

Primer of Adaptive Finite Element Methods

Ricardo H. Nochetto and Andreas Veeger

Abstract Adaptive finite element methods (AFEM) are a fundamental numerical instrument in science and engineering to approximate partial differential equations. In the 1980s and 1990s a great deal of effort was devoted to the design of a posteriori error estimators, following the pioneering work of Babuška. These are computable quantities, depending on the discrete solution(s) and data, that can be used to assess the approximation quality and improve it adaptively. Despite their practical success, adaptive processes have been shown to converge, and to exhibit optimal cardinality, only recently for dimension $d > 1$ and for linear elliptic PDE. These series of lectures presents an up-to-date discussion of AFEM encompassing the derivation of upper and lower a posteriori error bounds for residual-type estimators, including a critical look at the role of oscillation, the design of AFEM and its basic properties, as well as a complete discussion of convergence, contraction property and quasi-optimal cardinality of AFEM.

Ricardo H. Nochetto

Department of Mathematics and Institute of Physical Science and Technology, University of Maryland, College Park, MD 20742, e-mail: rhn@math.umd.edu. Partially supported by NSF grant DMS-0807811.

Andreas Veeger

Dipartimento di Matematica, Università degli Studi di Milano, Via C. Saldini 50, I-20133 Milano, Italy, e-mail: andreas.veeger@unimi.it. Partially supported by Italian PRIN 2008 “Analisi e sviluppo di metodi numerici avanzati per EDP”.

1 Piecewise Polynomial Approximation

We start with a discussion of piecewise polynomial approximation in W_p^k Sobolev spaces and graded meshes in any dimension d . We first compare pointwise approximation over uniform and graded meshes for $d = 1$ in §1.1, which reveals the advantages of the latter over the former and sets the tone for the rest of the paper. We continue with the concept of Sobolev number in §1.2.

We explore the geometric aspects of mesh refinement for conforming meshes in §1.3 and nonconforming meshes in §1.7, but postpone a full discussion until §6. We include a statement about complexity of the refinement procedure, which turns out to be instrumental later.

We briefly discuss the construction of finite element spaces in §1.4, along with polynomial interpolation of functions in Sobolev spaces in §1.5. This provides local estimates adequate for comparison of quasi-uniform and graded meshes for $d > 1$. We exploit them in developing the so-called error equidistribution principle and the construction of suitably graded meshes via thresholding in §1.6. We conclude that graded meshes can deliver optimal interpolation rates for certain classes of singular functions, and thus supersede quasi-uniform refinement.

1.1 Classical vs Adaptive Pointwise Approximation

We start with a simple motivation in 1d for the use of adaptive procedures, due to DeVore [22]. Given $\Omega = (0, 1)$, a partition $\mathcal{T}_N = \{x_i\}_{i=0}^N$ of Ω

$$0 = x_0 < x_1 < \cdots < x_n < \cdots < x_N = 1$$

and a continuous function $u : \Omega \rightarrow \mathbb{R}$, we consider the problem of *interpolating* u by a *piecewise constant* function U_N over \mathcal{T}_N . To quantify the difference between u and U_N we resort to the *maximum norm* and study two cases depending on the regularity of u .

Case 1: W_∞^1 -Regularity. Suppose that u is Lipschitz in $[0, 1]$. We consider the approximation

$$U_N(x) := u(x_{n-1}) \quad \text{for all } x_{n-1} \leq x < x_n.$$

Since

$$|u(x) - U_N(x)| = |u(x) - u(x_{n-1})| = \left| \int_{x_{n-1}}^x u'(t) dt \right| \leq h_n \|u'\|_{L^\infty(x_{n-1}, x_n)}$$

we conclude that

$$\|u - U_N\|_{L^\infty(\Omega)} \leq \frac{1}{N} \|u'\|_{L^\infty(\Omega)}, \quad (1)$$

provided the local mesh-size h_n is about constant (*quasi-uniform* mesh), and so proportional to N^{-1} (the reciprocal of the number of degrees of freedom). Note that the same integrability is used on both sides of (1). A natural question arises: *Is it possible to achieve the same asymptotic decay rate N^{-1} with weaker regularity demands?*

Case 2: W_1^1 -Regularity. To answer this question, we suppose $\|u'\|_{L^1(\Omega)} = 1$ and consider the non-decreasing function

$$\phi(x) := \int_0^x |u'(t)| dt$$

which satisfies $\phi(0) = 0$ and $\phi(1) = 1$. Let $\mathcal{T}_N = \{x_i\}_{i=0}^N$ be the partition given by

$$\int_{x_{n-1}}^{x_n} |u'(t)| dt = \phi(x_n) - \phi(x_{n-1}) = \frac{1}{N}.$$

Then, for $x \in [x_{n-1}, x_n]$,

$$|u(x) - u(x_{n-1})| = \left| \int_{x_{n-1}}^x u'(t) dt \right| \leq \int_{x_{n-1}}^x |u'(t)| dt \leq \int_{x_{n-1}}^{x_n} |u'(t)| dt = \frac{1}{N},$$

whence

$$\|u - U_N\|_{L^\infty(\Omega)} \leq \frac{1}{N} \|u'\|_{L^1(\Omega)}. \quad (2)$$

We thus conclude that we could achieve the same rate of convergence N^{-1} for rougher functions with just $\|u'\|_{L^1(\Omega)} < \infty$. The following comments are in order for Case 2.

Remark 1 (Equidistribution). The optimal mesh \mathcal{T}_N *equidistributes* the max-error. This mesh is graded instead of uniform but, in contrast to a uniform mesh, such a partition may not be adequate for another function with the same basic regularity as u . It is instructive to consider the singular function $u(x) = x^\gamma$ with $\gamma = 0.1$ and error tolerance 10^{-2} to quantify the above computations: if N_1 and N_2 are the number of degrees of freedom with uniform and graded partitions, we obtain $N_1/N_2 = 10^{18}$.

Remark 2 (Nonlinear Approximation). The regularity of u in (2) is measured in $W_1^1(\Omega)$ instead of $W_\infty^1(\Omega)$ and, consequently, the fractional γ regularity measured in $L^\infty(\Omega)$ increases to one full derivative when expressed in $L^1(\Omega)$. This exchange of integrability between left and right-hand side of (2), and gain of differentiability, is at the heart of the matter and the very reason why suitably graded meshes achieve optimal asymptotic error decay for singular functions. By those we mean functions which are not in the usual linear Sobolev scale, say $W_\infty^1(\Omega)$ in this example, but rather in a nonlinear scale [22]. We will get back to this issue in §7.

1.2 The Sobolev Number: Scaling and Embedding

In order to make Remark 2 more precise, we introduce the Sobolev number. Let $\Omega \subset \mathbb{R}^d$ with $d > 1$ be a Lipschitz and bounded domain, and let $k \in \mathbb{N}$, $1 \leq p \leq \infty$. The Sobolev space $W_p^k(\Omega)$ is defined by

$$W_p^k(\Omega) := \{v : \Omega \rightarrow \mathbb{R} \mid D^\alpha v \in L^p(\Omega) \ \forall |\alpha| \leq k\}.$$

If $p = 2$ we set $H^k(\Omega) = W_2^k(\Omega)$ and note that this is a Hilbert space. The *Sobolev number* of $W_p^k(\Omega)$ is given by

$$\text{sob}(W_p^k) := k - \frac{d}{p}. \quad (3)$$

This number governs the scaling properties of the semi-norm

$$|v|_{W_p^k(\Omega)} := \left(\sum_{|\alpha|=k} \|D^\alpha v\|_{L^p(\Omega)}^p \right)^{1/p},$$

because rescaling variables $\hat{x} = \frac{1}{h}x$ for all $x \in \Omega$, transforms Ω into $\hat{\Omega}$ and v into \hat{v} , while the corresponding norms scale as

$$|\hat{v}|_{W_p^k(\hat{\Omega})} = h^{\text{sob}(W_p^k)} |v|_{W_p^k(\Omega)}.$$

In addition, we have the following *compact embedding*: if $m > k$ and $\text{sob}(W_q^m) > \text{sob}(W_p^k)$, then

$$W_q^m(\Omega) \subset W_p^k(\Omega).$$

We say that two Sobolev spaces are in the same nonlinear Sobolev scale if they have the same Sobolev number. We note that for compactness the space $W_q^m(\Omega)$ must be above the Sobolev scale of $W_p^k(\Omega)$. A relevant example for $d = 2$ are the pair $H^1(\Omega)$ and $L^\infty(\Omega)$ which have the same Sobolev number, in fact $\text{sob}(H^1) = \text{sob}(L^\infty) = 0$, but the former is not even contained in the latter: in fact

$$v(x) = \log \log \frac{|x|}{2} \in H^1(\Omega) \setminus L^\infty(\Omega)$$

in the unit ball. This is a source of difficulties for polynomial interpolation theory and the need for quasi-interpolation operators. This is discussed in §1.5.

We conclude with a comment about Remark 2. We see that $d = 1$ and $\text{sob}(W_1^1) = \text{sob}(L^\infty) = 0$ but $W_1^1(\Omega)$ is compactly embedded in $L^\infty(\Omega)$ in this case. This shows that these two spaces are in the same nonlinear Sobolev scale and that the above inequality between Sobolev numbers for a compact embedding is only sufficient.

1.3 Conforming Meshes: The Bisection Method

In order to approximate functions in $W_p^k(\Omega)$ by piecewise polynomials, we decompose Ω into simplices. We briefly discuss the *bisection* method, the most elegant and successful technique for subdividing Ω in any dimension into a conforming mesh. We also discuss briefly nonconforming meshes in §1.7. We present complete proofs, especially of the complexity of bisection, later in §6.

We focus on $d = 2$ and follow Binev, Dahmen, and DeVore [7], but the results carry over to any dimension $d > 1$ (see Stevenson [53]). We refer to Nochetto, Siebert, and Veiser [45] for a rather complete discussion for $d > 1$.

Let \mathcal{T} denote a *mesh* (triangulation or grid) made of simplices T , and let \mathcal{T} be *conforming* (edge-to-edge). Each element is labeled, namely it has an edge $E(T)$ assigned for refinement (and an opposite vertex $v(T)$ for $d = 2$); see Figure 1.

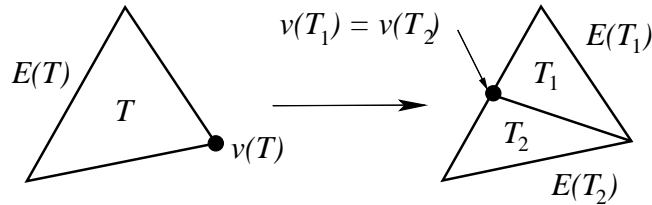


Fig. 1 Triangle $T \in \mathcal{T}$ with vertex $v(T)$ and opposite refinement edge $E(T)$. The bisection rule for $d = 2$ consists of connecting $v(T)$ with the midpoint of $E(T)$, thereby giving rise to children T_1, T_2 with common vertex $v(T_1) = v(T_2)$, the newly created vertex, and opposite refinement edges $E(T_1), E(T_2)$.

The bisection method consists of a suitable *labeling* of the initial mesh \mathcal{T}_0 and a rule to assign the refinement edge to the two children. For $d = 2$ we consider the *newest vertex bisection* as depicted in Figure 1. For $d > 2$ the situation is more complicated and one needs the concepts of type and vertex order [45, 53].

Bisection creates a *unique* master forest \mathbb{F} of binary trees with infinite depth, where each node is a simplex (triangle in 2d), its two successors are the two children created by bisection, and the roots of the binary trees are the elements of the initial conforming partition \mathcal{T}_0 . It is important to realize that, no matter how an element arises in the subdivision process, its associated newest vertex is unique and only depends on the labeling of \mathcal{T}_0 : so $v(T)$ and $E(T)$ are independent of the order of the subdivision process for all $T \in \mathbb{F}$; see Lemma 16 in §6. Therefore, \mathbb{F} is unique.

A finite subset $\mathcal{F} \subset \mathbb{F}$ is called a *forest* if $\mathcal{T}_0 \subset \mathcal{F}$ and the nodes of \mathcal{F} satisfy

- all nodes of $\mathcal{F} \setminus \mathcal{T}_0$ have a predecessor;
- all nodes in \mathcal{F} have either two successors or none.

Any node $T \in \mathcal{F}$ is thus uniquely connected with a node T_0 of the initial triangulation \mathcal{T}_0 , i.e. T belongs to the infinite tree $\mathbb{F}(T_0)$ emanating from T_0 . Furthermore, any forest may have *interior nodes*, i.e. nodes with successors, as well as *leaf nodes*, i.e.

nodes without successors. The set of leaves corresponds to a mesh (or triangulation, grid, partition) $\mathcal{T} = \mathcal{T}(\mathcal{F})$ of \mathcal{T}_0 which may not be conforming or edge-to-edge.

We thus introduce the set \mathbb{T} of all conforming refinements of \mathcal{T}_0 :

$$\mathbb{T} := \{ \mathcal{T} = \mathcal{T}(\mathcal{F}) \mid \mathcal{F} \subset \mathbb{F} \text{ is finite and } \mathcal{T}(\mathcal{F}) \text{ is conforming} \}.$$

If $\mathcal{T}_* = \mathcal{T}(\mathcal{F}_*) \in \mathbb{T}$ is a conforming refinement of $\mathcal{T} = \mathcal{T}(\mathcal{F}) \in \mathbb{T}$, we write $\mathcal{T}_* \geq \mathcal{T}$ and understand this inequality in the sense of trees, namely $\mathcal{F} \subset \mathcal{F}_*$.

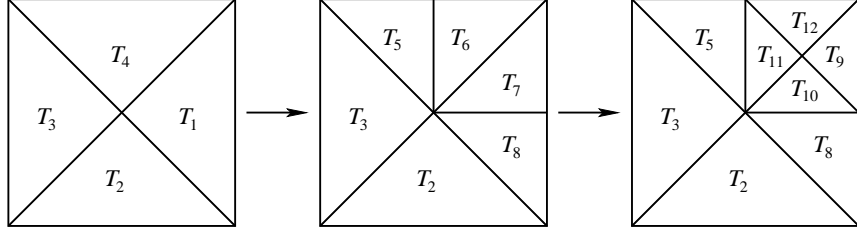


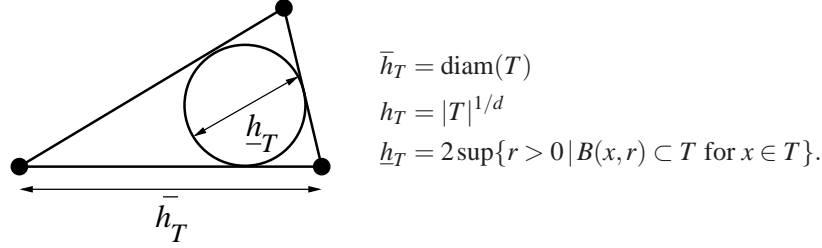
Fig. 2 Sequence of bisection meshes $\{\mathcal{T}_k\}_{k=0}^2$ starting from the initial mesh $\mathcal{T}_0 = \{T_i\}_{i=1}^4$ with longest edges labeled for bisection. Mesh \mathcal{T}_1 is created from \mathcal{T}_0 upon bisecting T_1 and T_4 , whereas mesh \mathcal{T}_2 arises from \mathcal{T}_1 upon refining T_6 and T_7 . The bisection rule is described in Figure 1.



Fig. 3 Forest \mathcal{F}_2 corresponding to the grid sequence $\{\mathcal{T}_k\}_{k=0}^2$ of Figure 2. The roots of \mathcal{F}_2 from the initial mesh \mathcal{T}_0 and the leaves of \mathcal{F}_2 constitute the conforming bisection mesh \mathcal{T}_2 . Moreover, each level of \mathcal{F}_2 corresponds to all elements with generation equal to the level.

Example: Consider $\mathcal{T}_0 = \{T_i\}_{i=1}^4$ and the longest edge to be the refinement edge. Figure 2 displays a sequence of conforming meshes $\mathcal{T}_k \in \mathbb{T}$ created by bisection. Each element T_i of \mathcal{T}_0 is a root of a finite tree emanating from T_i , which together form the forest \mathcal{F}_2 corresponding to mesh $\mathcal{T}_2 = \mathcal{T}(\mathcal{F}_2)$. Figure 3 displays \mathcal{F}_2 , whose leaf nodes are the elements of \mathcal{T}_2 .

Properties of Bisection. We now discuss several crucial geometric properties of bisection. We start with the concept of shape regularity. For any $T \in \mathcal{T}$, we define



Then

$$\underline{h}_T \leq h_T \leq \bar{h}_T \leq \sigma \underline{h}_T \quad \forall T \in \mathcal{T},$$

where $\sigma > 1$ is the shape regularity constant. We say that a sequence of meshes is *shape regular* if σ is uniformly bounded, or in other words that the element shape does not degenerate with refinement. The next lemma guarantees that bisection keeps σ bounded.

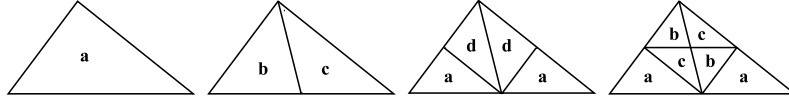


Fig. 4 Bisection produces at most 4 similarity classes for any triangle.

Lemma 1 (Shape Regularity). *The partitions \mathcal{T} generated by newest vertex bisection satisfy a uniform minimal angle condition, or equivalently σ is uniformly bounded, only depending on the initial partition \mathcal{T}_0 .*

Proof. Each $T \in \mathcal{T}_0$ gives rise to a fixed number of similarity classes, namely 4 for $d = 2$ according to Figure 4. This, combined with the fact that $\#\mathcal{T}_0$ is finite, yields the assertion. \square

We define the *generation (or level)* $g(T)$ of an element $T \in \mathcal{T}$ as the number of bisections needed to create T from its ancestor $T_0 \in \mathcal{T}_0$. Since bisection splits an element into two children with equal measure, we realize that

$$h_T = 2^{-g(T)/2} h_{T_0} \quad \forall T \in \mathcal{T}. \quad (4)$$

Referring to Figure 3 we observe that the leaf nodes $T_9, T_{10}, T_{11}, T_{12}$ have generation 2, whereas T_5, T_8 have generation 1 and T_2, T_3 have generation 0.

The following geometric property is a simple consequence of (4).

Lemma 2 (Element Size vs Generation). *There exist constants $0 < D_1 < D_2$, only depending on \mathcal{T}_0 , such that*

$$D_1 2^{-g(T)/2} \leq h_T < \bar{h}_T \leq D_2 2^{-g(T)/2} \quad \forall T \in \mathcal{T}. \quad (5)$$

Labeling and Bisection Rule. Whether the recursive application of bisection does not lead to inconsistencies depends on a suitable initial labeling of edges and a bisection rule. For $d = 2$ they are simple to state [7], but for $d > 2$ we refer to Condition (b) of Section 4 of [53]. Given $T \in \mathcal{T}$ with generation $g(T) = i$, we assign the label $(i + 1, i + 1, i)$ to T with i corresponding to the refinement edge $E(T)$. The following rule dictates how the labeling changes with refinement: the side i is bisected and both new sides as well as the bisector are labeled $i + 2$ whereas the remaining labels do not change. To guarantee that the label of an edge is independent of the elements sharing this edge, we need a special labeling for \mathcal{T}_0 [7]:

$$\text{edges of } \mathcal{T}_0 \text{ have labels 0 or 1 and all elements } T \in \mathcal{T} \text{ have exactly two edges with label 1 and one with label 0.} \quad (6)$$

It is not obvious that such a labeling exists, but if it does then all elements of \mathcal{T}_0 can be split into pairs of compatibly divisible elements. We refer to Figure 5 for an example of initial labeling of \mathcal{T}_0 satisfying (6) and the way it evolves for two successive refinements $\mathcal{T}_2 \geq \mathcal{T}_1 \geq \mathcal{T}_0$ corresponding to Figure 2.

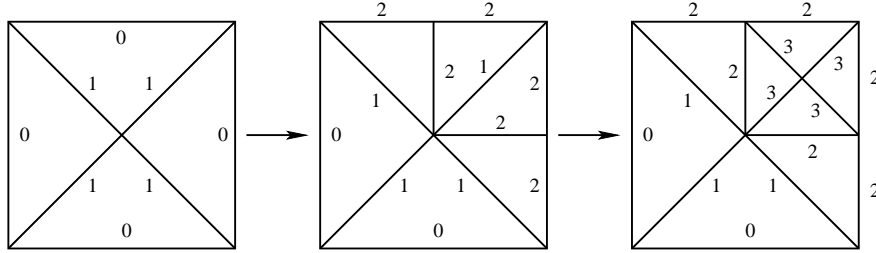


Fig. 5 Initial labeling and its evolution for the sequence of conforming refinements $\mathcal{T}_0 \leq \mathcal{T}_1 \leq \mathcal{T}_2$ of Figure 2.

To guarantee (6) we can proceed as follows: given a coarse mesh of elements T we can bisect twice each T and label the 4 grandchildren, as indicated in Figure 6 for the resulting mesh \mathcal{T}_0 to satisfy the initial labeling [7]. A similar, but much trickier,

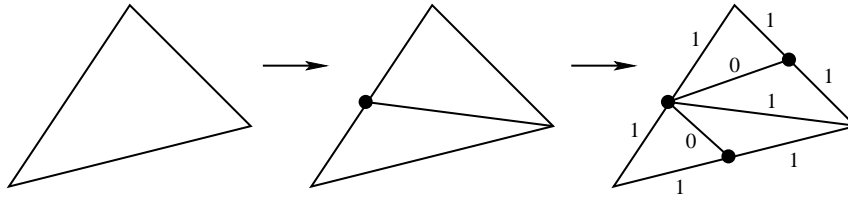


Fig. 6 Bisecting each triangle of \mathcal{T}_0 twice and labeling edges in such a way that all boundary edges have label 1 yields an initial mesh satisfying (6).

construction can be made in any dimension $d > 2$ (see Stevenson [53]). For $d = 3$

the number of elements increases by an order of magnitude, which indicates that (6) is a severe restriction in practice. Finding alternative, more practical, conditions is an open and important problem.

The Procedure REFINE. Given $\mathcal{T} \in \mathbb{T}$ and a subset $\mathcal{M} \subset \mathcal{T}$ of marked elements, the procedure

$$\mathcal{T}_* = \text{REFINE}(\mathcal{T}, \mathcal{M})$$

creates a new conforming refinement \mathcal{T}_* of \mathcal{T} by bisecting all elements of \mathcal{M} at least once and perhaps additional elements to keep conformity.

Conformity is a constraint in the refinement procedure that prevents it from being completely local. The propagation of refinement beyond the set of marked elements \mathcal{M} is a rather delicate matter, which we discuss later in §6. For instance, we show that a naive estimate of the form

$$\#\mathcal{T}_* - \#\mathcal{T} \leq \Lambda_0 \#\mathcal{M}$$

is *not* valid with an absolute constant Λ_0 independent of the refinement level. This can be repaired upon considering the cumulative effect for a sequence of conforming bisection meshes $\{\mathcal{T}_k\}_{k=0}^{\infty}$. This is expressed in the following crucial complexity result due to Binev, Dahmen, and DeVore [7] for $d = 2$ and Stevenson [53] for $d > 2$. We present a complete proof later in §6.

Theorem 1 (Complexity of REFINE). *If \mathcal{T}_0 satisfies the initial labeling (6) for $d = 2$, or that in [53, Section 4] for $d > 2$, then there exists a constant $\Lambda_0 > 0$ only depending on \mathcal{T}_0 and d such that for all $k \geq 1$*

$$\#\mathcal{T}_k - \#\mathcal{T}_0 \leq \Lambda_0 \sum_{j=0}^{k-1} \#\mathcal{M}_j.$$

If elements $T \in \mathcal{M}$ are to be bisected $b \geq 1$ times, then the procedure REFINE can be applied recursively, and Theorem 1 remains valid with Λ_0 also depending on b .

1.4 Finite Element Spaces

Given a conforming mesh $\mathcal{T} \in \mathbb{T}$ we define the finite element space of continuous piecewise polynomials of degree $n \geq 1$

$$\mathbb{S}^{n,0}(\mathcal{T}) := \{v \in C^0(\overline{\Omega}) \mid v|_T \in \mathbb{P}_n(T) \ \forall T \in \mathcal{T}\};$$

note that $\mathbb{S}^{n,0}(\mathcal{T}) \subset H^1(\Omega)$. We refer to Braess [10], Brenner-Scott [11], Ciarlet [19] and Siebert [50] for a discussion on the local construction of this space along with its properties.

We focus on the piecewise linear case $n = 1$ (Courant elements). Global continuity can be simply enforced by imposing continuity at the vertices z of \mathcal{T} , the so-called *nodal values*. We denote by \mathcal{N} the set of vertices z of \mathcal{T} .

However, the following local construction leads to global continuity. If T is a generic simplex of \mathcal{T} , namely the convex hull of $\{z_i\}_{i=0}^d$, then we associate to each vertex z_i a *barycentric coordinate* λ_i^T , which is the linear function in T with nodal value 1 at z_i and 0 at the other vertices of T . Upon pasting together the barycentric coordinates λ_z^T of all simplices T containing vertex $z \in \mathcal{N}$, we obtain a continuous piecewise linear function $\phi_z \in \mathbb{S}^{1,0}(\mathcal{T})$ as depicted in Figure 7 for $d = 2$: The set

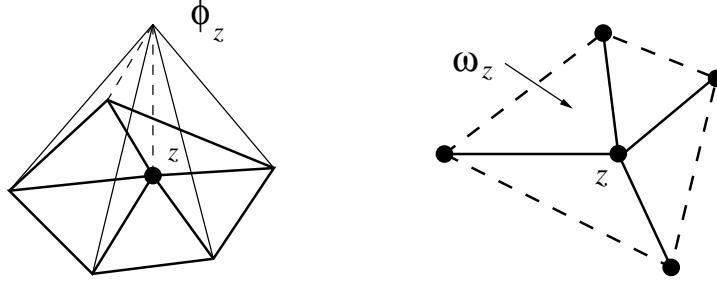


Fig. 7 Piecewise linear basis function ϕ_z corresponding to interior node z , support ω_z of ϕ_z and skeleton γ_z , the latter being composed of all sides within the interior of ω_z .

$\{\phi_z\}_{z \in \mathcal{N}}$ of all such functions is the nodal basis of $\mathbb{S}^{1,0}(\mathcal{T})$, or Courant basis. We denote by $\omega_z := \text{supp}(\phi_z)$ the support of ϕ_z , from now on called *star* associated to z , and by γ_z the skeleton of ω_z , namely all the sides containing z .

We denote functions in $\mathbb{S}^{n,0}(\mathcal{T})$ with capital letters. In view of the definition of ϕ_z , we have the following unique representation of any function $V \in \mathbb{S}^{n,0}(\mathcal{T})$

$$V(x) = \sum_{z \in \mathcal{N}} V(z) \phi_z(x).$$

If we further impose $V(z) = 0$ for all $z \in \partial\Omega \cap \mathcal{N}$, then $V \in H_0^1(\Omega)$. We denote by

$$\mathbb{V}(\mathcal{T}) := \mathbb{S}^{n,0}(\mathcal{T}) \cap H_0^1(\Omega)$$

the subspace of finite element functions which vanish on $\partial\Omega$. Note that we do not explicitly refer to the polynomial degree, which will be clear in each context.

For each simplex $T \in \mathcal{T}$, generated by vertices $\{z_i\}_{i=0}^d$, the *dual functions* $\{\lambda_i^*\}_{i=0}^d \subset \mathbb{P}_1(T)$ to the barycentric coordinates $\{\lambda_i\}_{i=0}^d$ satisfy the bi-orthogonality relation $\int_T \lambda_i^* \lambda_j = \delta_{ij}$, and are given by

$$\lambda_i^* = \frac{(1+d)^2}{|T|} \lambda_i - \frac{1+d}{|T|} \sum_{j \neq i} \lambda_j \quad \forall 0 \leq i \leq d.$$

The *Courant dual basis* $\phi_z^* \in \mathbb{S}^{n-1}(\mathcal{T})$ are the discontinuous piecewise linear functions over \mathcal{T} given by

$$\phi_z^* = \frac{1}{v_z} \sum_{T \ni z} (\lambda_z^T)^* \chi_T \quad \forall z \in \mathcal{N},$$

where $v_z \in \mathbb{N}$ is the valence of z (number of elements of \mathcal{T} containing z) and χ_T is the characteristic function of T . These functions have the same support ω_z as the nodal basis ϕ_z and satisfy the global bi-orthogonality relation

$$\int_{\Omega} \phi_z^* \phi_y = \delta_{zy} \quad \forall z, y \in \mathcal{N}.$$

1.5 Polynomial Interpolation in Sobolev Spaces

If $v \in C^0(\overline{\Omega})$ we define the *Lagrange interpolant* $I_{\mathcal{T}}v$ of v as follows:

$$I_{\mathcal{T}}v(x) = \sum_{z \in \mathcal{N}} v(z) \phi_z(x).$$

For functions without point values, such as functions in $H^1(\Omega)$ for $d > 1$, we need to determine nodal values by averaging. For any conforming refinement $\mathcal{T} \geq \mathcal{T}_0$ of \mathcal{T}_0 , the averaging process extends beyond nodes and so gives rise to the discrete neighborhood

$$N_{\mathcal{T}}(T) := \{T' \in \mathcal{T} \mid T' \cap T \neq \emptyset\}$$

for each element $T \in \mathcal{T}$ along with the *local quasi-uniformity* properties

$$\max_{T \in \mathcal{T}} \#N_{\mathcal{T}}(T) \leq C(\mathcal{T}_0), \quad \max_{T' \in N_{\mathcal{T}}(T)} \frac{|T|}{|T'|} \leq C(\mathcal{T}_0),$$

where $C(\mathcal{T}_0)$ depends only on the shape coefficient of \mathcal{T}_0 given by

$$\sigma(\mathcal{T}_0) := \max_{T \in \mathcal{T}_0} \frac{\bar{h}_T}{h_T}.$$

We introduce now one such operator $I_{\mathcal{T}}$ due to Scott-Zhang [11, 48], from now on called *quasi-interpolation operator*. We focus on polynomial degree $n = 1$, but the construction is valid for any n ; see [11, 48] for details. We recall that $\{\phi_z\}_{z \in \mathcal{N}}$ is the global Lagrange basis of $\mathbb{S}^{1,0}(\mathcal{T})$, $\{\phi_z^*\}_{z \in \mathcal{N}}$ is the global dual basis, and $\text{supp } \phi_z^* = \text{supp } \phi_z$ for all $z \in \mathcal{N}$. We thus define $I_{\mathcal{T}} : L^1(\Omega) \rightarrow \mathbb{S}^{1,0}(\mathcal{T})$ to be

$$I_{\mathcal{T}}v = \sum_{z \in \mathcal{N}} \langle v, \phi_z^* \rangle \phi_z,$$

If $0 \leq s \leq 2$ is a regularity index and $1 \leq p \leq \infty$ is an integrability index, then we would like to prove the *quasi-local error estimate*

$$\|D^t(v - I_{\mathcal{T}}v)\|_{L^q(T)} \lesssim h_T^{\text{sob}(W_p^s) - \text{sob}(W_q^t)} \|D^s v\|_{L^p(N_{\mathcal{T}}(T))} \quad (7)$$

for all $T \in \mathcal{T}$, provided $0 \leq t \leq s$, $1 \leq q \leq \infty$ are such that $\text{sob}(W_p^s) > \text{sob}(W_q^t)$.

We first observe that by construction $I_{\mathcal{T}}$ is invariant in $\mathbb{S}^{1,0}(\mathcal{T})$, namely,

$$I_{\mathcal{T}}P = P \quad \text{for all } P \in \mathbb{S}^{1,0}(\mathcal{T}).$$

Since the averaging process giving rise to the values of $I_{\mathcal{T}}v$ for each element $T \in \mathcal{T}$ takes place in the neighborhood $N_{\mathcal{T}}(T)$, we also deduce the local invariance

$$I_{\mathcal{T}}P|_T = P \quad \text{for all } P \in \mathbb{P}_1(N_{\mathcal{T}}(T))$$

as well as the local stability estimate for any $1 \leq q \leq \infty$

$$\|I_{\mathcal{T}}v\|_{L^q(T)} \lesssim \|v\|_{L^q(N_{\mathcal{T}}(T))}.$$

We thus may write

$$v - I_{\mathcal{T}}v|_T = (v - P) - I_{\mathcal{T}}(v - P)|_T \quad \text{for all } T \in \mathcal{T},$$

where $P \in \mathbb{P}_{s-1}$ is arbitrary ($P = 0$ if $s = 0$). It suffices now to prove (7) in the reference element \hat{T} and scale back and forth to T ; the definition (3) of Sobolev number accounts precisely for this scaling. We keep the notation T for \hat{T} , apply the inverse estimate for linear polynomials $\|D^t(I_{\mathcal{T}}v)\|_{L^q(T)} \lesssim \|I_{\mathcal{T}}v\|_{L^q(T)}$ to $v - P$ instead of v , and use the above local stability estimate, to infer that

$$\|D^t(v - I_{\mathcal{T}}v)\|_{L^q(T)} \lesssim \|v - P\|_{W_q^t(N_{\mathcal{T}}(T))} \lesssim \|v - P\|_{W_p^s(N_{\mathcal{T}}(T))}.$$

The last inequality is a consequence $W_p^s(N_{\mathcal{T}}(T)) \subset W_q^t(N_{\mathcal{T}}(T))$ because $\text{sob}(W_p^s) > \text{sob}(W_q^t)$ and $t \leq s$. Estimate (7) now follows from the Bramble-Hilbert lemma [11, Lemma 4.3.8], [19, Theorem 3.1.1]

$$\inf_{P \in \mathbb{P}_{s-1}(N_{\mathcal{T}}(T))} \|v - P\|_{W_p^s(N_{\mathcal{T}}(T))} \lesssim \|D^s v\|_{L^p(N_{\mathcal{T}}(T))}. \quad (8)$$

This proves (7) for $n = 1$. The construction of $I_{\mathcal{T}}$ and ensuing estimate (7) are still valid for any $n > 1$ [11, 48].

Proposition 1 (Quasi-Interpolant without Boundary Values). *Let s, t be regularity indices with $0 \leq t \leq s \leq n + 1$, and $1 \leq p, q \leq \infty$ be integrability indices so that $\text{sob}(W_p^s) > \text{sob}(W_q^t)$.*

There exists a quasi-interpolation operator $I_{\mathcal{T}} : L^1(\Omega) \rightarrow \mathbb{S}^{n,0}(\mathcal{T})$, which is invariant in $\mathbb{S}^{n,0}(\mathcal{T})$ and satisfies

$$\|D^t(v - I_{\mathcal{T}}v)\|_{L^q(T)} \lesssim h_T^{\text{sob}(W_p^s) - \text{sob}(W_q^t)} \|D^s v\|_{L^p(N_{\mathcal{T}}(T))} \quad \forall T \in \mathcal{T}. \quad (9)$$

The hidden constant in (7) depends on the shape coefficient of \mathcal{T}_0 and d .

To impose a vanishing trace on $I_{\mathcal{T}}v$ we may suitably modify the averaging process for boundary nodes. We thus define a set of dual functions with respect to an L^2 -scalar product over $(d-1)$ -subsimpllices contained on $\partial\Omega$; see again [11, 48] for details. This retains the invariance property of $I_{\mathcal{T}}$ on $\mathbb{S}^{n,0}(\mathcal{T})$ and guarantees that $I_{\mathcal{T}}v$ has a zero trace if $v \in W_1^1(\Omega)$ does. Hence, the above argument applies and (9) follows provided $s \geq 1$.

Proposition 2 (Quasi-Interpolant with Boundary Values). *Let s, t, p, q be as in Proposition 1. There exists a quasi-interpolation operator $I_{\mathcal{T}} : W_1^1(\Omega) \rightarrow \mathbb{S}^{n,0}(\mathcal{T})$ invariant in $\mathbb{S}^{n,0}(\mathcal{T})$ which satisfies (9) for $s \geq 1$ and preserves the boundary values of v provided they are piecewise polynomial of degree $\leq n$. In particular, if $v \in W_1^1(\Omega)$ has a vanishing trace on $\partial\Omega$, then so does $I_{\mathcal{T}}v$.*

Remark 3 (Fractional Regularity). We observe that (7) does not require the regularity indices t and s to be integer. The proof follows the same lines but replaces the polynomial degree $s-1$ by the greatest integer smaller than s ; the generalization of (8) can be taken from [26].

Remark 4 (Local Error Estimate for Lagrange Interpolant). Let the regularity index s and integrability index $1 \leq p \leq \infty$ satisfy $s - d/p > 0$. This implies that $\text{sob}(W_p^s) > \text{sob}(L^\infty)$, whence $W_p^s(\Omega) \subset C(\overline{\Omega})$ and the Lagrange interpolation operator $I_{\mathcal{T}} : W_p^s(\Omega) \rightarrow \mathbb{S}^{n,0}(\mathcal{T})$ is well defined and satisfies the *local error estimate*

$$\|D^t(v - I_{\mathcal{T}}v)\|_{L^q(T)} \lesssim h_T^{\text{sob}(W_p^s) - \text{sob}(W_q^t)} \|D^s v\|_{L^p(T)}, \quad (10)$$

provided $0 \leq t \leq s$, $1 \leq q \leq \infty$ are such that $\text{sob}(W_p^s) > \text{sob}(W_q^t)$. We point out that $N_{\mathcal{T}}(T)$ in (7) is now replaced by T in (10). We also remark that if v vanishes on $\partial\Omega$ so does $I_{\mathcal{T}}v$. The proof of (10) proceeds along the same lines as that of Proposition 1 except that the nodal evaluation does not extend beyond the element $T \in \mathcal{T}$ and the inverse and stability estimates over the reference element are replaced by

$$\|D^t I_{\mathcal{T}}v\|_{L^q(\hat{T})} \lesssim \|I_{\mathcal{T}}v\|_{L^q(\hat{T})} \lesssim \|v\|_{L^\infty(\hat{T})} \lesssim \|v\|_{W_p^s(\hat{T})}.$$

We are now in a position to derive a global interpolation error estimate. To this end, it is convenient to introduce the mesh-size function $h \in L^\infty(\Omega)$ given by

$$h|_T = h_T \quad \text{for all } T \in \mathcal{T}. \quad (11)$$

Notice that the following estimate encompasses the linear as well as the nonlinear Sobolev scales.

