ is the result of applying the quadrature scheme to the integral $F(v)$.

More precisely, let $D \in \mathcal{T}$ be arbitrary and $\Phi_D : S \to \overline{D}$ be an affine homeomorphism from the standard (closed) simplex $S$ onto $D$. It follows from a change of variables that

$$\int_D fv\mathrm{d}x = \int_S (\det J_{\Phi_D})(f \circ \Phi_D)(v \circ \Phi_D)\mathrm{d}\hat{x}$$

where without loss of generality $\det J_{\Phi_D}$ (the determinant of the Jacobian matrix of $\Phi$) is assumed to be a positive number. The integral on $S$ is evaluated by a quadrature scheme $I_S$ and the right-hand side of (157) is then given by

$$F_h(v) = \sum_{D \in \mathcal{T}} I_S\big((\det J_{\Phi_D})(f \circ \Phi_D)(v \circ \Phi_D)\big) \tag{158}$$

The error $u - u_\mathcal{T}$ can be estimated by the following analog of (140):

$$\|u - u_\mathcal{T}\|_a \le \inf_{v \in V_\mathcal{T}} \|u - v\|_a + \sup_{v \in V_\mathcal{T} \setminus \{0\}} \frac{a(u - u_\mathcal{T}, v)}{\|v\|_a} \tag{159}$$

The first term on the right-hand side of (159) can be estimated by $\|u - \Pi_\mathcal{T} u\|_a$ as in Section 4.1. The second term on the right-hand side of (159) measures the effect of numerical quadrature.

Below we give conditions on the quadrature scheme $w \mapsto I_S(w)$ and the function $f$ so that the magnitude of the quadrature error is identical with that of the optimal interpolation error for the finite element space.

We assume that the quadrature scheme $w \mapsto I_S(w)$ has the following properties:

$$|I_S(w)| \quad \leq \quad C_S \max_{\hat{x} \in S} |w(\hat{x})| \qquad \forall\, w \in C^0(S), \tag{160}$$

$$I_S(w) \quad = \quad \int_S w \mathrm{d}\hat{x} \qquad \forall\, w \in P_{2n-2} \tag{161}$$

We also assume that $f \in W_q^n(\Omega)$ such that $q \geq 2$ and $n > d/q$, which implies, in particular, by the Sobolev embedding theorem that $f \in C^0(\Omega)$ so that (158) makes sense.

Under these conditions it can be shown (Ciarlet, 1978) by using the Bramble-Hilbert lemma on $S$ that

$$\left| \int_D f v \mathrm{d}x - \mathrm{I_S}\big( (\det \mathrm{J}_{\Phi_\mathrm{D}})(\mathrm{f} \circ \Phi_\mathrm{D})(\mathrm{v} \circ \Phi_\mathrm{D}) \big) \right|$$
$$\leq C (\operatorname{diam} D)^n |D|^{(1/2)-(1/q)} \|f\|_{W_q^n(D)} |v|_{H^1(D)} \quad \forall\, v \in V_{\mathcal{T}} \tag{162}$$

where the positive constant $C$ depends only on the shape regularity of $D$. It then follows from (1), (158), (162) and Hölder's inequality that

$$|a(u - u_{\mathcal{T}}, v)| \quad = \quad \left| \int_\Omega f v \mathrm{d}x - \mathrm{F_h(v)} \right|$$
$$\leq \quad C h_{\mathcal{T}}^n \|f\|_{W_q^n(\Omega)} \|v\|_{H^1(\Omega)} \qquad \forall\, v \in V_{\mathcal{T}} \tag{163}$$

We conclude from (159) and (163) that

$$\|u - u_{\mathcal{T}}\|_a \leq \|u - \Pi_{\mathcal{T}} u\|_a + C h_{\mathcal{T}}^n \|f\|_{W_q^n(\Omega)} \tag{164}$$

We see by comparing (43) and (164) that the overall accuracy of the finite element method is not affected by the numerical integration.

We now consider a general symmetric positive definite elliptic boundary problem whose variational form is defined by

$$a(w, v) = \sum_{i,j=1}^d \int_\Omega a_{ij}(x) \frac{\partial w}{\partial x_i} \frac{\partial v}{\partial x_j} \mathrm{d}x + \int_\Omega b(x) w v \mathrm{d}x \tag{165}$$

where $a_{ij}, b \in W_\infty^2(\Omega)$, $b \geq 0$ on $\Omega$ and there exists a positive constant $c$ such that

$$\sum_{i,j=1}^d a_{ij}(x) \xi_i \xi_j \geq c |\xi|^2 \qquad \forall\, x \in \Omega \quad \text{and} \quad \xi \in \mathbb{R}^d \tag{166}$$

In this case, the bilinear form (165) must also be evaluated by numerical integration and the approximation solution $u_{\mathcal{T}} \in V_{\mathcal{T}}$ to (1) is defined by

$$a_{\mathcal{T}}(u_{\mathcal{T}}, v) = F_{\mathcal{T}}(v) \qquad \forall\, v \in V_{\mathcal{T}} \tag{167}$$

where $a_{\mathcal{T}}(w,v)$ is the result of applying the quadrature scheme $I_S$ to the pull-back of

$$\sum_{i,j=1}^{d} \int_D a_{ij}(x) \frac{\partial w}{\partial x_i} \frac{\partial v}{\partial x_j} \mathrm{dx} + \int_D b(x) wv \mathrm{dx}$$

on the standard simplex $S$ under the affine homeomorphism $\Phi_D$, over all $D \in \mathcal{T}$.

The error $u - u_{\mathcal{T}}$ can be estimated under (160), (161) and the additional condition that

$$g \geq 0 \text{ on } \bar{S} \implies I_S(g) \geq 0 \tag{168}$$

Indeed (161), (166), (168) and the sign of $b$ imply that

$$a_{\mathcal{T}}(v,v) \geq c|v|^2_{H^1(\Omega)} \qquad \forall v \in V_{\mathcal{T}} \tag{169}$$

and we have the following analog of (159)

$$|u - u_{\mathcal{T}}|_{H^1(\Omega)} \leq \inf_{v \in V_{\mathcal{T}}} |u - v|_{H^1(\Omega)}$$

$$+ \frac{1}{c} \left\{ \sup_{v \in V_{\mathcal{T}} \setminus \{0\}} \frac{a_{\mathcal{T}}(u,v) - a(u,v)}{|v|_{H^1(\Omega)}} + \sup_{v \in V_{\mathcal{T}} \setminus \{0\}} \frac{a(u,v) - a_{\mathcal{T}}(u_{\mathcal{T}},v)}{|v|_{H^1(\Omega)}} \right\} \tag{170}$$

The first term on the right-hand side of (170) is dominated by $|u - \Pi_{\mathcal{T}} u|_{H^1(\Omega)}$. Since $a(u,v) - a_{\mathcal{T}}(u_{\mathcal{T}},v) = \int_{\Omega} fv \mathrm{dx} - \mathrm{F}_{\mathcal{T}}(\mathrm{v})$, the third term is controlled by the estimate (163). The second term, which measures the effect of numerical quadrature on the bilinear form $a(\cdot, \cdot)$, is controlled by the estimate

$$\begin{aligned} |a_{\mathcal{T}}(u,v) - a(u,v)| &\leq C \sum_{D \in \mathcal{T}} (\operatorname{diam} D)^{\alpha(D)} \\ &\times |u|_{H^{1+\alpha(D)}(D)} |v|_{H^1(D)} \end{aligned} \tag{171}$$

provided the solution $u$ belongs to $H^{1+\alpha(D)}(D)$ for each $D \in \mathcal{T}$ and $1/2 < \alpha(D) \leq 1$. The estimate (171) follows from (160), (161) and the Bramble-Hilbert lemma, and the positive constant $C$ in (171) depends only on the $W^2_{\infty}$ norms of $a_{ij}$ and $b$ and the shape regularity of $\mathcal{T}$. Again we see by comparing (56), (164), and (171) that the overall accuracy of the finite element method is not affected by the numerical integration.

**Remark** 22. For problems that exhibit the phenomenon of locking, the choice of a lower order quadrature scheme in the evaluation of the stiffness matrix may help alleviate the effect of locking (Malkus and Hughes, 1978).

### 7.4. Curved domains

So far, we have restricted the discussion to polygonal (polyhedral) domains. In this section, we consider the second-order elliptic boundary value problem (5) on a domain $\Omega \subset \mathbb{R}^2$ with a curved boundary. For simplicity, we assume that $\Gamma = \partial \Omega$.
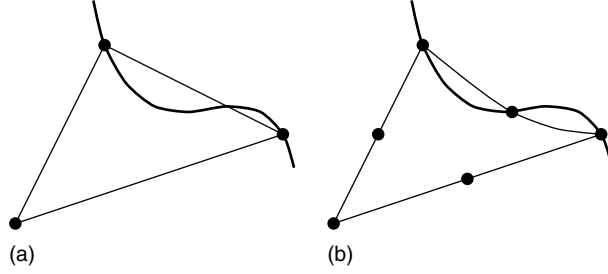
Figure 18: Triangulations for curved domains.