Theorem 2 (Global Interpolation Error Estimate). *Let $1 \leq s \leq n+1$ and $1 \leq p \leq 2$ satisfy $r := \text{sob}(W_p^s) - \text{sob}(H^1) > 0$. If $v \in W_p^s(\Omega)$, then*

$$\|\nabla(v - I_{\mathcal{T}}v)\|_{L^2(\Omega)} \lesssim \|h^r D^s v\|_{L^p(\Omega)}. \quad (12)$$

Proof. Use Proposition 1 along with the elementary property of series $\sum_n a_n \leq (\sum_n a_n^q)^{1/q}$ for $0 < q := p/2 \leq 1$. \square

Quasi-Uniform Meshes. We now apply Theorem 2 to quasi-uniform meshes, namely meshes $\mathcal{T} \in \mathbb{T}$ for which all its elements are of comparable size h regardless of the refinement level. In this case, we have

$$h \approx (\#\mathcal{T})^{-1/d}.$$

Corollary 1 (Quasi-Uniform Meshes). *Let $1 \leq s \leq n+1$ and $u \in H^s(\Omega)$. If $\mathcal{T} \in \mathbb{T}$ is quasi-uniform, then*

$$\|\nabla(v - I_{\mathcal{T}}v)\|_{L^2(\Omega)} \lesssim |v|_{H^s(\Omega)} (\#\mathcal{T})^{-(s-1)/d}. \quad (13)$$

Remark 5 (Optimal Rate). If $s = n+1$, and so v has the maximal regularity $v \in H^{n+1}(\Omega)$, then we obtain the optimal convergence rate in a linear Sobolev scale

$$\|\nabla(v - I_{\mathcal{T}}v)\|_{L^2(\Omega)} \lesssim |v|_{H^{n+1}(\Omega)} (\#\mathcal{T})^{-n/d}. \quad (14)$$

The order $-n/d$ is just dictated by the polynomial degree n and cannot be improved upon assuming either higher regularity than $H^{n+1}(\Omega)$ or a graded mesh \mathcal{T} .

Example (Corner Singularity in 2d). To explore the effect of a geometric singularity on (13), we let Ω be the L-shaped domain of Figure 1.5 and $v \in H^1(\Omega)$ be

$$u(r, \theta) = r^{2/3} \sin(2\theta/3) - r^2/4.$$

This function $v \in H^1(\Omega)$ exhibits the typical corner singularity of the solution of $-\Delta v = f$ with suitable Dirichlet boundary condition: $v \in H^s(\Omega)$ for $s < 5/3$. Table 1 displays the best approximation error for polynomial degree $n = 1, 2, 3$ and the sequence of *uniform* refinements depicted in Figure 1.5 in the seminorm $|\cdot|_{H^1(\Omega)}$. This gives a *lower* bound for the interpolation error in (13).

h	linear ($n = 1$)	quadratic ($n = 2$)	cubic ($n = 3$)
1/4	1.14	9.64	9.89
1/8	0.74	0.67	0.67
1/16	0.68	0.67	0.67
1/32	0.66	0.67	0.67
1/64	0.66	0.67	0.67
1/128	0.66	0.67	0.67

Table 1 The asymptotic rate of convergence in term of mesh-size h is about $h^{2/3}$, or equivalently $(\#\mathcal{T})^{-1/3}$, irrespective of the polynomial degree n . This provides a lower bound for $\|v - I_{\mathcal{T}}v\|_{L^2(\Omega)}$ and thus shows that (13) is sharp.

Even though s is fractional, the error estimate (13) is still valid as stated in Remark 3. In fact, for uniform refinement, (13) can be derived by space interpolation between $H^1(\Omega)$ and $H^{n+1}(\Omega)$. The asymptotic rate $(\#\mathcal{T})^{-1/3}$ reported in Table 1 is

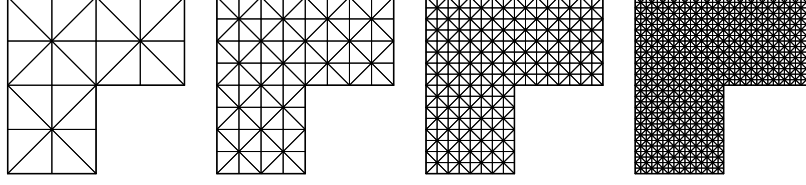


Fig. 8 Sequence of consecutive uniform meshes for L-shaped domain Ω created by 2 bisections.

consistent with (13) and independent of the polynomial degree n ; this shows that (13) is sharp. It is also suboptimal as compared with the optimal rate $(\#\mathcal{T})^{-n/2}$ of Remark 5.

The question arises whether the rate $(\#\mathcal{T})^{-1/3}$ in Table 1 is just a consequence of uniform refinement or unavoidable. It is important to realize that $v \notin H^s(\Omega)$ for $s \geq 5/3$ and thus (13) is not applicable. However, the problem is not that second order derivatives of v do not exist but rather that they are not square-integrable. In particular, it is true that $v \in W_p^2(\Omega)$ if $1 \leq p < 3/2$. We therefore may apply Theorem 2 with, e.g., $n = 1$, $s = 2$, and $p \in [1, 3/2)$ and then ask whether the structure of (12) can be exploited, e.g., by compensating the local behavior of $D^s u$ with the local mesh-size h . This enterprise naturally leads to *graded* meshes adapted to v .

1.6 Adaptive Approximation

Principle of Error Equidistribution. We investigate the relation between local mesh-size and regularity for the design of graded meshes adapted to $v \in H^1(\Omega)$ for $d = 2$. We formulate this as an optimization problem:

Given a function $v \in C^2(\Omega) \cap W_p^2(\Omega)$ and an integer $N > 0$ find conditions for a shape regular mesh \mathcal{T} to minimize the error $|v - I_{\mathcal{T}}v|_{H^1(\Omega)}$ subject to the constraint that the number of degrees of freedom $\#\mathcal{T} \leq N$.

We first convert this *discrete* optimization problem into a *continuous model*, following Babuška and Rheinboldt [5]. Let

$$\#\mathcal{T} = \int_{\Omega} \frac{dx}{h(x)^2}$$

be the number of elements of \mathcal{T} and let the Lagrange interpolation error

$$\|\nabla(v - I_{\mathcal{T}}v)\|_{L^2(\Omega)}^p = \int_{\Omega} h(x)^{2(p-1)} |D^2 v(x)|^p dx$$

be dictated by (12) with $s = 2$ and $1 < p \leq 2$; note that $r = \text{sob}(W_p^2) - \text{sob}(H^1) = 2 - 2/p$ whence $rp = 2(p-1)$ is the exponent of $h(x)$. We next propose the Lagrangian

$$\mathcal{L}[h, \lambda] = \int_{\Omega} \left(h(x)^{2(p-1)} |D^2 v(x)|^p - \frac{\lambda}{h(x)^2} \right) dx$$

with Lagrange multiplier $\lambda \in \mathbb{R}$. The optimality condition reads (Problem 4)

$$h(x)^{2(p-1)+2} |D^2 v(x)|^p = \Lambda \quad (15)$$

where $\Lambda > 0$ is a constant. In order to interpret this expression, we compute the interpolation error E_T incurred in element $T \in \mathcal{T}$. According to (10), E_T is given by

$$E_T^p \approx h_T^{2(p-1)} \int_T |D^2 v(x)|^p \approx \Lambda$$

provided $D^2 v(x)$ is about constant in T . Therefore we reach the heuristic, but insightful, conclusion that E_T is about constant, or equivalently

$$\text{A graded mesh is quasi-optimal if the local error is equidistributed.} \quad (16)$$

Corner Singularities. Meshes satisfying (16) have been constructed by Babuška et al [3] for corner singularities and $d = 2$; see also [30]. If the function v possess the typical behavior

$$v(x) \approx r(x)^\gamma, \quad 0 < \gamma < 1,$$

where $r(x)$ is the distance from $x \in \Omega$ to a reentrant corner of Ω , then (15) implies the mesh grading

$$h(x) = \Lambda^{\frac{1}{2p}} r(x)^{-\frac{1}{2}(\gamma-2)}$$

whence

$$\#\mathcal{T} = \int_{\Omega} h(x)^{-2} dx = \Lambda^{-\frac{1}{p}} \int_0^{\text{diam}(\Omega)} r^{\gamma-1} dr \approx \Lambda^{-\frac{1}{p}}.$$

This crucial relation is valid for any $\gamma > 0$ and $p > 1$; in fact the only condition on p is that $r = 2 - 2/p > 0$, or equivalently $\text{sob}(W_p^2) > \text{sob}(H^1)$. Therefore,

$$\|\nabla(v - I_{\mathcal{T}}v)\|_{L^2(\Omega)}^2 = \sum_{T \in \mathcal{T}} E_T^2 = \Lambda^{\frac{2}{p}} (\#\mathcal{T}) = (\#\mathcal{T})^{-1} \quad (17)$$

gives the optimal decay rate for $d = 2, n = 1$, according to Remark 5. We explore the case $d \geq 2$ and $n \geq 1$ in Problem 6. What this argument does not address is whether such meshes \mathcal{T} exist in general and, more importantly, whether they can actually be constructed upon bisecting the initial mesh \mathcal{T}_0 so that $\mathcal{T} \in \mathbb{T}$.

Thresholding. We now construct graded bisection meshes \mathcal{T} for $n = 1, d = 2$ that achieve the optimal decay rate $(\#\mathcal{T})^{-1/2}$ of (14) and (17) under the global regularity assumption

$$v \in W_p^2(\Omega), \quad p > 1. \quad (18)$$

Following the work of Binev et al. [8], we use a thresholding algorithm that is based on the knowledge of the element errors and on bisection. The algorithm hinges

on (16): if $\delta > 0$ is a given tolerance, the element error is equidistributed, that is $E_T \approx \delta^2$, and the global error decays with maximum rate $(\#\mathcal{T})^{-1/2}$, then

$$\delta^4 \#\mathcal{T} \approx \sum_{T \in \mathcal{T}} E_T^2 = |v - I_{\mathcal{T}}v|_{H^1(\Omega)}^2 \lesssim (\#\mathcal{T})^{-1}$$

that is $\#\mathcal{T} \lesssim \delta^{-2}$. With this in mind, we impose $E_T \leq \delta^2$ as a common threshold to stop refining and expect $\#\mathcal{T} \lesssim \delta^{-2}$. The following algorithm implements this idea.

Thresholding Algorithm. Given a tolerance $\delta > 0$ and a conforming mesh \mathcal{T}_0 , THRESHOLD finds a conforming refinement $\mathcal{T} \geq \mathcal{T}_0$ of \mathcal{T}_0 by bisection such that $E_T \leq \delta^2$ for all $T \in \mathcal{T}$: let $\mathcal{T} = \mathcal{T}_0$ and

```

THRESHOLD( $\mathcal{T}, \delta$ )
while  $\mathcal{M} := \{T \in \mathcal{T} \mid E_T > \delta^2\} \neq \emptyset$ 
   $\mathcal{T} := \text{REFINE}(\mathcal{T}, \mathcal{M})$ 
end while
return( $\mathcal{T}$ )

```

We get $W_p^2(\Omega) \subset C^0(\overline{\Omega})$, because $p > 1$, and can use the Lagrange interpolant and local estimate (10) with $r = \text{sob}(W_p^2) - \text{sob}(H^1) = 2 - 2/p > 0$. We deduce that

$$E_T \lesssim h_T^r \|D^2v\|_{L^p(T)}, \quad (19)$$

and that THRESHOLD *terminates* because h_T decreases monotonically to 0 with bisection. The quality of the resulting mesh is assessed next.

Theorem 3 (Thresholding). *If $v \in H_0^1(\Omega)$ verifies (18), then the output $\mathcal{T} \in \mathbb{T}$ of THRESHOLD satisfies*

$$|v - I_{\mathcal{T}}v|_{H^1(\Omega)} \leq \delta^2 (\#\mathcal{T})^{1/2}, \quad \#\mathcal{T} - \#\mathcal{T}_0 \lesssim \delta^{-2} |\Omega|^{1-1/p} \|D^2v\|_{L^p(\Omega)}.$$

Proof. Let $k \geq 1$ be the number of iterations of THRESHOLD before termination. Let $\mathcal{M} = \mathcal{M}_0 \cup \dots \cup \mathcal{M}_{k-1}$ be the set of marked elements. We organize the elements in \mathcal{M} by size in such a way that allows for a counting argument. Let \mathcal{P}_j be the set of elements T of \mathcal{M} with size

$$2^{-(j+1)} \leq |T| < 2^{-j} \quad \Rightarrow \quad 2^{-(j+1)/2} \leq h_T < 2^{-j/2}.$$

We proceed in several steps.

□ We first observe that all T 's in \mathcal{P}_j are *disjoint*. This is because if $T_1, T_2 \in \mathcal{P}_j$ and $T_1 \cap T_2 \neq \emptyset$, then one of them is contained in the other, say $T_1 \subset T_2$, due to the bisection procedure. Thus

$$|T_1| \leq \frac{1}{2} |T_2|$$

contradicting the definition of \mathcal{P}_j . This implies

$$2^{-(j+1)} \#\mathcal{P}_j \leq |\Omega| \quad \Rightarrow \quad \#\mathcal{P}_j \leq |\Omega| 2^{j+1}. \quad (20)$$

□² In light of (19), we have for $T \in \mathcal{P}_j$

$$\delta^2 \leq E_T \lesssim 2^{-(j/2)r} \|D^2 v\|_{L^p(T)}.$$

Therefore

$$\delta^{2p} \#\mathcal{P}_j \lesssim 2^{-(j/2)rp} \sum_{T \in \mathcal{P}_j} \|D^2 v\|_{L^p(T)}^p \leq 2^{-(j/2)rp} \|D^2 v\|_{L^p(\Omega)}^p$$

whence

$$\#\mathcal{P}_j \lesssim \delta^{-2p} 2^{-(j/2)rp} \|D^2 v\|_{L^p(\Omega)}^p. \quad (21)$$

□³ The two bounds for $\#\mathcal{P}$ in (20) and (21) are complementary. The first is good for j small whereas the second is suitable for j large (think of $\delta \ll 1$). The crossover takes place for j_0 such that

$$2^{j_0+1} |\Omega| = \delta^{-2p} 2^{-j_0(rp/2)} \|D^2 v\|_{L^p(\Omega)}^p \Rightarrow 2^{j_0} \approx \delta^{-2} \frac{\|D^2 v\|_{L^p(\Omega)}}{|\Omega|^{1/p}}.$$

□⁴ We now compute

$$\#\mathcal{M} = \sum_j \#\mathcal{P}_j \lesssim \sum_{j \leq j_0} 2^j |\Omega| + \delta^{-2p} \|D^2 v\|_{L^p(\Omega)}^p \sum_{j > j_0} (2^{-rp/2})^j.$$

Since

$$\sum_{j \leq j_0} 2^j \approx 2^{j_0}, \quad \sum_{j > j_0} (2^{-rp/2})^j \lesssim 2^{-(rp/2)j_0} = 2^{-(p-1)j_0}$$

we can write

$$\#\mathcal{M} \lesssim (\delta^{-2} + \delta^{-2p} \delta^{2(p-1)}) |\Omega|^{1-1/p} \|D^2 v\|_{L^p(\Omega)} \approx \delta^{-2} |\Omega|^{1-1/p} \|D^2 v\|_{L^p(\Omega)}.$$

We finally apply Theorem 1 to arrive at

$$\#\mathcal{T} - \#\mathcal{T}_0 \lesssim \#\mathcal{M} \lesssim \delta^{-2} |\Omega|^{1-1/p} \|D^2 v\|_{L^p(\Omega)}.$$

□⁵ It remains to estimate the energy error. We have, upon termination of THRESHOLD, that $E_T \leq \delta^2$ for all $T \in \mathcal{T}$. Then

$$|v - I_{\mathcal{T}} v|_{H^1(\Omega)}^2 = \sum_{T \in \mathcal{T}} E_T^2 \leq \delta^4 \#\mathcal{T}.$$

This concludes the Theorem. □

By relating the threshold value δ and the number of refinements N , we obtain a result about the convergence rate.

Corollary 2 (Convergence Rate). *Let $v \in H_0^1(\Omega)$ satisfy (18). Then for $N > \#\mathcal{T}_0$ integer there exists $\mathcal{T} \in \mathbb{T}$ such that*

$$|v - I_{\mathcal{T}}v|_{H^1(\Omega)} \lesssim |\Omega|^{1-1/p} \|D^2v\|_{L^p(\Omega)} N^{-1/2}, \quad \#\mathcal{T} - \#\mathcal{T}_0 \lesssim N.$$

Proof. Choose $\delta^2 = |\Omega|^{1-1/p} \|D^2v\|_{L^p(\Omega)} N^{-1}$ in Theorem 3. Then, there exists $\mathcal{T} \in \mathbb{T}$ such that $\#\mathcal{T} - \#\mathcal{T}_0 \lesssim N$ and

$$\begin{aligned} |v - I_{\mathcal{T}}v|_{H^1(\Omega)} &\lesssim |\Omega|^{1-1/p} \|D^2v\|_{L^p(\Omega)} N^{-1} (N + \#\mathcal{T}_0)^{1/2} \\ &\lesssim |\Omega|^{1-1/p} \|D^2v\|_{L^p(\Omega)} N^{-1/2} \end{aligned}$$

because $N > \#\mathcal{T}_0$. This finishes the Corollary. \square

Remark 6 (Piecewise smoothness). The global regularity (18) can be weakened to *piecewise* W_p^2 *regularity* over the initial mesh \mathcal{T}_0 , namely $W_p^2(\Omega; \mathcal{T}_0)$, and global $H_0^1(\Omega)$. This is because $W_p^2(T) \hookrightarrow C^0(\bar{T})$ for all $T \in \mathcal{T}_0$, whence $I_{\mathcal{T}}$ can be taken to be the Lagrange interpolation operator.

Remark 7 (Case $p < 1$). We consider now polynomial degree $n \geq 1$. The integrability p corresponding to differentiability $n+1$ results from equating Sobolev numbers:

$$n+1 - \frac{d}{p} = \text{sob}(H^1) = 1 - \frac{d}{2} \quad \Rightarrow \quad p = \frac{2d}{2n+d}.$$

Depending on $d \geq 2$ and $n \geq 1$, this may lead to $0 < p < 1$, in which case $W_p^{n+1}(\Omega)$ is to be replaced by the Besov space $B_{p,p}^{n+1}(\Omega)$ [22]; see Problem 6. The argument of Theorem 3 works provided we replace (19) by a modulus of regularity; in fact, $D^{n+1}v$ would not be locally integrable and so would fail to be a distribution.

Remark 8 (Isotropic vs anisotropic elements). Theorem 3 and Problem 5 show that isotropic graded meshes can always deal with geometric singularities for $d = 2$. This is no longer the case for $d > 2$ and is explored in Problem 6.

1.7 Nonconforming Meshes

More general subdivisions of Ω than those in §1.3 are used in practice. If the elements of \mathcal{T}_0 are quadrilaterals for $d = 2$, or their multidimensional variant for $d > 2$, then it is natural to allow for improper or *hanging nodes* for the resulting refinements \mathcal{T} to be graded; see Figure 9 (a). On the other hand, if \mathcal{T}_0 is made of triangles for $d = 2$, or simplices for $d > 2$, then red refinement without green completion also gives rise to graded meshes with hanging nodes; see Figure 9 (b). In both cases, the presence of hanging nodes is inevitable to enforce mesh grading. Finally, bisection may produce meshes with hanging nodes, as depicted in Figure 9 (c), if the completion process is incomplete. All three refinements maintain shape regularity, but for both practice and theory, they cannot be arbitrary: we need to restrict the level of incompatibility; see Problem 10. We discuss this next.

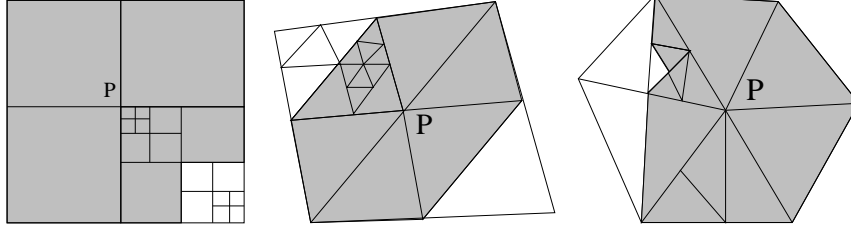


Fig. 9 Nonconforming meshes made of quadrilaterals (a), triangles with red refinement (b), and triangles with bisection (c). The shaded regions depict the domain of influence of a proper or conforming node P .

We start with the notion of domain of influence of a proper node, introduced by Babuška and Miller in the context of K -meshes [4]; see Figure 9. For simplicity, we restrict ourselves to polynomial degree $n = 1$. We say that a node P of \mathcal{T} is a *proper* (or *conforming*) node if it is a vertex of all elements containing P ; otherwise, we say that P is an *improper* (nonconforming or hanging) node. Since we only prescribe degrees of freedom at the proper nodes, it is natural to describe the canonical continuous piecewise linear basis functions ϕ_P associated with each proper node P .

We do this recursively. As in §1.3, the *generation* $g(T)$ of an element $T \in \mathcal{T}$ is the number of subdivisions needed to create T from its ancestor in the initial mesh \mathcal{T}_0 , hereafter assumed to be conforming. We first rearrange the elements in $\mathcal{T} = \{T_i\}_{i=1}^{\#\mathcal{T}}$ by generation: $g(T_i) \leq g(T_{i+1})$ for all $i \geq 0$. Suppose that ϕ_P has been already defined for each $T \in \mathcal{T}$ with $g(T) < i$. We proceed as follows to define ϕ_P at each vertex z of $T \in \mathcal{T}$ with $g(T) = i$:

- if z is a proper node, then we set $\phi_P(z) = 1$ if $z = P$ and $\phi_P(z) = 0$ otherwise;
- if z is a hanging node, then z belongs to an edge of another element $T' \in \mathcal{T}$ with $g(T') < i$ and set $\phi_P(z)|_T = \phi_P(z)|_{T'}$.

This definition is independent of the choice of T' since, by construction, ϕ_P is continuous across interelements of lower level. We also observe that $\{\phi_P\}_{P \in \mathcal{N}}$ is a basis of the finite element space $\mathbb{V}(\mathcal{T})$ of *continuous* piecewise linear functions, thus

$$V = \sum_{P \in \mathcal{N}} V(P)\phi_P \quad \forall V \in \mathbb{V}(\mathcal{T}).$$

The *domain of influence* of a proper node P is the support of ϕ_P :

$$\omega_{\mathcal{T}}(P) = \text{supp}(\phi_P).$$

We say that a sequence of nonconforming meshes $\{\mathcal{T}\}$ is *admissible* if there is a universal constant $\Lambda_* \leq 1$, independent of the refinement level and \mathcal{T} , such that

$$\text{diam}(\omega_{\mathcal{T}}(P)) \leq \Lambda_* h_T \quad \forall T \in \mathcal{T}. \quad (22)$$

An important example is quadrilaterals with *one* hanging node per edge. We observe, however, that (22) can neither be guaranteed with more than one hanging node per edge for quadrilaterals, nor for triangles with one hanging node per edge (see Problem 10).

Given an admissible grid \mathcal{T} , a subset \mathcal{M} of elements marked for refinement, and a desired number $\rho \geq 1$ of subdivisions to be performed in each marked element, the procedure

$$\mathcal{T}_* = \text{REFINE}(\mathcal{T}, \mathcal{M})$$

creates a minimal admissible mesh $\mathcal{T}_* \geq \mathcal{T}$ such that all the elements of \mathcal{M} are subdivided at least ρ times. In order for \mathcal{T}_* to be admissible, perhaps other elements not in \mathcal{M} must be partitioned. Despite the fact that admissibility is a constraint on the refinement procedure weaker than conformity, it cannot avoid the propagation of refinements beyond \mathcal{M} . The complexity of REFINE is again an issue which we discuss in §6.4: we show that Theorem 1 extends to this case.

Lemma 3 (REFINE for Nonconforming Meshes). *Let \mathcal{T}_0 be an arbitrary conforming partition of Ω , except for bisection in which case \mathcal{T}_0 satisfies the labeling (6) for $d = 2$ or its higher dimensional counterpart [53]. Then the estimate*

$$\#\mathcal{T}_k - \#\mathcal{T}_0 \leq \Lambda_0 \sum_{j=0}^{k-1} \#\mathcal{M}_j \quad \forall k \geq 1$$

holds with a constant Λ_0 depending on \mathcal{T}_0, d and ρ .

We conclude by emphasizing that the polynomial interpolation and adaptive approximation theories of §§1.5 and 1.6 extend to nonconforming meshes with fixed level of incompatibility as well.

1.8 Notes

The use of Sobolev numbers is not so common in the finite element literature, but allows as to write compact error estimates and speak about nonlinear Sobolev scale. The latter concept is quite natural in nonlinear approximation theory [22].

The discussion of bisection for $d = 2$ follows Binev, Dahmen, and DeVore [7]. Stevenson extended the theory to $d > 2$ [53]. We refer to the survey by Nochetto, Siebert and Veiser [45] for a rather complete discussion for $d > 1$, and to §6 for a proof of Theorem 1 for $d = 2$, which easily extends to $d > 2$.

The discussion of finite element spaces [10, 11, 19] and polynomial interpolation [11, 26, 48] is rather classical. In contrast, the material of adaptive approximation is much less documented. The principle of equidistribution goes back to Babuška and Rheinboldt [5] and the a priori design of optimal meshes for corner singularities for $d = 2$ is due to Babuška, Kellogg, and Pitkäranta [3]. The construction of optimal meshes via bisection using thresholding is extracted from Binev, Dahmen, DeVore, and Petrushev [8].

Finally the discussion of nonconforming meshes follows Bonito and Nochetto [9], and continues in §6 with the proof of Lemma 3.

1.9 Problems

Problem 1 (Nonconforming element). Given a d -simplex T in R^d with vertices z_0, \dots, z_d , construct a basis $\bar{\lambda}_0, \dots, \bar{\lambda}_d$ of $\mathbb{P}_1(T)$ such that $\bar{\lambda}_i(\bar{z}_j) = \delta_{ij}$ for all $i, j \in \{1, \dots, d\}$, where \bar{z}_j denotes the barycenter of the face opposite to the vertex z_j . Does this local basis also lead to a global one in $\mathbb{S}^{1,0}(\mathcal{T})$?

Problem 2 (Quadratic basis functions). Express the nodal basis of $\mathbb{P}_2(T)$ in terms of barycentric coordinates of $T \in \mathcal{T}$.

Problem 3 (Quadratic dual functions). Derive expressions for the dual functions of the quadratic local Lagrange basis of $\mathbb{P}_2(T)$ for each element $T \in \mathcal{T}$. Construct a global discontinuous dual basis $\phi_z^* \in \mathbb{S}^{2,-1}(\mathcal{T})$ of the global Lagrange basis $\phi_z \in \mathbb{S}^{2,0}(\mathcal{T})$ for all $z \in \mathcal{N}_2(\mathcal{T})$.

Problem 4 (Lagrangian). Let $h(x)$ be a smooth function locally equivalent to the mesh-size and $v \in C^2(\Omega) \cap W_p^2(\Omega)$. Prove that a stationary point of the Lagrangian

$$\mathcal{L}[h, \lambda] = \int_{\Omega} \left(h(x)^{2(p-1)} |D^2 v(x)|^p - \frac{\lambda}{h(x)^2} \right) dx$$

satisfies the optimality condition

$$h(x)^{2(p-1)+2} |D^2 v(x)|^p = \text{constant}.$$

Problem 5 (W_p^2 -regularity). Consider the function $v(r, \theta) = r^\gamma \phi(\theta)$ in polar coordinates (r, θ) for $d = 2$ with $\phi(\theta)$ smooth. Show that $v \in W_p^2(\Omega) \setminus H^2(\Omega)$ for $1 \leq p < 2/(2 - \gamma)$.

Problem 6 (Edge singularities). This problem explores *formally* the effect of edge singularities for dimension $d > 2$ and polynomial degree $n \geq 1$. Since edge (or line) singularities are two dimensional locally, away from corners, we assume the behavior $v(x) \approx r(x)^\gamma$ where $r(x)$ is the distance of $x \in \Omega$ to an edge of Ω and $\gamma > 0$.

(a) Use the Principle of Equidistribution with $p = 2$ to determine the mesh grading

$$h(x) \approx \Lambda^{\frac{1}{2n+d}} r(x)^{2d \frac{\gamma-(n+1)}{2n+d}}.$$

(b) Show the following relation between Λ and number of elements $\#\mathcal{T} = \int_{\Omega} h(x)^{-d}$

$$\gamma > \frac{(d-2)n}{d} \quad \Rightarrow \quad \#\mathcal{T} \approx \Lambda^{-\frac{d}{2n+d}}.$$

(c) If $\gamma > \frac{(d-2)n}{d}$, then deduce the optimal interpolation error decay

$$\|\nabla(v - I_{\mathcal{T}}v)\|_{L^2(\Omega)} \lesssim (\#\mathcal{T})^{-\frac{n}{d}}.$$

(d) Prove that $\gamma > \frac{(d-2)n}{d}$ is equivalent to the regularity $\int_{\Omega} |D^{n+1}v|^p < \infty$ for $p > \frac{2d}{2n+d}$. If $\tau := \frac{2d}{2n+d} \geq 1$, then this would mean $v \in W_p^{n+1}(\Omega)$. However, it is easy to find examples $d > 2$ or $n > 1$ for which $\tau < 1$, in which case the Sobolev space $W_p^{n+1}(\Omega)$ must be replaced by the Besov space $B_{p,p}^{n+1}(\Omega)$ [22]. Note that $p > \tau$ is precisely what yields the crucial relation between Sobolev numbers

$$\text{sob}(B_{p,p}^{n+1}) = n + 1 - \frac{d}{p} > \text{sob}(H^1) = 1 - \frac{d}{2}.$$

We observe that for $d = 2$ all singular exponents $\gamma > 0$ lead to optimal meshes, but this is not true for $d = 3$: $n = 1$ requires $\gamma > \frac{1}{3}$ whereas $n = 2$ needs $\gamma > \frac{2}{3}$. The latter corresponds to a dihedral angle $\omega > \frac{3\pi}{2}$ and can be easily checked computationally. We thus conclude that *isotropic* graded meshes are sufficient to deal with geometric singularities for $d = 2$ but not for $d > 2$, for which *anisotropic* graded meshes are the only ones which exhibit optimal behavior. Their adaptive construction is open.

Problem 7 (Local H^2 -regularity). Consider the function $v(x) \approx r(x)^\gamma$ where $r(x)$ is the distance to the origin and $d = 2$.

- (a) Examine the construction of a graded mesh via the thresholding algorithm.
- (b) Repeat the proof of Theorem 3 replacing the W_p^2 regularity by the corresponding local H^2 -regularity of v depending on the distance to the origin.

Problem 8 (Thresholding for $d > 2$). Let $d > 2$, $n = 1$, and $v \in W_p^2(\Omega)$ with $p > \frac{2d}{2+d}$. This implies that $v \in H^1(\Omega)$ but not necessarily in $C^0(\overline{\Omega})$. Use the quasi-interpolant $I_{\mathcal{T}}$ of Proposition 1 to define the local H^1 -error E_T for each element $T \in \mathcal{T}$ and use the thresholding algorithm to show Theorem 3 and Corollary 2.

Problem 9 (Reduced rate). Let $d \geq 2$, $n = 1$, and $v \in W_p^s(\Omega)$ with $1 < s < 2$ and $\text{sob}(W_p^s) > \text{sob}(H^1)$, namely $s - \frac{d}{p} > 1 - \frac{d}{2}$. Use the quasi-interpolant $I_{\mathcal{T}}$ of Proposition 1 to define the local H^1 -error E_T for each element $T \in \mathcal{T}$ and use the thresholding algorithm to show Corollary 2: given $N > \#\mathcal{T}_0$ there exists $\mathcal{T} \in \mathbb{T}$ with $\#\mathcal{T} - \#\mathcal{T}_0 \lesssim N$ such that

$$\|v - I_{\mathcal{T}}v\|_{H^1(\Omega)} \lesssim \|D^s v\|_{L^p(\Omega)} N^{-\frac{s-1}{d}}.$$

Problem 10 (Level of incompatibility). This problem shows that keeping the number of hanging nodes per side bounded does not guarantee a bounded level of incompatibility for $d = 2$. The situation is similar for $d > 2$.

- (a) *Square elements*: construct a selfsimilar quad-refinement of the unit square with only 2 hanging nodes per side and unbounded level of incompatibility.
- (b) *Triangular elements*: construct selfsimilar red-refinements and bisection refinements of the unit reference triangle with 1 hanging node per side and unbounded level of incompatibility.

Problem 11 (Quasi-interpolation of discontinuous functions). Let \mathcal{T} be an admissible nonconforming mesh. Let $\mathbb{V}(\mathcal{T})$ denote the space of discontinuous piecewise polynomials of degree $\leq n$ over \mathcal{T} , and $\mathbb{V}^0(\mathcal{T})$ be the subspace of continuous functions. Construct a local quasi-interpolation operator $I_{\mathcal{T}} : \mathbb{V}(\mathcal{T}) \rightarrow \mathbb{V}^0(\mathcal{T})$ with the following approximation property for all $V \in \mathbb{V}(\mathcal{T})$ and $|\alpha| = 0, 1$

$$\|D^\alpha(V - I_{\mathcal{T}}V)\|_{L^2(T)} \lesssim h_T^{\frac{1-|\alpha|}{2}} \|\llbracket V \rrbracket\|_{L^2(\Sigma_{\mathcal{T}}(T))} \quad \forall T \in \mathcal{T},$$

where $\Sigma_{\mathcal{T}}(T)$ stands for all sides within $N_{\mathcal{T}}(T)$ and $\llbracket V \rrbracket$ denotes the jump of V across sides.

2 Error Bounds for Finite Element Solutions

In §1 we have seen that approximating a given known function with meshes which are adapted to that function can impressively outperform the approximation with quasi-uniform meshes. In view of the fact that the solution of a boundary value problem is given only implicitly, it is not all clear if this is also true for its adaptive numerical solution. Considering a simple model problem and discretization, we now derive two upper bounds for the error of the finite element solution: an a priori one and an a posteriori one. The a priori bound reveals that an adaptive variant of the finite element method has the potential of a similar performance. The a posteriori bound is a first step to design such a variant, which has to face the complication that the target function is given only implicitly.

2.1 Model Boundary Value Problem

In order to minimize technicalities in the presentation, let us consider the following simple boundary value problem as a model problem: find a scalar function $u = u(x)$ such that

$$\begin{aligned} -\operatorname{div}(\mathbf{A}\nabla u) &= f && \text{in } \Omega, \\ u &= 0 && \text{on } \partial\Omega, \end{aligned} \quad (23)$$

where $\Omega \subset \mathbb{R}^d$ is a bounded domain with Lipschitz boundary $\partial\Omega$, $\mathbf{A} = \mathbf{A}(x)$ a map into the positive definite $d \times d$ matrices, and $f = f(x)$ a scalar load term. Introducing the Hilbert space

$$\mathbb{V} := H_0^1(\Omega) := \{v \in H^1(\Omega) \mid v|_{\partial\Omega} = 0\}, \quad \|v\|_{\mathbb{V}} := \left(\int_{\Omega} |\nabla v|^2 \right)^{1/2},$$

and the bilinear form

$$\mathcal{B}[v, w] := \int_{\Omega} \mathbf{A}\nabla v \cdot \nabla w, \quad v, w \in \mathbb{V},$$

the weak solution of (23) is characterized by

$$u \in \mathbb{V} : \quad \mathcal{B}[u, v] = \langle f, v \rangle \quad \text{for all } v \in \mathbb{V}. \quad (24)$$

Hereafter $\langle \cdot, \cdot \rangle$ stands for the $L^2(\Omega)$ -scalar product and also for a duality pairing. We assume that $f \in \mathbb{V}^* = H^{-1}(\Omega) := H_0^1(\Omega)^*$ and that there exist constants $0 < \alpha_1 \leq \alpha_2$ with

$$\forall x \in \Omega, \xi \in \mathbb{R}^d \quad \alpha_1 |\xi|^2 \leq \mathbf{A}(x)\xi \cdot \xi \text{ and } |\mathbf{A}(x)\xi| \leq \alpha_2 |\xi|. \quad (25)$$

The latter implies that the operator $-\operatorname{div}(\mathbf{A}\nabla \cdot)$ is uniformly elliptic. Moreover, the bilinear form \mathcal{B} is coercive and continuous with constants α_1 and α_2 , respectively.

Lax-Milgram Theorem and Poincaré-Friedrichs Inequality

$$\|v\|_{\Omega} \leq \text{diam}(\Omega) \|\nabla v\|_{\Omega} \quad \text{for all } v \in \mathbb{V} = H_0^1(\Omega) \quad (26)$$

thus ensure existence and uniqueness of the weak solution (24).

Note that \mathbf{A} is not assumed to be symmetric and so the bilinear form \mathcal{B} may be nonsymmetric. For the a posteriori upper bound, we will require some additional regularity on the data f and \mathbf{A} in §2.4.

2.2 Galerkin Solutions

Since \mathbb{V} has infinite dimension, problem (24) cannot be implemented on a computer and solved numerically. Given a subspace $\mathbb{S} \subset \mathbb{V}$, the corresponding Galerkin solution or approximation of (24) is given by

$$U \in \mathbb{S}: \quad \mathcal{B}[U, V] = \langle f, V \rangle \quad \text{for all } V \in \mathbb{S}. \quad (27)$$

We simply replaced each occurrence of \mathbb{V} in (24) by \mathbb{S} . If \mathbb{S} is finite-dimensional, we can choose a basis of \mathbb{S} and the coefficients of the expansion of U can be determined by solving a square linear system.

Residual. Associate the functional $\mathcal{R} \in \mathbb{V}^*$ given by

$$\langle \mathcal{R}, v \rangle := \langle f, v \rangle - \mathcal{B}[U, v],$$

to $U \in \mathbb{S}$. The functional \mathcal{R} is called the residual and depends only on the approximate solution U and data \mathbf{A} and f . Moreover, it has the following properties:

- It relates to the typically unknown error function $u - U$ in the following manner:

$$\langle \mathcal{R}, v \rangle = \mathcal{B}[u - U, v] \quad \text{for all } v \in \mathbb{V}. \quad (28)$$

This is a direct consequence of the characterization (24) of the exact solution.

- It vanishes for discrete test functions, which in the case of symmetric \mathbf{A} corresponds to the so-called Galerkin orthogonality:

$$\mathcal{B}[u - U, V] = \langle \mathcal{R}, V \rangle = 0 \quad \text{for all } V \in \mathbb{S}. \quad (29)$$

This immediately follows from (28) and the definition (27) of the Galerkin solution.

Quasi-Best Approximation. Property (25) of \mathbf{A} , Galerkin orthogonality (29) and Cauchy-Schwarz Inequality in $L^2(\Omega)$ imply

$$\begin{aligned} \alpha_1 \|u - U\|_{\mathbb{V}}^2 &\leq \mathcal{B}[u - U, u - U] = \mathcal{B}[u - U, u - V] \\ &\leq \alpha_2 \|u - U\|_{\mathbb{V}} \|u - V\|_{\mathbb{V}} \end{aligned}$$

for arbitrary $V \in \mathbb{S}$. This proves the famous

Theorem 4 (Céa Lemma). *The Galerkin solution is a quasi-best approximation from \mathbb{S} with respect to the \mathbb{V} -norm:*

$$\|u - U\|_{\mathbb{V}} \leq \frac{\alpha_2}{\alpha_1} \inf_{V \in \mathbb{S}} \|u - V\|_{\mathbb{V}}. \quad (30)$$

If the bilinear \mathcal{B} is also symmetric and one considers the error with respect to the so-called energy norm $\mathcal{B}[\cdot, \cdot]^{1/2}$, the Galerkin solution is even the best approximation from \mathbb{S} ; see Problem 12.

2.3 Finite Element Solutions and A Priori Bound

Problem (27) can be solved numerically on a computer, if we dispose of an implementable basis of \mathbb{S} . As an example of such space, let \mathcal{T} be a conforming triangulation of Ω into d -simplices (this imposes further conditions on Ω) and consider

$$\mathbb{S} = \mathbb{V}(\mathcal{T}) := \{V \in \mathbb{S}^{n,0}(\mathcal{T}) \mid V|_{\partial\Omega} = 0\}, \quad (31)$$

where, as in §1.4, $\mathbb{S}^{n,0}(\mathcal{T})$ the space of continuous functions that are piecewise polynomial up to degree n . This is in fact a subspace of $\mathbb{V} = H_0^1(\Omega)$ thanks to the continuity requirement and boundary condition for the functions in $\mathbb{V}(\mathcal{T})$. Moreover, the basis $\{\phi_z\}_{z \in \mathcal{N} \cap \Omega}$ from §1.4 can be easily constructed in the computer; see for example Siebert [50].

The space $\mathbb{V}(\mathcal{T})$ is a popular choice for \mathbb{S} in (27) and their combination may be viewed as a model finite element discretization.

In §1.5 we studied the approximation properties of $\mathbb{S}^{n,0}(\mathcal{T})$ with the help of (quasi-)interpolation operators $I_{\mathcal{T}}$. Since the right-hand side of (30) is bounded in terms of $\|u - I_{\mathcal{T}}u\|_{\mathbb{V}} = \|\nabla(u - I_{\mathcal{T}}u)\|_{L^2(\Omega)}$ with $I_{\mathcal{T}}$ as in Proposition 2, the discussion of §1.5 applies also to the error of the Galerkin solution $U_{\mathcal{T}}$ in $\mathbb{V}(\mathcal{T})$. In particular, the combination of the Céa Lemma and Theorem 2 yields the following upper bound. Since it does not involve the discrete solution, it is also comes with the adjective ‘a priori’.

Theorem 5 (A priori upper bound). *Assume that the exact solution u of (24) satisfies $u \in W^{s,p}(\Omega)$ with $1 \leq s \leq n+1$, $1 \leq p \leq 2$, and set*

$$r := \text{sob}(W_p^s(\Omega)) - \text{sob}(H^1(\Omega)) > 0.$$

Then the error of the finite element solution $U_{\mathcal{T}} \in \mathbb{S} = \mathbb{V}(\mathcal{T})$ of (27) satisfies the global a priori upper bound

$$\|u - U_{\mathcal{T}}\|_{\mathbb{V}} \lesssim \|h^r D^s u\|_{L^p(\Omega)}. \quad (32)$$

The discussion in §1.6 about adaptively graded meshes only partially carries over to the error of the finite element solution $U_{\mathcal{T}}$, from now on denoted U . In view of the Céa Lemma, §1.6 shows that there are sequences of meshes such that the error of U decays as $\#\mathcal{T}^{-1/2}$ if, for example, $d = 2$ and $u \in W^{2,p}(\Omega)$ with $p > 1$. Notice however that the thresholding algorithm utilizes the local errors $E_T = \|\nabla(u - I_{\mathcal{T}}u)\|_{L^2(T)}$, which are typically not computable. The construction of appropriate meshes when the target function is given only implicitly by a boundary value problem is much more subtle. A first step towards this goal is developed in the next section.

2.4 A Posteriori Upper Bound

The *a priori* upper bound (32) is not computable and essentially provides only asymptotic information, namely the asymptotic convergence rate. The goal of this section is to derive an alternative bound, so-called *a posteriori* bound, that provides information beyond asymptotics and is computable in terms of data and the approximate solution. It is worth noting that such bounds are useful not only for adaptivity but also for the quality assessment of the approximate solution.

Since in this section the grid \mathcal{T} is (arbitrary but) fixed, we simplify the notation by suppressing the subscript indicating the dependence on the grid in case of the approximate solution and similar quantities.

Error and Residual. Our starting point is the algebraic relationship (28) between the residual \mathcal{R} and the error function $u - U$. It implies (Problem 13)

$$\|u - U\|_{\mathbb{V}} \leq \frac{1}{\alpha_1} \|\mathcal{R}\|_{\mathbb{V}^*} \leq \frac{\alpha_2}{\alpha_1} \|u - U\|_{\mathbb{V}}, \quad (33)$$

which means that the dual norm

$$\|\ell\|_{\mathbb{V}^*} := \sup \{ \langle \ell, v \rangle \mid v \in \mathbb{V}, \|v\|_{\mathbb{V}} \leq 1 \} \quad (34)$$

is a good measure for the residual \mathcal{R} if we are interested in the error $\|u - U\|_{\mathbb{V}} = \|\nabla(u - U)\|_{L^2(\Omega)}$. However the evaluation of $\|\mathcal{R}\|_{\mathbb{V}^*} = \|\mathcal{R}\|_{H^{-1}(\Omega)}$ is impractical and, moreover, does not provide local information for guiding an adaptive mesh refinement. We therefore aim at a sharp upper bound of $\|\mathcal{R}\|_{H^{-1}(\Omega)}$ in terms of locally computable quantities.

Assumptions and Structure of Residual. For the derivation of a computable upper bound of the dual norm of the residual, we require that

$$f \in L^2(\Omega) \quad \text{and} \quad \mathbf{A} \in W^{1,\infty}(\Omega; \mathcal{T}) \quad (35)$$

where the latter means that \mathbf{A} is Lipschitz in each element of \mathcal{T} . Under these assumptions, we can write $\langle \mathcal{R}, v \rangle$ as integrals over elements $T \in \mathcal{T}$ and elementwise integration by parts yields the representation:

$$\begin{aligned}
\langle \mathcal{R}, v \rangle &= \int_{\Omega} f v - \mathbf{A} \nabla U \cdot \nabla v = \sum_{T \in \mathcal{T}} \int_T f v - \mathbf{A} \nabla U \cdot \nabla v \\
&= \sum_{T \in \mathcal{T}} \int_T r v + \sum_{S \in \mathcal{S}} \int_S j v,
\end{aligned} \tag{36}$$

where

$$\begin{aligned}
r &= f + \operatorname{div}(\mathbf{A} \nabla U) \quad \text{in any simplex } T \in \mathcal{T}, \\
j &= [[\mathbf{A} \nabla U]] \cdot \mathbf{n} = \mathbf{n}^+ \cdot \mathbf{A} \nabla U|_{T^+} + \mathbf{n}^- \cdot \mathbf{A} \nabla U|_{T^-} \quad \text{on any internal side } S \in \mathcal{S}
\end{aligned} \tag{37}$$

and \mathbf{n}^+ , \mathbf{n}^- are unit normals pointing towards T^+ , $T^- \in \mathcal{T}$. We see that the distribution \mathcal{R} consists of a regular part r , called *interior or element residual*, and a singular part j , called *jump or interelement residual*. The regular part is absolutely continuous w.r.t. the d -dimensional Lebesgue measure and is related to the strong form of the PDE. The singular part is supported on the skeleton $\Gamma = \bigcup_{S \in \mathcal{S}} S$ of \mathcal{T} and is absolutely continuous w.r.t. the $(d-1)$ -dimensional Hausdorff measure.

We point out that this structure of the residual is not special to the model problem and its discretization but rather arises from the weak formulation of the PDE and the piecewise construction of finite element spaces.

Scaled Integral Norms. In view of the structure of the residual \mathcal{R} , we make our goal precise as follows: we aim at a sharp upper bound for $\|\mathcal{R}\|_{H^{-1}(\Omega)}$ in terms of local Lebesgue norms of the element and interelement residuals r and j , which are considered to be computable because they can be easily approximated with numerical integration. This approach is usually called standard a posteriori error estimation.

The sharpness of these bounds crucially hinges on appropriate local scaling constants for the aforementioned Lebesgue norms, which depend on the local geometry of the mesh. For simplicity, we will explicitly trace only the dependence on the local mesh-size and write ' \lesssim ' instead of ' $\leq C$ ', where the constant C is bounded in terms of the shape coefficient $\sigma(\mathcal{T}) = \max_{T \in \mathcal{T}} \bar{h}_T / h_T$ of the triangulation \mathcal{T} and the dimension d .

Localization. As a first step, we decompose the residual \mathcal{R} into local contributions with the help of the Courant basis $\{\phi_z\}_{z \in \mathcal{V}}$ from §1.4. Hereafter \mathcal{V} stands for the set of vertices of \mathcal{T} , which coincide with the nodes of $\mathbb{S}^{1,0}(\mathcal{T})$. The Courant basis has the following properties:

- It provides a partition of unity:

$$\sum_{z \in \mathcal{V}} \phi_z = 1 \quad \text{in } \Omega. \tag{38}$$

- For each interior vertex z , the corresponding basis function ϕ_z is contained in $\mathbb{V}(\mathcal{T})$ and so the residual is orthogonal to the interior contributions of the partition of unity:

$$\langle \mathcal{R}, \phi_z \rangle = 0 \quad \text{for all } z \in \mathcal{V}^\circ := \mathcal{V} \cap \Omega. \tag{39}$$

The second property corresponds to the Galerkin orthogonality. Notice that the first property involves all vertices, while in the second one the boundary vertices are excluded.

Given any $v \in H_0^1(\Omega)$, we apply (38) and then (39) to write

$$\langle \mathcal{R}, v \rangle = \sum_{z \in \mathcal{V}} \langle \mathcal{R}, v \phi_z \rangle = \sum_{z \in \mathcal{V}} \langle \mathcal{R}, (v - c_z) \phi_z \rangle, \quad (40)$$

where $c_z \in \mathbb{R}$ and $c_z = 0$ whenever $z \in \partial\Omega$. Exploiting representation (36), $0 \leq \phi_z \leq 1$, and the fact that the ϕ_z are locally supported, we can bound each local contribution $\langle \mathcal{R}, (v - c_z) \phi_z \rangle$ in the following manner:

$$|\langle \mathcal{R}, (v - c_z) \phi_z \rangle| \leq \left| \int_{\omega_z} r(v - c_z) \phi_z \right| + \left| \int_{\gamma_z} j(v - c_z) \phi_z \right|, \quad (41)$$

where $\omega_z := \cup_{T \ni z} T$ is the star (or patch) around a vertex $z \in \mathcal{V}$ in \mathcal{T} and γ_z is the skeleton of ω_z , i.e. the union of all sides emanating from z ; note that r in (41) is computed elementwise. We examine the two terms on the right-hand side separately.

Bounding the Element Residual. We first consider the terms associated with the element residual r . The key tool for a sharp bound is the following local Poincaré-type inequality. Let

$$h_z := |\omega_z|^{1/d}$$

and notice that this quantity is, up to the shape coefficient $\sigma(\mathcal{T})$, equivalent to the diameter of ω_z , to $h_T = |T|^{1/d}$ if T is a d -simplex of ω_z and to $h_S := |S|^{1/(d-1)}$ if S is a side of γ_z .

Lemma 4 (Local Poincaré-type inequality). *For any $v \in H_0^1(\Omega)$ and $z \in \mathcal{V}$ there exists $c_z \in \mathbb{R}$ such that*

$$\|v - c_z\|_{L^2(\omega_z)} \lesssim h_z \|\nabla v\|_{L^2(\omega_z)}. \quad (42)$$

If $z \in \partial\Omega$ is a boundary vertex, then we can take $c_z = 0$.

We postpone the proof of Lemma 4. Combining the Cauchy-Schwarz inequality in $L^2(\omega_z)$ and Lemma 4 readily yields

$$\left| \int_{\omega_z} r(v - c_z) \phi_z \right| \leq \|r \phi_z^{1/2}\|_{L^2(\omega_z)} \|v - c_z\|_{L^2(\omega_z)} \lesssim h_z \|r \phi_z^{1/2}\|_{L^2(\omega_z)} \|\nabla v\|_{L^2(\omega_z)}. \quad (43)$$

Notice that the right-hand side consists of two factors: a computable one in the desired form and one that involves the test function in a local variant of the norm of the test space.

Bounding the Jump Residual. Next, we consider the terms associated to the jump residual j . Recall that j is supported on sides and so proceeding similarly as for the element residual will bring up traces of the test function. The following trace inequality exactly meets our needs.

Lemma 5 (Scaled trace inequality). *For any side S of a d -simplex T the following inequality holds:*

$$\|w\|_{L^2(S)} \lesssim h_S^{-1/2} \|w\|_{L^2(T)} + h_S^{1/2} \|\nabla w\|_{L^2(T)} \quad \text{for all } w \in H^1(T). \quad (44)$$

We again postpone the proof, now of Lemma 5. We apply first the Cauchy-Schwarz inequality in $L^2(\gamma_z)$, then Lemma 5 and finally Lemma 4 to obtain

$$\left| \int_{\gamma_z} j(v - c_z) \phi_z \right| \leq \|j \phi_z^{1/2}\|_{L^2(\gamma_z)} \|v - c_z\|_{L^2(\gamma_z)} \lesssim h_z^{1/2} \|j \phi_z^{1/2}\|_{L^2(\gamma_z)} \|\nabla v\|_{L^2(\omega_z)}, \quad (45)$$

where the right-hand side has the same structure as that of the element residual.

Upper Bound for Residual Norm. We collect the local estimates and sum them up in order to arrive at the desired bound for the dual norm of the residual. Inserting the estimates (43) and (45) for element and jump residuals into (41) gives

$$|\langle \mathcal{R}, v \phi_z \rangle| \lesssim \left(h_z \|r \phi_z^{1/2}\|_{L^2(\omega_z)} + h_z^{1/2} \|j \phi_z^{1/2}\|_{L^2(\gamma_z)} \right) \|\nabla v\|_{L^2(\omega_z)}.$$

Recalling the decomposition (40), we sum over $z \in \mathcal{V}$ and use Cauchy-Schwarz in $\mathbb{R}^{\#\mathcal{T}}$ to arrive at

$$|\langle \mathcal{R}, v \rangle| \lesssim \left(\sum_{z \in \mathcal{V}} h_z^2 \|r \phi_z^{1/2}\|_{L^2(\omega_z)}^2 + h_z \|j \phi_z^{1/2}\|_{L^2(\gamma_z)}^2 \right)^{1/2} \left(\sum_{z \in \mathcal{V}} \|\nabla v\|_{L^2(\omega_z)}^2 \right)^{1/2}.$$

For bounding the second factor, we resort to the finite overlapping property of stars, namely

$$\sum_{z \in \mathcal{V}} \chi_{\omega_z}(x) \leq d + 1,$$

and infer that

$$\sum_{z \in \mathcal{V}} \|\nabla v\|_{L^2(\omega_z)}^2 \lesssim \|\nabla v\|_{L^2(\Omega)}^2.$$

Since mesh refinement is typically based upon element subdivision, we regroup the terms within the first factor. To this end, denote by $h: \Omega \rightarrow \mathbb{R}^+$ the mesh-size function given by $h(x) := |S|^{1/k}$ if x belongs to the interior of the k -subsimplex S of \mathcal{T} with $k \in \{1, \dots, d\}$. Then for all $x \in \omega_z$ we have $h_z \lesssim h(x)$. Therefore employing (38) once more and recalling that Γ is the union of all interior sides of \mathcal{T} , we deduce

$$\begin{aligned} \sum_{z \in \mathcal{V}} h_z^2 \|r \phi_z^{1/2}\|_{L^2(\omega_z)}^2 + h_z \|j \phi_z^{1/2}\|_{L^2(\gamma_z)}^2 &\lesssim \sum_{z \in \mathcal{V}} \|h r \phi_z^{1/2}\|_{L^2(\omega_z)}^2 + \|h^{1/2} j \phi_z^{1/2}\|_{L^2(\Gamma)}^2 \\ &= \|h r\|_{L^2(\Omega)}^2 + \|h^{1/2} j\|_{L^2(\Gamma)}^2. \end{aligned}$$

Thus, introducing the *element indicators*

$$\mathcal{E}_{\mathcal{T}}^2(U, T) := h_T^2 \|r\|_{L^2(T)}^2 + h_T \|j\|_{L^2(\partial T \setminus \partial \Omega)}^2 \quad (46)$$

and the *error estimator*

$$\mathcal{E}_{\mathcal{T}}^2(U) = \sum_{T \in \mathcal{T}} \mathcal{E}_{\mathcal{T}}^2(U, T) \quad (47)$$

we arrive at the following upper bound for the dual norm of the residual:

$$\|\mathcal{R}\|_{H^{-1}(\Omega)} \lesssim \mathcal{E}_{\mathcal{T}}(U). \quad (48)$$

Hereafter, we write $\mathcal{E}_{\mathcal{T}}(U, \mathcal{M})$ to indicate that the estimator is computed over $\mathcal{M} \subset \mathcal{T}$, whereas $\mathcal{E}_{\mathcal{T}}(U, \mathcal{T}) = \mathcal{E}_{\mathcal{T}}(U)$ if no confusion arises.

Proofs of Poincaré-Type and Trace Inequalities. We now prove Lemmas 4 and 5. We start with a formula for the mean value of a trace, which follows from the Divergence Theorem.

Lemma 6 (Trace identity). *Let T be a d -simplex, S a side of T , and z the vertex opposite to S . Defining the vector field \mathbf{q}_S by*

$$\mathbf{q}_S(x) := x - z$$

the following equality holds

$$\frac{1}{|S|} \int_S w = \frac{1}{|T|} \int_T w + \frac{1}{d|T|} \int_T \mathbf{q}_S \cdot \nabla w \quad \text{for all } w \in W_1^1(T).$$

Proof. We start with properties of the vector field \mathbf{q}_S . Let S' be an arbitrary side of T and fix some $y \in S'$. We then see $\mathbf{q}_S(x) \cdot \mathbf{n}_T = \mathbf{q}_S(y) \cdot \mathbf{n}_T + (x - y) \cdot \mathbf{n}_T = \mathbf{q}_S(y) \cdot \mathbf{n}_T$ for any $x \in S'$ since $x - y$ is a tangent vector to S' . Therefore, on each side of T , the associated normal flux $\mathbf{q}_S \cdot \mathbf{n}_T$ is constant. In particular, we see $\mathbf{q}_S \cdot \mathbf{n}_T$ vanishes on $\partial T \setminus S$ by choosing $y = z$ for sides emanating from z . Moreover, $\operatorname{div} \mathbf{q}_S = d$. Thus, if $w \in C^1(\bar{T})$, the Divergence Theorem yields

$$\int_T \mathbf{q}_S \cdot \nabla w = -d \int_T w + (\mathbf{q}_S \cdot \mathbf{n}_T)|_S \int_S w.$$

Take $w = 1$ to show $(\mathbf{q}_S \cdot \mathbf{n}_T)|_S = d|T|/|S|$ and extend the result to $w \in W_1^1(T)$ by density. \square

Proof of Lemma 5. Apply Lemma 6 to $|w|^2$; for the details see Problem 17. \square

Proof of Lemma 4. \square For any $z \in \mathcal{V}$ the value

$$\bar{c}_z = \frac{1}{|\omega_z|} \int_{\omega_z} v$$

is an optimal choice and (42) follows from (8) with $c_z = \bar{c}_z$.

\square If $z \in \partial\Omega$, then we observe that there exists a side $S \subset \partial\omega_z \cap \partial\Omega$ such that $v = 0$ on S . We therefore can write

$$v = v - \frac{1}{|S|} \int_S v = (v - \bar{c}_z) - \frac{1}{|S|} \int_S (v - \bar{c}_z)$$

whence, using Lemma 5 and Step 1 for the second term,

$$\|v\|_{L^2(\omega_z)} \lesssim \|v - \bar{c}_z\|_{L^2(\omega_z)} + h_z \|\nabla v\|_{L^2(\omega_z)} \lesssim h_z \|\nabla v\|_{L^2(\omega_z)},$$

which establishes the supplement for boundary vertices. \square

Upper Bound for Error. Inserting the bound (48) for the dual norm of the residual in the first bound of (33), we obtain the main result of this section.

Theorem 6 (A posteriori upper bound). *Let u be the exact solution of the model problem (24) satisfying (25) and (35). The error of the finite element solution $U \in \mathbb{S} = \mathbb{V}(\mathcal{T})$ of (27) is bounded in terms of the estimator (47) as follows:*

$$\|u - U\|_{\mathbb{V}} \lesssim \frac{1}{\alpha_1} \mathcal{E}_{\mathcal{T}}(U), \quad (49)$$

where the hidden constant depends only on the shape coefficient $\sigma(\mathcal{T})$ of the triangulation \mathcal{T} and on the dimension d .

Notice that the a posteriori bound in Theorem 6 does not require additional regularity on the exact solution u as the a priori one in Theorem 5. On the other hand, the dependence of the estimator on the approximate solution prevents us from directly extracting information such as asymptotic decay rate of the error. The question thus arises how sharp the a posteriori bound in Theorem 6 is.

In this context it is worth noticing that if we did not exploit orthogonality and used a global Poincaré-type inequality instead of the local ones, the resulting scalings of the element and jump residuals would be, respectively, 1 and $h_T^{-1/2}$ and the corresponding upper bound would have a lower asymptotic decay rate. We will show in the next §3 that the upper bound in Theorem 6 is sharp in an asymptotic sense.

2.5 Notes

The discussion of the quasi-best approximation and the a priori upper bound of the error of the finite element solution are classical; see Braess [10], Brenner-Scott [11], and Ciarlet [19]. The core of the a posteriori upper bound is a bound of the dual norm of the residual in terms of scaled Lebesgue norms. This approach is usually called *standard a posteriori error estimation* and has been successfully used for a variety of problems and discretizations. For alternative approaches we refer to the monographs of Ainsworth and Oden [2] and Verfürth [58] on a posteriori error estimation.

Typically standard a posteriori error estimation is carried out with the help of error estimates for quasi-interpolation as in §1.5, which in turn rely on local Bramble-Hilbert lemmas. The above presentation invokes only the special case of Poincaré-type inequalities. It is a simplified version of derivation in Veiser and Verfürth [56], which has been influenced by Babuška and Rheinboldt [5], Carstensen and Funken

[12], and Morin, Nochetto and Siebert [42], and provides in particular constants that are explicit in terms of local Poincaré constants. It is worth mentioning that the ensuing constants are found in [56] for sample meshes and have values close to 1.

The setting and assumption of the model problem and discretization in this section avoids the following complications: numerical integration, approximation of boundary values, approximation of the domain, and inexact solution of the discrete system. While all these issues have been analyzed in an a priori context, only some of them have been considered in a posteriori error estimation; see Ainsworth and Kelly [1], Dörfler and Rumpf [25], Morin, Nochetto, and Siebert [42], Nochetto, Siebert, Schmidt and Veeseer [46], and Sacchi and Veeseer [47].

2.6 Problems

Problem 12 (Best approximation for symmetric problems). Consider the model problem (24), assume in addition to (25) that A is symmetric and denote the energy norm associated with the differential operator $-\operatorname{div}(A\nabla\cdot)$ by

$$\|v\|_{\Omega} := \left(\int_{\Omega} A\nabla v \cdot \nabla v \right)^{1/2}.$$

Prove that the Galerkin solution is the best approximation from $\mathbb{S} = \mathbb{V}(\mathcal{T})$ with respect to the energy norm:

$$\|u - U\|_{\Omega} = \min_{V \in \mathbb{S}} \|u - V\|_{\Omega}. \quad (50)$$

Derive from this that in this case (30) improves to

$$\|u - U\|_{\mathbb{V}} \leq \sqrt{\frac{\alpha_2}{\alpha_1}} \inf_{V \in \mathbb{S}} \|u - V\|_{\mathbb{V}}.$$

Problem 13 (Equivalence of error and residual norm). Prove the equivalence (33) between error and dual norm of the residual. Consider the model problem also with a symmetric A and derive a similar relationship for the energy norm error.

Problem 14 (Dominance of jump residual). Considering the model problem (24) and its discretization (27) with (31) and $n = 1$, show that, up to higher order terms, the jump residual

$$\eta_{\mathcal{T}}(U) = \left(\sum_{S \in \mathcal{T}} \|h^{1/2} j\|_{L^2(S)}^2 \right)^{1/2}$$

bounds $\|\mathcal{R}\|_{H^{-1}(\Omega)}$, which entails that the estimator $\mathcal{E}_{\mathcal{T}}(U)$ is dominated by $\eta_{\mathcal{T}}(U)$. To this end, revise the proof of the upper bound for $\|\mathcal{R}\|_{H^{-1}(\Omega)}$, use

$$c_z = \frac{1}{\int_{\omega_z} \phi_z} \int_{\omega_z} f \phi_z.$$

and rewrite $\int_{\omega_z} f(v - c_z)\phi_z$ by exploiting this weighted L^2 -orthogonality.

Problem 15 (A posteriori upper bound with quasi-interpolation). Consider the model problem (24) and its discretization (27) with space $\mathbb{S} = \mathbb{V}(\mathcal{T})$, and derive the upper a posteriori error bound without using the discrete partition of unity. To this end, use (36) and combine the scaled trace inequality (44) with the local interpolation error estimate (7). Show as an intermediate step the upper bound

$$|\langle \mathcal{R}, v \rangle| \lesssim \sum_{T \in \mathcal{T}} \mathcal{E}_{\mathcal{T}}(U, T) \|\nabla v\|_{L^2(N_{\mathcal{T}}(T))} \quad (51)$$

with $N_{\mathcal{T}}(T)$ from §1.5. This bound will be useful in §4.

Problem 16 (Upper bound for certain singular loads). Revising the proof of Theorem 6, derive an a posteriori upper bound in the case of right-hand sides of the form

$$\langle f, v \rangle = \int_{\Omega} g_0 v + \int_{\Gamma} g_1 v, \quad v \in \mathbb{V} = H_0^1(\Omega),$$

where $g_0 \in L^2(\Omega)$, $g_1 \in L^2(\Gamma)$, and Γ stands for the skeleton of the mesh \mathcal{T} .

Problem 17 (Scaled trace inequality). Work out the details of the proof of Lemma 5, taking into account that $h_T \approx |T|/|S| \approx h_S$.

Problem 18 (A posteriori upper bound for L^2 -error). Assuming that Ω is convex and applying a duality argument, establish a variant of (33) between the L^2 -error $\|u - U\|_{L^2(\Omega)}$ and a suitable dual norm of the residual. Use this to derive the a posteriori upper bound

$$\|u - U\|_{L^2(\Omega)} \lesssim \left(\sum_{T \in \mathcal{T}} h_T^2 \mathcal{E}_{\mathcal{T}}(U, T)^2 \right)^{1/2},$$

where the hidden constant depends in addition on the domain Ω .

3 Lower A Posteriori Bounds

The goal of this section is to assess the sharpness of the a posteriori upper bound for the model problem and discretization. We show not only that it is sharp in an asymptotic sense like the a priori bound but also in a local sense and, for certain data, in a non-asymptotic sense. Moreover, we verify that the latter cannot be true for all data and argue that this is the price to pay for the upper bound to be computable.

As in §2.4, ‘ \lesssim ’ stands for ‘ $\leq C$ ’, where the constant C is bounded in terms of the shape coefficient $\sigma(\mathcal{T})$ of the triangulation \mathcal{T} and the dimension d and, often, we do not indicate the dependence on the arbitrary but fixed triangulation.

3.1 Local Lower Bounds

The first step in the derivation of the upper bound (49) is that the error is bounded in terms of an appropriate dual norm of the residual. In the case of the model problem (24) this relies on the continuity of $[-\operatorname{div}(\mathbf{A}\nabla\cdot)]^{-1} : H^{-1}(\Omega) \rightarrow H_0^1(\Omega)$. Notice that the inverse is a global operator, while $-\operatorname{div}(\mathbf{A}\nabla\cdot)$ in the classical sense is a local one. One thus may suspect that an appropriate local dual norm of the residual is bounded in terms of the local error. Let us verify this for the model problem (24).

Local Dual Norms. Let ω be a subdomain of Ω and notice that $H^{-1}(\omega)$ is a good candidate for the local counterpart of $H^{-1}(\Omega)$. Given $v \in H_0^1(\omega)$ (and extending it by zero on $\Omega \setminus \omega$), the algebraic relationship (28), the Cauchy-Schwarz inequality in $L^2(\omega)$, and (25) readily yield

$$\langle \mathcal{R}, v \rangle = \mathcal{B}[u - U, v] = \int_{\omega} \mathbf{A}\nabla(u - U) \cdot \nabla v \leq \alpha_2 \|\nabla(u - U)\|_{L^2(\omega)} \|\nabla v\|_{L^2(\omega)}.$$

Consequently,

$$\|\mathcal{R}\|_{H^{-1}(\omega)} \leq \alpha_2 \|\nabla(u - U)\|_{L^2(\omega)}, \quad (52)$$

entailing that lower bounds for the local error $\|\nabla(u - U)\|_{L^2(\omega)}$ may be shown by bounding the local dual norm $\|\mathcal{R}\|_{H^{-1}(\omega)}$ from below.

Local Dual and Scaled Integral Norms. As for the a posteriori upper bound, we assume (35). If we take $\omega = T \in \mathcal{T}$ in the preceding paragraph, then there holds

$$\|\mathcal{R}\|_{H^{-1}(T)} = \sup_{\|\nabla v\|_{L^2(T)} \leq 1} \langle \mathcal{R}, v \rangle = \sup_{\|\nabla v\|_{L^2(T)} \leq 1} \int_T r v = \|r\|_{H^{-1}(T)} \quad (53)$$

thanks to the representation (36). Recall that the corresponding indicator $\mathcal{E}_{\mathcal{T}}(U, T)$ contains the term $h_T \|r\|_{L^2(T)}$ and therefore we wonder about the relationship of $\|r\|_{H^{-1}(T)}$ and $h_T \|r\|_{L^2(T)}$. Mimicking the local part in the derivation of the a posteriori upper bound in §2.4, we obtain

$$\int_T rv \leq \|r\|_{L^2(T)} \|v\|_{L^2(T)} \lesssim h_T \|r\|_{L^2(T)} \|\nabla v\|_{L^2(T)}$$

with the help of the Poincaré-Friedrichs inequality (26). Hence there holds

$$\|r\|_{H^{-1}(T)} \lesssim h_T \|r\|_{L^2(T)}. \quad (54)$$

Since $L^2(\Omega)$ is a proper subspace of $H^{-1}(\Omega)$, the inverse inequality cannot hold for arbitrary r . Consequently, $h_T \|r\|_{L^2(T)}$ may overestimate $\|r\|_{H^{-1}(T)}$. On the other hand, if $r \in \mathbb{R}$ is *constant* and $\eta = \eta_T$ denotes a non-negative function with properties

$$|T| \lesssim \int_T \eta, \quad \text{supp } \eta = T, \quad \|\nabla \eta\|_{L^\infty(T)} \lesssim h_T^{-1} \quad (55)$$

(postpone the question of existence until (59) below), we deduce

$$\begin{aligned} \|r\|_{L^2(T)}^2 &\lesssim \int_T r(r\eta) \leq \|r\|_{H^{-1}(T)} \|\nabla(r\eta)\|_{L^2(T)} \\ &\leq \|r\|_{H^{-1}(T)} \|r\|_{L^2(T)} \|\nabla \eta\|_{L^\infty(T)} \lesssim h_T^{-1} \|r\|_{H^{-1}(T)} \|r\|_{L^2(T)}, \end{aligned}$$

whence

$$h_T \|r\|_{L^2(T)} \lesssim \|r\|_{H^{-1}(T)}. \quad (56)$$

This shows that overestimation in (54) is caused by *oscillation* of r at a scale finer than the mesh-size. Notice that (56) is a so-called inverse estimate, where one norm is a dual norm. It is also valid for $r \in \mathbb{P}_n(T)$, but the constant deteriorates with the degree n ; see Problem 22.

Local Lower Bound with Element Residual. Motivated by the observations of the preceding paragraph, we expect that $h_T \|r\|_{L^2(T)}$ bounds asymptotically $\|\mathcal{R}\|_{H^{-1}(T)}$ from below and introduce the *oscillation of the interior residual* in T defined by

$$h_T \|r - \bar{r}_T\|_{L^2(T)},$$

where \bar{r}_T denotes the mean value of r in T . Replacing r by \bar{r}_T in (56) and by $r - \bar{r}_T$ in (54), as well as recalling (53), we derive

$$\begin{aligned} h_T \|r\|_{L^2(T)} &\leq h_T \|\bar{r}_T\|_{L^2(T)} + h_T \|r - \bar{r}_T\|_{L^2(T)} \\ &\lesssim \|\bar{r}_T\|_{H^{-1}(T)} + h_T \|r - \bar{r}_T\|_{L^2(T)} \\ &\lesssim \|r\|_{H^{-1}(T)} + \|r - \bar{r}_T\|_{H^{-1}(T)} + h_T \|r - \bar{r}_T\|_{L^2(T)} \\ &\lesssim \|\mathcal{R}\|_{H^{-1}(T)} + h_T \|r - \bar{r}_T\|_{L^2(T)}. \end{aligned} \quad (57)$$

This is the desired statement because the oscillation $h_T \|r - \bar{r}_T\|_{L^2(T)}$ is expected to converge faster than $h_T \|r\|_{L^2(T)}$ under refinement. In particular, if $n = 1$, then $r = f$ and the oscillation of the interior residual becomes *data oscillation*:

$$\text{osc}_{\mathcal{T}}(f, T) := \|h(f - \bar{f}_T)\|_{L^2(T)} \quad \text{for all } T \in \mathcal{T}. \quad (58)$$

Note that in this case there is one additional order of convergence if $f \in H^1(\Omega)$.

The inequality (57) holds also with \bar{r}_T chosen from $\mathbb{P}_{n_1}(T)$, with $n_1 \geq 1$, at the price of a larger constant hidden in \lesssim . We postpone the discussion of the higher order nature of the oscillation in this case after Theorem 7 below.

We conclude this paragraph by commenting on the choice of the cut-off function $\eta_T \in W_\infty^1(T)$ with (55). For example, we may take

$$\eta_T = (d+1)^{d+1} \prod_{z \in \mathcal{V} \cap T} \lambda_z, \quad (59)$$

where λ_z , $z \in \mathcal{V} \cap T$, are the barycentric coordinates of T from §1.4. This choice is due to Verfürth [57, 58]. Another choice, due to Dörfler [24], can be defined as follows: refine T such that there appears an interior node and take the corresponding Courant basis function on the virtual triangulation of T ; see Fig. 10 for the 2-dimensional case.

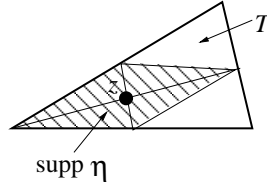


Fig. 10 Virtual refinement of a triangle for the Dörfler cut-off function.

The Dörfler cut-off function has the additional property that it is an element of a refined finite element space. This is not important here but useful when proving lower bounds for the differences of two discrete solutions; see e.g. Problem 23. Such estimates are therefore called *discrete lower bound* whereas the bound for the true error is called *continuous lower bound*.

Local Lower Bound with Jump Residual. We next strive for a local lower bound for the error in terms of the jump residual $h_S^{1/2} \|j\|_{L^2(S)}$, $S \in \mathcal{S}$, and use the local lower bound in terms of the element residual as guideline.

We first notice that $j = \llbracket \mathbf{A} \nabla U \rrbracket$ may not be constant on an interior side $S \in \mathcal{S}$ due to the presence of \mathbf{A} . We therefore introduce the *oscillation of the jump residual* in S ,

$$h_S^{1/2} \|j - \bar{j}_S\|_{L^2(S)},$$

where \bar{j}_S stands for the mean value of j on S , and write

$$h_S^{1/2} \|j\|_{L^2(S)} \leq h_S^{1/2} \|\bar{j}_S\|_{L^2(S)} + h_S^{1/2} \|j - \bar{j}_S\|_{L^2(S)}. \quad (60)$$

Notice that here the important question about the order of this oscillation is not obvious because, in contrast to the oscillation of the element residual in the case $n = 1$,

the approximate solution U is involved. We postpone the corresponding discussion until after Theorem 7 below.

To choose a counterpart of η_T , let ω_S denote the patch composed of the two elements of \mathcal{T} sharing S ; see Fig. 11 (left) for the 2-dimensional case. Obviously ω_S has a nonempty interior. Let $\eta_S \in W_\infty^1(\omega_S)$ be a cut-off function with the properties

$$|S| \lesssim \int_S \eta_S, \quad \text{supp } \eta_S = \omega_S, \quad \|\eta_S\|_{L^\infty(\omega_S)} = 1, \quad \|\nabla \eta_S\|_{L^\infty(\omega_S)} \lesssim h_S^{-1}. \quad (61)$$

Following Verfürth [57, 58] we may take η_S given by

$$\eta_S|_T = d^d \prod_{z \in \mathcal{V} \cap S} \lambda_z^T, \quad (62)$$

where $T \subset \omega_S$ and λ_z^T , $z \in \mathcal{V} \cap T$, are the barycentric coordinates of T . Also here



Fig. 11 Patch ω_S of triangles associated to interior side (left) and its refinement for Dörfler cut-off function (right).

Dörfler [24] proposed the following alternative: refine ω_S such that there appears an interior node of S and take the corresponding Courant basis function on the virtual triangulation of ω_S ; see Fig. 11 (right) for the 2-dimensional case.