First, we consider the $P_1$ Lagrange finite element method for a domain $\Omega$ with a $C^2$ boundary. We approximate $\Omega$ by a polygonal domain $\Omega_h$ on which a simplicial triangulation $\mathcal{T}_h$ of mesh size $h$ is imposed. We assume that the vertices of $\mathcal{T}_h$ belong to the closure of $\Omega$. A typical triangle in $\mathcal{T}_h$ near $\partial\Omega$ is depicted in Figure 18(a).

Let $V_h \subset H_0^1(\Omega_h)$ be the $P_1$ finite element space associated with $\mathcal{T}_h$. The approximate solution $u_h \in V_h$ for (5) is then defined by

$$a_h(u_h, v) = F_h(v) \qquad \forall\, v \in V_h \tag{172}$$

where

$$a_h(w, v) = \int_{\Omega_h} \nabla w \cdot \nabla v \mathrm{d}x \qquad \forall\, w, v \in H^1(\Omega_h) \tag{173}$$

and $F_h(v)$ represented the result of applying a numerical quadrature scheme to the integral $\int_{\Omega_h} \tilde{f} v \mathrm{d}x$ (cf. (158) with $f$ and $\mathcal{T}$ replaced by $\tilde{f}$ and $\mathcal{T}_h$). Here $\tilde{f}$ is an extension of $f$ to $\mathbb{R}^2$. We assume that the numerical scheme uses only the values of $\tilde{f}$ at the nodes of $\mathcal{T}_h$ and hence the discrete problem (172) is independent of the extension $\tilde{f}$.

We assume that $f \in W_q^1(\Omega)$ for $2 < q < \infty$ and hence $u \in W_q^3(\Omega)$ by elliptic regularity. Let $\tilde{u} \in W_q^3(\mathbb{R}^2)$ be an extension of $u$ such that

$$\|\tilde{u}\|_{W_q^3(\mathbb{R}^2)} \le C_\Omega \|u\|_{W_q^3(\Omega)} \le C_\Omega \|f\|_{W_q^1(\Omega)} \tag{174}$$

We can then take $\tilde{f} = -\Delta \tilde{u} \in W_q^1(\mathbb{R}^2)$ to be the extension appearing in the definition of $F_h$.

The error $\tilde{u} - u_h$ over $\Omega_h$ can be estimated by the following analog of (140):

$$\|\tilde{u} - u_h\|_{a_h} \le \inf_{v \in V_h} \|\tilde{u} - v\|_{a_h} + \sup_{v \in V_h \backslash \{0\}} \frac{a_h(\tilde{u} - u_h, v)}{\|v\|_{a_h}} \tag{175}$$

The first term on the right-hand side is dominated by $\|\tilde{u} - \Pi_h \tilde{u}\|_{a_h}$ where $\Pi_h$ is the nodal interpolation operator. The second term is controlled by

$$\begin{aligned} \left| a_h(\tilde{u} - u_h, v) \right| &= \left| \int_{\Omega_h} \tilde{f} v \mathrm{d}x - F_h(v) \right| \\ &\le Ch \|\tilde{f}\|_{W_q^1(\Omega_h)} \|v\|_{H^1(\Omega_h)} \end{aligned} \tag{176}$$

which is a special case of (163), provided that the conditions (160) and (161) (with $n = 1$) on the numerical quadrature scheme are satisfied.

Combining (43) and (174)-(176) we see that

$$|\tilde{u} - u_h|_{H^1(\Omega_h)} = \|\tilde{u} - u_h\|_{a_h} \le Ch\|f\|_{W_q^1(\Omega)} \tag{177}$$

that is, the $P_1$ finite element method for the curved domain retains the optimal $O(h)$ accuracy.

The approximation of $\Omega_h$ to $\Omega$ can be improved if we replace straight-edge triangles by triangles with a curved edge (Figure 18(b)). This can be achieved by the *isoparametric* finite element methods. We will illustrate the idea using the $P_2$ Lagrange element.

Let $\Omega_h$, an approximation of $\Omega$, be the union of straight-edge triangles (in the interior of $\Omega_h$) and triangles with one curved edge (at the boundary of $\Omega_h$), which form a triangulation $\mathcal{T}_h$ of $\Omega_h$. The finite element associated with the interior triangles is the standard $P_2$ Lagrange element. For a triangle $D$ at the boundary, we assume that there is a homeomorphism $\Phi_D$ from the standard simplex $S$ onto $D$ such that $\Phi_D(\hat{x}) = (\Phi_{D,1}(\hat{x}), \Phi_{D,2}(\hat{x}))$ where $\Phi_{D,1}(\hat{x})$ and $\Phi_{D,2}(\hat{x})$ are quadratic polynomials in $\hat{x} = (\hat{x}_1, \hat{x}_2)$. The space of shape functions $\mathcal{P}_{\overline{D}}$ is then defined by

$$\mathcal{P}_{\overline{D}} = \{v \in C^\infty(D) : v \circ \Phi_D \in P_2(S)\} \tag{178}$$

that is, the functions in $\mathcal{P}_{\overline{D}}$ are quadratic polynomials in the *curvilinear* coordinates on $D$ induced by $\Phi_D^{-1}$. The set $\mathcal{N}_{\overline{D}}$ of nodal variables consist of pointwise evaluations of the shape functions at the nodes corresponding to the nodes of the $P_2$ element on $S$ (cf. Figures 2 and 18) under the map $\Phi_D$. We assume that all such nodes belong to $\overline{\Omega}$ and the nodes on the curved edge of $D$ belong to $\partial\Omega$ (cf. Figure 18).

In other words, the finite element $(\overline{D}, \mathcal{P}_{\overline{D}}, \mathcal{N}_{\overline{D}})$ is pulled back to the $P_2$ Lagrange finite element on $\bar{S}$ under $\Phi_D$. It is called an isoparametric element because the components of the parameterization map $\Phi_D$ are shape functions of the $P_2$ element on $\bar{S}$. The corresponding finite element space defined by (32) (with $\mathcal{T}$ replaced by $\mathcal{T}_h$) is a subspace of $H^1(\Omega_h)$ that contains all the continuous piecewise linear polynomials with respect to $\mathcal{T}_h$. By setting the nodal values on $\partial\Omega_h$ to be zero, we have a finite element space $V_h \subset H_0^1(\Omega_h)$. The discrete problem for $u_h \in V_h$ is then defined by

$$\tilde{a}_h(u_h, v) = F_h(v) \tag{179}$$

where the numerical quadrature scheme in the definition of $F_h$ involves only the nodes of the finite element space so that the discrete problem is independent of the choice of the extension of $f$ and the variational form $\tilde{a}_h(\cdot, \cdot)$ is obtained from $a_h(\cdot, \cdot)$ by the numerical quadrature scheme.

We assume that $f \in W_q^2(\Omega)$ for $1 < q < \infty$ and hence $u \in W_q^4(\Omega)$ by elliptic regularity (assuming that $\Omega$ has a $C^3$ boundary). Let $\tilde{u} \in W_q^4(\mathbb{R}^2)$ be an extension of $u$ such that

$$\|\tilde{u}\|_{W_q^4(\mathbb{R}^2)} \le C_\Omega\|u\|_{W_q^4(\Omega)} \le C_\Omega\|f\|_{W_q^2(\Omega)} \tag{180}$$

Under the condition (168), we have

$$\tilde{a}_h(v, v) \ge c|v|_{H^1(\Omega_h)}^2 \qquad \forall\, v \in V_h \tag{181}$$

and the error of $\tilde{u} - u_h$ over $\Omega_h$ can be estimated by the following analog of (170)

$$|\tilde{u} - u_h|_{H^1(\Omega_h)} \le \inf_{v \in V_h} |\tilde{u} - v|_{H^1(\Omega_h)}$$

$$+ \frac{1}{c} \left\{ \sup_{v \in V_h \setminus \{0\}} \frac{\tilde{a}_h(\tilde{u}, v) - a_h(\tilde{u}, v)}{|v|_{H^1(\Omega_h)}} + \sup_{v \in V_h \setminus \{0\}} \frac{a_h(\tilde{u}, v) - \tilde{a}_h(u_h, v)}{|v|_{H^1(\Omega_h)}} \right\} \tag{182}$$

The analysis of the terms on the right-hand side of (182) involves the shape regularity of a curved triangle, which can be defined as follows. Let $\mathcal{A}_D$ be the affine map that agrees with $\Phi_D$ at the vertices of the standard simplex $S$, and $\widetilde{D}$ be the image of $S$ under $\mathcal{A}_D$. ($\widetilde{D}$ is the triangle in Figure 18(a), while $D$ is the curved triangle in (b).) The shape regularity of the curved triangle $D$ is measured by the aspect ratio $\gamma(\widetilde{D})$ (cf. (29)) of the straight-edged triangle $\widetilde{D}$ and the parameter $\kappa(D)$ defined by

$$\kappa(D) = \max \left\{ h^{-1} |\Phi_D \circ \mathcal{A}_D^{-1}|_{W_\infty^1(D)}, h^{-2} |\Phi_D \circ \mathcal{A}_D^{-1}|_{W_\infty^2(D)} \right\} \tag{183}$$

and we can take the aspect ratio $\gamma(D)$ to be the maximum of $\gamma(\tilde{D})$ and $\kappa(D)$. Note that in the case where $D = \widetilde{D}$ the parameter $\kappa(D) = 0$ and $\gamma(D) = \gamma(\widetilde{D})$.