After these preparations we are ready to derive a counterpart of (57). In view of the properties of η_S , we have

$$\|\bar{j}_S\|_{L^2(S)}^2 \lesssim \int_S \bar{j}_S (\bar{j}_S \eta_S) = \int_S j v_S + \int_S (\bar{j}_S - j) v_S \quad (63)$$

with $v_S = \bar{j}_S \eta_S$. We rewrite the first term on the right-hand side with the representation formula (36) as follows:

$$\int_S j v_S = - \int_{\omega_S} r v_S + \langle \mathcal{R}, v_S \rangle;$$

in contrast to (53), the jump residual couples with the element residual. Hence

$$\left| \int_{\omega_S} j v_S \right| \leq \|r\|_{L^2(\omega_S)} \|v_S\|_{L^2(\omega_S)} + \|\mathcal{R}\|_{H^{-1}(\omega_S)} \|\nabla v_S\|_{L^2(\omega_S)}.$$

In view of the Poincaré-Friedrichs inequality (26), $|\omega_S| \lesssim h_S |S|$ and (61), we have

$$\|v_S\|_{L^2(\omega_S)} \lesssim h_S \|\nabla v_S\|_{L^2(\omega_S)} \leq h_S \|\bar{j}_S\|_{L^2(\omega_S)} \|\nabla \eta_S\|_{L^\infty(\omega_S)} \lesssim h_S^{1/2} \|\bar{j}_S\|_{L^2(S)}.$$

We thus infer that

$$\left| \int_{\omega_S} j v_S \right| \lesssim \left(h_S^{1/2} \|r\|_{L^2(\omega_S)} + h_S^{-1/2} \|\mathcal{R}\|_{H^{-1}(\omega_S)} \right) \|\bar{j}_S\|_{L^2(S)}$$

and, using (44),

$$\left| \int_S (\bar{j}_S - j) v_S \right| \leq \|\bar{j}_S - j\|_{L^2(S)} \|v_S\|_{L^2(S)} \lesssim \|\bar{j}_S - j\|_{L^2(S)} \|\bar{j}_S\|_{L^2(S)}.$$

Inserting these estimates into (63) yields

$$\|\bar{j}_S\|_{L^2(S)}^2 \lesssim \left(h_S^{1/2} \|r\|_{L^2(\omega_S)} + h_S^{-1/2} \|\mathcal{R}\|_{H^{-1}(\omega_S)} + \|\bar{j}_S - j\|_{L^2(S)} \right) \|\bar{j}_S\|_{L^2(S)}$$

whence, recalling (60),

$$h_S^{1/2} \|j\|_{L^2(S)} \lesssim \|\mathcal{R}\|_{H^{-1}(\omega_S)} + h_S \|r\|_{L^2(\omega_S)} + h_S^{1/2} \|\bar{j}_S - j\|_{L^2(S)}. \quad (64)$$

This estimate also holds if $\bar{j}_S \in \mathbb{P}_{n_2}(S)$ is a polynomial of degree $\leq n_2$ (Problem 27).

Local Lower Bound with Indicator. We combine the two results on interior and jump residual and exploit also the local relationship between residual and error in order to obtain a local lower bound in terms of a single indicator.

To this end, we introduce the following notation for the oscillations. Recall the mesh-size function h from §1.5 and let

$$\bar{r} = P_{2n-2}r, \quad \bar{j} = P_{2n-1}j,$$

where $P_{2n-2}r|_T$ and $P_{2n-1}j|_S$ are the L^2 -orthogonal projections of r and j onto $\mathbb{P}_{2n-2}(T)$ and $\mathbb{P}_{2n-1}(S)$, respectively. The choice of the polynomial degrees arise from the desire that the oscillations are of higher order. Details are discussed after Theorem 7. Moreover, we associate with each simplex $T \in \mathcal{T}$ the patch

$$\omega_T := \bigcup_{S \subset \partial T \setminus \partial \Omega} \omega_S$$

(see Fig. 12 for the 2-dimensional case), and define the oscillation in ω_T by

$$\text{osc}_{\mathcal{T}}(U, \omega_T) = \|h(r - \bar{r})\|_{L^2(\omega_T)} + \|h^{1/2}(j - \bar{j})\|_{L^2(\partial T \setminus \partial \Omega)}. \quad (65)$$

In general, as indicated by the notation, the oscillation depends on the approximation U . However, in certain situations, it may be independent of the approximation U and then becomes *data* oscillation (58); see also Problem 19.

Theorem 7 (Local lower bound). *Let u be the exact solution of the model problem (24) satisfying (25) and (35). Each element indicator of (46) bounds, up to oscilla-*

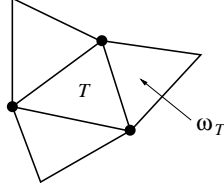


Fig. 12 Patch associated to a triangle in the local lower bound.

tion, the local error of an approximation $U \in \mathbb{V}(\mathcal{T})$ from below:

$$\mathcal{E}_{\mathcal{T}}(U, T) \lesssim \alpha_2 \|\nabla(u - U)\|_{L^2(\omega_T)} + \text{osc}_{\mathcal{T}}(U, \omega_T) \quad \text{for all } T \in \mathcal{T}, \quad (66)$$

where the hidden constant depends only on the shape coefficients of the simplices in ω_T , the dimension d and the polynomial degree n .

Proof. We combine (57) and (64), where \bar{r} and \bar{j} are piecewise polynomial of degree $2n - 2$ and $2n - 1$, respectively. Noting $h_S \approx h_T$ for all interior sides $S \in \mathcal{S}$ with $S \subset \partial T$ and $\|\mathcal{R}\|_{H^{-1}(T')} \leq \|\mathcal{R}\|_{H^{-1}(\omega_T)}$ for $T' \subset T$, we thus derive

$$\mathcal{E}_{\mathcal{T}}(U, T) \lesssim \|\mathcal{R}\|_{H^{-1}(\omega_T)} + \|h(r - \bar{r})\|_{L^2(\omega_T)} + \|h^{1/2}(\bar{j} - j)\|_{L^2(\partial T \setminus \partial \Omega)}.$$

Thus, the special case

$$\|\mathcal{R}\|_{H^{-1}(\omega_T)} \leq \alpha_2 \|\nabla(u - U)\|_{L^2(\omega_T)}$$

of (52) finishes the proof. \square

A discussion of the significance of local lower bound in Theorem 7 is in order. To this end, we first consider the decay properties of the oscillation terms, which are crucial for the relevance of the aforementioned bound. Then we remark about the importance of the fact that the bound in Theorem 7 is local. Finally, in the next section, we provide a global lower bound as corollary and discuss its relationship with the upper bound in Theorem 6.

Higher Order Nature of Oscillation. In some sense the oscillation pollutes the local lower bound in Theorem 7. It is therefore important that the oscillation is or gets small relative to the local error. We therefore compare the convergence order of the oscillation (65) with that of the local error.

To this end, let us first observe that the choices of the polynomial degrees in the oscillation allow us to derive the following upper bound of the oscillation (see Problem 29):

$$\begin{aligned} \text{osc}_{\mathcal{T}}(U, \omega_T) &\lesssim \|h(f - P_{2n-2}f)\|_{L^2(\omega_T)} \\ &\quad + \left(\|h(\text{div} \mathbf{A} - P_{n-1}(\text{div} \mathbf{A}))\|_{L^\infty(\omega_T)} + \|\mathbf{A} - P_n \mathbf{A}\|_{L^\infty(\omega_T)} \right) \|\nabla U\|_{L^2(\omega_T)}. \end{aligned} \quad (67)$$

If f and \mathbf{A} are smooth, one expects that the local error vanishes like

$$\|\nabla(u - U)\|_{L^2(T)} = \mathcal{O}(h_T^{d/2+n})$$

and, in view of (67), oscillation like

$$\text{osc}_{\mathcal{T}}(U, \omega_T) = \mathcal{O}(h_T^{d/2+n+1}).$$

See also Problem 30 for a stronger result for the jump residual.

The oscillation $\text{osc}_{\mathcal{T}}(U, \omega_T)$ is therefore expected to be of higher order as $h_T \downarrow 0$. However, as Problem 32 below illustrates, it may be relevant on relatively coarse triangulations \mathcal{T} .

Local Lower Bound and Marking. In contrast to the upper bound in Theorem 6, the lower bound in Theorem 7 is local. This is very welcome in a context of adaptivity. In fact, if $\text{osc}_{\mathcal{T}}(U, \omega_T) \ll \|\nabla(u - U)\|_{L^2(\omega_T)}$, as we expect asymptotically, then (66) translates into

$$\mathcal{E}_{\mathcal{T}}(U, T) \lesssim \alpha_2 \|\nabla(u - U)\|_{L^2(\omega_T)}. \quad (68)$$

This means that an element T with relatively large error indicator contains a large portion of the error. To improve the solution U effectively, such T must be split giving rise to a procedure that tries to equidistribute errors. This is consistent with the discussion of adaptive approximation of §1.1 for $d = 1$ and of §1.6 for $d > 1$.

3.2 Global Lower Bound

We derive a global lower bound from the local one in Theorem 7 and discuss its relationship with the global upper bound in Theorem 6.

The global counterpart of $\text{osc}_{\mathcal{T}}(U, \omega_T)$ from (65) is given by

$$\text{osc}_{\mathcal{T}}(U) = \|h(r - \bar{r})\|_{L^2(\Omega)} + \|h^{1/2}(j - \bar{j})\|_{L^2(\Gamma)}, \quad (69)$$

where r is computed elementwise over \mathcal{T} and Γ is the interior skeleton of \mathcal{T} .

Corollary 3 (Global lower bound). *Let u be the exact solution of the model problem (24) satisfying (25) and (35). The estimator (47) bounds, up to oscillation, the error of an approximation $U \in \mathbb{V}(\mathcal{T})$ from below:*

$$\mathcal{E}_{\mathcal{T}}(U) \lesssim \alpha_2 \|u - U\|_{\mathbb{V}} + \text{osc}_{\mathcal{T}}(U) \quad (70)$$

where the hidden constant depends on the shape coefficient of \mathcal{T} , the dimension d , and the polynomial degree n .

Proof. Sum (66) over all $T \in \mathcal{T}$ and take into account that each element is contained in at most by $d + 2$ patches ω_T . \square

Supposing that the approximation U is the Galerkin solution (27) with (31), the upper and lower a posteriori bounds in Theorem 6 and Corollary 3 imply

$$\|u - U\|_{\mathbb{V}} \lesssim \frac{1}{\alpha_1} \mathcal{E}_{\mathcal{T}}(U) \lesssim \frac{\alpha_2}{\alpha_1} \|u - U\|_{\mathbb{V}} + \frac{1}{\alpha_1} \text{osc}_{\mathcal{T}}(U). \quad (71)$$

In other words, the error and estimator are equivalent up to oscillation.

In Problem 32 we present an example for which the ratio $\|u - U\|_{\mathbb{V}}/\mathcal{E}_{\mathcal{T}}(U)$ can be made arbitrarily small. Consequently, a lower bound without pollution and a perfect equivalence of error and estimator cannot be true in general. Moreover, for that example there holds $\mathcal{E}_{\mathcal{T}}(U) = \text{osc}_{\mathcal{T}}(U)$, indicating that $\text{osc}_{\mathcal{T}}(U)$ is a good measure to account for the discrepancy.

We see that $\text{osc}_{\mathcal{T}}(U)$ intervenes in the relationship of error and estimator and, therefore, cannot be ignored in an analysis of an adaptive algorithm using the estimator $\mathcal{E}_{\mathcal{T}}(U)$; we will come back to this in §7. The case of data oscillation will be simpler than the general case in which $\text{osc}_{\mathcal{T}}(U)$ depends on the approximation U ; the latter dependence creates a nonlinear interaction in the adaptive algorithm.

The presence of oscillation is also consistent with our previous comparison of local dual norms and scaled integral norms. Since we invoked scaled integral norm in order to have an (almost) computable upper bound, this suggests that, at least for standard a posteriori error estimation, oscillation is a price that we have to pay for computability.

Fortunately, as we have illustrated in §3.1, oscillation is typically of higher order and then the a posteriori upper bound in Theorem 6 is asymptotically sharp in that its decay rate coincides with the one of the error, as the a priori bound of Theorem 5. Notice however the lower bound in Corollary 3 provides information beyond asymptotics: for example, if we consider the linear finite element method, that is $n = 1$, then $\text{osc}_{\mathcal{T}}(U)$ vanishes for all triangulations on which f and \mathbf{A} are piecewise constant and in this class of meshes error and estimator are thus equivalent:

$$\|\nabla(u - U)\|_{L^2(\Omega)} \approx \mathcal{E}_{\mathcal{T}}(U).$$

In summary: the estimator $\mathcal{E}_{\mathcal{T}}(U)$ from (47) is computable, it may be used to quantify the error and, in view of the local properties in §3.1, its indicators may be employed to provide the problem-specific information for local refinement.

3.3 Notes

Local lower bounds first appear in the work of Babuška and Miller [4]. Their derivation with continuous bubble functions is due to Verfürth [57], while the discrete lower bounds are due to Dörfler [24].

The discussion of the relationship between local dual norms and scaled integral norms as the reason for oscillation is an elaborated version of Sacchi-Veeser's one [47, Remark 3.1]. It is worth mentioning that there the indicators associated with

the approximation of the Dirichlet boundary values do not need to invoke scaled Lebesgue norms and are overestimation-free. Binev, Dahmen and DeVore [7] and Stevenson [52] arrange the a posteriori analysis such that oscillation is measured in $H^{-1}(\Omega)$. This avoids overestimation but brings back the question how to (approximately) evaluate the $H^{-1}(\Omega)$ -norm at acceptable cost. This question is open.

One may think that the issue of oscillation is specific to standard a posteriori error estimation. However all estimators we are aware of suffer from oscillations of the data that are finer than the mesh-size. For example, in the case of hierarchical estimators $\eta_{\mathcal{T}}(U)$ [2, 55, 58], as well as those based upon local discrete problems [2, 12, 42] or on gradient recovery [2, 27], the oscillation arises in the upper but not in the lower bounds and so creates a similar gap as that discussed here, namely

$$\eta_{\mathcal{T}}(U) \lesssim \|\nabla(u - U)\|_{L^2(\Omega)} \lesssim \eta_{\mathcal{T}}(U) + \text{osc}_{\mathcal{T}}(U). \quad (72)$$

3.4 Problems

Problem 19 (Data oscillation). Check that $\text{osc}_{\mathcal{T}}(U, \omega_T)$ in (65) does not depend on the approximation U if U is piecewise affine and \mathbf{A} is piecewise constant, and is given by

$$\text{osc}(f, \omega_T) = \|h(f - \bar{f})\|_{L^2(\omega_T)},$$

which corresponds to element data oscillation in (58).

Problem 20 (Energy norm case). Consider model problem (24) and discretization (27) with $\mathbb{S} = \mathbb{V}(\mathcal{T})$ and \mathbf{A} symmetric. Derive the counterparts of (66) and (71) for the energy norm and discuss the difference to the case presented here.

Problem 21 (Cut-off functions for simplices). Verify that a suitable multiple of the Verfürth cut-off function (59) satisfies the properties (55). To this end, exploit affine equivalence of T to a fixed reference simplex and shape regularity. Repeat for the corresponding Dörfler cut-off function.

Problem 22 (Inverse estimate for general polynomials). (Try this problem after Problem 21.) Show that the choice (59) for η_T verifies, for all $p \in \mathbb{P}_n(T)$,

$$\int_T p^2 \lesssim \int_T p^2 \eta_T, \quad \|\nabla(p\eta_T)\|_{L^2(T)} \lesssim h_T^{-1} \|p\|_{L^2(T)}$$

with constants depending on n and the shape coefficient of T . To this end, recall the equivalence of norms in finite-dimensional spaces. Derive the estimate

$$h_T \|r\|_{L^2(T)} \lesssim \|r\|_{H^{-1}(T)}$$

for $r \in \mathbb{P}_n(T)$.

Problem 23 (Lower bound for correction). Consider the model problem and its discretization for $d = 2$ and $n = 1$. Let U_1 be the solution over a triangulation \mathcal{T}_1

and U_2 the solution over \mathcal{T}_2 , where \mathcal{T}_2 has been obtained by applying at least 3 bisections to every triangle of \mathcal{T}_1 . Moreover, suppose that f is piecewise constant over \mathcal{T}_1 . Show

$$\|\nabla(U_2 - U_1)\|_{L^2(\Omega)} \geq \|h_1 f\|_{L^2(\Omega)},$$

where h_1 is the mesh-size function of \mathcal{T}_1 .

Problem 24 (Cut-off functions for sides). Verify that a suitable multiple of the Verfürth cut-off function (62) satisfies the properties (61). Repeat for the corresponding Dörfler cut-off function.

Problem 25 (Polynomial extension). Let S be a side of a simplex T . Show that for each $q \in \mathbb{P}_n(S)$ there exists a $p \in \mathbb{P}_n(T)$ such that

$$p = q \text{ on } S \quad \text{and} \quad \|p\|_{L^2(T)} \lesssim h_T^{1/2} \|q\|_{L^2(S)}.$$

Problem 26 (Norm equivalences with cut-off functions of sides). Let S be a side of a simplex T . Show that the choice (62) for η_S verifies, for all $q \in \mathbb{P}_n(S)$ and all $p \in \mathbb{P}_m(T)$,

$$\int_S q^2 \lesssim \int_S q^2 \eta_S, \quad \|\nabla(p\eta_S)\|_{L^2(T)} \lesssim h_T^{-1} \|p\|_{L^2(T)}$$

with constants depending on m, n , and the shape coefficient of T .

Problem 27 (Lower bound with jump residual and general oscillation). Exploit the claims in Problems 25 and 26, to rederive the estimate (64) but this time with \bar{r} and \bar{j} piecewise polynomials of degree $\leq n_1$ and n_2 .

Problem 28 (Best approximation of a product). Let K be either a d or a $(d-1)$ -simplex. For $\ell \in \mathbb{N}$ denote by $P_m^\ell: L^p(K, \mathbb{R}^\ell) \rightarrow \mathbb{P}_m(K, \mathbb{R}^\ell)$ the operator of best L^p -approximation in K . Prove that, for all $v \in L^\infty(K, \mathbb{R}^\ell)$, $V \in \mathbb{P}_n(K, \mathbb{R}^\ell)$ and $m \geq n$,

$$\|vV - P_m^2(vV)\|_{L^2(K)} \leq \|v - P_{m-n}^\infty v\|_{L^\infty(K)} \|V\|_{L^2(K)}.$$

Problem 29 (Upper bound for oscillation). Verify the upper bound (67) for the oscillation by exploiting Problem 28.

Problem 30 (Superconvergence of jump residual oscillation). Show that if \mathbf{A} is smooth across interelement boundaries, then the oscillation of the jump residual is superconvergent in that

$$\|j - \bar{j}_S\|_{L^2(S)} = \mathcal{O}(h_S^n) \|j\|_{L^2(S)} \quad \text{as } h_S \searrow 0.$$

Problem 31 (Simplified bound of oscillation). Using (67), show that (35) implies

$$\text{osc}_{\mathcal{T}}(U, \omega_T) \lesssim h_T \left(\|f\|_{L^2(\omega_T)} + \|\nabla U\|_{L^2(\omega_T)} \right), \quad (73)$$

where the hidden constant depends also on \mathbf{A} .

Problem 32 (Necessity of oscillation). Let $\varepsilon = 2^{-K}$ for K integer and extend the function $\frac{1}{2}x(\varepsilon - |x|)$ defined on $(-\varepsilon, \varepsilon)$ to a 2ε -periodic C^1 function u_ε on $\Omega = (-1, 1)$. Moreover, let the forcing function be $f_\varepsilon = -u''$, which is 2ε -periodic and piecewise constant with values ± 1 that change at multiples of ε ; see Fig. 13. Let \mathcal{T}_ε be a uniform mesh with mesh-size $h = 2^{-k}$, with $k \ll K$. We consider piece-

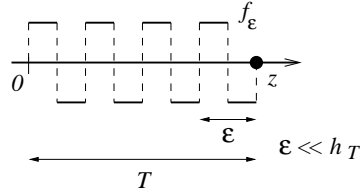


Fig. 13 An strongly oscillating forcing function.

wise linear finite elements $\mathbb{V}(\mathcal{T}_\varepsilon)$ and corresponding Galerkin solution $U_\varepsilon \in \mathbb{V}(\mathcal{T}_\varepsilon)$. Observing that f_ε is $L^2(\Omega)$ -orthogonal to both the space of piecewise constants and linears over \mathcal{T}_ε , show that

$$\begin{aligned} \|u'_\varepsilon - U'_\varepsilon\|_{L^2(\Omega)} &= \|u'_\varepsilon\|_{L^2(\Omega)} = \frac{\varepsilon}{\sqrt{6}} = \frac{2^{-K}}{\sqrt{6}} \\ &\ll 2^{-k} = h = \|hf_\varepsilon\|_{L^2(\Omega)} = \text{osc}_{\mathcal{T}_\varepsilon}(U_\varepsilon) = \mathcal{E}_{\mathcal{T}_\varepsilon}(U_\varepsilon). \end{aligned}$$

Extend this 1d example via a checkerboard pattern to any dimension.

4 Convergence of AFEM

The purpose of this section is to formulate an adaptive finite element method (AFEM) and to prove that it generates a sequence of approximate solutions converging to the exact one. The method consists in the following main steps:

SOLVE \rightarrow ESTIMATE \rightarrow MARK \rightarrow REFINES.

By their nature, adaptive algorithms define the sequence of approximate solutions as well as associated meshes and spaces only implicitly. This fact requires an approach that differs from ‘classical’ convergence proofs. In particular, a proof of convergence will hinge on results of an a posteriori analysis as in §§2.4 and 3.

The approach presented in this section covers wide classes of problems, discrete spaces, estimators and marking strategies. Here we do not strive for such generality but instead, in order to minimize technicalities, illustrate the main arguments only in a model case and then hint on possible generalizations.

It is worth noticing that, conceptually, the following convergence proof does not suppose any additional regularity of the exact solution. Consequently, it does not (and cannot) provide any information about the convergence speed. This important issue will be the concern of §7 for smaller classes of problems and algorithms.

4.1 A Model Adaptive Algorithm

We first present an AFEM for the model problem (24), which is an example of a *standard* iterative process that is often used in practice. In §4.2 we then prove its convergence and, finally, in §4.3 we comment on generalizations still covered by the given approach.

AFEM. The main structure of the adaptive finite element method is as follows: given an initial grid \mathcal{T}_0 , set $k = 0$ and iterate

1. $U_k = \text{SOLVE}(\mathcal{T}_k)$;
2. $\{\mathcal{E}_k(U_k, T)\}_{T \in \mathcal{T}_k} = \text{ESTIMATE}(U_k, \mathcal{T}_k)$;
3. $\mathcal{M}_k = \text{MARK}(\{\mathcal{E}_k(U_k, T)\}_{T \in \mathcal{T}_k}, \mathcal{T}_k)$;
4. $\mathcal{T}_{k+1} = \text{REFINE}(\mathcal{M}_k, \mathcal{T}_k)$; $k \leftarrow k + 1$.

Thus, the algorithm produces sequences $(\mathcal{T}_k)_{k=0}^\infty$ of meshes, $(U_k)_{k=0}^\infty$ of approximate solutions, and, implicitly, $(\mathbb{V}_k)_{k=0}^\infty$ of discrete spaces.

We next state our main assumptions and define the aforementioned modules for the problem at hand in detail.

Assumptions on continuous problem. We assume that the model problem (24) satisfies (25) and (35) so that the a posteriori error bounds of §§2 and 3 are available.

Initial grid. Assume that \mathcal{T}_0 is some initial triangulation of Ω such that \mathbf{A} is piecewise Lipschitz over \mathcal{T}_0 .

Solve. Let

$$\mathbb{V}_k := \mathbb{V}(\mathcal{T}_k) := \{V \in \mathbb{S}^{n,0}(\mathcal{T}_k) \mid V|_{\partial\Omega} = 0\}$$

be the space of continuous functions that are piecewise polynomial of degree $\leq n$ over \mathcal{T}_k , and compute the Galerkin solution U_k in \mathbb{V}_k given by

$$U_k \in \mathbb{V}_k : \quad \int_{\Omega} \mathbf{A} \nabla U_k \cdot \nabla V = \int_{\Omega} f V \quad \text{for all } V \in \mathbb{V}_k.$$

Estimate. Compute the error estimator $\{\mathcal{E}_k(U_k, T)\}_{T \in \mathcal{T}_k}$ given by

$$\mathcal{E}_k(U_k, T) := \left(h_T^2 \|r\|_{L^2(T)}^2 + h_T \|j\|_{L^2(\partial T \setminus \partial\Omega)}^2 \right)^{1/2}$$

where $h_T = |T|^{1/d}$, r and j are the element and jump residuals from (37) associated to the approximate solution U_k .

Mark. Collect a subset $\mathcal{M}_k \subset \mathcal{T}_k$ of marked elements with the following property:

$$\forall T \in \mathcal{T}_k \quad \mathcal{E}_k(U_k, T) = \mathcal{E}_{k,\max} > 0 \quad \implies \quad T \in \mathcal{M}_k \quad (74)$$

with $\mathcal{E}_{k,\max} := \max_{T \in \mathcal{T}_k} \mathcal{E}_k(U_k, T)$.

Refine. Refine \mathcal{T}_k into \mathcal{T}_{k+1} using bisection, as explained in §1.3, in such a way that each element in \mathcal{M}_k is bisected at least once and, finally, increment k .

Classical convergence proofs consider the case of uniform, or ‘non-adaptive’, refinement, which is included in the above class of algorithms by choosing $\mathcal{M}_k = \mathcal{T}_k$, thereby ignoring the information provided by the estimator. These convergence proofs rely on the fact that the maximum mesh-size decreases to 0 and therefore $\bigcup_{k=0}^{\infty} \mathbb{V}_k = H_0^1(\Omega)$. The above algorithm does not require this property, neither explicitly nor implicitly in general. In fact this property is not desirable in an adaptive context, since (30) reveals that it is sufficient to approximate only one function of $H_0^1(\Omega)$, namely the exact solution u of (24). In the next section we see that the above algorithm ensures just convergence to u by a subtle combination of properties of estimator and marking strategy.

4.2 Convergence

The goal of this section is to prove the convergence of the AFEM in §4.1. More precisely, we show that the sequence $(U_k)_{k=0}^{\infty}$ of approximate solutions converges to the exact solution u of the model problem (24).

Throughout this section ‘ \lesssim ’ stands for ‘ $\leq C$ ’, where the constant is independent of the iteration number k in the adaptive algorithm.

Convergence to Some Function. We expect the Galerkin solutions $(U_k)_{k=0}^\infty$ to approximate the exact solution u in $\mathbb{V} = H_0^1(\Omega)$. In any event, we may regard them as approximations to the Galerkin solution U_∞ in the limit

$$\mathbb{V}_\infty := \overline{\bigcup_{k=0}^{\infty} \mathbb{V}_k}$$

of the discrete spaces. Notice that \mathbb{V}_∞ is a subspace of \mathbb{V} , which may not coincide with \mathbb{V} (see below). In the next lemma we adopt this viewpoint and show that $(U_k)_{k=0}^\infty$ converges to U_∞ .

Lemma 7 (Limit of approximate solutions). *The finite element solutions $(U_k)_{k=0}^\infty$ converge in \mathbb{V} to the Galerkin solution $U_\infty \in \mathbb{V}_\infty$ given by*

$$\int_{\Omega} \mathbf{A} \nabla U_\infty \cdot \nabla V = \int_{\Omega} f V \quad \text{for all } V \in \mathbb{V}_\infty.$$

Proof. Since the sequence of $(\mathbb{V}_k)_{k=0}^\infty$ is nested (see Problem 33), the set \mathbb{V}_∞ is a closed linear subspace of \mathbb{V} . Hence \mathbb{V}_∞ is a Hilbert space and the bilinear form \mathcal{B} is coercive and continuous on \mathbb{V}_∞ . The Lax-Milgram Theorem therefore ensures existence and uniqueness of U_∞ .

Let $k \in \mathbb{N}_0$ and note that $\mathbb{V}_k \subset \mathbb{V}_\infty$. We can therefore replace \mathbb{V} by \mathbb{V}_∞ in Theorem 4 and obtain

$$\|U_\infty - U_k\|_{\mathbb{V}} \leq \inf_{V \in \mathbb{V}_k} \|U_\infty - V\|_{\mathbb{V}}.$$

Sending $k \rightarrow \infty$ then finishes the proof, because the right-hand side decreases to 0 by the very definition of \mathbb{V}_∞ . \square

Lemma 7 reduces our task to showing that $U_\infty = u$. Notice that this is equivalent to the condition $u \in \mathbb{V}_\infty$, illustrating that in general there is no need for $\mathbb{V}_\infty = \mathbb{V}$.

The identity $U_\infty = u$ hinges on the design of the adaptive algorithm. To illustrate this point, let us consider two extreme examples:

- It may happen that all indicators vanish in iteration k^* . Then $\mathcal{E}_{k^*, \max} = 0$ and (74) is compatible with $\mathcal{M}_k = \emptyset$ and $\mathbb{V}_\infty = \mathbb{V}_k$ for all $k \geq k^*$. In this case, $U_\infty = U_{k^*}$ and convergence is only ensured if a vanishing estimator implies a vanishing error. The latter is given in particular if the estimator bounds the error from above.
- It may happen that only the simplices containing a fixed point are bisected in each iteration, but the exact solution u has a more complex structure so that $u \notin \mathbb{V}_\infty$. Since $u \neq U_\infty$, and uniform refinement is not enforced, the adaptive procedure must depend on the unknown function u .

Convergence therefore will require that the module **ESTIMATE** extracts enough relevant information about the error, the module **MARK** uses this information correctly, and the module **REFINE** reduces the mesh-size where requested.

Mesh-Size before Bisection. The module **REFINE** bisects elements and so halves their volume. This implies the following useful property of elements to be bisected, which include the marked elements.

Lemma 8 (Sequences of elements to be bisected). *For any sequence $(T_k)_{k=0}^\infty$ of elements with $T_k \in \mathcal{T}_k \setminus \mathcal{T}_{k+1}$ there holds $\lim_{k \rightarrow \infty} |T_k| = 0$.*

Proof. Suppose that $\limsup_{k \rightarrow \infty} |T_k| \geq c > 0$, that is there exists a infinite subsequence $(T_{k_\ell})_\ell$ such that $\lim_{\ell \rightarrow \infty} |T_{k_\ell}| \geq c$. Recall that the children of a bisection have half the volume of the parent. Consequently, only a finite number of children of any generation of each T_{k_ℓ} can appear in the sequence $(T_{k_\ell})_\ell$. Eliminating inductively these children, we obtain an infinite sequence of simplices with disjoint interiors and volume greater than $c > 0$. This however contradicts the boundedness of Ω , whence $\limsup_{k \rightarrow \infty} |T_k| \leq 0$, which is equivalent to the assertion. \square

It is instructive and convenient to reformulate Lemma 8 in terms of mesh-size functions.

Lemma 9 (Mesh-size of elements to be bisected). *If χ_k denotes the characteristic function of the union $\cup_{T \in \mathcal{T}_k \setminus \mathcal{T}_{k+1}} T$ of elements to be bisected and h_k is the mesh-size function of \mathcal{T}_k , then*

$$\lim_{k \rightarrow \infty} \|h_k \chi_k\|_{L^\infty(\Omega)} = 0$$

Proof. We may assume that $\mathcal{T}_k \setminus \mathcal{T}_{k+1} \neq \emptyset$ for all $k \in \mathbb{N}_0$ without loss of generality. Choose $(T_k)_{k=0}^\infty$ such that $T_k \in \mathcal{T}_k \setminus \mathcal{T}_{k+1}$ and $h_{T_k} = \max_{T \in \mathcal{T}_k \setminus \mathcal{T}_{k+1}} h_T$ and, recalling that $h_T = |T|^{1/d}$, use Lemma 8 to deduce the assertion. \square

Lemma 9 may be viewed as a generalization of $\lim_{k \rightarrow \infty} \|h_k\|_{L^\infty(\Omega)} = 0$ in the case of uniform refinement. It may be proven also by invoking the limiting mesh-size h_∞ ; see Problems 34 and 35. The limiting mesh-size describes the local structure of \mathbb{V}_∞ and may differ from the zero function.

Convergence to Exact Solution. In order to achieve $U_\infty = u$, we may investigate the residual of U_∞ , which is related to the residuals of the finite element solutions U_k . The latter are in turn controlled by the element indicators $\mathcal{E}_k(U_k, T)$, $T \in \mathcal{T}_k$, which are employed in the step MARK. The fact that indicators with maximum value are marked yields the following property of the largest element indicator $\mathcal{E}_{k,\max}$.

Lemma 10 (Convergence of maximum indicator). *There holds*

$$\lim_{k \rightarrow \infty} \mathcal{E}_{k,\max} = 0.$$

Proof. We may assume that $\mathcal{M}_k \neq \emptyset$ for all $k \in \mathbb{N}_0$ without loss of generality. Choose a sequence $(T_k)_{k=0}^\infty$ of elements such that $T_k \in \mathcal{T}_k$ and $\mathcal{E}_k(U_k, T_k) = \mathcal{E}_{k,\max}$. Thanks to (74), we have $T_k \in \mathcal{M}_k$ and so Lemma 8 and module REFINE yield $\lim_{k \rightarrow \infty} |T_k| = 0$. Exploiting the local lower bound in Theorem 7 and the simplified upper bound for the local oscillation (73), we derive the following estimate for any indicator for $T \in \mathcal{T}_k$:

$$\begin{aligned} \mathcal{E}_k(U_k, T) &\lesssim \|\nabla(U_k - U_\infty)\|_{L^2(\omega_T)} + \|\nabla(U_\infty - u)\|_{L^2(\omega_T)} \\ &\quad + h_T \left(\|f\|_{L^2(\omega_T)} + \|\nabla U_k\|_{L^2(\omega_T)} \right). \end{aligned} \tag{75}$$

Taking $T = T_k$, we obtain

$$\begin{aligned} \mathcal{E}_{k,\max}^e &= \mathcal{E}_k(U_k, T_k) \lesssim \|U_k - U_\infty\|_{\mathbb{V}} + \|\nabla(U_\infty - u)\|_{L^2(\omega_k)} \\ &\quad + |T_k|^{1/d} \left(\|f\|_{L^2(\omega_k)} + \|U_k\|_{\mathbb{V}} \right) \end{aligned}$$

with $\omega_k := \omega_{T_k}$. Consequently, Lemma 7 and $\lim_{k \rightarrow \infty} |T_k| = 0$, which also entails $\lim_{k \rightarrow \infty} |\omega_k| = 0$, prove the assertion. \square

With these preparations we are ready for the first main result of this section.

Theorem 8 (Convergence of approximate solutions). *Let u be the exact solution of the model problem (24) satisfying (25) and (35). The finite element solutions $(U_k)_{k=0}^\infty$ of the AFEM of §4.1 converge to the exact one in \mathbb{V} :*

$$U_k \rightarrow u \text{ in } \mathbb{V} \text{ as } k \rightarrow \infty.$$

Proof. \square In view of Lemma 7, it remains to show that $U_\infty = u$. This is equivalent to

$$0 = \langle \mathcal{R}_\infty, v \rangle := \int_{\Omega} f v - \int_{\Omega} A \nabla U_\infty \cdot \nabla v \quad \text{for all } v \in \mathbb{V} = H_0^1(\Omega). \quad (76)$$

Here we can take the test functions from $C_0^\infty(\Omega)$, because $C_0^\infty(\Omega)$ is a dense subset of the Hilbert space $H_0^1(\Omega)$. Lemma 7 therefore ensures that (76) follows from

$$0 = \lim_{k \rightarrow \infty} \langle \mathcal{R}_k, \varphi \rangle \quad \forall \varphi \in C_0^\infty(\Omega), \quad (77)$$

where $\mathcal{R}_k \in \mathbb{V}^*$ is the residual of U_k given by

$$\langle \mathcal{R}_k, v \rangle := \int_{\Omega} f v - \int_{\Omega} A \nabla U_k \cdot \nabla v.$$

\square In order to show (77), let $\varphi \in C_0^\infty(\Omega)$ and introduce the set

$$\mathcal{T}_\ell^* := \bigcap_{m \geq \ell} \mathcal{T}_m$$

of elements in \mathcal{T}_ℓ that will no longer be bisected; note that if $\mathcal{T}_\ell^* \neq \emptyset$, then $\mathbb{V} \neq \mathbb{V}_\infty$. Given $\ell \leq k$, $\mathbb{V}_\ell \subset \mathbb{V}_k$ and (51) imply

$$\langle \mathcal{R}_k, \varphi \rangle = \langle \mathcal{R}_k, \varphi - I_\ell \varphi \rangle \lesssim S_{\ell,k} + S_{\ell,k}^*, \quad (78)$$

where we expect that

$$S_{\ell,k} := \sum_{T \in \mathcal{T}_k \setminus \mathcal{T}_\ell^*} \mathcal{E}_k(U_k, T) \|\nabla(\varphi - I_\ell \varphi)\|_{L^2(N_k(T))}$$

gets small because of decreasing mesh-size whereas

$$S_{\ell,k}^* := \sum_{T \in \mathcal{T}_\ell^*} \mathcal{E}_k(U_k, T) \|\nabla(\varphi - I_\ell \varphi)\|_{L^2(N_k(T))}$$

gets small because of properties of the adaptive algorithm.

□ We first deal with $S_{\ell,k}$. The Cauchy-Schwarz inequality in some \mathbb{R}^N yields

$$S_{\ell,k} \leq \mathcal{E}_k(U_k, \mathcal{T}_k \setminus \mathcal{T}_\ell^*) \left(\sum_{T \in \mathcal{T}_k \setminus \mathcal{T}_\ell^*} \|\nabla(\varphi - I_\ell \varphi)\|_{L^2(N_k(T))}^2 \right)^{1/2},$$

where the first factor

$$\begin{aligned} \mathcal{E}_k(U_k, \mathcal{T}_k \setminus \mathcal{T}_\ell^*) &\lesssim \|U_k - U_\infty\|_{\mathbb{V}} + \|U_\infty - u\|_{\mathbb{V}} \\ &\quad + \|h_k \chi_\ell\|_{L^\infty(\Omega)} (\|f\|_{L^2(\Omega)} + \|U_k\|_{\mathbb{V}}) \lesssim 1 \end{aligned} \quad (79)$$

is uniformly bounded thanks to (75) and the second factor satisfies

$$\begin{aligned} \left(\sum_{T \in \mathcal{T}_k \setminus \mathcal{T}_\ell^*} \|\nabla(\varphi - I_\ell \varphi)\|_{L^2(N_k(T))}^2 \right)^{1/2} &\lesssim \left(\sum_{T \in \mathcal{T}_\ell \setminus \mathcal{T}_\ell^*} \|\nabla(\varphi - I_\ell \varphi)\|_{L^2(N_\ell(T))}^2 \right)^{1/2} \\ &\lesssim \|h_\ell \chi_\ell\|_{L^\infty(\Omega)}^n \|D^{n+1} \varphi\|_{L^2(\Omega)} \end{aligned}$$

because of $\mathcal{T}_k \geq \mathcal{T}_\ell$ and Proposition 2. Hence Lemma 9 implies

$$S_{\ell,k} \rightarrow 0 \text{ as } \ell \rightarrow \infty \text{ uniformly in } k. \quad (80)$$

□ Next, we deal with $S_{\ell,k}^*$. Here the Cauchy-Schwarz inequality yields

$$S_{\ell,k}^* \leq \mathcal{E}_k(U_k, \mathcal{T}_\ell^*) \left(\sum_{T \in \mathcal{T}_\ell^*} \|\nabla(\varphi - I_\ell \varphi)\|_{L^2(N_k(T))}^2 \right)^{1/2},$$

where the first factor satisfies

$$\mathcal{E}_k(U_k, \mathcal{T}_\ell^*) \leq \#\mathcal{T}_\ell \mathcal{E}_{k,\max} \quad (81)$$

and the second factor

$$\left(\sum_{T \in \mathcal{T}_\ell^*} \|\nabla(\varphi - I_\ell \varphi)\|_{L^2(N_k(T))}^2 \right)^{1/2} \lesssim \|h_\ell\|_{L^\infty(\Omega)}^n \|D^{n+1} \varphi\|_{L^2(\Omega)} \lesssim 1$$

is uniformly bounded. Lemma 10 therefore implies

$$S_{\ell,k}^* \rightarrow 0 \text{ as } k \rightarrow \infty \text{ for any fixed } \ell. \quad (82)$$

□ Given $\varepsilon > 0$, we exploit (80) and (82) by first choosing ℓ so that $S_{\ell,k} \leq \varepsilon/2$ and next $k \geq \ell$ so that $S_{\ell,k}^* \leq \varepsilon/2$. Inserting this into (78) yields the desired convergence (77) and finishes the proof. □

Convergence of Estimator. Theorem 8 ensures convergence of the finite element solutions U_k but says nothing about the behavior of the estimators

$$\mathcal{E}_k(U_k) = \left(\sum_{T \in \mathcal{T}_k} \mathcal{E}_k(U_k, T)^2 \right)^{1/2},$$

which enables one to monitor that convergence. The convergence of the estimators is ensured by the following theorem. Notice that this is not a simple consequence of Theorem 8 and Corollary 3 because of the presence of the oscillation $\text{osc}_k(U_k)$ in the global lower bound; see also Problem 36.