The first term on the right-hand side of (182) is dominated by $\|\tilde{u} - \Pi_h \tilde{u}\|_{a_h}$, where $\Pi_h$ is the nodal interpolation operator. Note that, by using the Bramble-Hilbert lemma on $S$ and scaling, we have the following generalization of (43):

$$|\tilde{u} - \Pi_D \tilde{u}|_{H^1(D)} \le C (\operatorname{diam} D)^2 \|\tilde{u}\|_{H^3(D)} \tag{184}$$

where $\Pi_D$ is the element nodal interpolation operator and the constant $C$ depends only on an upper bound of $\gamma(D)$, and hence

$$\|\tilde{u} - \Pi_h \tilde{u}_h\|_{H^1(\Omega_h)} \le C h^2 \|\tilde{u}\|_{H^3(\Omega_h)} \tag{185}$$

In order to analyze the third term on the right-hand side of (182), we take $\tilde{f} = -\Delta \tilde{u}$ and impose the conditions (160) and (161) (with $n = 2$) on the numerical quadrature scheme. We then have the following special case of (163):

$$\begin{aligned} |a_h(\tilde{u}, v) - \tilde{a}_h(u_h, v)| &= \left| \int_{\Omega_h} \tilde{f} v \mathrm{d}x - F_h(v) \right| \\ &\le C h^2 \|\tilde{f}\|_{W_q^2(\Omega_h)} |v|_{H^1(\Omega_h)} \end{aligned} \tag{186}$$

Similarly, the second term on the right-hand side of (182), which measures the effect of numerical integration on the variational form $a_h(\cdot, \cdot)$, is controlled by the estimate

$$|\tilde{a}_h(\tilde{u}, v) - a_h(\tilde{u}, v)| \le C h^2 |\tilde{u}|_{H^3(\Omega_h)} |v|_{H^1(\Omega_h)} \qquad \forall\, v \in V_h \tag{187}$$

Combining (182) and (185)-(187) we have

$$\|\tilde{u} - u_h\| \le C h^2 \|f\|_{W_q^2(\Omega)} \tag{188}$$

where $C$ depends only on an upper bound of $\{\gamma(D) : D \in \mathcal{T}_h\}$ and the constants in (180). Therefore, the $P_2$ isoparametric finite element method retains the optimal $O(h^2)$ accuracy. On the other hand, if only straight-edged triangles are used in the construction of $\Omega_h$, then the accuracy of the $P_2$ Lagrange finite element method is only of order $O(h^{3/2})$ (Strang and Berger, 1971).

The discussion above can be generalized to higher-order isoparametric finite element methods, higher dimensions, and elliptic problems with variable coefficients (Ciarlet, 1978).

**Remark** 23. Estimates such as (175) and (177) are useful only when a sequence of domains $\Omega_{h_i}$ with corresponding triangulations $\mathcal{T}_{h_i}$ can be constructed so that $h_i \downarrow 0$ and the aspect ratios of all the triangles (straight or curved) in the triangulations remain bounded. We refer the readers to Scott (1973) for 2-D constructions and to Lenoir (1986) for the 3-D case.

**Remark** 24. Other finite element methods for curved domains can be found in Zlámal (1973, 1974), Scott (1975), and Bernardi (1989).

**Remark** 25. Let $\Omega_i$ be a sequence of convex polygons approaching the unit disc. The displacement of the simply supported plate on $\Omega_i$ with unit loading does not converge to the displacement of the simply supported plate on the unit disc (also with unit loading) as $i \to \infty$. This is known as Babuška's plate paradox (Babuška and Pitkäranta, 1990). It shows that numerical solutions obtained by approximating a curved domain with polygonal domains, in general, do not converge to the solution of a fourth-order problem defined on the curved domain. We refer the readers to Mansfield (1978) for the construction of finite element spaces that are subspaces of $H^2(\Omega)$.

### 7.5. Pointwise estimates

Besides the estimates in $L_2$-based Sobolev spaces discussed in Section 4, there also exist a priori error estimates for finite element methods in $L_p$-based Sobolev spaces with $p \neq 2$. In particular, error estimates in the $L_\infty$-based Sobolev spaces can provide pointwise error estimates. Below we describe some results for second-order elliptic boundary value problems with homogeneous Dirichlet boundary conditions.

In the one-dimensional case (Wheeler, 1973) where $\Omega$ is an interval, the finite element solution $u_\mathcal{T}$ for a given triangulation $\mathcal{T}$ with mesh size $h_\mathcal{T}$ satisfies

$$\|u - u_\mathcal{T}\|_{L_\infty(\Omega)} \leq C h_\mathcal{T}^n |u|_{W_\infty^n(\Omega)} \tag{189}$$

provided the solution $u$ of (5) belongs to $W_\infty^n(\Omega)$ and the finite element space contains all the piecewise polynomial functions of degree $\leq n - 1$. The estimate (189) also holds in higher dimensions (Douglas, Dupont and Wheeler, 1974) in the case where $\Omega$ is a product of intervals and $u_\mathcal{T}$ is the solution in the $Q_{n-1}$ finite element space of Example 7.

For a two-dimensional convex polygonal domain (Natterer, 1975; Scott, 1976; Nitsche, 1977), the estimate (189) holds in the case where $n \geq 3$ and $u_\mathcal{T}$ is the $P_{n-1}$ triangular finite element

solution for a general triangulation $\mathcal{T}$. In the case where $u_{\mathcal{T}}$ is the $P_1$ finite element solution, (189) is replaced by

$$\|u - u_{\mathcal{T}}\|_{L_\infty(\Omega)} \leq C h_{\mathcal{T}}^2 |\ln h_{\mathcal{T}}| \, |u|_{W_\infty^2(\Omega)} \tag{190}$$

$L_\infty$ estimates for general triangulations on polygonal domains with reentrant corners and higher-dimensional domains can be found in Schatz and Wahlbin (1978, 1979, 1982) and Schatz (1998). The estimate (190) was also established in Gastaldi and Nochetto (1987) for the Crouzeix-Raviart nonconforming $P_1$ element of Example 16.

It is also known (Rannacher and Scott, 1982; Brenner and Scott, 2002) that

$$|u - u_{\mathcal{T}}|_{W_\infty^1(\Omega)} \leq C \inf_{v \in V_{\mathcal{T}}} C |u - v|_{W_\infty^1(\Omega)} \tag{191}$$

where $\Omega$ is a convex polygonal domain in $\mathbb{R}^2$ and $u_{\mathcal{T}}$ is the $P_n$ $(n \geq 1)$ triangular finite element solution obtained from a general triangulation $\mathcal{T}$ of $\Omega$. Optimal order estimates for $|u - u_{\mathcal{T}}|_{W_\infty^1(\Omega)}$ can be derived immediately from (191). Extension of (191) to higher dimensions can be found in Schatz and Wahlbin (1995).

### 7.6. Interior estimates and pollution effects

Let $\Omega$ be the $L$-shaped polygon in Figure 1. The solution $u$ of the Poisson problem (5) on $\Omega$ with homogeneous Dirichlet boundary condition is singular near the reentrant corner and $u \notin H^2(\Omega)$. Consequently, the error estimate $|u - u_{\mathcal{T}}|_{H^1(\Omega)} \leq C h_{\mathcal{T}} \|f\|_{L_2(\Omega)}$ does not hold for the $P_1$ triangular finite element solution $u_{\mathcal{T}}$ associated with a quasi-uniform triangulation $\mathcal{T}$ of mesh size $h_{\mathcal{T}}$.

However, $u$ does belong to $H^2(\Omega_\delta)$ where $\Omega_\delta$ is the subset of the points of $\Omega$ whose distances to the reentrant corner are strictly greater than the positive number $\delta$. Therefore, it is possible that

$$|u - u_{\mathcal{T}}|_{H^1(\Omega_\delta)} \leq C h_{\mathcal{T}} \|f\|_{L_2(\Omega)} \tag{192}$$

That the estimate (192) indeed holds is a consequence of the following *interior estimate* (Nitsche and Schatz, 1974):

$$|u - u_{\mathcal{T}}|_{H^1(\Omega_\delta)} \leq C \left( \inf_{v \in V_{\mathcal{T}}} |u - v|_{H^1(\Omega_{\delta/2})} + \|u - u_{\mathcal{T}}\|_{L_2(\Omega_{\delta/2})} \right) \tag{193}$$

where $V_{\mathcal{T}} \subset H_0^1(\Omega)$ is the $P_1$ triangular finite element space. Interior estimates in various Sobolev norms can be established for subdomains of general $\Omega$ in $\mathbb{R}^d$ and general finite elements. We refer the readers to Wahlbin (1991) for a survey of such results and to Schatz (2000) for some recent developments.

On the other hand, since $u_{\mathcal{T}}$ is obtained by solving a global system that involves the nodal values near the reentrant corner of the $L$-shaped domain, the effect of the singularity at the reentrant corner can propagate into other parts of $\Omega$. This is known as the *pollution effect* and is reflected, for example, by the following estimate (Wahlbin, 1984):

$$\|u - u_{\mathcal{T}}\|_{L_2(\Omega)} \geq C h_{\mathcal{T}}^{2\beta} \tag{194}$$

where $\beta = \pi/(3\pi/2) = 2/3$. Similar estimates can also be established for other Sobolev norms.

### 7.7. Superconvergence

Let $u_\mathcal{T}$ be the finite element solution of a second-order elliptic boundary value problem. Suppose that the space of shape functions on each element contains all the polynomials of degree $\leq n$ but not all the polynomials of degree $n+1$. Then the $L_\infty$ norm of the error $u - u_\mathcal{T}$ is at most of order $h^{n+1}$, even if the solution $u$ is smooth. However, the absolute value of $u - u_\mathcal{T}$ at certain points can be of order $h^{n+1+\sigma}$ for some $\sigma > 0$. This is known as the phenomenon of *superconvergence* and such points are the superconvergence points for $u_\mathcal{T}$. Similarly, a point where the absolute value of a derivative of $u - u_\mathcal{T}$ is of order $h^{n+\sigma}$ is a superconvergence point for the derivative of $u_\mathcal{T}$.

The division points of a partition $\mathcal{T}$ for a two point boundary value problem with smooth coefficients provides the simplest example of superconvergence points. Let $u_\mathcal{T}$ be the finite element solution from the $P_n$ Lagrange finite element space. Since the Green's function $G_p$ associated with a division point $p$ is continuous in the interval $\Omega$ and smooth on the two subintervals divided by $p$, we have (Douglas and Dupont, 1974)

$$\begin{aligned}
|(u - u_\mathcal{T})(p)| &= |a(u - u_\mathcal{T}, G_p)| \\
&= |a(u - u_\mathcal{T}, G_p - \Pi_\mathcal{T}^N G_p)| \\
&\leq C\|u - u_\mathcal{T}\|_{H^1(\Omega)}\|G_p - \Pi_\mathcal{T}^N G_p)\|_{H^1(\Omega)} \leq Ch^{2n}
\end{aligned}$$

provided that $u$ is sufficiently smooth. Therefore, $p$ is a superconvergence point for $u_\mathcal{T}$ if $n \geq 2$.

For general superconvergence results in various dimensions, we refer the readers to Křížek and Neittaanmäki (1987), Chen and Huang (1995), Wahlbin (1995), Lin and Yan (1996), Schatz, Sloan and Wahlbin (1996), Křížek, Neittaanmäki and Stenberg (1998), Chen (1999), and Babuška and Strouboulis (2001).

### 7.8. Finite element program in 8 lines of MATLAB

It is the purpose of this section to introduce a short (two-dimensional) $P_1$ finite element program.

The data for a given triangulation $\mathcal{T} = \{T_1, \ldots, T_m\}$ into triangles with a set of nodes $\mathcal{N} = \{z_1, \ldots, z_n\}$ are described in user-specified matrices called `c4n` and `n4e`. Figure 19 displays a triangulation with $m$ triangles and $n$ nodes as well as a fixed enumeration and the corresponding data. The coordinates of the nodes $z_k = (x_k, y_k)$ ($d$ real components in general) are stored in the $k$th row of the two-dimensional matrix `c4n`. Each element $T_j = \text{conv}\{z_k, z_\ell, z_m\}$ is represented by the labels of its vertices $(k, \ell, m)$ stored in the $j$th row of the two-dimensional matrix `n4e`. The chosen permutation of $(k, \ell, m)$ describes the element in a counterclockwise orientation. Homogeneous Dirichlet conditions are prescribed on the
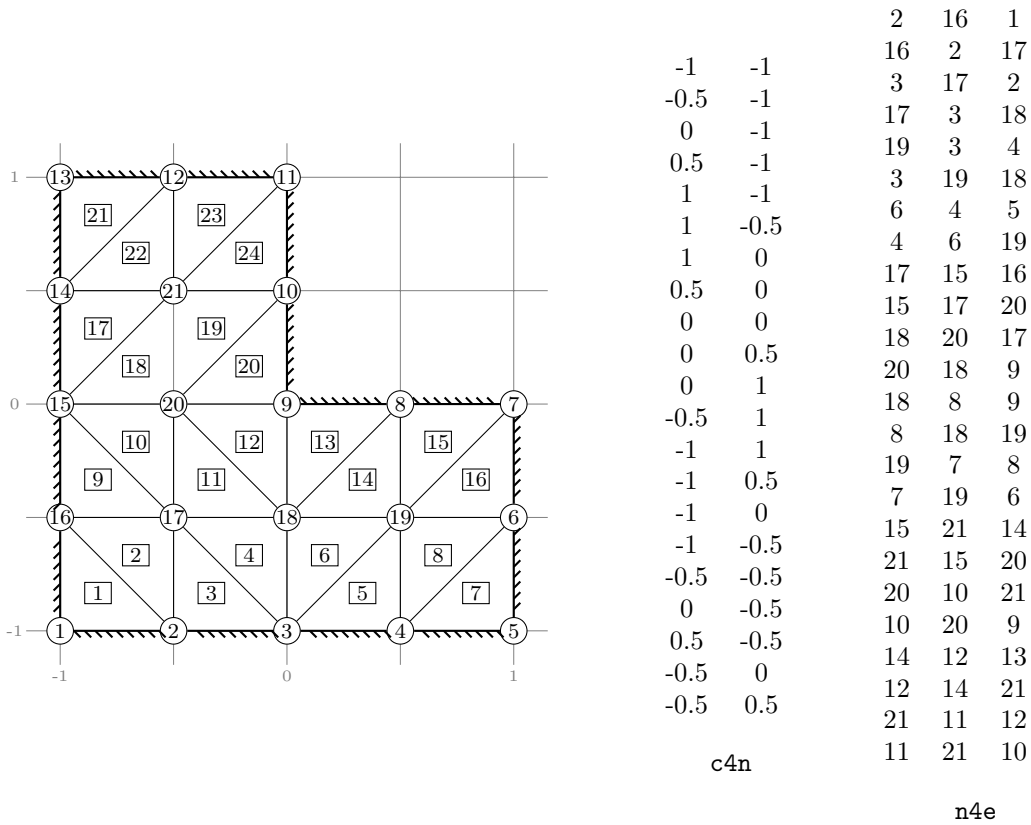
| c4n | |
|---|---|
| -1 | -1 |
| -0.5 | -1 |
| 0 | -1 |
| 0.5 | -1 |
| 1 | -1 |
| 1 | -0.5 |
| 1 | 0 |
| 0.5 | 0 |
| 0 | 0 |
| 0 | 0.5 |
| 0 | 1 |
| -0.5 | 1 |
| -1 | 1 |
| -1 | 0.5 |
| -1 | 0 |
| -1 | -0.5 |
| -0.5 | -0.5 |
| 0 | -0.5 |
| 0.5 | -0.5 |
| -0.5 | 0 |
| -0.5 | 0.5 |

| n4e | | |
|---|---|---|
| 2 | 16 | 1 |
| 16 | 2 | 17 |
| 3 | 17 | 2 |
| 17 | 3 | 18 |
| 19 | 3 | 4 |
| 3 | 19 | 18 |
| 6 | 4 | 5 |
| 4 | 6 | 19 |
| 17 | 15 | 16 |
| 15 | 17 | 20 |
| 18 | 20 | 17 |
| 20 | 18 | 9 |
| 18 | 8 | 9 |
| 8 | 18 | 19 |
| 19 | 7 | 8 |
| 7 | 19 | 6 |
| 15 | 21 | 14 |
| 21 | 15 | 20 |
| 20 | 10 | 21 |
| 10 | 20 | 9 |
| 14 | 12 | 13 |
| 12 | 14 | 21 |
| 21 | 11 | 12 |
| 11 | 21 | 10 |

Figure 19: Picture of a triangulation $\mathcal{T} = \text{conv}\{(-1/2, -1), (-1, -1/2), (-1, -1)\}$, $\text{conv}\{(-1, -1/2), (-1/2, -1), (-1/2, -1/2)\}, \ldots, \text{conv}\{(1/2, 1), (1, 1/2), (1, 1)\}$ with $m = 24$ triangles and $n = 21$ nodes (a). The picture indicates an enumeration of nodes (numbers in circles) and elements (numbers in boxes) given in the matrices c4n (b) and n4e (c). The Dirichlet boundary conditions on the exterior nodes are included in the vector dirichlet $= (1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16)$ of the labels in a counterclockwise enumeration. The data c4n, n4e, and dirichlet are the input of the finite element program to compute a displacement vector x as its output.

boundary specified by an input vector dirichlet of all fixed nodes at the outer boundary; cf. Figure 19.