Corollary 4 (Estimator convergence). *Assume again that the model problem (24) satisfies (25) and (35). The estimators $(\mathcal{E}_k(U_k))_{k=0}^\infty$ of AFEM in §4.1 converge to 0:*

$$\lim_{k \rightarrow \infty} \mathcal{E}_k(U_k) = 0.$$

Proof. Theorem 8 implies $U_\infty = u$. Using this and $h_k \leq h_\ell$ for $\ell \leq k$, along with $\|U_k - U_\infty\|_{\mathbb{V}} \lesssim \|U_\ell - U_\infty\|_{\mathbb{V}}$, after the first inequality of (79) yields

$$\mathcal{E}_k(U_k, \mathcal{T}_k \setminus \mathcal{T}_\ell^*) \rightarrow 0 \text{ as } \ell \rightarrow \infty \text{ uniformly in } k \quad (83)$$

with the help of Lemmas 7 and 9. In view of

$$\mathcal{E}_k^2(U_k) = \mathcal{E}_k^2(U_k, \mathcal{T}_k \setminus \mathcal{T}_\ell^*) + \mathcal{E}_k^2(U_k, \mathcal{T}_\ell^*),$$

we realize that (81), (83), and Lemma 10 complete the proof. \square

We conclude this section with a few remarks about variants of Theorem 8 and Corollary 4 for general estimators. Theorem 8 holds for any estimator that provides an upper bound of the form

$$|\langle \mathcal{R}_k, v \rangle| \lesssim \sum_{T \in \mathcal{T}_k} \mathcal{E}_k(U_k, T) \|\nabla v\|_{L^2(N_k(T))} \quad \text{for all } v \in \mathbb{V}, \quad (84)$$

which is locally stable in the sense

$$\mathcal{E}_k(U_k, T) \lesssim h_T \|f\|_{L^2(\omega_T)} + \|\nabla U_k\|_{L^2(\omega_T)} \quad \text{for all } T \in \mathcal{T}_k; \quad (85)$$

see Problem 37. While the first assumption (84) appears natural, and is in fact crucial in view of the first example after Lemma 7, the second assumption (85) may appear artificial. However, Problem 38 reveals that is also crucial and, thus, the suggested variant of Theorem 8 is ‘sharp’. Problem 39 proposes the construction of an estimator verifying the two assumptions (84) and (85) which, however, does not decrease to 0. On the other hand, Corollary 4 hinges on the local lower bound (75), which is a sort of minimal requirement of efficiency if the finite element solutions U_k converge. Roughly speaking, convergence of U_k relies on reliability and stability of the estimator, while the convergence of the estimator depends on the efficiency of the estimator. This shows that the assumptions on the estimator for Theorem 8 and

Corollary 4 are of different nature. In particular, we see that convergence of U_k can be achieved even with estimators that are too poor to quantify the error.

4.3 Notes

The convergence proof in §4.2 is a simplified version of Siebert [49], which unifies the work of Morin, Siebert, and Veeseer [44] with the standard a priori convergence theory based on (global) density. In order to further discuss the underlying assumptions of the approach in 4.2, we now compare these two works in more detail.

Solve. Both works [44] and [49] consider well-posed linear problems and invoke a generalization of Lemma 7 that follows from a discrete inf-sup condition on the discretization. The latter assumption appears natural since it is necessary for convergence in the particular case of uniform refinement; see [10, Problem 3.9]. In the case of a problem with potential or ‘energy’, the explicit construction of U_∞ can be replaced by a convergent sequence of approximate energy minima. Examples are the convergence analyses for the p -Laplacian by Veeseer [55] and for the obstacle problem by Siebert and Veeseer [51], which are the first steps in the terrain of nonlinear and nonsmooth problems and are predecessors of [44] and [49].

Estimate and Mark. Paper [44] differs from [49] on the assumptions on estimators and marking strategy. More precisely, [44] assumes that the estimator provides a discrete local lower bound and that the marking strategy essentially ensures

$$\mathcal{E}_k(U_k, T) \leq \left(\sum_{T' \in \mathcal{M}_k} \mathcal{E}_k(U_k, T')^2 \right)^{1/2} \quad \text{for all } T \in \mathcal{T}_k \setminus \mathcal{M}_k, \quad (86)$$

whereas [49] essentially assumes (84), (85), and (74). Thus, the assumptions on the estimator are weaker in [49], while those on the marking strategy are weaker in [44]; see also Problem 40. Since both works verify that their assumptions on the marking strategy are necessary, this shows that (minimal) assumptions on the estimator and marking strategy are coupled.

Refine. Both [44, 49] consider the same framework for REFINE. This does not only include bisection for conforming meshes (see §1.3), but also nonconforming meshes (see §1.7) and other manners of subdividing elements. Moreover, [44, 49] assume the minimal requirement of subdividing the marked elements, as in §4.1.

Further Variants and Generalizations. These approaches can be further developed in several directions:

- Morin, Siebert, and Veeseer [43] give a proof of convergence of a variant of the AFEM in §4.1 when the estimator provides upper and local lower bounds for the error in ‘weak’ norms, e.g. similar to the L^2 -norm in Problem 18.
- Demlow [20] proves convergence of a variant of the AFEM in §4.1 with estimators for local energy norm errors.

- Garau, Morin, and Zuppa [28] show convergence of a variant of the AFEM in §4.1 for symmetric eigenvalue problems.
- Holst, Tsogtgerel, and Zhu [31] extend [44] to nonlinear partial differential equations, the linearization of which are well-posed.

4.4 Problems

Problem 33 (Nesting of spaces). Let \mathcal{T}_1 and \mathcal{T}_2 be triangulations such that $\mathcal{T}_1 \leq \mathcal{T}_2$, that is \mathcal{T}_2 is a refinement by bisection of \mathcal{T}_1 . Show that the corresponding Lagrange finite element spaces from (31) are nested, i. e., $\mathbb{V}(\mathcal{T}_1) \subset \mathbb{V}(\mathcal{T}_2)$.

Problem 34 (Limiting mesh-size function). Prove that there exists a limiting mesh-size function $h_\infty \in L^\infty(\Omega)$ such that

$$\|h_k - h_\infty\|_{L^\infty(\Omega)} \rightarrow 0 \text{ as } k \rightarrow \infty.$$

Can you construct an example with $h_\infty \neq 0$?

Problem 35 (Alternative proof of Lemma 9). For any iteration k , let χ_k be the characteristic function of the union $\cup_{T \in \mathcal{T}_k \setminus \mathcal{T}_{k+1}} T$ of elements to be bisected and h_k the mesh-size function of \mathcal{T}_k . Show

$$\lim_{k \rightarrow \infty} \|h_k \chi_k\|_{L^\infty(\Omega)} = 0$$

by means of Problem 34 and the fact that bisection reduces the mesh-size.

Problem 36 (Persistence of oscillation). Choosing appropriately the data of the model problem (24), provide an example where the exact solution is (locally) piecewise affine and the (local) oscillation does not vanish.

Problem 37 (Convergence for general estimators). Check that Lemma 10 and Theorem 8 hold for any estimator $\{\mathcal{E}_k(U_k, T)\}_{T \in \mathcal{T}_k}$ that is reliable in the sense of (84) and locally stable in the sense of (85).

Problem 38 ('Necessity' of local estimator stability). Construct an estimator that satisfies (84) and its indicators are always largest around a fixed point, entailing that (74) is compatible with refinement only around that fixed point, irrespective of the exact solution u .

Problem 39 (No estimator convergence). Assuming that the exact solution u of the model problem (24) does not vanish, construct an estimator satisfying (84) and (85) which does not decrease to 0.

Problem 40 (Assumptions for marking strategies). Check that (86) is weaker than (74) by considering the bulk-chasing strategy (90).

5 Contraction Property of AFEM

This section discusses the contraction property for a special AFEM for the *model problem* (23), which we rewrite for convenience:

$$\begin{aligned} -\operatorname{div}(\mathbf{A}(x)\nabla u) &= f && \text{in } \Omega \\ u &= 0 && \text{on } \partial\Omega, \end{aligned} \tag{87}$$

with piecewise smooth coefficient matrix \mathbf{A} on \mathcal{T}_0 . The matrix \mathbf{A} is assumed to be (uniformly) SPD so that the problem is *coercive*, and *symmetric*. We consider a loop of the form

SOLVE \rightarrow ESTIMATE \rightarrow MARK \rightarrow REFINES

that produces a sequence $(\mathcal{T}_k, \mathbb{V}_k, U_k)_{k=0}^\infty$ of conforming meshes \mathcal{T}_k , spaces of conforming elements \mathbb{V}_k (typically C^0 piecewise linears $n = 1$), and Galerkin solutions $U_k \in \mathbb{V}_k$.

The desired contraction property hinges on *error monotonicity*. Since this is closely related to a minimization principle, it is natural to consider the coercive problem (87). We cannot expect a similar theory for problems governed by an inf-sup condition; this is an important open problem.

We next follow Cascón, Kreuzer, Nochetto and Siebert [14]. We refer to [7, 9, 16, 17, 23, 24, 37, 40, 41, 42] for other approaches and to §5.6 for a discussion.

5.1 Modules of AFEM for the Model Problem

We present further properties of the four basic modules of AFEM for (87). The main additional restrictions with respect to §4 are symmetry and coercivity of the bilinear form and the marking strategy.

Module SOLVE. If $\mathcal{T} \in \mathbb{T}$ is a conforming refinement of \mathcal{T}_0 and $\mathbb{V} = \mathbb{V}(\mathcal{T})$ is the finite element space of C^0 piecewise polynomials of degree $\leq n$, then

$$U = \text{SOLVE}(\mathcal{T})$$

determines the Galerkin solution *exactly*, namely,

$$U \in \mathbb{V} : \int_{\Omega} \mathbf{A}\nabla U \cdot \nabla V = \int_{\Omega} fV \quad \text{for all } V \in \mathbb{V}. \tag{88}$$

Module ESTIMATE. Given a conforming mesh $\mathcal{T} \in \mathbb{T}$ and the Galerkin solution $U \in \mathbb{V}(\mathcal{T})$, the output of

$$\{\mathcal{E}_{\mathcal{T}}(U, T)\}_{T \in \mathcal{T}} = \text{ESTIMATE}(U, \mathcal{T}).$$

are the element indicators defined in (46). For convenience, we recall the definitions (37) of *interior* and *jump residuals*

$$\begin{aligned} r(V)|_T &= f + \operatorname{div}(\mathbf{A}\nabla V) && \text{for all } T \in \mathcal{T} \\ j(V)|_S &= \llbracket \mathbf{A}\nabla V \rrbracket \cdot \mathbf{n}|_S && \text{for all } S \in \mathcal{S} \quad (\text{internal sides of } \mathcal{T}), \end{aligned}$$

and $j(V)|_S = 0$ for boundary sides $S \in \mathcal{S}$, as well as the element indicator

$$\mathcal{E}_{\mathcal{T}}^2(V, T) = h_T^2 \|r(V)\|_{L^2(T)}^2 + h_T \|j(V)\|_{L^2(\partial T)}^2 \quad \text{for all } T \in \mathcal{T}. \quad (89)$$

We observe that we now write explicitly the argument V in both r and j because this dependence is relevant for the present discussion.

Module MARK. Given $\mathcal{T} \in \mathbb{T}$, the Galerkin solution $U \in \mathbb{V}(\mathcal{T})$, and element indicators $\{\mathcal{E}_{\mathcal{T}}(U, T)\}_{T \in \mathcal{T}}$, the module MARK selects elements for refinement using *Dörfler Marking* (or bulk chasing), i. e., using a fixed parameter $\theta \in (0, 1]$ the output

$$\mathcal{M} = \text{MARK}(\{\mathcal{E}_{\mathcal{T}}(U, T)\}_{T \in \mathcal{T}}, \mathcal{T})$$

satisfies

$$\mathcal{E}_{\mathcal{T}}(U, \mathcal{M}) \geq \theta \mathcal{E}_{\mathcal{T}}(U, \mathcal{T}). \quad (90)$$

This marking guarantees that \mathcal{M} contains a substantial part of the total (or bulk), thus its name. This is a crucial property in our arguments. The choice of \mathcal{M} does not have to be minimal at this stage, that is, the marked elements $T \in \mathcal{M}$ do not necessarily must be those with largest indicators. However, minimality of \mathcal{M} will be crucial to derive rates of convergence in §7.

Module REFINE. Let $b \in \mathbb{N}$ be the number of desired bisections per marked element. Given $\mathcal{T} \in \mathbb{T}$ and a subset \mathcal{M} of marked elements, the output $\mathcal{T}_* \in \mathbb{T}$ of

$$\mathcal{T}_* = \text{REFINE}(\mathcal{T}, \mathcal{M})$$

is the smallest refinement \mathcal{T}_* of \mathcal{T} so that all elements of \mathcal{M} are at least bisected b times. Therefore, we have $h_{\mathcal{T}_*} \leq h_{\mathcal{T}}$ and the strict reduction property

$$h_{\mathcal{T}_*}|_T \leq 2^{-b/d} h_{\mathcal{T}}|_T \quad \text{for all } T \in \mathcal{M}. \quad (91)$$

We finally let $\mathcal{R}_{\mathcal{T} \rightarrow \mathcal{T}_*}$ be the subset of refined elements of \mathcal{T} and note that

$$\mathcal{M} \subset \mathcal{R}_{\mathcal{T} \rightarrow \mathcal{T}_*}.$$

AFEM. The following procedure is identical to that of §4.1 except for the module MARK, which uses Dörfler marking with parameter $0 < \theta \leq 1$: given an initial grid \mathcal{T}_0 , set $k = 0$ and iterate

1. $U_k = \text{SOLVE}(\mathcal{T}_k)$;
2. $\{\mathcal{E}_k(U_k, T)\}_{T \in \mathcal{T}_k} = \text{ESTIMATE}(U_k, \mathcal{T}_k)$;
3. $\mathcal{M}_k = \text{MARK}(\{\mathcal{E}_k(U_k, T)\}_{T \in \mathcal{T}_k}, \mathcal{T}_k)$;
4. $\mathcal{T}_{k+1} = \text{REFINE}(\mathcal{T}_k, \mathcal{M}_k)$; $k \leftarrow k + 1$.

5.2 Basic Properties of AFEM

We next summarize some basic properties of AFEM that emanate from the symmetry of the differential operator (i.e. of \mathbf{A}) and features of the modules. In doing this, any explicit constant or hidden constant in \lesssim will only depend on the uniform shape-regularity of \mathbb{T} , the dimension d , the polynomial degree n , and the (global) eigenvalues of \mathbf{A} , but not on a specific grid $\mathcal{T} \in \mathbb{T}$, except if explicitly stated. Furthermore, u will always be the weak solution of (24).

The following property relies on the fact that the bilinear form \mathcal{B} is coercive and symmetric, and so induces a scalar product in \mathbb{V} equivalent to the H_0^1 -scalar product.

Lemma 11 (Pythagoras). *Let $\mathcal{T}, \mathcal{T}_* \in \mathbb{T}$ such that $\mathcal{T} \leq \mathcal{T}_*$. The corresponding Galerkin solutions $U \in \mathbb{V}(\mathcal{T})$ and $U_* \in \mathbb{V}(\mathcal{T}_*)$ satisfy the following orthogonality property in the energy norm $\|\cdot\|_{\Omega}$*

$$\|u - U\|_{\Omega}^2 = \|u - U_*\|_{\Omega}^2 + \|U_* - U\|_{\Omega}^2. \quad (92)$$

Proof. See Problem 41. \square

Property (92) is valid for (87) for the energy norm exclusively. This restricts the subsequent analysis to the energy norm, or equivalent norms, but does not extend to other, perhaps more practical, norms such as the maximum norm. This is an important open problem and a serious limitation of this theory.

We now recall the concept of oscillation from §3.1. In view of (65), we denote by $\text{osc}_{\mathcal{T}}(V, T)$ the *element oscillation* for any $V \in \mathbb{V}$

$$\text{osc}_{\mathcal{T}}(V, T) = \|h(r(V) - \overline{r(V)})\|_{L^2(T)} + \|h^{1/2}(j(V) - \overline{j(V)})\|_{L^2(\partial T \cap \Omega)}, \quad (93)$$

where $\overline{r(V)} = P_{2n-2}r(V)$ and $\overline{j(V)} = P_{2n-1}j(V)$ stand for L^2 -projections of the residuals $r(V)$ and $j(V)$ onto the polynomials $\mathbb{P}_{2n-2}(T)$ and $\mathbb{P}_{2n-1}(S)$ defined on the element T or side $S \subset \partial T$, respectively. For variable \mathbf{A} , $\text{osc}_{\mathcal{T}}(V, T)$ depends on the discrete function $V \in \mathbb{V}$, and its study is more involved than for piecewise constant \mathbf{A} . In the latter case, $\text{osc}_{\mathcal{T}}(V, T)$ is given by (58) and is called *data oscillation*; see also Problem 19.

We now rewrite the a posteriori error estimates of Theorems 6 and 7 in the energy norm.

Lemma 12 (A posteriori error estimates). *There exist constants $0 < C_2 \leq C_1$, such that for any $\mathcal{T} \in \mathbb{T}$ and the corresponding Galerkin solution $U \in \mathbb{V}(\mathcal{T})$ there holds*

$$\|u - U\|_{\Omega}^2 \leq C_1 \mathcal{E}_{\mathcal{T}}^2(U) \quad (94a)$$

$$C_2 \mathcal{E}_{\mathcal{T}}^2(U) \leq \|u - U\|_{\Omega}^2 + \text{osc}_{\mathcal{T}}^2(U). \quad (94b)$$

The constants C_1 and C_2 depend on the smallest and largest global eigenvalues of \mathbf{A} . This dependence can be improved if the a posteriori analysis is carried out directly in the energy norm instead of the H_0^1 -norm; see Problem 20. The definitions of $\overline{r(V)}$ and $\overline{j(V)}$, as well as the lower bound (94b), are immaterial for deriving a contraction property. However, they will be important for proving convergence rates in §7.

One serious difficulty in dealing with AFEM is that one has access to the energy error $\|u - U\|_{\Omega}$ only through the estimator $\mathcal{E}_{\mathcal{T}}(U)$. The latter, however, fails to be monotone because it depends on the discrete solution $U \in \mathbb{V}(\mathcal{T})$ that changes with the mesh. We first show that $\mathcal{E}_{\mathcal{T}}(V)$ decreases strictly provided V does not change (Lemma 13) and next we account for the effect of changing V but keeping the mesh (Lemma 14). Combining these two lemmas we get Proposition 3. In formulating these results we rely on the following notation: given $\mathcal{T} \in \mathbb{T}$ let $\mathcal{M} \subset \mathcal{T}$ denote a set of elements that are bisected $b \geq 1$ times at least, let $\mathcal{T}_* \geq \mathcal{T}$ be a conforming refinement of \mathcal{T} that contains the bisected elements of \mathcal{M} , and let

$$\lambda = 1 - 2^{-b/d}.$$

Lemma 13 (Reduction of $\mathcal{E}_{\mathcal{T}}(V)$ wrt \mathcal{T}). *For any $V \in \mathbb{V}(\mathcal{T})$, we have*

$$\mathcal{E}_{\mathcal{T}_*}^2(V, \mathcal{T}_*) \leq \mathcal{E}_{\mathcal{T}}^2(V, \mathcal{T}) - \lambda \mathcal{E}_{\mathcal{T}}^2(V, \mathcal{M}). \quad (95)$$

Proof. We decompose $\mathcal{E}_{\mathcal{T}_*}^2(V, \mathcal{T}_*)$ over elements $T \in \mathcal{T}$, and distinguish whether or not $T \in \mathcal{M}$. If $T \in \mathcal{M}$, then T is bisected at least b times and so T can be written as the union of elements $T' \in \mathcal{T}_*$. We denote this set of elements $\mathcal{T}_*(T)$ and observe that, according with (91), $h_{T'} \leq 2^{-b/d} h_T$ for all $T' \in \mathcal{T}_*(T)$. Therefore

$$\sum_{T' \in \mathcal{T}_*(T)} h_{T'}^2 \|r(V)\|_{L^2(T')}^2 \leq 2^{-2b/d} h_T^2 \|r(V)\|_{L^2(T)}^2$$

and

$$\sum_{T' \in \mathcal{T}_*(T)} h_{T'} \|j(V)\|_{L^2(\partial T' \cap \Omega)}^2 \leq 2^{-b/d} h_T \|j(V)\|_{L^2(\partial T \cap \Omega)}^2,$$

because $V \in \mathbb{V}(\mathcal{T})$ only jumps across the boundary of T . This implies

$$\mathcal{E}_{\mathcal{T}_*}^2(V, T) \leq 2^{-b/d} \mathcal{E}_{\mathcal{T}}^2(V, T) \quad \text{for all } T \in \mathcal{M}.$$

For the remaining elements $T \in \mathcal{T} \setminus \mathcal{M}$ we only know that mesh-size does not increased because $\mathcal{T} \leq \mathcal{T}_*$, whence

$$\mathcal{E}_{\mathcal{T}_*}^2(V, T) \leq \mathcal{E}_{\mathcal{T}}^2(V, T) \quad \text{for all } T \in \mathcal{T} \setminus \mathcal{M}.$$

Combining the two estimates we see that

$$\begin{aligned}\mathcal{E}_{\mathcal{T}_*}^2(V, \mathcal{T}_*) &\leq 2^{-b/d} \mathcal{E}_{\mathcal{T}}^2(V, \mathcal{M}) + \mathcal{E}_{\mathcal{T}}^2(V, \mathcal{T} \setminus \mathcal{M}) \\ &= \mathcal{E}_{\mathcal{T}}^2(V, \mathcal{T}) - (1 - 2^{-b/d}) \mathcal{E}_{\mathcal{T}}^2(V, \mathcal{M}),\end{aligned}$$

which, in light of the definition of λ , is the desired estimate. \square

Lemma 14 (Lipschitz property of $\mathcal{E}_{\mathcal{T}}(V, T)$ wrt V). *For all $T \in \mathcal{T}$, let ω_T denote the union of elements sharing a side with T , $\operatorname{div} \mathbf{A} \in L^\infty(\Omega; \mathbb{R}^d)$ be the divergence of \mathbf{A} computed by rows, and*

$$\eta_{\mathcal{T}}(\mathbf{A}, T) := h_T \|\operatorname{div} \mathbf{A}\|_{L^\infty(T)} + \|\mathbf{A}\|_{L^\infty(\omega_T)}.$$

Then the following estimate is valid

$$|\mathcal{E}_{\mathcal{T}}(V, T) - \mathcal{E}_{\mathcal{T}}(W, T)| \lesssim \eta_{\mathcal{T}}(\mathbf{A}, T) \|\nabla(V - W)\|_{L^2(\omega_T)} \quad \text{for all } V, W \in \mathbb{V}(\mathcal{T}).$$

Proof. Recalling the definition of the indicators, the triangle inequality yields

$$|\mathcal{E}_{\mathcal{T}}(V, T) - \mathcal{E}_{\mathcal{T}}(W, T)| \leq h_T \|r(V) - r(W)\|_{L^2(T)} + h_T^{1/2} \|j(V) - j(W)\|_{L^2(\partial T)}.$$

We set $E := V - W \in \mathbb{V}(\mathcal{T})$, and observe that

$$r(V) - r(W) = \operatorname{div}(\mathbf{A} \nabla E) = \operatorname{div} \mathbf{A} \cdot \nabla E + \mathbf{A} : D^2 E,$$

where $D^2 E$ is the Hessian of E . Since E is a polynomial of degree $\leq n$ in T , applying the inverse estimate $\|D^2 E\|_{L^2(T)} \lesssim h_T^{-1} \|\nabla E\|_{L^2(T)}$, we deduce

$$h_T \|r(V) - r(W)\|_{L^2(T)} \lesssim \eta_{\mathcal{T}}(\mathbf{A}, T) \|\nabla E\|_{L^2(T)}.$$

On the other hand, for any $S \subset \partial T$ applying the inverse estimate of Problem 43 gives

$$\|j(V) - j(W)\|_{L^2(S)} = \|j(E)\|_{L^2(S)} = \|\llbracket \mathbf{A} \nabla E \rrbracket\|_{L^2(S)} \lesssim h_T^{-1/2} \|\nabla E\|_{L^2(\omega_T)}$$

where the hidden constant is proportional to $\eta_{\mathcal{T}}(\mathbf{A}, T)$. This finishes the proof. \square

Proposition 3 (Estimator reduction). *Given $\mathcal{T} \in \mathbb{T}$ and a subset $\mathcal{M} \subset \mathcal{T}$ of marked elements, let $\mathcal{T}_* = \operatorname{REFINE}(\mathcal{T}, \mathcal{M})$. Then there exists a constant $\Lambda > 0$, such that for all $V \in \mathbb{V}(\mathcal{T})$, $V_* \in \mathbb{V}_*(\mathcal{T}_*)$ and any $\delta > 0$ we have*

$$\begin{aligned}\mathcal{E}_{\mathcal{T}_*}^2(V_*, \mathcal{T}_*) &\leq (1 + \delta) (\mathcal{E}_{\mathcal{T}}^2(V, \mathcal{T}) - \lambda \mathcal{E}_{\mathcal{T}}^2(V, \mathcal{M})) \\ &\quad + (1 + \delta^{-1}) \Lambda \eta_{\mathcal{T}}^2(\mathbf{A}, \mathcal{T}) \|V_* - V\|_{\Omega}^2.\end{aligned}$$

Proof. Apply Lemma 14 to $V, V_* \in \mathbb{V}(\mathcal{T}_*)$ in conjunction with Lemma 13 for V (see Problem 44). \square

5.3 Contraction Property of AFEM

A key question to ask is what is (are) the quantity(ies) that AFEM may contract. In light of (92), an obvious candidate is the energy error $\|u - U_k\|_\Omega$. We first show, in the simplest scenario of piecewise constant data \mathbf{A} and f , that this is in fact the case provided an interior node property holds; see Lemma 15. However, the energy error may not contract in general unless REFINE enforces several levels of refinement; see Example 1. We then present a more general approach that eliminates the interior node property at the expense of a more complicated contractive quantity, the quasi error; see Theorem 9.

Piecewise Constant Data. We now assume that both f and \mathbf{A} are piecewise constant in the initial mesh \mathcal{T}_0 , so that $\text{osc}_k(U_k) = 0$ for all $k \geq 0$. The following property was introduced by Morin, Nochetto, and Siebert [40].

Definition 1 (Interior node property). The refinement $\mathcal{T}_{k+1} \geq \mathcal{T}_k$ satisfies an interior node property with respect to \mathcal{T}_k if each element $T \in \mathcal{M}_k$ contains at least one node of \mathcal{T}_{k+1} in the interiors of T and of each side of T .

This property is valid upon enforcing a fixed number b_* of bisections ($b_* = 3, 6$ for $d = 2, 3$). An immediate consequence of this property, proved in [40, 41], is the following *discrete* lower a posteriori bound:

$$C_2 \mathcal{E}_k^2(U_k, \mathcal{M}_k) \leq \|U_k - U_{k+1}\|_\Omega^2 + \text{osc}_k^2(U_k); \quad (96)$$

see also Problem 23 for a related result.

Lemma 15 (Contraction property for piecewise constant data). *Let \mathbf{A}, f be piecewise constant in the initial mesh \mathcal{T}_0 . If \mathcal{T}_{k+1} satisfies an interior node property with respect to \mathcal{T}_k , then for $\alpha := (1 - \theta^2 \frac{C_2}{C_1})^{1/2} < 1$ there holds*

$$\|u - U_{k+1}\|_\Omega \leq \alpha \|u - U_k\|_\Omega, \quad (97)$$

where $0 < \theta < 1$ is the parameter in (90) and $C_1 \geq C_2$ are the constants in (94).

Proof. For convenience, we use the notation

$$e_k = \|u - U_k\|_\Omega, \quad E_k = \|U_{k+1} - U_k\|_\Omega, \quad \mathcal{E}_k = \mathcal{E}_k(U_k, \mathcal{T}_k), \quad \mathcal{E}_k(\mathcal{M}_k) = \mathcal{E}_k(U_k, \mathcal{M}_k).$$

The key idea is to use the Pythagoras equality (11)

$$e_{k+1}^2 = e_k^2 - E_k^2,$$

and show that E_k is a significant portion of e_k . Since (96) with $\text{osc}_k(U_k) = 0$ implies

$$C_2 \mathcal{E}_k^2(\mathcal{M}_k) \leq E_k^2, \quad (98)$$

applying Dörfler marking (90) and the upper bound (94a), we deduce

$$E_k^2 \geq C_2 \theta^2 \mathcal{E}_k^2 \geq \frac{C_2}{C_1} \theta^2 e_k^2.$$

This is the desired property of E_k and leads to (97). \square

Example 1 (Strict monotonicity). Let $\Omega = (0, 1)^2$, $\mathbf{A} = \mathbf{I}$, $f = 1$ (constant data), and consider the following sequences of meshes depicted in Figure 14. If ϕ_0 denotes the basis function associated with the only interior node of the initial mesh \mathcal{T}_0 , then

$$U_0 = U_1 = \frac{1}{12} \phi_0, \quad U_2 \neq U_1.$$

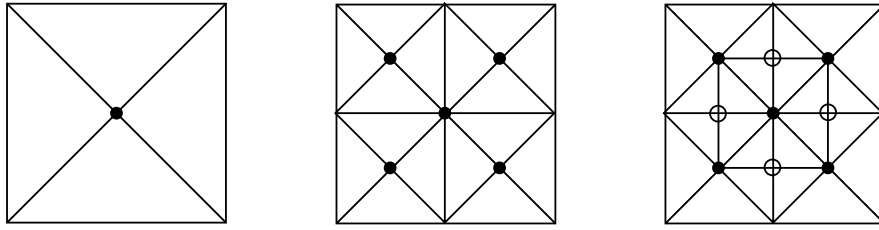


Fig. 14 Grids \mathcal{T}_0 , \mathcal{T}_1 , and \mathcal{T}_2 of Example 1. The mesh \mathcal{T}_1 has nodes in the middle of sides of \mathcal{T}_0 , but only \mathcal{T}_2 has nodes in the interior of elements of \mathcal{T}_0 . Hence, \mathcal{T}_2 satisfies the interior node property of Definition 1 with respect to \mathcal{T}_0 .

The mesh $\mathcal{T}_1 \geq \mathcal{T}_0$ is produced by a standard 2-step bisection ($b = 2$) in $2d$. Since $U_0 = U_1$ we conclude that the energy error may not change

$$\|u - U_0\|_{\Omega} = \|u - U_1\|_{\Omega}$$

between two consecutive steps of AFEM for $b = d = 2$. This is no longer true provided an interior node in each marked element is created, as in Definition 1, because then Lemma 15 holds. This example appeared first in [40, 41], and was used to justify the interior node property.

General Data. If $\text{osc}_k(U_k) \neq 0$, then the contraction property of AFEM becomes trickier because the energy error and estimator are no longer equivalent regardless of the interior node property. The first question to ask is what quantity replaces the energy error in the analysis. We explore this next and remove the interior node property.

Heuristics. According to (92), the energy error is monotone

$$\|u - U_{k+1}\|_{\Omega} \leq \|u - U_k\|_{\Omega},$$

but the previous example shows that strict inequality may fail. However, if $U_{k+1} = U_k$, estimate (95) reveals a strict estimator reduction $\mathcal{E}_{k+1}(U_k) < \mathcal{E}_k(U_k)$. We thus

expect that, for a suitable scaling factor $\gamma > 0$, the so-called *quasi error*

$$\|u - U_k\|_{\Omega}^2 + \gamma \mathcal{E}_k^2(U_k) \quad (99)$$

may be contractive. This heuristics illustrates a distinct aspect of AFEM theory, the interplay between continuous quantities such the energy error $\|u - U_k\|_{\Omega}$ and discrete ones such as the estimator $\mathcal{E}_k(U_k)$: no one alone has the requisite properties to yield a contraction between consecutive adaptive steps.

Theorem 9 (Contraction property). *Let $\theta \in (0, 1]$ be the Dörfler Marking parameter, and $\{\mathcal{T}_k, \mathbb{V}_k, U_k\}_{k=0}^{\infty}$ be a sequence of conforming meshes, finite element spaces and discrete solutions created by AFEM for the model problem (87).*

Then there exist constants $\gamma > 0$ and $0 < \alpha < 1$, additionally depending on the number $b \geq 1$ of bisections and θ , such that for all $k \geq 0$

$$\|u - U_{k+1}\|_{\Omega}^2 + \gamma \mathcal{E}_{k+1}^2(U_{k+1}) \leq \alpha^2 \left(\|u - U_k\|_{\Omega}^2 + \gamma \mathcal{E}_k^2(U_k) \right). \quad (100)$$

Proof. We split the proof into four steps and use the notation in Lemma 15.

1 The error orthogonality (92) reads

$$e_{k+1}^2 = e_k^2 - E_k^2. \quad (101)$$

Employing Proposition 3 with $\mathcal{T} = \mathcal{T}_k$, $\mathcal{T}_* = \mathcal{T}_{k+1}$, $V = U_k$ and $V_* = U_{k+1}$ gives

$$\mathcal{E}_{k+1}^2 \leq (1 + \delta) (\mathcal{E}_k^2 - \lambda \mathcal{E}_k^2(\mathcal{M}_k)) + (1 + \delta^{-1}) \Lambda_0 E_k^2, \quad (102)$$

where $\Lambda_0 = \Lambda \eta_{\mathcal{T}_0}^2(\mathbf{A}, \mathcal{T}_0) \geq \Lambda \eta_{\mathcal{T}_k}^2(\mathbf{A}, \mathcal{T}_k)$. After multiplying (102) by $\gamma > 0$, to be determined later, we add (101) and (102) to obtain

$$e_{k+1}^2 + \gamma \mathcal{E}_{k+1}^2 \leq e_k^2 + (\gamma(1 + \delta^{-1}) \Lambda_0 - 1) E_k^2 + \gamma(1 + \delta) (\mathcal{E}_k^2 - \lambda \mathcal{E}_k^2(\mathcal{M}_k)).$$

2 We now choose the parameters δ, γ , the former so that

$$(1 + \delta)(1 - \lambda \theta^2) = 1 - \frac{\lambda \theta^2}{2},$$

and the latter to verify

$$\gamma(1 + \delta^{-1}) \Lambda_0 = 1.$$

Note that this choice of γ yields

$$e_{k+1}^2 + \gamma \mathcal{E}_{k+1}^2 \leq e_k^2 + \gamma(1 + \delta) (\mathcal{E}_k^2 - \lambda \mathcal{E}_k^2(\mathcal{M}_k)).$$

3 We next employ Dörfler Marking, namely $\mathcal{E}_k(\mathcal{M}_k) \geq \theta \mathcal{E}_k$, to deduce

$$e_{k+1}^2 + \gamma \mathcal{E}_{k+1}^2 \leq e_k^2 + \gamma(1 + \delta)(1 - \lambda \theta^2) \mathcal{E}_k^2$$

which, in conjunction with the choice of δ , gives

$$e_{k+1}^2 + \gamma \mathcal{E}_{k+1}^2 \leq e_k^2 + \gamma \left(1 - \frac{\lambda \theta^2}{2}\right) \mathcal{E}_k^2 = e_k^2 - \frac{\gamma \lambda \theta^2}{4} \mathcal{E}_k^2 + \gamma \left(1 - \frac{\lambda \theta^2}{4}\right) \mathcal{E}_k^2.$$

□ Finally, the upper bound (94a), namely $e_k^2 \leq C_1 \mathcal{E}_k^2$, implies that

$$e_{k+1}^2 + \gamma \mathcal{E}_{k+1}^2 \leq \left(1 - \frac{\gamma \lambda \theta^2}{4C_1}\right) e_k^2 + \gamma \left(1 - \frac{\lambda \theta^2}{4}\right) \mathcal{E}_k^2.$$

This in turn leads to

$$e_{k+1}^2 + \gamma \mathcal{E}_{k+1}^2 \leq \alpha^2 (e_k^2 + \gamma \mathcal{E}_k^2),$$

with

$$\alpha^2 := \max \left\{ 1 - \frac{\gamma \lambda \theta^2}{4C_1}, 1 - \frac{\lambda \theta^2}{4} \right\},$$

and proves the theorem because $\alpha^2 < 1$. □

Remark 9 (Ingredients). The basic ingredients of this proof are: Dörfler marking; coercivity and symmetry of \mathcal{B} and nesting of spaces, which imply the Pythagoras identity (Lemma 11); the a posteriori upper bound (Lemma 12); and the estimator reduction property (Proposition 3). It does not use the lower bound (94b) and does not require marking by oscillation, as previous proofs do [17, 37, 40, 41, 42].

Remark 10 (Separate marking). MARK is driven by \mathcal{E}_k exclusively, as it happens in all practical AFEM. Previous proofs in [17, 37, 40, 41, 42] require separate marking by estimator and oscillation. It is shown in [14] that separate marking may lead to suboptimal convergence rates. On the other hand, we will prove in §7 that the present AFEM yields quasi-optimal convergence rates.