Given the aforementioned data in the model Dirichlet problem with right-hand side $f = 1$, the $P_1$ finite element space $\widetilde{V} := \text{span}\{\varphi_j : z_j \in \mathcal{K}\}$ is formed by the nodal basis functions $\varphi_j$ of each free node $z_k$; the set $\mathcal{K}$ of free nodes, the interior nodes, is represented in the $N$ vector freenodes, the vector of labels in $1 : n$ without dirichlet.

The resulting discrete equation is the $N \times N$ linear system of equations $Ax = b$ with the

positive definite symmetric stiffness matrix $A$ and right-hand side $b$. Their components are defined (as a subset of)

$$A_{jk} \quad := \quad \int_{\Omega} \nabla \varphi_j \cdot \nabla \varphi_k \mathrm{dx}$$

and

$$b_j \quad := \quad \int_{\Omega} f \varphi_j \mathrm{dx} \quad \text{for} \quad j, k = 1, \dots, n$$

The computation of the entries $A_{jk}$ and $b_j$ is performed elementwise for the additivity of the integral and since $\mathcal{T}$ is a partition of the domain $\Omega$. Given the triangle $T_j$ number $j$, the MATLAB command `c4n(n4e(j,:),:)` returns the $3 \times 2$ matrix $(P_1, P_2, P_3)^T$ of its vertices. Then, the *local stiffness matrix* reads

$$\mathrm{STIMA}(T_j)_{\alpha\beta} := \int_{T_j} \nabla \varphi_k \cdot \nabla \varphi_\ell \mathrm{dx} \quad \text{for } \alpha, \beta = 1, 2, 3$$

for those numbers $k$ and $\ell$ of two vertices $z_k = P_\alpha$ and $z_\ell = P_\beta$ of $T_j$. The correspondence of global and local indices, i.e. the numbers of vertices in $(z_k, z_\ell, z_m) = (P_1, P_2, P_3)$, of $T_j$ can be formalized by

$$I(T_j) = \{ (\alpha, k) \in \{1, 2, 3\} \times \{1, \dots, n\} : P_\alpha = z_k \in \mathcal{N} \}$$

The local stiffness matrix is in fact

$$\mathrm{STIMA}(T_j) = \det \frac{P}{2} \left( QQ^T \right) \text{ with } P := \begin{bmatrix} 1 & 1 & 1 \\ P_1 & P_2 & P_3 \end{bmatrix}$$

$$\text{and } Q := P^{-1} \begin{bmatrix} 0 & 0 \\ 1 & 0 \\ 0 & 1 \end{bmatrix}$$

This formula allows a compact programming in MATLAB as shown (for any dimension $d$)

```
function stima=stima(vertices)
P=[ones(1,size(vertices,2)+1);vertices'];
Q=P\[zeros(1,size(vertices,2));...
  eye(size(vertices,2))];
stima=det(P)*Q*Q'/prod(1:size(vertices,2));
```

Utilizing the index sets $I$, the assembling of all local stiffness matrices reads

$$\mathrm{STIMA} = \sum_{T_j \in \mathcal{T}} \sum_{(\alpha, k) \in I(T_j)} \sum_{(\beta, \ell) \in I(T_j)} \mathrm{STIMA}(T_j)_{\alpha\beta} \, e_k \otimes e_\ell$$

($e_k$ is the $k$th canonical unit vector with the $\ell$th component equal to the Kronecker delta $\delta_{k\ell}$ and $\otimes$ is the dyadic product.) The implementation of each summation is realized by adding $\mathrm{STIMA}(T_j)$ to the $3 \times 3$ submatrix of the rows and columns corresponding to $k$, $\ell$, $m$; see the MATLAB program below.

```
function [x,A]=FEM(c4n,n4e,Dirichlet)
N=size(c4n,1);d=size(c4n,2);A=sparse(N,N);b=zeros(N,1);x=b;
```
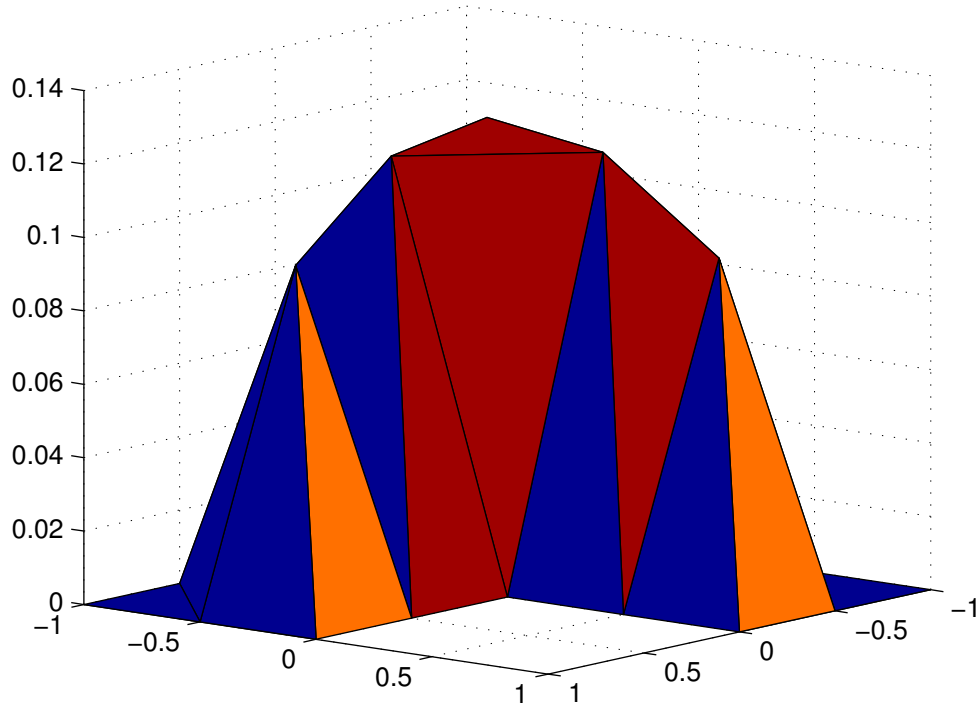
Figure 20: Discrete solution of $-\Delta u = 1$ with homogeneous Dirichlet boundary data based on the triangulation of Figure 19.

```
for j=1:size(n4e,1)
area=abs(det([ones(1,d+1);c4n(n4e(j,:),:)'])/factorial(d));
grads=[ones(1,d+1);c4n(n4e(j,:),:)']\[zeros(1,d);eye(d)];
A(n4e(j,:),n4e(j,:))=A(n4e(j,:),n4e(j,:))+area*(grads*grads');
b(n4e(j,:))=b(n4e(j,:))+ones(d+1,1)*area/(d+1);end
dof=setdiff(1:N,Dirichlet(:));x(dof)=A(dof,dof)_(dof); end
```

Given the output vector x, a plot of the discrete solution

$$\tilde{u} = \sum_{j=1}^{n} x_j \, \varphi_n$$

is generated by the command `trisurf(n4e,c4n(:,1),c4n(:,2),x)` and displayed in Figure 20.

For alternative programs with numerical examples and full documentation, the interested

readers are referred to Alberty, Carstensen and Funken (1999) and Alberty *et al* (2002). The closest more commercial finite element package might be FEMLAB. The internet provides over 200 000 entries under the search for 'Finite Element Method Program'. Amongst public domain software are the programs DEAL II and FREEFEM and details of implementation of 2D adaptive mesh-refining close to the implementation can be found in Funken, Praetorius and Wissgott (2011).

## REFERENCES

Adams RA. *Sobolev Spaces.* Academic Press, New York, 1995.

Agmon S. *Lectures on Elliptic Boundary Value Problems.* Van Nostrand, Princeton, 1965.

Ainsworth M and Oden JT. *A Posteriori Error Estimation in Finite Element Analysis.* Wiley-Interscience, New York, 2000.

Alberty J, Carstensen C and Funken S. Remarks around 50 lines of Matlab: Short finite element implementation. *Numer. Algorithms* 1999; **20**:117-137.

Alberty J, Carstensen C, Funken S and Klose R. Matlab implementation of the finite element method in elasticity. *Computing* 2002; **60**:239-263.

Apel T. *Anisotropic Finite Elements: Local Estimates and Applications.* Teubner Verlag, Stuttgart, 1999.

Apel T and Dobrowolski M. Anisotropic interpolation with applications to the finite element method. *Computing* 1992; **47**:277-293.

Arnold DN and Brezzi F. Mixed and nonconforming finite element methods: implementation, postprocessing and error estimates. *RAIRO Modél. Math. Anal. Numér.* 1982; **19**:7-32.

Arnold DN, Boffi D and Falk RS. Approximation by quadrilateral finite elements. *Math. Comp.* 2002; **71**:909-922.

Arnold DN, Mukherjee A and Pouly L. Locally adapted tetrahedral meshes using bisection. *SIAM J. Sci. Comput.* 2000; **22**:431-448.

Aziz AK (ed.). *The Mathematical Foundations of the Finite Element Method with Applications to Partial Differential Equations.* Academic Press, New York, 1972.

Babuška I. Courant element: before and after. *Lecture Notes in Pure and Applied Mathematics*, vol. 164. Marcel Dekker: New York, 1994; 37-51.

Babuška I and Aziz AK. Survey lectures on the mathematical foundations of the finite element method. In *The Mathematical Foundations of the Finite Element Method with Applications to Partial Differential Equations*, Aziz AK (ed.). Academic Press, New York, 1972; 3-359.

Babuška I and Aziz AK. On the angle condition in the finite element method. *SIAM J. Numer. Anal.* 1976; **13**:214-226.

Babuška I and Guo B. The $h, p$, and $h-p$ versions of the finite element methods in 1 dimension. *Numer. Math.* 1986; **49**:613-657.

Babuška I and Kellogg RB. Nonuniform error estimates for the finite element method. *SIAM J. Numer. Anal.* 1975; **12**:868-875.

Babuška I and Miller A. A feedback finite element method with a posteriori error estimation. I. The finite element method and some properties of the a posteriori estimator. *Comput. Methods Appl. Mech. Eng.* 1987; **61**:1-40.

Babuška I and Osborn J. Eigenvalue problems. In *Handbook of Numerical Analysis*, vol. II, Ciarlet PG and Lions JL (eds). North Holland: Amsterdam, 1991; 641-787.

Babuška I and Pitkäranta J. The plate paradox for hard and soft simple support. *SIAM J. Math. Anal.* 1990; **21**:551-576.

Babuška I and Rheinboldt WC. Error estimates for adaptive finite element computations. *SIAM J. Numer. Anal.* 1978; **15**:736-754.

Babuška I and Rheinboldt WC. Analysis of optimal finite-element meshes in $\mathbb{R}^1$. *Math. Comp.* 1979; **33**:435-463.

Babuška I and Strouboulis T. *The Finite Element Method and its Reliability*. Oxford University Press, New York, 2001.

Babuška I and Suri M. On locking and robustness in the finite element method. *SIAM J. Numer. Anal.* 1992; **29**:1261-1293.

Babuška I and Vogelius R. Feedback and adaptive finite element solution of one-dimensional boundary value problems. *Numer. Math.* 1984; **44**:75-102.

Bangerth W and Rannacher R. *Adaptive Finite Element Methods for Differential Equations*, Lectures in Mathematics ETH Zürich. Birkhäuser Verlag, Basel, 2003.

Bank RE and Weiser A. Some a posteriori error estimators for elliptic differential equations. *Math. Comp.* 1985; **44**:283-301.

Bank RE and Xu J. Asymptotically Exact a Posteriori Error Estimators, Part I: Grids with Superconvergence. *SIAM J. Numer. Anal.* 2003; **41**:2294-2312.

Bänsch E. Local mesh refinement in 2 and 3 dimensions. *IMPACT Comput. Sci. Eng.* 1991; **3**:181-191.

Bartels S and Carstensen C. Each averaging technique yields reliable a posteriori error control in FEM on unstructured grids. II. Higher order FEM. *Math. Comp.* 2002; **71**:971-994.

Bathe K-J. *Finite Element Procedures*. Prentice Hall, Upper Saddle River, 1996.

Becker EB, Carey GF and Oden JT. *Finite Elements. An Introduction.* Prentice Hall, Englewood Cliffs, 1981.

Becker R , Mao S, and Shi Z-C. A convergent nonconforming adaptive finite element method with quasi-optimal complexity. *SIAM J. Numer. Anal.* 2010; **47**:4639-4659.

Becker R and Rannacher R. A feed-back approach to error control in finite element methods: basic analysis and examples. *East-West J. Numer. Math.* 1996; **4**:237-264.

Becker R and Rannacher R. An optimal control approach to a posteriori error estimation in finite element methods. *Acta Numer.* 2001; **10**:1-102.

Ben Belgacem F and Brenner SC. Some nonstandard finite element estimates with applications to 3D Poisson and Signorini problems. *Electron. Trans. Numer. Anal.* 2001; **12**:134-148.

Berger A, Scott R and Strang G. Approximate boundary conditions in the finite element method. *Symposia Mathematica,* vol. X (Convegno di Analisi Numerica, INDAM, Rome, 1972). Academic Press: London, 1972; 295-313.

Bernardi C. Optimal finite-element interpolation on curved domains. *SIAM J. Numer. Anal.* 1989; **26**:1212-1240.

Bernardi C and Girault V. A local regularization operator for triangular and quadrilateral finite elements. *SIAM J. Numer. Anal.* 1998; **35**:1893-1916.

Bey J. Tetragonal grid refinement. *Computing* 1995; **55**:355-378.

Binev P, Dahmen W and DeVore R. Adaptive finite element methods with convergence rates. *Numer. Math.* 2004; **97**: 219-268.

Braess D. *Finite Elements* (3nd edn). Cambridge University Press, Cambridge, 2007.

Bramble JH and Hilbert AH. Estimation of linear functionals on Sobolev spaces with applications to Fourier transforms and spline interpolation. *SIAM J. Numer. Anal.* 1970; **7**:113-124.

Bramble JH, Pasciak JE and Schatz AH. The construction of preconditioners for elliptic problems by substructuring, I. *Math. Comp.* 1986; **47**:103-134.

Bramble JH, Pasciak JE and Steinbach O. On the stability of the $L_2$-projection in $H^1(\Omega)$. *Math. Comp.* 2002; **71**:147-156.

Brenner SC.Two-level additive Schwarz preconditioners for nonconforming finite elements. In *Domain Decomposition Methods in Scientific and Engineering Computing*, Keyes DE and Xu J (eds). American Mathematical Society: Providence, 1994; 9-14.

Brenner SC. Poincaré-Friedrichs inequalities for piecewise $H^1$ functions. *SIAM J. Numer. Anal.* 2003; **41**:306-324.

Brenner SC. Korn's inequalities for piecewise $H^1$ vector fields. *Math. Comp.* 2004; **73**:1067-1087.

Brenner SC. Forty years of the Crouzeix-Raviart element. *Numer. Methods Partial Differential Equations.* 2015; **31**:367-396.

Brenner SC and Scott LR. *The Mathematical Theory of Finite Element Methods* (2nd edn). Springer-Verlag, New York, 2002.

Brenner SC and Sung LY. Discrete Sobolev and Poincaré inequalities via Fourier series. *East-West J. Numer. Math.* 2000; **8**:83-92.

Brenner SC and Sung LY. Piecewise $H^1$ functions and vector fields associated with meshes generated by independent refinements. *Math. Comp.* 2015; **84**:1017-1036.

Brenner SC, Wang K and Zhao J. Poincaré-Friedrichs inequalities for piecewise $H^2$ functions. *Numer. Funct. Anal. Optim.* 2004; **25**:463-478.

Carstensen C. Quasi-interpolation and a posteriori error analysis in finite element methods. *M2AN Math. Model. Numer. Anal.* 1999; **33**:1187-1202.

Carstensen C. Merging the Bramble-Pasciak-Steinbach and the Crouzeix-Thomée criterion for $H^1$-stability of the $L^2$-projection onto finite element spaces. *Math. Comp.* 2002; **71**:157-163.

Carstensen C. All first-order averaging techniques for a posteriori finite element error control on unstructured grids are efficient and reliable. *Math. Comp.* 2004; **73**:1153-1165.

Carstensen C. An adaptive mesh-refining algorithm allowing for an $H^1$-stable $L^2$-projection onto Courant finite element spaces. *Constructive Approximation Theory.*

Carstensen C. Some remarks on the history and future of averaging techniques in a posteriori finite element error analysis. *ZAMM* 2004; **84**:3-21.

Carstensen C and Alberty J. Averaging techniques for reliable a posteriori FE-error control in elastoplasticity with hardening. *Comput. Methods Appl. Mech. Eng.* 2003; **192**:1435-1450.

Carstensen C and Bartels S. Each averaging technique yields reliable a posteriori error control in FEM on unstructured grids. I. Low order conforming, nonconforming and Mixed FEM. *Math. Comp.* 2002; **71**:945-969.

Carstensen C , Eigel M , Hoppe RHW and Löbhard C. A review of unified a posteriori finite element error control. *Numer. Math. Theory Methods Appl.* 2012; **5**:509-558.

Carstensen C, Feischl M, Page M and Praetorius C. Axioms of adaptivity. *Comput. Methods Appl. Math.* 2014; **67**:1195-1253.

Carstensen C and Funken SA. Fully reliable localised error control in the FEM. *SIAM J. Sci. Comput.* 1999/00; **21**:1465-1484.

Carstensen C and Funken SA. Constants in Clément-interpolation error and residual based a posteriori estimates in finite element methods. *East-West J. Numer. Math.* 2000; **8**:153-175.

Carstensen C and Funken SA. A posteriori error control in low-order finite element discretizations of incompressible stationary flow problems. *Math. Comp.* 2001a; **70**:1353-1381.