5.4 Example: Discontinuous Coefficients

We invoke the formulas derived by Kellogg [34] to construct an exact solution of an elliptic problem with piecewise constant coefficients and vanishing right-hand side f ; data oscillation is thus immaterial. We now write these formulas in the particular case $\Omega = (-1, 1)^2$, $\mathbf{A} = a_1 \mathbf{I}$ in the first and third quadrants, and $\mathbf{A} = a_2 \mathbf{I}$ in the second and fourth quadrants. An exact weak solution u of the model problem (87) for $f \equiv 0$ is given in polar coordinates by $u(r, \theta) = r^\gamma \mu(\theta)$ (see Figure 15), where

$$\mu(\theta) = \begin{cases} \cos((\pi/2 - \sigma)\gamma) \cdot \cos((\theta - \pi/2 + \rho)\gamma) & \text{if } 0 \leq \theta \leq \pi/2, \\ \cos(\rho\gamma) \cdot \cos((\theta - \pi + \sigma)\gamma) & \text{if } \pi/2 \leq \theta \leq \pi, \\ \cos(\sigma\gamma) \cdot \cos((\theta - \pi - \rho)\gamma) & \text{if } \pi \leq \theta < 3\pi/2, \\ \cos((\pi/2 - \rho)\gamma) \cdot \cos((\theta - 3\pi/2 - \sigma)\gamma) & \text{if } 3\pi/2 \leq \theta \leq 2\pi, \end{cases}$$

and the numbers γ, ρ, σ satisfy the nonlinear relations

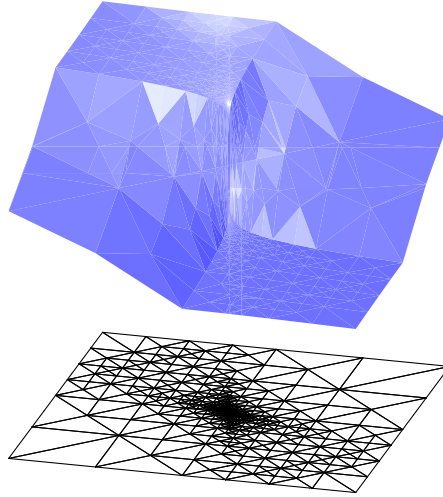


Fig. 15 Discontinuous coefficients in checkerboard pattern: Graph of the discrete solution, which is $u \approx r^{0.1}$, and underlying strongly graded grid. Notice the singularity of u at the origin.

$$\begin{cases} R := a_1/a_2 = -\tan((\pi/2 - \sigma)\gamma) \cdot \cot(\rho\gamma), \\ 1/R = -\tan(\rho\gamma) \cdot \cot(\sigma\gamma), \\ R = -\tan(\sigma\gamma) \cdot \cot((\pi/2 - \rho)\gamma), \\ 0 < \gamma < 2, \\ \max\{0, \pi\gamma - \pi\} < 2\gamma\rho < \min\{\pi\gamma, \pi\}, \\ \max\{0, \pi - \pi\gamma\} < -2\gamma\sigma < \min\{\pi, 2\pi - \pi\gamma\}. \end{cases} \quad (103)$$

Since we want to test the algorithm AFEM in a worst case scenario, we choose $\gamma = 0.1$, which produces a very singular solution u that is barely in H^1 ; in fact $u \in H^s(\Omega)$ for $s < 1.1$ and piecewise in $W_p^2(\Omega)$ for $p > 1$. We then solve (103) for R , ρ , and σ using Newton's method to obtain

$$R = a_1/a_2 \cong 161.4476387975881, \quad \rho = \pi/4, \quad \sigma \cong -14.92256510455152,$$

and finally choose $a_1 = R$ and $a_2 = 1$. A smaller γ would lead to a larger ratio R , but in principle γ may be as close to 0 as desired.

We realize from Figure 16 that AFEM attains optimal decay rate for the energy norm. This is consistent with adaptive approximation for functions piecewise in $W_p^2(\Omega)$ (see §1.6), but nonobvious for AFEM which does not have direct access to u ; this is the topic of §7. We also notice from Figure 17 that a graded mesh with mesh-size of order 10^{-10} at the origin is achieved with about 2×10^3 elements. To reach a similar resolution with a uniform mesh we would need $N \approx 10^{20}$ elements! This example clearly reveals the advantages and potentials of adaptivity for the FEM even with modest computational resources.

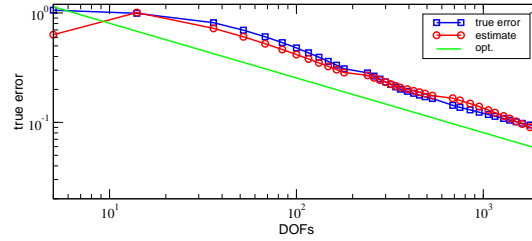


Fig. 16 Quasi-optimality of AFEM for discontinuous coefficients: estimate and true error. The optimal decay for piecewise linear elements in 2d is indicated by the line with slope $-1/2$.

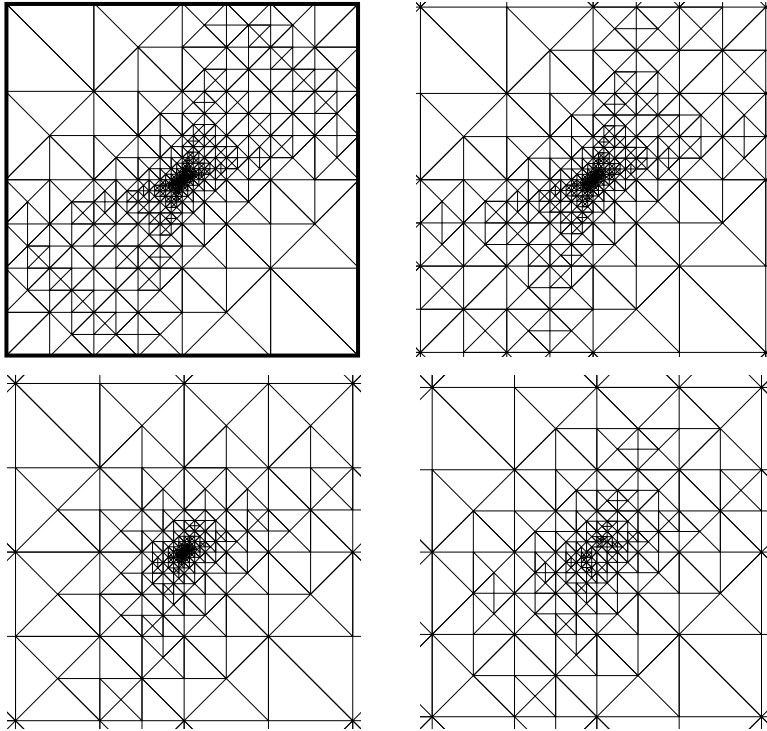


Fig. 17 Discontinuous coefficients in checkerboard pattern: Final grid (full grid with < 2000 nodes) (top left), zooms to $(-10^{-3}, 10^{-3})^2$ (top right), $(-10^{-6}, 10^{-6})^2$ (bottom left), and $(-10^{-9}, 10^{-9})^2$ (bottom right). The grid is highly graded towards the origin. For a similar resolution, a uniform grid would require $N \approx 10^{20}$ elements.

What is missing is an explanation of the recovery of optimal error decay $N^{-1/2}$ through mesh grading. This is the subject of §7, where we have to deal with the interplay between continuous and discrete quantities as already alluded to in the heuristics.

5.5 Extensions and Restrictions

It is important to take a critical look at the theory just developed and wonder about its applicability. Below we list a few extensions of the theory and acknowledge some restrictions.

Nonconforming Meshes. Theorem 9 easily extends to non-conforming meshes since conformity plays no role. This is reported in Bonito and Nochetto [9].

Non-Residual Estimators. The contraction property (100) has been derived for residual estimators $\mathcal{E}_k(U_k)$. This is because the estimator reduction property (95) is not known to hold for other estimators, such as hierarchical, Zienkiewicz-Zhu, and Braess-Schoerbel estimators, as well as those based on the solution of local problems. A common feature of these estimators $\eta_{\mathcal{T}}(U)$ is the lack of reliability in the preasymptotic regime, in which oscillation $\text{osc}_{\mathcal{T}}(U)$ may dominate. In fact, we recall the upper a posteriori bound from (72)

$$\|u - U\|_{\Omega}^2 \leq C_1 \left(\eta_{\mathcal{T}}^2(U) + \text{osc}_{\mathcal{T}}^2(U) \right) =: \mathcal{E}_{\mathcal{T}}^2(U),$$

which gives rise to Dörfler marking for the total estimator $\mathcal{E}_{\mathcal{T}}(U)$. Cascón and Nochetto [15] have recently extended Theorem 9 for $n = 1$ upon allowing an interior node property after a fixed number of adaptive loops and combining Lemma 15 with Theorem 9; this is easy to implement within ALBERTA [50]. At the same time, using the local equivalence of the above estimators with the residual one, Kreuzer and Siebert have proved an error reduction property after several adaptive loops [35].

Elliptic PDE on Manifolds. Meckhay, Morin and Nochetto extended this theory to the Laplace-Beltrami operator [38]. In this case, an additional geometric error due to piecewise polynomial approximation of the surface must be accounted for.

Discontinuous Galerkin Methods (dG). The convergence results available in the literature are for the *interior penalty* method [9, 32, 33]. The simplest contraction property (97) for a right-hand side f in the finite element space and the Laplace operator was first derived by Karakashian and Pascal [33], and later improved by Hoppe, Kanschat, and Warburton [32] for $f \in L^2$ and just one bisection per marked element. In both cases, the theory is developed for $d = 2$. The most general result, valid for $d \geq 2$, operators with discontinuous variable coefficients, and L^2 data, has been developed by Bonito and Nochetto [9]. The theory in [9] deals with non-conforming meshes made of quadrilaterals or triangles, or their multidimensional generalizations, which are natural in the dG context. A key theoretical issue is the control of the jump term, which is not monotone with refinement [33, 9].

Saddle Point Problems. The contraction properties (97) and (100) rely crucially on the Pythagoras orthogonality property (92) and does not extend to saddle point problems. However, a modified AFEM based on an inexact Uzawa iteration and separate marking was shown to converge by Bänsch, Morin, and Nochetto for the Stokes equation [6]. The situation is somewhat simpler for mixed FEM for scalar

second order elliptic PDE, and has been tackled directly for $d = 2$ by Carstensen and Hoppe for the lowest order Raviart-Thomas element [13], and by Chen, Holst, and Xu for any order [18]. They exploit the underlying special structure: the flux error is L^2 -orthogonal to the discrete divergence free subspace, whereas the nonvanishing divergence component of the flux error can be bounded by data oscillation. This is not valid for the Stokes system, which remains open.

Beyond the Energy Framework. The contraction properties (97) and (100) may fail also for other norms of practical interest. An example is the maximum norm, for which there is no convergence result known yet of AFEM. Demlow proved a contraction property for local energy errors [20], and Demlow and Stevenson [21] showed a contraction property for the L^2 norm provided that the mesh grading is sufficiently mild.

5.6 Notes

The theory for conforming meshes in dimension $d > 1$ started with Dörfler [24], who introduced the crucial marking (90), the so-called *Dörfler marking*, and proved strict energy error reduction for the Laplacian provided the initial mesh \mathcal{T}_0 satisfies a fineness assumption. This marking plays an essential role in the present discussion, which does not seem to extend to other marking strategies such as those in §4. Morin, Nochetto, and Siebert [40, 41] showed that such strict energy error reduction does not hold in general even for the Laplacian. They introduced the concept of data oscillation and the interior node property, and proved convergence of the AFEM without restrictions on \mathcal{T}_0 . The latter result, however, is valid only for \mathbf{A} in (23) piecewise constant on \mathcal{T}_0 . Inspired by the work of Chen and Feng [17], Mekchay and Nochetto [37] proved a contraction property for the *total error*, namely the sum of the energy error plus oscillation for \mathbf{A} piecewise smooth. The total error will reappear in the study of convergence rates in §7.

Diening and Kreuzer proved a similar contraction property for the p -Laplacian replacing the energy norm by a so-called quasi-norm [23]. They were able to avoid marking for oscillation by using the fact that oscillation is dominated by the estimator. Most results for nonlinear problems utilize the equivalence of the energy error and error in the associated (nonlinear) energy; compare with Problem 42. This equivalence was first used by Veeseer in a convergence analysis for the p -Laplacian [55] and later on by Siebert and Veeseer for the obstacle problem [51].

The result of Diening and Kreuzer inspired the work by Cascón, Kreuzer, Nochetto, and Siebert [14]. This approach hinges solely on a strict reduction of the mesh-size within refined elements, the upper a posteriori error bound, an orthogonality property natural for (87) in nested approximation spaces, and Dörfler marking. This appears to be the simplest approach currently available.

5.7 Problems

Problem 41 (Pythagoras). Let $\mathbb{V}_1 \subset \mathbb{V}_2 \subset \mathbb{V} = H_0^1(\Omega)$ be nested, conforming and closed subspaces. Let $u \in \mathbb{V}$ be the weak solution to (87), $U_1 \in \mathbb{V}_1$ and $U_2 \in \mathbb{V}_2$ the respective Ritz-Galerkin approximations to u . Prove the orthogonality property

$$\|u - U_1\|_{\Omega}^2 = \|u - U_2\|_{\Omega}^2 + \|U_2 - U_1\|_{\Omega}^2. \quad (104)$$

Problem 42 (Error in energy). Let $\mathbb{V}_1 \subset \mathbb{V}_2 \subset \mathbb{V}$ and U_1, U_2, u be as in Problem 41. Recall that u, U_1, U_2 are the unique minimizers of the quadratic energy

$$I[v] := \frac{1}{2} \mathcal{B}[v, v] - \langle f, v \rangle$$

in $\mathbb{V}, \mathbb{V}_1, \mathbb{V}_2$ respectively. Show that (104) is equivalent to the identity

$$I[U_1] - I[u] = (I[U_2] - I[u]) + (I[U_1] - I[U_2]).$$

To this end prove

$$I[U_i] - I[u] = \frac{1}{2} \|U_i - u\|_{\Omega}^2 \quad \text{and} \quad I[U_1] - I[U_2] = \frac{1}{2} \|U_1 - U_2\|_{\Omega}^2.$$

Problem 43 (Inverse estimate). Let $S \in \mathcal{S}$ be an interior side of $T \in \mathcal{T}$, and let $\mathbf{A} \in L^\infty(S)$. Make use of a scaling argument to the reference element to show

$$\|\mathbf{A} \nabla V\|_S \lesssim h_S^{-1/2} \|\nabla V\|_T \quad \text{for all } V \in \mathbb{V}(\mathcal{T}),$$

where the hidden constant depends on the shape coefficient of \mathcal{T} , the dimension d , and $\|\mathbf{A}\|_{L^\infty(S)}$.

Problem 44 (Proposition 3). Complete the proof of Proposition 3 upon using Young inequality

$$(a + b)^2 \leq (1 + \delta)a^2 + (1 + \delta^{-1})b^2 \quad \text{for all } a, b \in \mathbb{R}.$$

Problem 45 (Quasi-local Lipschitz property). Let $\mathbf{A} \in W_\infty^1(T)$ for all $T \in \mathcal{T}$. Prove

$$|\text{osc}_{\mathcal{T}}(V, T) - \text{osc}_{\mathcal{T}}(W, T)| \lesssim \text{osc}_{\mathcal{T}}(\mathbf{A}, T) \|V - W\|_{H^1(\omega_T)} \quad \text{for all } V, W \in \mathbb{V},$$

where $\text{osc}_{\mathcal{T}}(\mathbf{A}, T) = h_T \|\text{div} \mathbf{A} - P_{n-1}^\infty(\text{div} \mathbf{A})\|_{L^\infty(T)} + \|\mathbf{A} - P_n^\infty \mathbf{A}\|_{L^\infty(\omega_T)}$. Proceed as in the proof of Lemma 14 and use Problem 28.

Problem 46 (Perturbation). Let $\mathcal{T}, \mathcal{T}_* \in \mathbb{T}$, with $\mathcal{T} \leq \mathcal{T}_*$. Use Problem 45 to prove that, for all $V \in \mathbb{V}(\mathcal{T})$ and $V_* \in \mathbb{V}(\mathcal{T}_*)$, there is a constant $\Lambda_1 > 0$ such that

$$\text{osc}_{\mathcal{T}}^2(V, \mathcal{T} \cap \mathcal{T}_*) \leq 2 \text{osc}_{\mathcal{T}_*}^2(V_*, \mathcal{T} \cap \mathcal{T}_*) + \Lambda_1 \text{osc}_{\mathcal{T}_0}^2(\mathbf{A}, \mathcal{T}_0)^2 \|V - V_*\|_{\Omega}^2.$$

6 Complexity of Refinement

This section is devoted to proving Theorem 1 for conforming meshes and Lemma 3 for nonconforming meshes. The results of Sections 6.1 and 6.2 are valid for $d = 2$ but the proofs of Theorem 1 in Section 6.3 and Lemma 3 in Section 6.4 extend easily to $d > 2$. We refer to the survey [45] for a full discussion for $d \geq 2$.

6.1 Chains and Labeling for $d = 2$

In order to study nonlocal effects of bisection for $d = 2$ we introduce now the concept of chain [7]; this concept is not adequate for $d > 2$ [45, 53]. Recall that $E(T)$ denotes the edge of T assigned for refinement. To each $T \in \mathcal{T}$ we associate the element $F(T) \in \mathcal{T}$ sharing the edge $E(T)$ if $E(T)$ is interior and $F(T) = \emptyset$ if $E(T)$ is on $\partial\Omega$. A *chain* $\mathcal{C}(T, \mathcal{T})$, with starting element $T \in \mathcal{T}$, is a sequence $\{T, F(T), \dots, F^m(T)\}$ with no repetitions of elements and with

$$F^{m+1}(T) = F^k(T) \text{ for } k \in \{0, \dots, m-1\} \text{ or } F^{m+1}(T) = \emptyset;$$

see Figure 18. We observe that if an element T belongs to two different grids, then the corresponding chains may be different as well. Two adjacent elements $T, T' =$

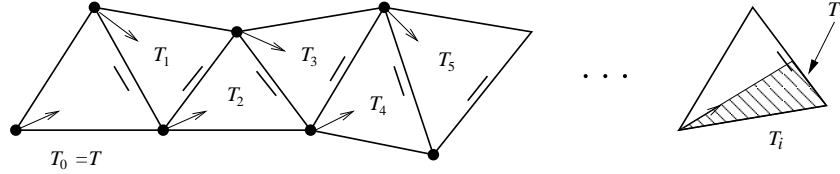


Fig. 18 Typical chain $\mathcal{C}(T, \mathcal{T}) = \{T_j\}_{j=0}^i$ emanating from $T = T_0 \in \mathcal{T}$ with $T_j = F(T_{j-1}), j \geq 1$.

$F(T)$ are *compatibly divisible* (or equivalently T, T' form a *compatible bisection patch*) if $F(T') = T$. Hence, $\mathcal{C}(T, \mathcal{T}) = \{T, T'\}$ and a bisection of either T or T' does not propagate outside the patch.

Example (Chains): Let $\mathcal{F} = \{T_i\}_{i=1}^{12}$ be the forest of Figure 3. Then $\mathcal{C}(T_6, \mathcal{F}) = \{T_6, T_7\}$, $\mathcal{C}(T_9, \mathcal{F}) = \{T_9\}$, and $\mathcal{C}(T_{10}, \mathcal{F}) = \{T_{10}, T_8, T_2\}$ are chains, but only $\mathcal{C}(T_6, \mathcal{F})$ is a compatible bisection patch.

To study the structure of chains we rely on the initial labeling (6) and the bisection rule of Section 1.3 (see Figure 5):

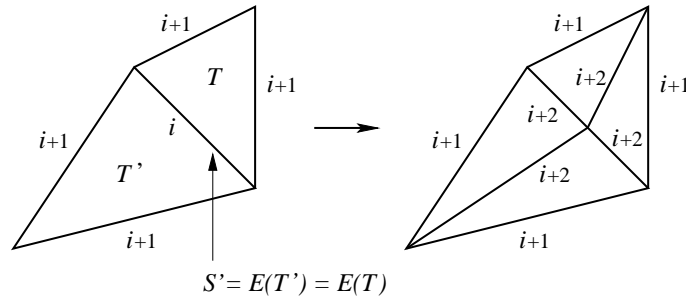
every triangle $T \in \mathcal{T}$ with generation $g(T) = i$ receives the label $(i+1, i+1, i)$ with i corresponding to the refinement edge $E(T)$, its side i is bisected and both new sides as well as the bisector are labeled $i+2$ whereas the remaining labels do not change. (105)

We first show that once the initial labeling and bisection rule are set, the resulting master forest \mathbb{F} is uniquely determined: the label of an edge is independent of the elements sharing this edge and no ambiguity arises in the recursion process.

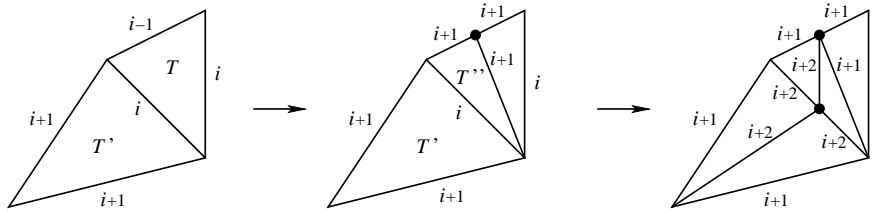
Lemma 16 (Labeling). *Let the initial labeling (6) for \mathcal{T}_0 and above bisection rule be enforced. If $\mathcal{T}_0 \leq \mathcal{T}_1 \leq \dots \leq \mathcal{T}_n$ are generated according to (105), then each side in \mathcal{T}_k has a unique label independent of the two triangles sharing this edge.*

Proof. We argue by induction over \mathcal{T}_k . For $k = 0$ the assertion is valid due to the initial labeling. Suppose the statement is true for \mathcal{T}_k . An edge S in \mathcal{T}_{k+1} can be obtained in two ways. The first is that S is a bisector, and so a new edge, in which case there is nothing to prove about its label being unique. The second possibility is that S was obtained by bisecting an edge $S' \in \mathcal{T}_k$. Let $T, T' \in \mathcal{T}_k$ be the elements sharing S' , and let us assume that $E(T') = S'$. Let $(i + 1, i + 1, i)$ be the label of T' , which means that S is assigned the label $i + 2$. By induction assumption over \mathcal{T}_k , the label of S' as an edge of T is also i . There are two possible cases for the label of T :

- Label $(i + 1, i + 1, i)$: this situation is symmetric, $E(T) = S'$, and S' is bisected with both halves getting label $i + 2$. This is depicted in the figure below.



- Label $(i, i, i - 1)$: a bisection of side $E(T)$ with label $i - 1$ creates a children T'' with label $(i + 1, i + 1, i)$ that is compatibly divisible with T' . Joining the new node of T with the midpoint of S' creates a conforming partition with level $i + 2$ assigned to S . This is depicted in the figure below.



Therefore, in both cases the label $i + 2$ assigned to S is the same from both sides, as asserted. \square

The two possible configurations displayed in the two figures above lead readily to the following statement about generations.

Corollary 5 (Generation of Consecutive Elements). *For any $\mathcal{T} \in \mathbb{T}$ and $T, T' = F(T) \in \mathcal{T}$ we either have:*

- (a) $g(T) = g(T')$ and T, T' are compatibly divisible, or
- (b) $g(T') = g(T) - 1$ and T is compatibly divisible with a child of T' .

Corollary 6 (Generations within a Chain). *For all $\mathcal{T} \in \mathbb{T}$ and $T \in \mathcal{T}$, its chain $\mathcal{C}(T, \mathcal{T}) = \{T_k\}_{k=0}^m$ with $T_k = F^k(T)$ have the property*

$$g(T_k) = g(T) - k \quad 0 \leq k \leq m - 1$$

and $T_m = F^m(T)$ has generation $g(T_m) = g(T_{m-1})$ or it is a boundary element with lowest labeled edge on $\partial\Omega$. In the first case, T_{m-1} and T_m are compatibly divisible.

Proof. Apply Corollary 5 repeatedly to consecutive elements of $\mathcal{C}(T, \mathcal{T})$. \square

6.2 Recursive Bisection

Given an element $T \in \mathcal{M}$ to be refined, the routine `REFINE_RECURSIVE` (\mathcal{T}, T) recursively refines the chain $\mathcal{C}(T, \mathcal{T})$ of T , from the end back to T , and creates a minimal conforming partition $\mathcal{T}_* \geq \mathcal{T}$ such that T is bisected once. This procedure reads as follows:

```

REFINE_RECURSIVE ( $\mathcal{T}, T$ )
  if  $g(F(T)) < g(T)$ 
     $\mathcal{T} := \text{REFINE\_RECURSIVE} (\mathcal{T}, F(T));$ 
  else
    bisect the compatible bisection patch  $\mathcal{C}(T, \mathcal{T})$ ;
    update  $\mathcal{T}$ ;
  end if
  return ( $\mathcal{T}$ )

```

We denote by $\mathcal{C}_*(T, \mathcal{T}) \subset \mathcal{T}_*$ the recursive refinement of $\mathcal{C}(T, \mathcal{T})$ (or completion of $\mathcal{C}(T, \mathcal{T})$) caused by bisection of T . Since `REFINE_RECURSIVE` refines solely compatible bisection patches, intermediate meshes are always conforming.

We refer to Figure 19 for an example of recursive bisection $\mathcal{C}_*(T_{10}, \mathcal{T})$ of $\mathcal{C}(T_{10}, \mathcal{T}) = \{T_{10}, T_8, T_2\}$ in Figure 2: `REFINE_RECURSIVE` starts bisecting from the end of $\mathcal{C}(T_{10}, \mathcal{T})$, namely T_2 , which is a boundary element, and goes back the chain bisecting elements twice until it gets to T_{10} .

We now establish a fundamental property of `REFINE_RECURSIVE` (\mathcal{T}, T) relating the generation of elements within $\mathcal{C}_*(T, \mathcal{T})$.

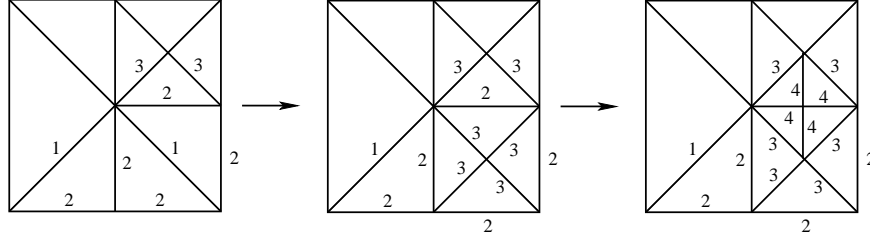


Fig. 19 Recursive refinement of $T_{10} \in \mathcal{T}$ in Figure 2 by `REFINE_RECURSIVE`. This entails refining the chain $\mathcal{C}(T_{10}, \mathcal{T}) = \{T_{10}, T_8, T_2\}$, starting from the last element $T_2 \in \mathcal{T}$, which form alone a compatible bisection patch because its refinement edge is on the boundary, and continuing with $T_8 \in \mathcal{T}$ and finally $T_{10} \in \mathcal{T}$. Note that the successive meshes are always conforming and that `REFINE_RECURSIVE` bisects elements in $\mathcal{C}(T_{10}, \mathcal{T})$ twice before getting back to T_{10} .

Lemma 17 (Recursive Refinement). *Let \mathcal{T}_0 satisfy the labeling (6), and let $\mathcal{T} \in \mathbb{T}$ be a conforming refinement of \mathcal{T}_0 . A call to `REFINE_RECURSIVE` (\mathcal{T}, T) terminates, for all $T \in \mathcal{M}$, and outputs the smallest conforming refinement \mathcal{T}_* of \mathcal{T} such that T is bisected. In addition, all newly created $T' \in \mathcal{C}_*(T, \mathcal{T})$ satisfy*

$$g(T') \leq g(T) + 1. \quad (106)$$

Proof. We first observe that T has maximal generation within $\mathcal{C}(T, \mathcal{T})$. So recursion is applied to elements with generation $\leq g(T)$, whence the recursion terminates. We also note that this procedure creates children of T and either children or grandchildren of triangles $T_k \in \mathcal{C}(T, \mathcal{T}) = \{T_i\}_{i=0}^m$ with $k \geq 1$. If T' is a child of T there is nothing to prove. If not, we consider first $m = 1$, in which case T' is a child of T_1 because T_0 and T_1 are compatibly divisible and so have the same generation; thus $g(T') = g(T_1) + 1 = g(T_0) + 1$. Finally, if $m > 1$, then $g(T_k) < g(T)$ and we apply Corollary 6 to deduce

$$g(T') \leq g(T_k) + 2 \leq g(T) + 1,$$

as asserted. \square

The following crucial lemma links generation and distance between T and $T' \in \mathcal{C}_*(T, \mathcal{T})$, the latter being defined as

$$\text{dist}(T', T) := \inf_{x' \in T', x \in T} |x' - x|.$$

Lemma 18 (Distance and Generation). *Let $T \in \mathcal{M}$. Any newly created $T' \in \mathcal{C}_*(T, \mathcal{T})$ by `REFINE_RECURSIVE` (\mathcal{T}, T) satisfies*

$$\text{dist}(T', T) \leq D_2 \frac{2}{\sqrt{2}-1} 2^{-g(T')/2}, \quad (107)$$

where $D_2 > 0$ is the constant in (5).

Proof. Suppose $T' \subset T_i \in \mathcal{C}(T, \mathcal{T})$ have been created by subdividing T_i (see Figure 18). If $i \leq 1$ then $\text{dist}(T', T) = 0$ and there is nothing to prove. If $i > 1$, then we observe that $\text{dist}(T', T_{i-1}) = 0$, whence

$$\begin{aligned} \text{dist}(T', T) &\leq \text{dist}(T_{i-1}, T) + \text{diam}(T_{i-1}) \leq \sum_{k=1}^{i-1} \text{diam}(T_k) \\ &\leq D_2 \sum_{k=1}^{i-1} 2^{-g(T_k)/2} < D_2 \frac{1}{1-2^{-1/2}} 2^{-g(T_{i-1})/2}, \end{aligned}$$

because the generations decrease exactly by 1 along the chain $\mathcal{C}(T)$ according to Corollary 5(b). Since T' is a child or grandchild of T_i , we deduce

$$g(T') \leq g(T_i) + 2 = g(T_{i-1}) + 1,$$

whence

$$\text{dist}(T', T) < D_2 \frac{2^{1/2}}{1-2^{-1/2}} 2^{-g(T')/2}.$$

This is the desired estimate. \square

The recursive procedure `REFINE_RECURSIVE` is the core of the routine `REFINE` of Section 1.3: given a conforming mesh $\mathcal{T} \in \mathbb{T}$ and a subset $\mathcal{M} \subset \mathcal{T}$ of marked elements, `REFINE` creates a conforming refinement $\mathcal{T}_* \geq \mathcal{T}$ of \mathcal{T} such that all elements of \mathcal{M} are bisected at least once:

```

REFINE ( $\mathcal{T}, \mathcal{M}$ )
for all  $T \in \mathcal{M} \cap \mathcal{T}$  do
   $\mathcal{T} := \text{REFINE\_RECURSIVE}(\mathcal{T}, T)$ ;
end
return ( $\mathcal{T}$ )

```

It may happen that an element $T' \in \mathcal{M}$ is scheduled prior to T for refinement and $T \in \mathcal{C}(T', \mathcal{T})$. Since the call `REFINE_RECURSIVE` (\mathcal{T}, T') bisects T , its two children replace T in \mathcal{T} . This implies that $T \notin \mathcal{M} \cap \mathcal{T}$, which prevents further refinement of T .

In practice, one often likes to bisect selected elements several times, for instance each marked element is scheduled for $b \geq 1$ bisections. This can be done by assigning the number $b(T) = b$ of bisections that have to be executed for each marked element T . If T is bisected then we assign $b(T) - 1$ as the number of pending bisections to its children and the set of marked elements is $\mathcal{M} := \{T \in \mathcal{T} \mid b(T) > 0\}$.

6.3 Conforming Meshes: Proof of Theorem 1

Figure 19 reveals that the issue of propagation of mesh refinement to keep conformity is rather delicate. In particular, an estimate of the form

$$\#\mathcal{T}_k - \#\mathcal{T}_{k-1} \leq \Lambda \#\mathcal{M}_{k-1}$$

is not valid with a constant Λ independent of k ; in fact the constant can be proportional to k according to Figure 19.

Binev, Dahmen, and DeVore [7] for $d = 2$ and Stevenson [53] for $d > 2$ show that control of the propagation of refinement by bisection is possible when considering the collective effect:

$$\#\mathcal{T}_k - \#\mathcal{T}_0 \leq \Lambda_0 \sum_{j=0}^{k-1} \#\mathcal{M}_j.$$

This can be heuristically motivated as follows. Consider the set $\mathcal{M} := \bigcup_{j=0}^{k-1} \mathcal{M}_j$ used to generate the sequence $\mathcal{T}_0 \leq \mathcal{T}_1 \leq \dots \leq \mathcal{T}_k =: \mathcal{T}$. Suppose that each element $T_* \in \mathcal{M}$ is assigned a fixed amount C_1 of money to spend on refined elements in \mathcal{T} , i. e., on $T \in \mathcal{T} \setminus \mathcal{T}_0$. Assume further that $\lambda(T, T_*)$ is the portion of money spent by T_* on T . Then it must hold

$$\sum_{T \in \mathcal{T} \setminus \mathcal{T}_0} \lambda(T, T_*) \leq C_1 \quad \text{for all } T_* \in \mathcal{M}. \quad (108a)$$

In addition, we suppose that the investment of all elements in \mathcal{M} is fair in the sense that each $T \in \mathcal{T} \setminus \mathcal{T}_0$ gets at least a fixed amount C_2 , whence

$$\sum_{T_* \in \mathcal{M}} \lambda(T, T_*) \geq C_2 \quad \text{for all } T \in \mathcal{T} \setminus \mathcal{T}_0. \quad (108b)$$

Therefore, summing up (108b) and using the upper bound (108a) we readily obtain

$$C_2(\#\mathcal{T} - \#\mathcal{T}_0) \leq \sum_{T \in \mathcal{T} \setminus \mathcal{T}_0} \sum_{T_* \in \mathcal{M}} \lambda(T, T_*) = \sum_{T_* \in \mathcal{M}} \sum_{T \in \mathcal{T} \setminus \mathcal{T}_0} \lambda(T, T_*) \leq C_1 \#\mathcal{M},$$

which proves Theorem 1 for \mathcal{T} and \mathcal{M} . In the remainder of this section we design such an allocation function $\lambda: \mathcal{T} \times \mathcal{M} \rightarrow \mathbb{R}^+$ in several steps and prove that recurrent refinement by bisection yields (108) provided \mathcal{T}_0 satisfies (6).

Construction of the Allocation Function. The function $\lambda(T, T_*)$ is defined with the help of two sequences $(a(\ell))_{\ell=-1}^{\infty}, (b(\ell))_{\ell=0}^{\infty} \subset \mathbb{R}^+$ of positive numbers satisfying

$$\sum_{\ell \geq -1} a(\ell) = A < \infty, \quad \sum_{\ell \geq 0} 2^{-\ell/2} b(\ell) = B < \infty, \quad \inf_{\ell \geq 1} b(\ell) a(\ell) = c_* > 0,$$

and $b(0) \geq 1$. Valid instances are $a(\ell) = (\ell + 2)^{-2}$ and $b(\ell) = 2^{\ell/3}$.

With these settings we are prepared to define $\lambda: \mathcal{T} \times \mathcal{M} \rightarrow \mathbb{R}^+$ by

$$\lambda(T, T_*) := \begin{cases} a(g(T_*) - g(T)), & \text{dist}(T, T_*) < D_3 B 2^{-g(T)/d} \text{ and } g(T) \leq g(T_*) + 1 \\ 0, & \text{else,} \end{cases}$$

where $D_3 := D_2(1 + 2(\sqrt{2} - 1)^{-1})$. Therefore, the investment of money by $T_* \in \mathcal{M}$ is restricted to cells T that are sufficiently close and are of generation $g(T) \leq g(T_*) + 1$. Only elements of such generation can be created during refinement of T_* according to Lemma 17. We stress that except for the definition of B , this construction is multidimensional and we refer to [45, 53] for details.

The following lemma shows that the total amount of money spend by the allocation function $\lambda(T, T_*)$ per marked element T_* is bounded.

Lemma 19 (Upper Bound). *There exists a constant $C_1 > 0$ only depending on \mathcal{T}_0 such that λ satisfies (108a), i. e.,*

$$\sum_{T \in \mathcal{T} \setminus \mathcal{T}_0} \lambda(T, T_*) \leq C_1 \quad \text{for all } T_* \in \mathcal{M}.$$

Proof. \square Given $T_* \in \mathcal{M}$ we set $g_* = g(T_*)$ and we let $0 \leq g \leq g_* + 1$ be a generation of interest in the definition of λ . We claim that for such g the cardinality of the set

$$\mathcal{T}(T_*, g) = \{T \in \mathcal{T} \mid \text{dist}(T, T_*) < D_3 B 2^{-g/2} \text{ and } g(T) = g\}$$

is uniformly bounded, i. e., $\#\mathcal{T}(T_*, g) \leq C$ with C solely depending on D_1, D_2, D_3, B .

From (5) we learn that $\text{diam}(T_*) \leq D_2 2^{-g_*/2} \leq 2D_2 2^{-(g_*+1)/2} \leq 2D_2 2^{-g/2}$ as well as $\text{diam}(T) \leq D_2 2^{-g/2}$ for any $T \in \mathcal{T}(T_*, g)$. Hence, all elements of the set $\mathcal{T}(T_*, g)$ lie inside a ball centered at the barycenter of T_* with radius $(D_3 B + 3D_2) 2^{-g/2}$. Again relying on (5) we thus conclude

$$\#\mathcal{T}(T_*, g) D_1 2^{-g} \leq \sum_{T \in \mathcal{T}(T_*, g)} |T| \leq c(D_3 B + 3D_2)^2 2^{-g},$$

whence $\#\mathcal{T}(T_*, g) \leq c D_1^{-1} (D_3 B + 3D_2)^2 =: C$.

\square Accounting only for non-zero contributions $\lambda(T, T_*)$ we deduce

$$\sum_{T \in \mathcal{T} \setminus \mathcal{T}_0} \lambda(T, T_*) = \sum_{g=0}^{g_*+1} \sum_{T \in \mathcal{T}(T_*, g)} a(g_* - g) \leq C \sum_{\ell=-1}^{\infty} a(\ell) = CA =: C_1,$$

which is the desired upper bound. \square

The definition of λ also implies that each refined element receives a fixed amount of money. We show this next.

Lemma 20 (Lower Bound). *There exists a constant $C_2 > 0$ only depending on \mathcal{T}_0 such that λ satisfies (108b), i. e.,*

$$\sum_{T_* \in \mathcal{M}} \lambda(T, T_*) \geq C_2 \quad \text{for all } T \in \mathcal{T} \setminus \mathcal{T}_0.$$

Proof. \square Fix an arbitrary $T_0 \in \mathcal{T} \setminus \mathcal{T}_0$. Then there is an iteration count $1 \leq k_0 \leq k$ such that $T_0 \in \mathcal{T}_{k_0}$ and $T_0 \notin \mathcal{T}_{k_0-1}$. Therefore there exists an $T_1 \in \mathcal{M}_{k_0-1} \subset \mathcal{M}$

such that T_0 is generated during `REFINE_RECURSIVE` (\mathcal{T}_{k_0-1}, T_1). Iterating this process we construct a sequence $\{T_j\}_{j=1}^J \subset \mathcal{M}$ with corresponding iteration counts $\{k_j\}_{j=1}^J$ such that T_j is created by `REFINE_RECURSIVE` ($\mathcal{T}_{k_j-1}, T_{j+1}$). The sequence is finite since the iteration counts are strictly decreasing and thus $k_J = 0$ for some $J > 0$, or equivalently $T_J \in \mathcal{T}_0$.

Since T_j is created during refinement of T_{j+1} we infer from (106) that

$$g(T_{j+1}) \geq g(T_j) - 1.$$

Accordingly, $g(T_{j+1})$ can decrease the previous value of $g(T_j)$ at most by 1. Since $g(T_j) = 0$ there exists a smallest value s such that $g(T_s) = g(T_0) - 1$. Note that for $j = 1, \dots, s$ we have $\lambda(T_0, T_j) > 0$ if $\text{dist}(T_0, T_j) \leq D_3 B g^{-g(T_0)/d}$.

\square We next estimate the distance $\text{dist}(T_0, T_j)$. For $1 \leq j \leq s$ and $\ell \geq 0$ we define the set

$$\mathcal{T}(T_0, \ell, j) := \{T \in \{T_0, \dots, T_{j-1}\} \mid g(T) = g(T_0) + \ell\}$$

and denote by $m(\ell, j)$ its cardinality. The triangle inequality combined with an induction argument yields

$$\begin{aligned} \text{dist}(T_0, T_j) &\leq \text{dist}(T_0, T_1) + \text{diam}(T_1) + \text{dist}(T_1, T_j) \\ &\leq \sum_{i=1}^j \text{dist}(T_{i-1}, T_i) + \sum_{i=1}^{j-1} \text{diam}(T_i). \end{aligned}$$

We apply (107) for the terms of the first sum and (5) for the terms of the second sum to obtain

$$\begin{aligned} \text{dist}(T_0, T_j) &< D_2 \frac{2}{\sqrt{2}-1} \sum_{i=1}^j 2^{-g(T_{i-1})/2} + D_2 \sum_{i=1}^{j-1} 2^{-g(T_i)/2} \\ &\leq D_2 \left(1 + \frac{2}{\sqrt{2}-1}\right) \sum_{i=0}^{j-1} 2^{-g(T_i)/2} \\ &= D_3 \sum_{\ell=0}^{\infty} m(\ell, j) 2^{-(g(T_0)+\ell)/2} \\ &= D_3 2^{-g(T_0)/2} \sum_{\ell=0}^{\infty} m(\ell, j) 2^{-\ell/2}. \end{aligned}$$

For establishing the lower bound we distinguish two cases depending on the size of $m(\ell, s)$. This is done next.

\square *Case 1:* $m(\ell, s) \leq b(\ell)$ for all $\ell \geq 0$. From this we conclude

$$\text{dist}(T_0, T_s) < D_3 2^{-g(T_0)/2} \sum_{\ell=0}^{\infty} b(\ell) 2^{-\ell/2} = D_3 B 2^{-g(T_0)/2}$$

and the definition of λ then readily implies

$$\sum_{T_* \in \mathcal{M}} \lambda(T_0, T_*) \geq \lambda(T_0, T_s) = a(g(T_s) - g(T_0)) = a(-1) > 0.$$

□ *Case 2:* There exists $\ell \geq 0$ such that $m(\ell, s) > b(\ell)$. For each of these ℓ 's there exists a smallest $j = j(\ell)$ such that $m(\ell, j(\ell)) > b(\ell)$. We let ℓ^* be the index ℓ that gives rise to the smallest $j(\ell)$, and set $j^* = j(\ell^*)$. Consequently

$$m(\ell, j^* - 1) \leq b(\ell) \quad \text{for all } \ell \geq 0 \quad \text{and} \quad m(\ell^*, j^*) > b(\ell^*).$$

As in Case 1 we see $\text{dist}(T_0, T_i) < D_3 B 2^{-g(T_0)/2}$ for all $i \leq j^* - 1$, or equivalently

$$\text{dist}(T_0, T_i) < D_3 B 2^{-g(T_0)/2} \quad \text{for all } T_i \in \mathcal{T}(T_0, \ell, j^*).$$

We next show that the elements in $\mathcal{T}(T_0, \ell^*, j^*)$ spend enough money on T_0 . We first consider $\ell^* = 0$ and note that $T_0 \in \mathcal{T}(T_0, 0, j^*)$. Since $m(0, j^*) > b(0) \geq 1$ we discover $j^* \geq 2$. Hence, there is an $T_i \in \mathcal{T}(T_0, 0, j^*) \cap \mathcal{M}$, which yields the estimate

$$\sum_{T_* \in \mathcal{M}} \lambda(T_0, T_*) \geq \lambda(T_0, T_i) = a(g(T_i) - g(T_0)) = a(0) > 0.$$

For $\ell^* > 0$ we see that $T_0 \notin \mathcal{T}(T_0, \ell^*, j^*)$, whence $\mathcal{T}(T_0, \ell^*, j^*) \subset \mathcal{M}$. In addition, $\lambda(T_0, T_i) = a(\ell^*)$ for all $T_i \in \mathcal{T}(T_0, \ell^*, j^*)$. From this we conclude

$$\begin{aligned} \sum_{T_* \in \mathcal{M}} \lambda(T_0, T_*) &\geq \sum_{T_* \in \mathcal{T}(T_0, \ell^*, j^*)} \lambda(T_0, T_*) = m(\ell^*, j^*) a(\ell^*) \\ &> b(\ell^*) a(\ell^*) \geq \inf_{\ell \geq 1} b(\ell) a(\ell) = c_* > 0. \end{aligned}$$

□ In summary we have proved the assertion since for any $T_0 \in \mathcal{T} \setminus \mathcal{T}_0$

$$\sum_{T_* \in \mathcal{M}} \lambda(T_0, T_*) \geq \min\{a(-1), a(0), c_*\} =: C_2 > 0. \quad \square$$

Remark 11 (Complexity with $b > 1$ Bisections). To show the complexity estimate when REFINE performs $b > 1$ bisections, the set \mathcal{M}_k is to be understood as a sequence of *single* bisections recorded in sets $\{\mathcal{M}_k(j)\}_{j=1}^b$, which belong to intermediate triangulations between \mathcal{T}_k and \mathcal{T}_{k+1} with $\#\mathcal{M}_k(j) \leq 2^{j-1} \#\mathcal{M}_k$, $j = 1, \dots, b$. Then we also obtain Theorem 1 because

$$\sum_{j=1}^b \#\mathcal{M}_k(j) \leq \sum_{j=1}^b 2^{j-1} \#\mathcal{M}_k = (2^b - 1) \#\mathcal{M}_k.$$

In practice, it is customary to take $b = d$ [50].

6.4 Nonconforming Meshes: Proof of Lemma 3

We now examine briefly the refinement process for quadrilaterals with one hanging node per edge, which gives rise to the so-called *1-meshes*. The refinement of $T \in \mathcal{T}$ might affect four elements of \mathcal{T} for $d = 2$ (or 2^d elements for any dimension $d \geq 2$), all contained in the *refinement patch* $R(T, \mathcal{T})$ of T in \mathcal{T} . The latter is defined as

$$R(T, \mathcal{T}) := \{T' \in \mathcal{T} \mid T' \text{ and } T \text{ share an edge and } g(T') \leq g(T)\},$$

and is called *compatible* provided $g(T') = g(T)$ for all $T' \in R(T, \mathcal{T})$. The generation gap between elements sharing an edge, in particular those in $R(T, \mathcal{T})$, is always ≤ 1 for 1-meshes, and is 0 if $R(T, \mathcal{T})$ is compatible. The element size satisfies

$$h_T = 2^{-g(T)} h_{T_0} \quad \forall T \in \mathcal{T}$$

where $T_0 \in \mathcal{T}_0$ is the ancestor of T in the initial mesh \mathcal{T}_0 . Lemma 2 is thus valid

$$h_T < \bar{h}_T \leq D_2 2^{-g(T)} \quad \forall T \in \mathcal{T}. \quad (109)$$

Given an element $T \in \mathcal{M}$ to be refined, the routine `REFINE_RECURSIVE` (\mathcal{T}, T) refines recursively $R(T, \mathcal{T})$ in such a way that the intermediate meshes are always 1-meshes, and reads as follows:

```

REFINE_RECURSIVE ( $\mathcal{T}, T$ )
if  $g = \min\{g(T'') : T'' \in R(T, \mathcal{T})\} < g(T)$ 
  let  $T' \in R(T, \mathcal{T})$  satisfy  $g(T') = g$ 
   $\mathcal{T} := \text{REFINE\_RECURSIVE}(\mathcal{T}, T')$ ;
else
  subdivide  $T$ ;
  update  $\mathcal{T}$  upon replacing  $T$  by its children;
end if
return ( $\mathcal{T}$ )

```

The conditional prevents the generation gap within $R(T, \mathcal{T})$ from getting larger than 1. If it fails, then the refinement patch $R(T, \mathcal{T})$ is compatible and refining T increases the generation gap from 0 to 1 without violating the 1-mesh structure. This implies Lemma 17: for all newly created elements $T' \in \mathcal{T}_*$

$$g(T') \leq g(T) + 1. \quad (110)$$

In addition, `REFINE_RECURSIVE` (\mathcal{T}, T) creates a minimal 1-mesh $\mathcal{T}_* \geq \mathcal{T}$ refinement of \mathcal{T} so that T is subdivided only *once*. This yields Lemma 18: there exist a geometric constant $D > 0$ such that for all newly created elements $T' \in \mathcal{T}_*$

$$\text{dist}(T, T') \leq D 2^{g(T')}. \quad (111)$$

The procedure `REFINE_RECURSIVE` is the core of `REFINE`, which is conceptually identical to that in Section 6.2. Suppose that each marked element $T \in \mathcal{M}$ is to be subdivided $\rho \geq 1$ times. We assign a flag $q(T)$ to each element T which is initialized $q(T) = \rho$ if $T \in \mathcal{M}$ and $q(T) = 0$ otherwise. The marked set \mathcal{M} is then the set of elements T with $q(T) > 0$, and every time T is subdivided it is removed from \mathcal{T} and replaced by its children, which inherit the flag $q(T) - 1$. This avoids the conflict of subdividing again an element that has been previously refined by `REFINE_RECURSIVE`. The procedure `REFINE` (\mathcal{T}, \mathcal{M}) reads

```

REFINE ( $\mathcal{T}, \mathcal{M}$ )
for all  $T \in \mathcal{M} \cap \mathcal{T}$  do
     $\mathcal{T} := \text{REFINE\_RECURSIVE} (\mathcal{T}, T)$ ;
end
return ( $\mathcal{T}$ )

```

and its output is a minimal 1-mesh $\mathcal{T}_* \geq \mathcal{T}$, refinement of \mathcal{T} , so that all marked elements of \mathcal{M} are refined at least ρ times. Since \mathcal{T}_* has one hanging node per side it is thus admissible in the sense (22). However, the refinement may spread outside \mathcal{M} and the issue of complexity of `REFINE` again becomes non-trivial.

With the above ingredients in place, the proof of Lemma 3 follows along the lines of Section 6.3; see Problem 50.

6.5 Notes

The complexity theory for bisection hinges on the initial labeling (6) for $d = 2$. That such a labeling exists is due to Mitchell [39, Theorem 2.9] and Binev, Dahmen, and DeVore [7, Lemma 2.1], but the proofs are not constructive. A couple of global bisections of \mathcal{T}_0 , as depicted in Figure 6, guarantee (6) over the ensuing mesh. For $d > 2$ the corresponding initial labeling is due to Stevenson [53, Section 4 - Condition (b)], who in turn improves upon Maubach [36] and Traxler [54] and shows how to impose it upon further refining each element of \mathcal{T}_0 . We refer to the survey [45] for a discussion of this condition: a key consequence is that every uniform refinement of \mathcal{T}_0 gives a conforming bisection mesh.

The fundamental properties of chains, especially Lemmas 17 and 18, along with the clever ideas of Section 6.3 are due to Binev, Dahmen, and DeVore [7] for $d = 2$, and Stevenson for $d > 2$; see [45]. Bonito and Nochetto [9] observed, in the context of dG methods, that such properties extend to admissible nonconforming meshes.

6.6 Problems

Problem 47 (Largest number of bisections). Show that `REFINE_RECURSIVE` (\mathcal{T}, T) for $d = 2$ bisects T exactly once and all the elements in the chain $\mathcal{C}(T, \mathcal{T})$ at

most twice. This property extends to $d > 2$ provided the initial labeling of Stevenson [53, Section 4 - Condition (b)] is enforced.

Problem 48 (Properties of generation for quad-refinement). Prove (110) and (111).

Problem 49 (Largest number of subdivisions for quads). Show that the procedure REFINE_RECURSIVE (\mathcal{T}, T) subdivides T exactly once and never subdivides any other quadrilateral of \mathcal{T} more than once.

Problem 50 (Lemma 3). Combine (110) and (111) to prove Lemma 3 for any $\rho \geq 1$.

7 Convergence Rates

We have already realized in §1.6 that we can a priori accommodate the degrees of freedom in such a way that the finite element approximation retains optimal energy error decay for a class of singular functions. This presumes knowledge of the exact solution u . At the same time, we have seen numerical evidence in §5.4 that the standard AFEM of §5.1, achieves such a performance without direct access to the exact solution u . Practical experience strongly suggests that this is even true for a much larger class of problems and adaptive methods. The challenge ahead is to reconcile these two distinct aspects of AFEM.

A crucial insight in such a connection for the simplest scenario, the Laplacian and piecewise constant forcing f , is due to Stevenson [52]:

any marking strategy that reduces the energy error relative to the current value must contain a substantial portion of $\mathcal{E}_{\mathcal{T}}(U)$, and so it can be related to Dörfler Marking. (112)

This allows one to compare meshes produced by AFEM with optimal ones and to conclude a quasi-optimal error decay. We discuss this issue in §7.3. However, this is not enough to handle the model problem (87) with variable \mathbf{A} and f .

The objective of this section is to study (87) for general data \mathbf{A} and f . This study hinges on the total error and its relation with the quasi error, which is contracted by AFEM. This approach allows us to improve upon and extend Stevenson [52] to variable data. In doing so, we follow closely Cascón, Kreuzer, Nochetto, and Siebert [14]. The present theory, however, does not extend to noncoercive problems and marking strategies other than Dörfler's. These remain important open questions.

As in §5, u will always be the weak solution of (87) and, except when stated otherwise, any explicit constant or hidden constant in \lesssim may depend on the uniform shape-regularity of \mathbb{T} , the dimension d , the polynomial degree n , the (global) eigenvalues of \mathbf{A} , and the oscillation $\text{osc}_{\mathcal{T}_0}(\mathbf{A})$ of \mathbf{A} on the initial mesh \mathcal{T}_0 , but not on a specific grid $\mathcal{T} \in \mathbb{T}$.

7.1 The Total Error

We first introduce the concept of *total error* for the Galerkin function $U \in \mathbb{V}(\mathcal{T})$

$$\|u - U\|_{\Omega}^2 + \text{osc}_{\mathcal{T}}^2(U), \quad (113)$$

(see Mekchay and Nochetto [37]), and next assert its equivalence to the quasi error (99). In fact, in view of the upper and lower a posteriori error bounds (94), and

$$\text{osc}_{\mathcal{T}}^2(U) \leq \mathcal{E}_{\mathcal{T}}^2(U),$$

we have

$$\begin{aligned} C_2 \mathcal{E}_{\mathcal{T}}^2(U) &\leq \|u - U\|_{\Omega}^2 + \text{osc}_{\mathcal{T}}^2(U) \\ &\leq \|u - U\|_{\Omega}^2 + \mathcal{E}_{\mathcal{T}}^2(U) \leq (1 + C_1) \mathcal{E}_{\mathcal{T}}^2(U), \end{aligned}$$

whence

$$\mathcal{E}_{\mathcal{T}}^2(U) \approx \|u - U\|_{\Omega}^2 + \text{osc}_{\mathcal{T}}^2(U). \quad (114)$$

Since AFEM selects elements for refinement based on information extracted exclusively from the error indicators $\{\mathcal{E}_{\mathcal{T}}(U, T)\}_{T \in \mathcal{T}}$, we realize that the decay rate of AFEM must be characterized by the total error. Moreover, on invoking the upper bound (94a) again, we also see that the total error is equivalent to the quasi error

$$\|u - U\|_{\Omega}^2 + \text{osc}_{\mathcal{T}}^2(U) \approx \|u - U\|_{\Omega}^2 + \mathcal{E}_{\mathcal{T}}^2(U).$$

The latter is the quantity being strictly reduced by AFEM (Theorem 9). Finally, the total error satisfies the following Cea's type-lemma, or equivalently AFEM is quasi-optimal regarding the total error. In fact, if the oscillation vanishes, then this is Cea's Lemma stated in Theorem 4; see also Problem 12.

Lemma 21 (Quasi-optimality of total error). *There exists an explicit constant Λ_2 , which depends on \mathbf{A} , \mathcal{T}_0 , n and d , such that for any $\mathcal{T} \in \mathbb{T}$ and the corresponding Galerkin solution $U \in \mathbb{V}(\mathcal{T})$ there holds*

$$\|u - U\|_{\Omega}^2 + \text{osc}_{\mathcal{T}}^2(U) \leq \Lambda_2 \inf_{V \in \mathbb{V}(\mathcal{T})} \left(\|u - V\|_{\Omega}^2 + \text{osc}_{\mathcal{T}}^2(V) \right).$$

Proof. For $\varepsilon > 0$ choose $V_{\varepsilon} \in \mathbb{V}(\mathcal{T})$, with

$$\|u - V_{\varepsilon}\|_{\Omega}^2 + \text{osc}_{\mathcal{T}}^2(V_{\varepsilon}) \leq (1 + \varepsilon) \inf_{V \in \mathbb{V}(\mathcal{T})} \left(\|u - V\|_{\Omega}^2 + \text{osc}_{\mathcal{T}}^2(V) \right).$$

Applying Problem 46 with $\mathcal{T}_* = \mathcal{T}$, $V = U$, and $V_* = V_{\varepsilon}$ yields

$$\text{osc}_{\mathcal{T}}^2(U) \leq 2 \text{osc}_{\mathcal{T}}^2(V_{\varepsilon}) + C_3 \|U - V_{\varepsilon}\|_{\Omega}^2,$$

with

$$C_3 := \Lambda_1 \text{osc}_{\mathcal{T}_0}(\mathbf{A})^2.$$

Since $U \in \mathbb{V}(\mathcal{T})$ is the Galerkin solution, $U - V_{\varepsilon} \in \mathbb{V}(\mathcal{T})$ is orthogonal to $u - U$ in the energy norm, whence $\|u - U\|_{\Omega}^2 + \|U - V_{\varepsilon}\|_{\Omega}^2 = \|u - V_{\varepsilon}\|_{\Omega}^2$ and

$$\begin{aligned} \|u - U\|_{\Omega}^2 + \text{osc}_{\mathcal{T}}^2(U) &\leq (1 + C_3) \|u - V_{\varepsilon}\|_{\Omega}^2 + 2 \text{osc}_{\mathcal{T}}^2(V_{\varepsilon}) \\ &\leq (1 + \varepsilon) \Lambda_2 \inf_{V \in \mathbb{V}(\mathcal{T})} \left(\|u - U\|_{\Omega}^2 + \text{osc}_{\mathcal{T}}^2(V) \right), \end{aligned}$$

with $\Lambda_2 = \max\{2, 1 + C_3\}$. The assertion follows upon taking $\varepsilon \rightarrow 0$. \square

7.2 Approximation Classes

In view of (114) and Lemma 21, the definition of approximation class \mathbb{A}_s hinges on the concept of best total error:

$$\inf_{V \in \mathbb{V}(\mathcal{T})} \left(\|u - V\|_{\Omega}^2 + \text{osc}_{\mathcal{T}}^2(V) \right).$$

We first let $\mathbb{T}_N \subset \mathbb{T}$ be the set of all possible conforming refinements of \mathcal{T}_0 with at most N elements more than \mathcal{T}_0 , i. e.,

$$\mathbb{T}_N = \{ \mathcal{T} \in \mathbb{T} \mid \#\mathcal{T} - \#\mathcal{T}_0 \leq N \}.$$

The quality of the best approximation in \mathbb{T}_N with respect to the total error is characterized by

$$\sigma(N; u, f, \mathbf{A}) := \inf_{\mathcal{T} \in \mathbb{T}_N} \inf_{V \in \mathbb{V}(\mathcal{T})} \left(\|u - V\|_{\Omega}^2 + \text{osc}_{\mathcal{T}}^2(V) \right)^{1/2},$$

and the approximation class \mathbb{A}_s for $s > 0$ is defined by

$$\mathbb{A}_s := \left\{ (v, f, \mathbf{A}) \mid |v, f, \mathbf{A}|_s := \sup_{N > 0} (N^s \sigma(N; v, f, \mathbf{A})) < \infty \right\}.$$

Therefore, if $(v, f, \mathbf{A}) \in \mathbb{A}_s$, then $\sigma(N; v, f, \mathbf{A}) \lesssim N^{-s}$ decays with rate N^{-s} . We point out the upper bound $s \leq n/d$ for polynomial degree $n \geq 1$; this can be seen with full regularity and uniform refinement (see (14)). Note that if $(v, f, \mathbf{A}) \in \mathbb{A}_s$ then for all $\varepsilon > 0$ there exist $\mathcal{T}_{\varepsilon} \geq \mathcal{T}_0$ conforming and $V_{\varepsilon} \in \mathbb{V}(\mathcal{T}_{\varepsilon})$ such that (see Problem 51)

$$\|v - V_{\varepsilon}\|_{\Omega}^2 + \text{osc}_{\mathcal{T}_{\varepsilon}}^2(V_{\varepsilon}) \leq \varepsilon^2 \quad \text{and} \quad \#\mathcal{T}_{\varepsilon} - \#\mathcal{T}_0 \leq |v, f, \mathbf{A}|_s^{1/s} \varepsilon^{-1/s}. \quad (115)$$

In addition, thanks to Lemma 21, the solution u with data (f, \mathbf{A}) satisfies

$$\sigma(N; u, f, \mathbf{A}) \approx \inf_{\mathcal{T} \in \mathbb{T}_N} \left\{ \mathcal{E}_{\mathcal{T}}(U, \mathcal{T}) \mid U = \text{SOLVE}(\mathbb{V}(\mathcal{T})) \right\}. \quad (116)$$

This establishes a direct connection between \mathbb{A}_s and AFEM.

Mesh Overlay. For the subsequent discussion it will be convenient to merge two conforming meshes $\mathcal{T}_1, \mathcal{T}_2 \in \mathbb{T}$. Given the corresponding forests $\mathcal{F}_1, \mathcal{F}_2 \in \mathbb{F}$ we consider the set $\mathcal{F}_1 \cup \mathcal{F}_2 \in \mathbb{F}$, which satisfies $\mathcal{T}_0 \subset \mathcal{F}_1 \cup \mathcal{F}_2$. Then $\mathcal{F}_1 \cup \mathcal{F}_2$ is a forest and its leaves are called the *overlay* of \mathcal{F}_1 and \mathcal{F}_2 :

$$\mathcal{T}_1 \oplus \mathcal{T}_2 = \mathcal{T}(\mathcal{F}_1 \cup \mathcal{F}_2).$$

We next bound the cardinality of $\mathcal{T}_1 \oplus \mathcal{T}_2$ in terms of that of \mathcal{T}_1 and \mathcal{T}_2 ; see [14, 52].

Lemma 22 (Overlay). *The overlay $\mathcal{T} = \mathcal{T}_1 \oplus \mathcal{T}_2$ is conforming and*

$$\#\mathcal{T} \leq \#\mathcal{T}_1 + \#\mathcal{T}_2 - \#\mathcal{T}_0. \quad (117)$$

Proof. See Problem 52. \square

Discussion of \mathbb{A}_s . We now would like to show a few examples of membership in \mathbb{A}_s and highlight some important open questions. We first investigate the class \mathbb{A}_s for piecewise polynomial coefficient matrix \mathbf{A} of degree $\leq n$ over \mathcal{T}_0 . In this simplified scenario, the oscillation $\text{osc}_{\mathcal{T}}(U)$ reduces to *data oscillation* (see (58) and (93)):

$$\text{osc}_{\mathcal{T}}(U) = \text{osc}_{\mathcal{T}}(f) := \|h(f - P_{2n-2}f)\|_{L^2(\Omega)}.$$

We then have the following characterization of \mathbb{A}_s in terms of the standard approximation classes [7, 8, 52]:

$$\begin{aligned} \mathcal{A}_s &:= \left\{ v \in \mathbb{V} \mid |v|_{\mathcal{A}_s} := \sup_{N>0} \left(N^s \inf_{\mathcal{T} \in \mathbb{T}_N} \inf_{V \in \mathbb{V}(\mathcal{T})} \|v - V\|_{\Omega} \right) < \infty \right\}, \\ \mathcal{A}_s^{\vec{}} &:= \left\{ g \in L^2(\Omega) \mid |g|_{\mathcal{A}_s^{\vec{}}} := \sup_{N>0} \left(N^s \inf_{\mathcal{T} \in \mathbb{T}_N} \text{osc}_{\mathcal{T}}(g) \right) < \infty \right\}. \end{aligned}$$

Lemma 23 (Equivalence of classes). *Let \mathbf{A} be piecewise polynomial of degree $\leq n$ over \mathcal{T}_0 . Then $(u, f, \mathbf{A}) \in \mathbb{A}_s$ if and only if $(u, f) \in \mathcal{A}_s \times \mathcal{A}_s^{\vec{}}$ and*

$$|u, f, \mathbf{A}|_s \approx |u|_{\mathcal{A}_s} + |f|_{\mathcal{A}_s^{\vec{}}}. \quad (118)$$

Proof. It is obvious that $(u, f, \mathbf{A}) \in \mathbb{A}_s$ implies $(u, f) \in \mathcal{A}_s \times \mathcal{A}_s^{\vec{}}$ as well as the bound $|u|_{\mathcal{A}_s} + |f|_{\mathcal{A}_s^{\vec{}}} \lesssim |u, f, \mathbf{A}|_s$.

In order to prove the reverse inequality, let $(u, f) \in \mathcal{A}_s \times \mathcal{A}_s^{\vec{}}$. Then there exist $\mathcal{T}_1, \mathcal{T}_2 \in \mathbb{T}_N$ so that $\|u - U_{\mathcal{T}_1}\|_{\Omega} \leq |u|_{\mathcal{A}_s} N^{-s}$ where $U_{\mathcal{T}_1} \in \mathbb{V}(\mathcal{T}_1)$ is the best approximation and $\text{osc}_{\mathcal{T}_2}(f, \mathcal{T}_2) \leq |f|_{\mathcal{A}_s^{\vec{}}} N^{-s}$.

The overlay $\mathcal{T} = \mathcal{T}_1 \oplus \mathcal{T}_2 \in \mathbb{T}_{2N}$ according to (117), and

$$\|u - U_{\mathcal{T}}\|_{\Omega} + \text{osc}_{\mathcal{T}}(f) \leq \|u - U_{\mathcal{T}_1}\|_{\Omega} + \text{osc}_{\mathcal{T}_2}(f) \leq 2^s (|u|_{\mathcal{A}_s} + |f|_{\mathcal{A}_s^{\vec{}}}) (2N)^{-s}.$$

This yields $(u, f, \mathbf{A}) \in \mathbb{A}_s$ together with the bound $|u, f, \mathbf{A}|_s \lesssim |u|_{\mathcal{A}_s} + |f|_{\mathcal{A}_s^{\vec{}}}$. \square

Corollary 7 (Membership in $\mathbb{A}_{1/2}$ with piecewise linear \mathbf{A}). *Let $d = 2$, $n = 1$, and $u \in H_0^1(\Omega)$ be the solution of the model problem with piecewise linear \mathbf{A} and $f \in L^2(\Omega)$. If $u \in W_p^2(\Omega; \mathcal{T}_0)$ is piecewise W_p^2 over the initial grid \mathcal{T}_0 and $p > 1$, then $(u, f, \mathbf{A}) \in \mathbb{A}_{1/2}$ and*

$$|u, f, \mathbf{A}|_{1/2} \lesssim \|D^2 u\|_{L^p(\Omega; \mathcal{T}_0)} + \|f\|_{L^2(\Omega)}. \quad (119)$$

Proof. Since $f \in L^2(\Omega)$, we realize that for all uniform refinements $\mathcal{T} \in \mathbb{T}$ we have

$$\text{osc}_{\mathcal{T}}(f) = \|h(f - P_0 f)\|_{L^2(\Omega)} \leq h_{\max}(\mathcal{T}) \|f\|_{L^2(\Omega)} \lesssim (\#\mathcal{T})^{-1/2} \|f\|_{L^2(\Omega)},$$

This implies $f \in \mathcal{A}_{1/2}^{\vec{}}$ with $|f|_{\mathcal{A}_{1/2}^{\vec{}}} \lesssim \|f\|_{L^2(\Omega)}$. Moreover, for $u \in W_p^2(\Omega; \mathcal{T}_0)$ we learn from Corollary 2 and Remark 6 of §1.6 that $u \in \mathcal{A}_{1/2}$ and $|u|_{\mathcal{A}_{1/2}} \lesssim \|D^2 u\|_{L^2(\Omega; \mathcal{T}_0)}$. The assertion then follows from Lemma 23. \square

Corollary 8 (Membership in $\mathbb{A}_{1/2}$ with variable \mathbf{A}). *Let $d = 2$, $n = 1$, $p > 1$, $f \in L^2(\Omega)$. Let $\mathbf{A} \in W_\infty^1(\Omega, \mathcal{T}_0)$ be piecewise Lipschitz and $u \in W_p^2(\Omega; \mathcal{T}_0) \cap H_0^1(\Omega)$ be piecewise W_p^2 over the initial mesh \mathcal{T}_0 . Then $(u, f, \mathbf{A}) \in \mathbb{A}_{1/2}$ and*

$$\|u, f, \mathbf{A}\|_{1/2} \lesssim \|D^2 u\|_{L^p(\Omega; \mathcal{T}_0)} + \|f\|_{L^2(\Omega)} + \|\mathbf{A}\|_{W_\infty^1(\Omega; \mathcal{T}_0)}. \quad (120)$$

Proof. Combine Problem 55 with Corollary 2. \square

Corollary 9 (Membership in \mathbb{A}_s with $s < 1/d$). *Let $d \geq 2$, $n = 1$, $1 < t < 2$, $p > 1$, and $f \in L^2(\Omega)$. Let $\mathbf{A} \in W_\infty^1(\Omega, \mathcal{T}_0)$ be piecewise Lipschitz and $u \in W_p^t(\Omega; \mathcal{T}_0) \cap H_0^1(\Omega)$ be piecewise W_p^t over the initial mesh \mathcal{T}_0 with $t - \frac{d}{p} > 1 - \frac{d}{2}$. Then $(u, f, \mathbf{A}) \in \mathbb{A}_{(t-1)/d}$ and*

$$\|u, f, \mathbf{A}\|_{(t-1)/d} \lesssim \|D^t u\|_{L^p(\Omega; \mathcal{T}_0)} + \|f\|_{L^2(\Omega)} + \|\mathbf{A}\|_{W_\infty^1(\Omega; \mathcal{T}_0)}. \quad (121)$$

Proof. Combine Problem 9 with Problem 55. \square

Example 2 (Pre-asymptotics). Corollary 7 shows that oscillation decays with rate $1/2$ for $f \in L^2(\Omega)$. Since the decay rate of the total error is $s \leq 1/2$, oscillation can be ignored asymptotically; this is verified in Problems 56, 57, and 58. However, oscillation may dominate the total error, or equivalently the class \mathbb{A}_s may fail to describe the behavior of $\|u - U_k\|_\Omega$, in the early stages of adaptivity. In fact, we recall from Problem 32 that the discrete solution $U_k = 0$, and $\|u - U_k\|_\Omega \approx 2^{-K}$ is constant for as many steps $k \leq K$ as desired. In contrast, $\mathcal{E}_k(U_k) = \text{osc}_k(U_k) = \|h(f - \bar{f})\|_{L^2(\Omega)} = \|hf\|_{L^2(\Omega)}$ reduces strictly for $k \leq K$ but overestimates $\|u - U_k\|_\Omega$. The fact that the preasymptotic regime $k \leq K$ for the energy error could be made arbitrarily long would be problematic if we were to focus exclusively on $\|u - U_k\|_\Omega$.

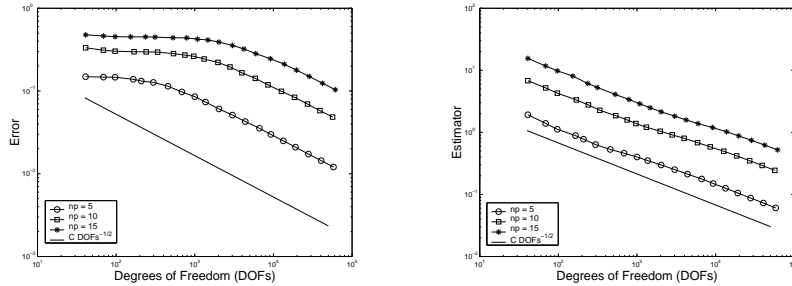


Fig. 20 Decay of the energy error (left) and the estimator (right) for the smooth solution u_S of (122) with frequencies $\kappa = 5, 10$, and 15 . The energy error exhibits a frequency-dependent plateau in the preasymptotic regime and later an optimal decay. This behavior is allowed by \mathbb{A}_s .

In practice, this effect is typically less dramatic because f is not orthogonal to $\mathbb{V}(\mathcal{T}_k)$. Figure 20 displays the behavior of the AFEM for the smooth solution u_S

$$u_S(x, y) = 10^{-2} a_i^{-1} (x^2 + y^2) \sin^2(\kappa\pi x) \sin^2(\kappa\pi y), \quad 1 \leq i \leq 4, \quad (122)$$

of the model problem (87) with discontinuous coefficients $\{a_i\}_{i=1}^4$ in checkerboard pattern as in §5.4 and frequencies $\kappa = 5, 10$, and 15 . We can see that the error exhibits a frequency-dependent plateau in the preasymptotic regime and later an optimal decay. In contrast, the estimator decays always with the optimal rate. Since all decisions of the AFEM are based on the estimator, this behavior has to be expected and is consistent with our notion of approximation class \mathbb{A}_s , which can be characterized just by the estimator according to (116).

7.3 Quasi-Optimal Cardinality: Vanishing Oscillation

In this section we follow the ideas of Stevenson [52] for the simplest scenario with vanishing oscillation $\text{osc}_{\mathcal{T}}(U) = 0$, and thereby explore the insight (112). We recall that in this case the a posteriori error estimates (94) become

$$C_2 \mathcal{E}_{\mathcal{T}}^2(U) \leq \| \|u - U\|_{\Omega}^2 \leq C_1 \mathcal{E}_{\mathcal{T}}^2(U). \quad (123)$$

It is then evident that the ratio $C_2/C_1 \leq 1$, between the *reliability* constant C_1 and the *efficiency* constant C_2 , is a quality measure of the estimator $\mathcal{E}_{\mathcal{T}}(U)$: the closer to 1 the better! This ratio is usually closer to 1 for non-residual estimators, such as those discussed in §5.5, but their theory is a bit more cumbersome.

Assumptions for Optimal Decay Rate. The following are further restrictions on AFEM to achieve optimal error decay, as predicted by the approximation class \mathcal{A}_s .

Assumption 1 (Marking parameter: vanishing oscillation). *The parameter θ of Dörfler marking satisfies $\theta \in (0, \theta_*)$ with*

$$\theta_* := \sqrt{\frac{C_2}{C_1}} \quad (124)$$

Assumption 2 (Cardinality of \mathcal{M}). *MARK selects a set \mathcal{M} with minimal cardinality.*

Assumption 3 (Initial labeling). *The labeling of the initial mesh \mathcal{T}_0 satisfies (6) for $d = 2$ or its multidimensional counterpart for $d > 2$ [52, 45].*

A few comments about these assumptions are now in order.

Remark 12 (Threshold $\theta_ < 1$).* It is reasonable to be cautious in making marking decisions if the constants C_1 and C_2 are very disparate, and thus the ratio C_2/C_1 is far from 1. This justifies the upper bound $\theta_* \leq 1$ in Assumption 1.

Remark 13 (Minimal \mathcal{M}). According to the equidistribution principle (16) and the local lower bound (68) without oscillation, it is natural to mark elements with largest error indicators. This leads to a minimal set \mathcal{M} , as stated in Assumption 2, and turns out to be crucial to link AFEM with optimal meshes and approximation classes.

Remark 14 (Initial triangulation). Assumption 3 guarantees the complexity estimate of module REFINED stated in Theorem 1 and proved in §6.3:

$$\#\mathcal{T}_k - \#\mathcal{T}_0 \leq \Lambda_0 \sum_{j=0}^{k-1} \#\mathcal{M}_j.$$

Assumption 3 is rather restrictive for dimension $d > 2$. Any other refinement giving the same complexity estimate can replace REFINED together with Assumption 3.

Even though we cannot expect local upper bounds between the continuous and discrete solution, the following crucial result shows that this is not the case between discrete solutions on nested meshes $\mathcal{T}_* \geq \mathcal{T}$: what matters is the set of elements of \mathcal{T} which are no longer in \mathcal{T}_* .

Lemma 24 (Localized upper bound). *Let $\mathcal{T}, \mathcal{T}_* \in \mathbb{T}$ satisfy $\mathcal{T}_* \geq \mathcal{T}$ and let $\mathcal{R} := \mathcal{R}_{\mathcal{T} \rightarrow \mathcal{T}_*}$ be the refined set. If $U \in \mathbb{V}$, $U_* \in \mathbb{V}_*$ are the corresponding Galerkin solutions, then*

$$\|U_* - U\|_{\Omega}^2 \leq C_1 \mathcal{E}_{\mathcal{T}}^2(U, \mathcal{R}). \quad (125)$$

Proof. See Problem 53. \square

We are now ready to explore Stevenson's insight (112) for the simplest scenario.

Lemma 25 (Dörfler marking: vanishing oscillation). *Let θ satisfy Assumption 1 and set $\mu := 1 - \theta^2/\theta_*^2 > 0$. Let $\mathcal{T}_* \geq \mathcal{T}$ and the corresponding Galerkin solution $U_* \in \mathbb{V}(\mathcal{T}_*)$ satisfy*

$$\|u - U_*\|_{\Omega}^2 \leq \mu \|u - U\|_{\Omega}^2. \quad (126)$$

Then the refined set $\mathcal{R} = \mathcal{R}_{\mathcal{T} \rightarrow \mathcal{T}_}$ satisfies the Dörfler property*

$$\mathcal{E}_{\mathcal{T}}(U, \mathcal{R}) \geq \theta \mathcal{E}_{\mathcal{T}}(U, \mathcal{T}). \quad (127)$$

Proof. Since $\mu < 1$ we use the lower bound in (123), in conjunction with (126) and Pythagoras equality (92), to derive

$$\begin{aligned} (1 - \mu) C_2 \mathcal{E}_{\mathcal{T}}^2(U, \mathcal{T}) &\leq (1 - \mu) \|u - U\|_{\Omega}^2 \\ &\leq \|u - U\|_{\Omega}^2 - \|u - U_*\|_{\Omega}^2 = \|U - U_*\|_{\Omega}^2. \end{aligned}$$

In view of Lemma 24, we thus deduce

$$(1 - \mu) C_2 \mathcal{E}_{\mathcal{T}}^2(U, \mathcal{T}) \leq C_1 \mathcal{E}_{\mathcal{T}}^2(U, \mathcal{R}),$$

which is the assertion in disguise. \square

To examine the cardinality of \mathcal{M}_k in terms of $\|u - U_k\|_{\Omega}$ we must relate AFEM with the approximation class \mathcal{A}_s . Even though this might appear like an undoable task, the key to unravel this connection is given by Lemma 25. We show this now.