Carstensen C and Funken SA. Averaging technique for FE-a posteriori error control in elasticity. I: Conforming FEM. II: λ-independent estimates. III: Locking-free nonconforming FEM. *Comput. Methods Appl. Mech. Eng.* 2001b; **190**:2483-2498; **190**:4663-4675; **191**:861-877.

Carstensen C, Gallistl D and Schedensack M. Discrete Reliability for Crouzeix-Raviart FEMs. *SIAM J. Numer. Anal.* 2013; **51**:2935-2955.

Carstensen C, Gallistl D and Gedicke J. Justification of the saturation assumption. *Numer. Math.* 2016; to appear.

Carstensen C and Hoppe RHW. Convergence analysis of an adaptive nonconforming finite element method. *Numer. Math.* 2006; **103**:251-266.

Carstensen C and Merdon C. Effective postprocessing for equilibration a posteriori error estimators. *Comput. Methods Appl. Math.* 2014; **14**:35-54.

Carstensen C and Merdon C. Computational survey on a posteriori error estimators for the Crouzeix-Raviart nonconforming finite element method for the Stokes problem. *Numer. Math.* 2013; **123**:425-459.

Carstensen C and Merdon C. Refined fully explicit a posteriori residual-based error control. *SIAM J. Numer. Anal.* 2014; **52**:1709-1728.

Carstensen C and Verfürth R. Edge residuals dominate a posteriori error estimates for low order finite element methods. *SIAM J. Numer. Anal.* 1999; **36**:1571-1587.

Carstensen C, Bartels S and Klose R. An experimental survey of a posteriori Courant finite element error control for the Poisson equation. *Adv. Comput. Math.* 2001; **15**:79-106.

Carstensen C and Rabus H Axioms of adaptivity for separate marking. *preprint* 2016

Chen C. Superconvergence for triangular finite elements. *Sci. China (Ser. A)* 1999; **42**:917-924.

Chen CM and Huang YQ. *High Accuracy Theory of Finite Element Methods*. Hunan Scientific and Technical Publisher, Changsha, 1995 (in Chinese).

Ciarlet PG. *The Finite Element Method for Elliptic Problems*. North Holland, Amsterdam, 1978 (Reprinted in the Classics in Applied Mathematics Series, SIAM, Philadelphia, 2002).

Ciarlet PG. *Mathematical Elasticity, Volume I: Three-dimensional Elasticity*. North Holland, Amsterdam, 1988.

Ciarlet PG. Basic error estimates for elliptic problems. In *Handbook of Numerical Analysis*, vol. II, Ciarlet PG and Lions JL (eds). North Holland: Amsterdam, 1991; 17-351.

Ciarlet PG. *Mathematical Elasticity, Volume II: Theory of Plates*. North Holland, Amsterdam, 1997.

Clément P. Approximation by finite element functions using local regularization. *RAIRO Modél. Math. Anal. Numér.* 1975; **9**:77-84.

Courant R. Variational methods for the solution of problems of equilibrium and vibration. *Bull. Am. Math. Soc.* 1943; **49**:1-23.

Crouzeix M and Raviart PA. Conforming and nonconforming finite element methods for solving the stationary Stokes equations I. *RAIRO Modél. Math. Anal. Numér.* 1973; **7**:33-75.

Crouzeix M and Thomée V. The stability in $L^p$ and $W^{1,p}$ of the $L^2$-projection onto finite element function paces. *Math. Comp.* 1987; **48**:521-532.

Dari E , Duran R, Padra C and Vampa C. A posteriori error estimators for nonconforming finite element methods. *RAIRO Model. Math. Anal. Numer.* 1996; **30**:385-400.

Dauge M. *Elliptic Boundary Value Problems on Corner Domains*. Lecture Notes in Mathematics 1341. Springer-Verlag, Berlin, 1988.

Davis PJ and Rabinowitz P. *Methods of Numerical Integration*. Academic Press, Orlando, 1984.

Demkowicz, L. *hp*-adaptive finite elements for time-harmonic Maxwell equations. *Topics in Computational Wave Propagation*, Lecture Notes in Computational Science and Engineering, pp. 163-199, Ainsworth M, Davies P, Duncan D, Martin P and Rynne B (eds). Springer-Verlag, Berlin, 2003.

Dörfler W. A convergent adaptive algorithm for Poison's equation. *SIAM J. Numer. Anal.* 1996; **33**:1106-1124.

Dörfler W and Nochetto RH. Small data oscillation implies the saturation assumption. *Numer. Math.* 2002; **91**:1-12.

Douglas Jr J and Dupont T. Galerkin approximations for the two-point boundary value problem using continuous, piecewise polynomial spaces. *Numer. Math.* 1974; **22**:99-109.

Douglas Jr J, Dupont T and Wheeler MF. An $L^\infty$ estimate and a superconvergence result for a Galerkin method for elliptic equations based on tensor products of piecewise polynomials. *RAIRO Modél. Math. Anal. Numér.* 1974; **8**:61-66.

Douglas Jr J, Dupont T, Percell P and Scott R. A family of $C^1$ finite elements with optimal approximation properties for various Galerkin methods for 2nd and 4th order problems. *RAIRO Modél. Math. Anal. Numér.* 1979; **13**:227-255.

Dupont T and Scott R. Polynomial approximation of functions in Sobolev spaces. *Math. Comp.* 1980; **34**:441-463.

Duvaut G and Lions JL. *Inequalities in Mechanics and Physics.* Springer-Verlag, Berlin, 1976.

Eriksson K and Johnson C. Adaptive finite element methods for parabolic problems. I. A linear model problem. *SIAM J. Numer. Anal.* 1991; **28**:43-77.

Eriksson K, Estep D, Hansbo P and Johnson C. Introduction to adaptive methods for differential equations. *Acta Numer.* 1995; **4**:105-158.

Friedrichs KO. On the boundary value problems of the theory of elasticity and Korn's inequality. *Ann. Math.* 1947; **48**:441-471.

Funken S and Praetorius D and Wissgott P Efficient implementation of adaptive P1-FEM in Matlab. *Comput. Methods Appl. Math.* 2011; **11**:460-490.

Gallistl D and Schedensack M and Stevenson RA remark on newest vertex bisection in any space dimension. *Comput. Methods Appl. Math.* 2014; **14**:317-320.

Gastaldi L and Nochetto RH. Optimal $L^\infty$-error estimates for nonconforming and mixed finite element methods of lowest order. *Numer. Math.* 1987; **50**:587-611.

Gilbarg D and Trudinger NS. *Elliptic Partial Differential Equations of Second Order* (2nd edn). Springer-Verlag, Berlin, 1983.

Girault V and Scott LR. Hermite interpolation of nonsmooth functions preserving boundary conditions. *Math. Comput.* 2002; **71**:1043-1074.

Grisvard P. *Elliptic Problems in Nonsmooth Domains.* Pitman, Boston, 1985.

Gudi T. A new error analysis for discontinuous finite element methods for linear elliptic problems. *Math. Comp.* 2010; **79**:2169-2189.

Hu J , Shi Z-C and Xu J. Convergence and optimality of the adaptive Morley element method. *Numer. Math.* 2012; **121**:731-752.

Hughes TJR. *The Finite Element Method. Linear Static and Dynamic Finite Element Analysis.* Prentice Hall, Englewood Cliffs, 1987 (Reprinted by Dover Publications, New York, 2000).

Jamet P. Estimations d'erreur pour des éléments finis droits presque dégénérés. *RAIRO Anal. Numér.* 1976; **10**:43-61.

Mardal KA and Winther R. An observation on Korn's inequality for nonconforming finite element methods . *Math. Comp.* 2006; **75**:1-6.

Kozlov VA, Maz'ya VG and Rossman J. *Elliptic Boundary Value Problems in Domains with Point Singularities.* American Mathematical Society, 1997.

Kozlov VA, Maz'ya VG and Rossman J. *Spectral Problems Associated with Corner Singularities of Solutions to Elliptic Problems.* American Mathematical Society, 2001.

Křížek M and Neittaanmäki P. On superconvergence techniques. *Acta Appl. Math.* 1987; **9**:175-198.

Křížek M, Neittaanmäki P and Stenberg R (eds.). *Finite Element Methods: Superconvergence, Post-processing and A Posteriori Estimates*, Proc. Conf. Univ. of Jyväskylä, 1996, Lecture Notes in Pure and Applied Mathematics 196. Marcel Dekker, New York, 1998.

Ladeveze P and Leguillon D. Error estimate procedure in the finite element method and applications. *SIAM J. Numer. Anal.* 1983; **20**:485-509.

Lenoir M. Optimal isoparametric finite elements and error estimates for domains involving curved boundaries. *SIAM J. Numer. Anal.* 1986; **23**:562-580.

Lin Q and Yan NN. *Construction and Analysis of Efficient Finite Element Methods*. Hebei University Press, Baoding, 1996 (in Chinese).

Mansfield L. Approximation of the boundary in the finite element solution of fourth order problems. *SIAM J. Numer. Anal.* 1978; **15**:568-579.

Malkus DS and Hughes TJR. Mixed finite element methods-reduced and selective integration techniques: A unification of concepts. *Comput. Methods Appl. Mech. Eng.* 1978; **15**:63-81.

Maubach JM. Local bisection refinement for $n$-simplicial grids generated by reflection. *SIAM J. Sci. Comput.* 1995; **16**:210-227.

Morin P, Nochetto RH and Siebert KG. Local problems on stars: a posteriori error estimation, convergence, and performance. *Math. Comp.* 2003a; **72**:1067-1097.

Morin P, Nochetto RH and Siebert KG. Convergence of adaptive finite element methods. *SIAM Rev.* 2003b; **44**:631-658.

Morley LSD. The triangular equilibrium problem in the solution of plate bending problems . *Aero.Quart.* 1968; **19**:149-169.

Natterer F. Über die punktweise Konvergenz finiter Elemente. *Numer. Math.* 1975; **25**:67-77.

Nazarov SA and Plamenevsky BA. *Elliptic Problems in Domains with Piecewise Smooth Boundaries.* Walter De Gruyter, Berlin, 1994.

Nečas J. *Les Méthodes Directes en Théorie des Équations Elliptiques.* Masson, Paris, 1967.

Nicaise S. *Polygonal Interface Problems.* Peter D Lang Verlag, Frankfurt am Main, 1993.

Nitsche JA. $L_\infty$-convergence of finite element approximations. In *Mathematical Aspects of Finite Element Methods*, Lecture Notes in Mathematics 606. Springer-Verlag: New York, 1977; 261-274.

Nitsche JA. On Korn's second inequality. *RAIRO Anal. Numér.* 1981; **15**:237-248.

Nitsche JA and Schatz AH. Interior estimates for Ritz-Galerkin methods. *Math. Comp.* 1974; **28**:937-958.

Nochetto RH. Removing the saturation assumption in a posteriori error analysis. *Istit. Lombardo Accad. Sci. Lett. Rend. A* 1993; **127**:67-82.

Nochetto RH and Wahlbin LB. Positivity preserving finite element approximation. *Math. Comp.* 2002; **71**:1405-1419.

Oden JT. Finite elements: an introduction. In *Handbook of Numerical Analysis*, vol. II, Ciarlet PG and Lions JL (eds). North Holland: Amsterdam, 1991; 3-15.

Oden JT and Demkowicz LF. *Applied Functional Analysis.* CRC Press, Boca Raton, 1996.

Oden JT and Reddy JN. *An Introduction to the Mathematical Theory of Finite Elements.* Wiley, New York, 1976.

Rannacher R and Scott R. Some optimal error estimates for piecewise linear finite element approximations. *Math. Comp.* 1982; **38**:437-445.

Rannacher R and Turek S. Simple nonconforming quadrilateral Stokes element. *Numer. Methods Partial Differential Equations* 1992; **8**:97-111.

Reddy JN. *Applied Functional Analysis and Variational Methods in Engineering.* McGraw-Hill, New York, 1986.

Repin S. *A Posteriori Estimates for Partial Differential Equations.* de Gruyter, Berlin, 2008.

Rivara MC. Algorithms for refining triangular grids suitable for adaptive and multigrid techniques. *Int. J. Numer. Methods Eng.* 1984; **20**:745-756.

Rodriguez R. Some remarks on Zienkiewicz-Zhu estimator. *Numer. Methods Partial Differential Equations* 1994a; **10**:625-635.

Rodriguez R. A posteriori error analysis in the finite element method. In *Lecture Notes in Pure and Applied Mathematics*, vol. 164. Marcel Dekker: New York, 1994b; 389-397.

Ruas V. A quadratic finite element method for solving biharmonic problems in $R^n$. *Numer. Math.* 1988; **52**:33-43.

Schatz AH. An observation concerning Ritz-Galerkin methods with indefinite bilinear forms. *Math. Comp.* 1974; **28**:959-962.

Schatz AH. Pointwise error estimates and asymptotic error expansion inequalities for the finite element method on irregular grids: Part I. Global estimates. *Math. Comp.* 1998; **67**:877-899.

Schatz AH. Pointwise error estimates and asymptotic error expansion inequalities for the finite element method on irregular grids: Part II. Interior estimates. *SIAM J. Numer. Anal.* 2000; **38**:1269-1293.

Schatz AH and Wahlbin LB. Maximum norm estimates in the finite element method on plane polygonal domains. Part 1. *Math. Comp.* 1978; **32**:73-109.

Schatz AH and Wahlbin LB. Maximum norm estimates in the finite element method on plane polygonal domains. Part 2. *Math. Comp.* 1979; **33**:465-492.

Schatz AH and Wahlbin LB. On the quasi-optimality in $L_\infty$ of the $\mathring{H}^1$-projection into finite element spaces. *Math. Comp.* 1982; **38**:1-22.

Schatz AH and Wahlbin LB. Interior maximum-norm estimates for finite element methods, Part II. *Math. Comp.* 1995; **64**:907-928.

Schatz AH, Sloan IH and Wahlbin LB. Superconvergence in finite element methods and meshes that are locally symmetric with respect to a point. *SIAM J. Numer. Anal.* 1996; **33**:505-521.

Schatz AH, Thomée V and Wendland WL. *Mathematical Theory of Finite and Boundary Element Methods.* Birkhäuser Verlag, Basel, 1990.

Scott R. *Finite Element Techniques for Curved Boundaries.* Doctoral thesis, Massachusetts Institute of Technology, 1973.

Scott R. Interpolated boundary conditions in the finite element method. *SIAM J. Numer. Anal.* 1975; **12**:404-427.

Scott R. Optimal $L^\infty$ estimates for the finite element method on irregular meshes. *Math. Comp.* 1976; **30**:681-697.

Scott LR and Zhang S. Finite element interpolation of non-smooth functions satisfying boundary conditions. *Math. Comp.* 1990; **54**:483-493.

Shi Z-C. On the convergence of the incomplete biquadratic nonconforming plate element . *Math. Numer. Sinica* 1986; **8**:53-62.

Shi Z-C. Error estimates of Morley element. *Numer. Math. & Appl.* 1990; **12**:9-15.

Stevenson R. The completion of locally refined simplicial partitions created by bisection. *Math. Comp.* 2008; **77**:227-241.

Stevenson R. Optimality of a standard adaptive finite element method. *Found. Comput. Math.* 2007; **7**:245-269.

Strang G and Berger A. The change in solution due to change in domain. *Proceedings AMS Symposium on Partial Differential Equations.* American Mathematical Society: Providence, 1971; 199-205.

Strang G. and Fix GJ. *An Analysis of the Finite Element Method.* Prentice Hall, Englewood Cliffs, 1973 (Reprinted by Wellesley-Cambridge Press, Wellesley, 1988).

Szabó BA and Babuška I. *Finite Element Analysis.* John Wiley & Sons, New York, 1991.

Traxler CT. An algorithm for adaptive mesh refinement in $n$ dimensions. *Computing* 1997; **59**:115-137.

Triebel H. *Interpolation Theory, Function Spaces, Differential Operators.* North Holland, Amsterdam, 1978.

Verfürth R. *A Review of A Posteriori Error Estimation and Adaptive Mesh-Refinement Techniques.* Wiley-Teubner, New York, 1996.

Verfürth R. A note on polynomial approximation in Sobolev spaces. *Modél. Math. Anal. Numér.* 1999; **33**:715-719.

Verfürth R. *A Posteriori Error Estimation Techniques for Finite Element Methods.* Oxford University Press, Oxford, 2013.

Wahlbin LB. On the sharpness of certain local estimates for $\mathring{H}^1$ projections into finite element spaces: Influence of a reentrant corner. *Math. Comp.* 1984; **42**:1-8.

Wahlbin LB. Local behavior in finite element methods. In *Handbook of Numerical Analysis*, vol. II, Ciarlet PG, Lions JL (eds). North Holland: Amsterdam, 1991; 355-522.

Wahlbin LB. *Superconvergence in Galerkin Finite Element Methods*, Lecture Notes in Mathematics 1605. Springer-Verlag, Berlin, 1995.

Wang M and Xu J. The Morley element for fourth order elliptic equations in any dimensions. *Numer. Math.* 2006; **103**:155-169.

Wheeler MF. An optimal $L_\infty$ error estimate for Galerkin approximations to solutions of two-point boundary value problems. *SIAM J. Numer. Anal.* 1973; **1**:914-917.

Wloka J. *Partial Differential Equations.* Cambridge University Press, Cambridge, 1987.

Yosida K. *Functional Analysis*, Classics in Mathematics. Springer-Verlag, Berlin, 1995.

Zienkiewicz OC and Taylor RL. *The Finite Element Method* (5th edn). Butterworth-Heinemann, Oxford, 2000.

Zlámal M. Curved elements in the finite element method. I. *SIAM J. Numer. Anal.* 1973; **10**:229-240.

Zlámal M. Curved elements in the finite element method. II. *SIAM J. Numer. Anal.* 1974; **11**:347-362.

Ženíšek A. Maximum-angle condition and triangular finite element of Hermite type. *Math. Comp.* 1995; **64**:929-941.