Lemma 26 (Cardinality of \mathcal{M}_k). *Let Assumptions 1 and 2 hold. If $u \in \mathcal{A}_s$ then*

$$\#\mathcal{M}_k \lesssim |u|_s^{1/s} \|u - U_k\|_{\Omega}^{-1/s} \quad \forall k \geq 0. \quad (128)$$

Proof. We invoke that $u \in \mathcal{A}_s$, together with Problem 51 with $\varepsilon^2 = \mu \|u - U_k\|_{\Omega}^2$, to find a mesh $\mathcal{T}_{\varepsilon} \in \mathbb{T}$ and the Galerkin solution $U_{\varepsilon} \in \mathbb{V}(\mathcal{T}_{\varepsilon})$ so that

$$\|u - U_{\varepsilon}\|_{\Omega}^2 \leq \varepsilon^2, \quad \#\mathcal{T}_{\varepsilon} - \#\mathcal{T}_0 \lesssim |u|_s^{1/s} \varepsilon^{-1/s}.$$

Since $\mathcal{T}_{\varepsilon}$ may be totally unrelated to \mathcal{T}_k , we introduce the overlay $\mathcal{T}_* = \mathcal{T}_{\varepsilon} \oplus \mathcal{T}_k$. We exploit the property $\mathcal{T}_* \geq \mathcal{T}_{\varepsilon}$ to conclude that the Galerkin solution $U_* \in \mathbb{V}(\mathcal{T}_*)$ satisfies (127):

$$\|u - U_*\|_{\Omega}^2 \leq \|u - U_{\varepsilon}\|_{\Omega}^2 \leq \varepsilon^2 = \mu \|u - U\|_{\Omega}^2.$$

Therefore, Lemma 25 implies that the refined set $\mathcal{R} = \mathcal{R}_{\mathcal{T} \rightarrow \mathcal{T}_*}$ satisfies a Dörfler marking with parameter $\theta < \theta_*$. But MARK delivers a minimal set \mathcal{M}_k with this property, according to Assumption 2, whence

$$\#\mathcal{M}_k \leq \#\mathcal{R} \leq \#\mathcal{T}_* - \#\mathcal{T}_k \leq \#\mathcal{T}_{\varepsilon} - \#\mathcal{T}_0 \lesssim |u|_s^{1/s} \varepsilon^{-1/s},$$

where we use Lemma 22 to account for the overlay. The proof is complete. \square

Proposition 4 (Quasi-optimality: vanishing oscillation). *Let Assumptions 1, 2, and 3 hold. If $u \in \mathcal{A}_s$, then AFEM gives rise to a sequence $(\mathcal{T}_k, \mathbb{V}_k, U_k)_{k=0}^{\infty}$ such that*

$$\|u - U_k\|_{\Omega} \lesssim |u|_s (\#\mathcal{T}_k - \#\mathcal{T}_0)^{-s} \quad \forall k \geq 1.$$

Proof. We make use of Assumption 3, along with Theorem 1, to infer that

$$\#\mathcal{T}_k - \#\mathcal{T}_0 \leq \Lambda_0 \sum_{j=0}^{k-1} \#\mathcal{M}_j \lesssim |u|_s^{1/s} \sum_{j=0}^{k-1} \|u - U_j\|_{\Omega}^{-1/s}.$$

We now use the contraction property (97) of Lemma 15

$$\|u - U_k\|_{\Omega} \leq \alpha^{k-j} \|u - U_j\|_{\Omega}$$

to replace the sum above by

$$\sum_{j=0}^{k-1} \|u - U_j\|_{\Omega}^{-1/s} \leq \|u - U_k\|_{\Omega}^{-1/s} \sum_{j=0}^{k-1} \alpha^{k-j/s} < \frac{\alpha^{1/s}}{1 - \alpha^{1/s}} \|u - U_k\|_{\Omega}^{-1/s},$$

because $\alpha < 1$ and the series is summable. This completes the proof. \square

7.4 Quasi-Optimal Cardinality: General Data

In this section we remove the restriction $\text{osc}_{\mathcal{T}}(U) = 0$, and thereby make use of the basic ingredients developed in §7.1 and §7.2. Therefore, we replace the energy error by the total error and the linear approximation class \mathcal{A}_s for u by the nonlinear class \mathbb{A}_s for the triple (u, f, \mathbf{A}) . To account for the presence of general f and \mathbf{A} , we need to make an even more stringent assumption on the threshold θ_* .

Assumption 4 (Marking parameter: general data). Let $C_3 = \Lambda_1 \text{osc}_{\mathcal{T}_0}^2(\mathbf{A})$ be the constant in Problem 46 and Lemma 21. The marking parameter θ satisfies $\theta \in (0, \theta_*)$ with

$$\theta_* = \sqrt{\frac{C_2}{1 + C_1(1 + C_3)}}. \quad (129)$$

We now proceed along the same lines as those of §7.3.

Lemma 27 (Dörfler marking: general data). Let Assumption 4 hold and set $\mu := \frac{1}{2}(1 - \frac{\theta^2}{\theta_*^2}) > 0$. If $\mathcal{T}_* \geq \mathcal{T}$ and the corresponding Galerkin solution $U_* \in \mathbb{V}(\mathcal{T}_*)$ satisfy

$$\|u - U_*\|_{\Omega}^2 + \text{osc}_{\mathcal{T}_*}^2(U_*) \leq \mu (\|u - U\|_{\Omega}^2 + \text{osc}_{\mathcal{T}}^2(U)), \quad (130)$$

then the refined set $\mathcal{R} = \mathcal{R}_{\mathcal{T} \rightarrow \mathcal{T}_*}$ satisfies the Dörfler property

$$\mathcal{E}_{\mathcal{T}}(U, \mathcal{R}) \geq \theta \mathcal{E}_{\mathcal{T}}(U, \mathcal{T}). \quad (131)$$

Proof. We split the proof into four steps.

□ In view of the global lower bound (94b)

$$C_2 \mathcal{E}_{\mathcal{T}}^2(U) \leq \|u - U\|_{\Omega}^2 + \text{osc}_{\mathcal{T}}^2(U)$$

and (130), we can write

$$\begin{aligned} (1 - 2\mu) C_2 \mathcal{E}_{\mathcal{T}}^2(U) &\leq (1 - 2\mu) (\|u - U\|_{\Omega}^2 + \text{osc}_{\mathcal{T}}^2(U)) \\ &\leq (\|u - U\|_{\Omega}^2 - 2\|u - U_*\|_{\Omega}^2) + (\text{osc}_{\mathcal{T}}^2(U) - 2\text{osc}_{\mathcal{T}_*}^2(U_*)). \end{aligned}$$

□ Combining the Pythagoras orthogonality relation (92)

$$\|u - U\|_{\Omega}^2 - \|u - U_*\|_{\Omega}^2 = \|U - U_*\|_{\Omega}^2.$$

with the localized upper bound Lemma 24 yields

$$\|u - U\|_{\Omega}^2 - 2\|u - U_*\|_{\Omega}^2 \leq \|U - U_*\|_{\Omega}^2 \leq C_1 \mathcal{E}_{\mathcal{T}}^2(U, \mathcal{R}).$$

□ To deal with oscillation we decompose the elements of \mathcal{T} into two disjoint sets: \mathcal{R} and $\mathcal{T} \setminus \mathcal{R}$. In the former case, we have

$$\text{osc}_{\mathcal{T}}^2(U, \mathcal{R}) - 2\text{osc}_{\mathcal{T}_*}^2(U_*, \mathcal{R}) \leq \text{osc}_{\mathcal{T}}^2(U, \mathcal{R}) \leq \mathcal{E}_{\mathcal{T}}^2(U, \mathcal{R}),$$

because $\text{osc}_{\mathcal{T}}(U, T) \leq \mathcal{E}_{\mathcal{T}}(U, T)$ for all $T \in \mathcal{T}$. On the other hand, we use that $\mathcal{T} \setminus \mathcal{R} = \mathcal{T} \cap \mathcal{T}_*$ and apply Problem 46 in conjunction with Lemma 24 to arrive at

$$\text{osc}_{\mathcal{T}}^2(U, \mathcal{T} \setminus \mathcal{R}) - 2 \text{osc}_{\mathcal{T}_*}^2(U_*, \mathcal{T} \setminus \mathcal{R}) \leq C_3 \|U - U_*\|_{\Omega}^2 \leq C_1 C_3 \mathcal{E}_{\mathcal{T}}^2(U, \mathcal{R}).$$

Adding these two estimates gives

$$\text{osc}_{\mathcal{T}}^2(U) - 2 \text{osc}_{\mathcal{T}_*}^2(U_*) \leq (1 + C_1 C_3) \mathcal{E}_{\mathcal{T}}^2(U, \mathcal{R}).$$

□ Returning to □ we realize that

$$(1 - 2\mu) C_2 \mathcal{E}_{\mathcal{T}}^2(U, \mathcal{T}) \leq (1 + C_1(1 + C_3)) \mathcal{E}_{\mathcal{T}}^2(U, \mathcal{R}),$$

which is the asserted estimate (131) in disguise. □

Lemma 28 (Cardinality of \mathcal{M}_k : general data). *Let Assumptions 2 and 4 hold. If $(u, f, \mathbf{A}) \in \mathbb{A}_s$, then*

$$\#\mathcal{M}_k \lesssim |u, f, \mathbf{A}|_s^{1/s} (\|u - U_k\|_{\Omega} + \text{osc}_k(U_k))^{-1/s} \quad \forall k \geq 0. \quad (132)$$

Proof. We split the proof into three steps.

□ We set $\varepsilon^2 := \mu \Lambda_2^{-1} (\|u - U_k\|_{\Omega}^2 + \text{osc}_k^2(U_k))$ with $\mu = \frac{1}{2} (1 - \frac{\theta^2}{\theta_*^2}) > 0$ as in Lemma 27 and Λ_2 given Lemma 21. Since $(u, f, \mathbf{A}) \in \mathbb{A}_s$, in view of Problem 51 there exists $\mathcal{T}_{\varepsilon} \in \mathbb{T}$ and $U_{\varepsilon} \in \mathbb{V}(\mathcal{T}_{\varepsilon})$ such that

$$\|u - U_{\varepsilon}\|_{\Omega}^2 + \text{osc}_{\varepsilon}^2(U_{\varepsilon}) \leq \varepsilon^2 \quad \text{and} \quad \#\mathcal{T}_{\varepsilon} - \#\mathcal{T}_0 \lesssim |u, f, \mathbf{A}|_s^{1/2} \varepsilon^{-1/s}.$$

Since $\mathcal{T}_{\varepsilon}$ may be totally unrelated to \mathcal{T}_k we introduce the overlay $\mathcal{T}_* = \mathcal{T}_k \oplus \mathcal{T}_{\varepsilon}$.

□ We claim that the total error over \mathcal{T}_* reduces by a factor μ relative to that one over \mathcal{T}_k . In fact, since $\mathcal{T}_* \geq \mathcal{T}_{\varepsilon}$ and so $\mathbb{V}(\mathcal{T}_*) \supset \mathbb{V}(\mathcal{T}_{\varepsilon})$, we use Lemma 21 to obtain

$$\begin{aligned} \|u - U_*\|_{\Omega}^2 + \text{osc}_{\mathcal{T}_*}^2(U_*) &\leq \Lambda_2 \left(\|u - U_{\varepsilon}\|_{\Omega}^2 + \text{osc}_{\varepsilon}^2(U_{\varepsilon}) \right) \\ &\leq \Lambda_2 \varepsilon^2 = \mu \left(\|u - U_k\|_{\Omega}^2 + \text{osc}_k^2(U_k) \right). \end{aligned}$$

Upon applying Lemma 27 we conclude that the set $\mathcal{R} = \mathcal{R}_{\mathcal{T}_k \rightarrow \mathcal{T}_*}$ of refined elements satisfies a Dörfler marking (131) with parameter $\theta < \theta_*$.

□ According to Assumption 2, MARK selects a minimal set \mathcal{M}_k satisfying this property. Therefore, we deduce

$$\#\mathcal{M}_k \leq \#\mathcal{R} \leq \#\mathcal{T}_* - \#\mathcal{T}_k \leq \#\mathcal{T}_{\varepsilon} - \#\mathcal{T}_0 \lesssim |u, f, \mathbf{A}|_s^{1/s} \varepsilon^{-1/s},$$

where we have employed Lemma 22 to account for the cardinality of the overlay. Finally, recalling the definition of ε we end up with the asserted estimate (132). □

Remark 15 (Blow-up of constant). The constant hidden in (132) blows up as $\theta \uparrow \theta_*$ because $\mu \downarrow 0$; see Problem 54.

We are ready to prove the main result of this section, which combines Theorem 9 and Lemma 28.

Theorem 10 (Quasi-optimality: general data). *Let Assumptions 2, 3 and 4 hold. If $(u, f, \mathbf{A}) \in \mathbb{A}_s$, then AFEM gives rise to a sequence $(\mathcal{T}_k, \mathbb{V}_k, U_k)_{k=0}^\infty$ such that*

$$\|u - U_k\|_\Omega + \text{osc}_k(U_k) \lesssim |u, f, \mathbf{A}|_s (\#\mathcal{T}_k - \#\mathcal{T}_0)^{-s} \quad \forall k \geq 1.$$

Proof. \square Since no confusion arises, we use the notation $\text{osc}_j = \text{osc}_j(U_j)$ and $\mathcal{E}_j = \mathcal{E}_j(U_j)$. In light of Assumption 3, which yields Theorem 1, and (132) we have

$$\#\mathcal{T}_k - \#\mathcal{T}_0 \lesssim \sum_{j=0}^{k-1} \#\mathcal{M}_j \lesssim |u, f, \mathbf{A}|_s^{1/s} \sum_{j=0}^{k-1} (\|u - U_j\|_\Omega^2 + \text{osc}_j^2)^{-1/(2s)}.$$

\square Let $\gamma > 0$ be the scaling factor in the (contraction) Theorem 9. The lower bound (94b) along with $\text{osc}_j \leq \mathcal{E}_j$ implies

$$\|u - U_j\|_\Omega^2 + \gamma \text{osc}_j^2 \leq \|u - U_j\|_\Omega^2 + \gamma \mathcal{E}_j^2 \leq \left(1 + \frac{\gamma}{C_2}\right) (\|u - U_j\|_\Omega^2 + \text{osc}_j^2).$$

\square Theorem 9 yields for $0 \leq j < k$

$$\|u - U_k\|_\Omega^2 + \gamma \mathcal{E}_k^2 \leq \alpha^{2(k-j)} (\|u - U_j\|_\Omega^2 + \gamma \mathcal{E}_j^2),$$

whence

$$\#\mathcal{T}_k - \#\mathcal{T}_0 \lesssim |u, f, \mathbf{A}|_s^{1/s} (\|u - U_k\|_\Omega^2 + \gamma \mathcal{E}_k^2)^{-1/(2s)} \sum_{j=0}^{k-1} \alpha^{(k-j)/s}.$$

Since $\sum_{j=0}^{k-1} \alpha^{(k-j)/s} = \sum_{j=1}^k \alpha^{j/s} < \sum_{j=1}^\infty \alpha^{j/s} < \infty$ because $\alpha < 1$, the assertion follows immediately. \square

We conclude this section with several applications of Theorem 10.

Corollary 10 (Estimator decay). *Let Assumptions 2, 3 and 4 be satisfied. If $(u, f, \mathbf{A}) \in \mathbb{A}_s$ then the estimator $\mathcal{E}_k(U_k)$ satisfies*

$$\mathcal{E}_k(U_k) \lesssim |u, f, \mathbf{A}|_s^{1/s} (\#\mathcal{T}_k - \#\mathcal{T}_0)^{-s} \quad \forall k \geq 1.$$

Proof. Use (114) and Theorem 10. \square

Corollary 11 (W_p^2 -regularity with piecewise linear \mathbf{A}). *Let $d = 2$, the polynomial degree $n = 1$, $f \in L^2(\Omega)$, and let \mathbf{A} be piecewise linear over \mathcal{T}_0 . If $u \in W_p^2(\Omega; \mathcal{T}_0)$ for $p > 1$, then AFEM gives rise to a sequence $\{\mathcal{T}_k, \mathbb{V}_k, U_k\}_{k=0}^\infty$ satisfying $\text{osc}_k(U_k) = \|h_k(f - P_0 f)\|_{L^2(\Omega)}$ and for all $k \geq 1$*

$$\|u - U_k\|_\Omega + \text{osc}_k(U_k) \lesssim \left(\|D^2 u\|_{L^p(\Omega; \mathcal{T}_0)} + \|f\|_{L^2(\Omega)} \right) (\#\mathcal{T}_k - \#\mathcal{T}_0)^{-1/2}.$$

Proof. Combine Corollary 7 with Theorem 10. \square

Corollary 12 (W_p^2 -regularity with variable \mathbf{A}). *Assume the setting of Corollary 11, but let \mathbf{A} be piecewise Lipschitz over the initial grid \mathcal{T}_0 . Then AFEM gives rise to a sequence $\{\mathcal{T}_k, \nabla_k, U_k\}_{k=0}^\infty$ satisfying for all $k \geq 1$*

$$\begin{aligned} & \|u - U_k\|_\Omega + \text{osc}_k(U_k) \\ & \lesssim \left(\|D^2 u\|_{L^p(\Omega; \mathcal{T}_0)} + \|f\|_{L^2(\Omega)} + \|\mathbf{A}\|_{W_\infty^1(\Omega; \mathcal{T}_0)} \right) (\#\mathcal{T}_k - \#\mathcal{T}_0)^{-1/2}. \end{aligned}$$

Proof. Combine Corollary 8 with Theorem 10. \square

Corollary 13 (W_p^s -regularity with $s < 1/d$). *Let $d \geq 2$, $n = 1$, $1 < t < 2$, $p > 1$, $f \in L^2(\Omega)$, and $\mathbf{A} \in W_\infty^1(\Omega, \mathcal{T}_0)$ be piecewise Lipschitz. If $u \in W_p^t(\Omega; \mathcal{T}_0) \cap H_0^1(\Omega)$ is piecewise W_p^t over the initial mesh \mathcal{T}_0 with $t - \frac{d}{p} > 1 - \frac{d}{2}$, then AFEM gives rise to a sequence $\{\mathcal{T}_k, \nabla_k, U_k\}_{k=0}^\infty$ satisfying for all $k \geq 1$*

$$\begin{aligned} & \|u - U_k\|_\Omega + \text{osc}_k(U_k) \\ & \lesssim \left(\|D^t u\|_{L^p(\Omega; \mathcal{T}_0)} + \|f\|_{L^2(\Omega)} + \|\mathbf{A}\|_{W_\infty^1(\Omega; \mathcal{T}_0)} \right) (\#\mathcal{T}_k - \#\mathcal{T}_0)^{-(t-1)/d}. \end{aligned}$$

Proof. Combine Corollary 9 with Theorem 10. \square

7.5 Extensions and Restrictions

We conclude with a brief discussion of extensions of the theory and some of its restrictions.

Optimal Complexity: Inexact Solvers, Quadrature, and Storage. We point out that we have never mentioned the notion of *complexity* so far. This is because complexity estimates entail crucial issues that we have ignored: inexact solvers to approximate the Galerkin solution; quadrature; and optimal storage. We comment on them now.

Multilevel solvers are known to deliver an approximate solution with cost proportional to the number of degrees of freedom. Even though the theory is well developed for uniform refinement, it is much less understood for adaptive refinement. This is due to the fact that the adaptive bisection meshes do not satisfy the so-called nested refinement assumption. Recently, Xu, Chen, and Nochetto [60] have bridged the gap between graded and quasi-uniform grids exploiting the geometric structure of bisection grids and a resulting new space decomposition. They designed and analyzed optimal additive and multiplicative multilevel methods for any dimension $d \geq 2$ and polynomial degree $n \geq 1$, thereby improving upon Wu and Chen [59]. The theories of §5 and §7 can be suitably modified to account for optimal iterative solvers; we refer to Stevenson [52].

Quadrature is a very delicate issue in a purely a posteriori context, that is without a priori knowledge of the functions involved. Even if we were to replace both data f and \mathbf{A} by piecewise polynomials so that quadrature would be simple, we would need to account for the discrepancy in adequate norms between exact and approximate data, again a rather delicate matter. This issue is to a large extent open.

Optimal storage is an essential, but often disregarded, aspect of a complexity analysis. For instance, ALBERTA is an excellent library for AFEM but does not have optimal storage capabilities [50].

Non-Residual Estimators. The cardinality analysis of this section extends to estimators other than the residual; we refer to Cascón and Nochetto [15] and Kreuzer and Siebert [35]. They include the hierarchical, Zienkiewicz-Zhu [2, 27, 55, 58], and Braess-Schoerbel estimators, as well as those based on the solution of local problems [12, 42]. Even though the contraction property of Theorem 9 is no longer valid between consecutive iterates, it is true after a fixed number of iterations, which is enough for the arguments in Proposition 4 and Theorem 10 to apply. The resulting error estimates possess constants proportional to this gap.

Nonconforming Meshes. Since REFINE exhibits optimal complexity for admissible nonconforming meshes, according to §6.4, and this is the only ingredient where nonconformity might play a role, the theory of this section extends. We refer to Bonito and Nochetto [9].

Discontinuous Galerkin Methods (dG). The study of cardinality for adaptive dG methods is rather technical. This is in part due to the fact that key Lemmas 26 and 28 hinge on mesh overlay, which in turn does not provide control of the level of refinement. This makes it difficult to compare broken energy norms

$$\|v\|_{\mathcal{T}}^2 = \|\mathbf{A}^{1/2} \nabla v\|_{L^2(\Omega; \mathcal{T})}^2 + \|h^{-1/2} \llbracket v \rrbracket\|_{L^2(\Sigma)}^2,$$

which contain jump terms with negative powers of the mesh-size over the skeleton Σ of \mathcal{T} . Consequently, the monotonicity of energy norms used in Lemmas 26 and 28 is no longer true!

To circumvent this difficulty, Bonito and Nochetto [9] resorted to continuous finite elements $\mathbb{V}^0(\mathcal{T})$ over the (admissible nonconforming) mesh \mathcal{T} , which have the same degree as their discontinuous counterpart $\mathbb{V}(\mathcal{T})$. This leads to a cardinality theory very much in the spirit of this section. However, it raises the question whether discontinuous elements deliver a better asymptotic rate over admissible nonconforming meshes. Since this result is of intrinsic interest, we report it now.

Lemma 29 (Equivalence of classes). *Let \mathbb{A}_s be the approximation class using discontinuous elements of degree $\leq n$ and \mathbb{A}_s^0 be the continuous counterpart. Then, for $0 < s \leq n/d$, total errors are equivalent on the same mesh, whence $\mathbb{A}_s = \mathbb{A}_s^0$.*

Proof. We use the notation of Problem 11. Since $\mathbb{V}^0(\mathcal{T}) \subset \mathbb{V}(\mathcal{T})$, the inclusion $\mathbb{A}_s^0 \subset \mathbb{A}_s$ is obvious. To prove the converse, we let $(u, f, \mathbf{A}) \in \mathbb{A}_s$ and, for $N > \#\mathcal{T}_0$, let $\mathcal{T}_* \in \mathbb{T}_N$ be an admissible nonconforming grid and $U_* \in \mathbb{V}(\mathcal{T}_*)$ be so that

$$\|u - U_*\|_{\mathcal{T}_*} + \text{osc}_{\mathcal{T}_*}(U_*) = \inf_{\mathcal{T} \in \mathbb{T}_N} \inf_{V \in \mathbb{V}(\mathcal{T})} \left(\|u - V\|_{\mathcal{T}} + \text{osc}_{\mathcal{T}}(V) \right) \lesssim N^{-s}.$$

Let $I_{\mathcal{T}} : \mathbb{V}(\mathcal{T}) \rightarrow \mathbb{V}^0(\mathcal{T})$ be the interpolation operator of Problem 11. Since $I_{\mathcal{T}_*} U_* \in \mathbb{V}^0(\mathcal{T}_*)$, if we were able to prove

$$\|u - I_{\mathcal{T}_*} U_*\|_{\mathcal{T}_*} + \text{osc}_{\mathcal{T}_*}(I_{\mathcal{T}_*} U_*) \lesssim N^{-s},$$

then $(u, f, \mathbf{A}) \in \mathbb{A}_s^0$. Using the triangle inequality, we get

$$\|u - I_{\mathcal{T}_*} U_*\|_{\mathcal{T}_*} \leq \|\mathbf{A}^{1/2} \nabla(u - U_*)\|_{L^2(\Omega; \mathcal{T}_*)} + \|\mathbf{A}^{1/2} \nabla(U_* - I_{\mathcal{T}_*} U_*)\|_{L^2(\Omega; \mathcal{T}_*)},$$

because $[u - I_{\mathcal{T}_*} U_*]$ vanish on Σ . Problem 11 implies the estimate

$$\|\mathbf{A}^{1/2} \nabla(U_* - I_{\mathcal{T}_*} U_*)\|_{L^2(\Omega; \mathcal{T}_*)} \lesssim \|h^{-1/2} [U_*]\|_{L^2(\Sigma_*)} \leq \|u - U_*\|_{\mathcal{T}_*},$$

whence

$$\|u - I_{\mathcal{T}_*} U_*\|_{\mathcal{T}_*} \lesssim \|u - U_*\|_{\mathcal{T}_*}.$$

Since $\|\mathbf{A}^{1/2} \nabla(U_* - I_{\mathcal{T}_*} U_*)\|_{L^2(\Omega; \mathcal{T}_*)} \lesssim \|U_* - I_{\mathcal{T}_*} U_*\|_{\mathcal{T}_*}$, the oscillation term can be treated similarly. In fact, Problem 46 adapted to discontinuous functions yields

$$\text{osc}_{\mathcal{T}_*}(I_{\mathcal{T}_*} U_*) \lesssim \text{osc}_{\mathcal{T}_*}(U_*) + \|u - U_*\|_{\mathcal{T}_*}.$$

Coupling the two estimates above, we end up with

$$\|u - I_{\mathcal{T}_*} U_*\|_{\mathcal{T}_*} + \text{osc}_{\mathcal{T}_*}(I_{\mathcal{T}_*} U_*) \lesssim \|u - U_*\|_{\mathcal{T}_*} + \text{osc}_{\mathcal{T}_*}(U_*) \lesssim N^{-s}.$$

Therefore, $(u, f, \mathbf{A}) \in \mathbb{A}_s^0$ as desired. \square

7.6 Notes

The theory presented in this section is rather recent. It started with the breakthrough (112) by Stevenson [52] for vanishing oscillation. If f is variable and \mathbf{A} is piecewise constant, then Stevenson extended this idea upon adding an inner loop to handle data oscillation to the usual AFEM. This idea does not extend to the model problem (87) with variable \mathbf{A} , because the oscillation then depends on the Galerkin solution.

The next crucial step was made by Cascón, Kreuzer, Nochetto, and Siebert [14], who dealt with the notion of total error of §7.1, as previously done by Mekchay and Nochetto [37], and introduced the nonlinear approximation class \mathbb{A}_s of §7.2. They derived the convergence rates of §7.4.

The analysis for nonconforming meshes is due to Bonito and Nochetto [9], who developed this theory in the context of dG methods for which they also derived convergence rates. The study of non-residual estimators is due to Kreuzer and Siebert [35] and Cascón and Nochetto [15].

The theory is almost exclusively devoted to the energy norm, except for the L^2 -analysis of Demlow and Stevenson [21], who proved an optimal convergence rate for mildly varying graded meshes. Convergence rates have been proved for Raviart-Thomas mixed FEM by Chen, Holst, and Xu [18].

7.7 Problems

Problem 51 (Alternative definition of \mathbb{A}_s). Show that $(v, f, \mathbf{A}) \in \mathbb{A}_s$ if and only there exists a constant $\Lambda > 0$ such that for all $\varepsilon > 0$ there exist $\mathcal{T}_\varepsilon \geq \mathcal{T}_0$ conforming and $V_\varepsilon \in \mathbb{V}(\mathcal{T}_\varepsilon)$ such that

$$\|v - V_\varepsilon\|_\Omega^2 + \text{osc}_{\mathcal{T}_\varepsilon}^2(V_\varepsilon) \leq \varepsilon^2 \quad \text{and} \quad \#\mathcal{T}_\varepsilon - \#\mathcal{T}_0 \leq \Lambda^{1/s} \varepsilon^{-1/s};$$

in this case $|v, f, \mathbf{A}|_s \leq \Lambda$. Hint: Let \mathcal{T}_ε be minimal for $\|v - V_\varepsilon\|_\Omega^2 + \text{osc}_{\mathcal{T}_\varepsilon}^2(V_\varepsilon) \leq \varepsilon^2$. This means that for all $\mathcal{T} \in \mathbb{T}$ such that $\#\mathcal{T} = \#\mathcal{T}_\varepsilon - 1$ we have $\|v - V_\varepsilon\|_\Omega^2 + \text{osc}_{\mathcal{T}_\varepsilon}^2(V_\varepsilon) > \varepsilon$.

Problem 52 (Lemma 22). Prove that the overlay $\mathcal{T} = \mathcal{T}_1 \oplus \mathcal{T}_2$ is conforming and

$$\#\mathcal{T} \leq \#\mathcal{T}_1 + \#\mathcal{T}_2 - \#\mathcal{T}_0.$$

Hint: for each $T \in \mathcal{T}_0$, consider two cases $\mathcal{T}_1(T) \cap \mathcal{T}_2(T) \neq \emptyset$ and $\mathcal{T}_1(T) \cap \mathcal{T}_2(T) = \emptyset$, where $\mathcal{T}_i(T)$ is the portion of the mesh \mathcal{T}_i contained in T .

Problem 53 (Lemma 24). Prove that if $\mathcal{T}, \mathcal{T}_* \in \mathbb{T}$ satisfy $\mathcal{T}_* \geq \mathcal{T}$, $\mathcal{R} := \mathcal{R}_{\mathcal{T} \rightarrow \mathcal{T}_*}$ is the refined set to go from \mathcal{T} to \mathcal{T}_* , and $U \in \mathbb{V}$, $U_* \in \mathbb{V}_*$ are the corresponding Galerkin solutions, then

$$\|U_* - U\|_\Omega^2 \leq C_1 \mathcal{E}_{\mathcal{T}}^2(U, \mathcal{R}).$$

To this end, write the equation fulfilled by $U_* - U \in \mathbb{V}_*$ and use as a test function the local quasi-interpolant $I_{\mathcal{T}}(U_* - U)$ of $U_* - U$ introduced in Proposition 2.

Problem 54 (Explicit dependence on θ and s). Trace the dependence on θ and s , as $\theta \rightarrow \theta_*$ and $s \rightarrow 0$, in the hidden constants in Lemma 28 and Theorem 10.

Problem 55 (Asymptotic decay of oscillation). Let $\mathbf{A} \in W_\infty^1(\Omega; \mathcal{T}_0)$ be piecewise Lipschitz over the initial grid \mathcal{T}_0 and $f \in L^2(\Omega)$. Show that

$$\inf_{\mathcal{T} \in \mathbb{T}_N} \text{osc}_{\mathcal{T}}(U) \lesssim \left(\|f\|_{L^2(\Omega)} + \|\mathbf{A}\|_{W_\infty^1(\Omega; \mathcal{T}_0)} \right) N^{-1/d}$$

is attained with uniform meshes.

Problem 56 (Faster decay of data oscillation). Let $d = 2$ and $n = 1$. Let f be piecewise W_1^1 over the initial mesh \mathcal{T}_0 , namely $f \in W_1^1(\Omega; \mathcal{T}_0)$. Show that

$$\inf_{\mathcal{T} \in \mathbb{T}_N} \|h_{\mathcal{T}}(f - P_0 f)\|_{L^2(\Omega)} \lesssim \|f\|_{W_1^1(\Omega; \mathcal{T}_0)} N^{-1},$$

using the thresholding algorithm of §1.6. Therefore, data oscillation decays twice as fast as the energy error asymptotically on suitably graded meshes.

Problem 57 (Faster decay of coefficient oscillation). Consider the coefficient oscillation weighted locally by the energy of the discrete solution U :

$$\eta_{\mathcal{T}}^2(\mathbf{A}, U) = \sum_{T \in \mathcal{T}} \text{osc}_{\mathcal{T}}^2(\mathbf{A}, T) \|\nabla U\|_{L^2(\omega_T)}^2,$$

where $\text{osc}_{\mathcal{T}}(\mathbf{A}, T)$ is defined in Problem 45. Let $d = 2, n = 1, p > 2$, and $\mathbf{A} \in W_p^2(\Omega; \mathcal{T}_0)$ be piecewise in W_p^2 over the initial grid \mathcal{T}_0 . Use the thresholding algorithm of §1.6 to show that $\eta_{\mathcal{T}}(\mathbf{A}, U)$ decays with a rate twice as fast as the energy error:

$$\inf_{\mathcal{T} \in \mathbb{T}_N} \eta_{\mathcal{T}}(\mathbf{A}, U) \lesssim \|\mathbf{A}\|_{W_p^2(\Omega; \mathcal{T}_0)} \|\nabla U\|_{L^2(\Omega)} N^{-1}.$$

Problem 58 (Faster decay of oscillation). Combine Problems 29, 56 and 57 for $d = 2, n = 1$ and $p > 2$ to prove that if $f \in W_1^1(\Omega; \mathcal{T}_0)$ and $\mathbf{A} \in W_p^2(\Omega; \mathcal{T}_0)$, then the oscillation $\text{osc}_{\mathcal{T}}(U, \mathcal{T})$ decays with a rate twice as fast as the energy error:

$$\inf_{\mathcal{T} \in \mathbb{T}_N} \text{osc}_{\mathcal{T}}(U) \lesssim \left(\|f\|_{W_1^1(\Omega; \mathcal{T}_0)} + \|\mathbf{A}\|_{W_p^2(\Omega; \mathcal{T}_0)} \right) N^{-1}.$$

References

1. M. AINSWORTH AND D. W. KELLY, *A posteriori error estimators and adaptivity for finite element approximation of the non-homogeneous Dirichlet problem*, Adv. Comput. Math., 15 (2001), pp. 3–23.
2. M. AINSWORTH AND J.T. ODEN, *A Posteriori Error Estimation in Finite Element Analysis*, Pure and Applied Mathematics (New York), Wiley-Interscience, 2000.
3. I. BABUŠKA, R.B. KELLOGG, AND J. PITKÄRANTA, *Direct and inverse error estimates for finite elements with mesh refinements*, Numer. Math., 33 (1979), pp. 447–471.
4. I. BABUŠKA AND A. MILLER, *A feedback finite element method with a posteriori error estimation. I. The finite element method and some basic properties of the a posteriori error estimator*, Comput. Methods Appl. Mech. Engrg. 61 (1) (1987), pp. 1–40.
5. I. BABUŠKA AND W. RHEINBOLDT, *Error estimates for adaptive finite element computations* SIAM J. Numer. Anal., 15 (1978), pp. 736–754.
6. E. BÄNSCH, P. MORIN, AND R. H. NOCHETTO, *An adaptive Uzawa FEM for the Stokes problem: convergence without the inf-sup condition*, SIAM J. Numer. Anal., 40 (2002), pp. 1207–1229 (electronic).
7. P. BINEV, W. DAHMEN, AND R. DEVORE, *Adaptive finite element methods with convergence rates*, Numer. Math., 97 (2004), pp. 219–268.
8. P. BINEV, W. DAHMEN, R. DEVORE, AND P. PETRUSHEV, *Approximation classes for adaptive methods*, Serdica Math. J., 28 (2002), pp. 391–416. Dedicated to the memory of Vassil Popov on the occasion of his 60th birthday.
9. A. BONITO AND R.H. NOCHETTO, *Quasi-optimal convergence rate for an adaptive discontinuous Galerkin method*, (submitted).
10. D. BRAESS, *Finite Elements. Theory, fast solvers, and applications in solid mechanics*, 2nd edition. Cambridge University Press (2001).
11. S. BRENNER AND L.R. SCOTT, *The Mathematical Theory of Finite Element Methods*, Springer Texts in Applied Mathematics 15 (2008).
12. C. CARSTENSEN AND S. A. FUNKEN, *Fully reliable localized error control in the FEM*, SIAM J. Sci. Comput., 21 (1999), pp. 1465–1484.
13. C. CARSTENSEN AND R.H.W. HOPPE, *Error reduction and convergence for an adaptive mixed finite element method*, Math. Comp. 75 (2006), 1033–1042.
14. J. M. CASCÓN, C. KREUZER, R. H. NOCHETTO, AND K. G. SIEBERT, *Quasi-optimal convergence rate for an adaptive finite element method*, SIAM J. Numer. Anal., 46 (2008), pp. 2524–2550.
15. J. M. CASCÓN AND R. H. NOCHETTO, *Convergence and quasi-optimality for AFEM based on non-residual a posteriori error estimators*, (in preparation).
16. J. M. CASCÓN, R. H. NOCHETTO, AND K. G. SIEBERT, *Design and convergence of AFEM in $H(\text{div})$* , Math. Models Methods Appl. Sci., 17 (2007), pp. 1849–1881.
17. Z. CHEN AND J. FENG, *An adaptive finite element algorithm with reliable and efficient error control for linear parabolic problems*, Math. Comp., 73 (2006), pp. 1167–1042.
18. L. CHEN, M. HOLST, AND J. XU, *Convergence and optimality of adaptive mixed finite element methods*, Math. Comp. 78 (2009), 35–53.
19. P.G. CIARLET, *The Finite Element Method for Elliptic Problems*, Classics in Applied Mathematics, 40, SIAM (2002).
20. A. DEMLOW, *Convergence of an adaptive finite element method for controlling local energy errors*, (submitted).
21. A. DEMLOW AND R.P. STEVENSON, *Convergence and quasi-optimality of an adaptive finite element method for controlling L_2 errors*, (submitted).
22. R.A. DEVORE, *Nonlinear approximation*, A. Iserles (ed.) Acta Numerica, vol. 7, pp. 51–150. Cambridge University Press (1998).
23. L. DIENING AND CH. KREUZER, *Convergence of an adaptive finite element method for the p -Laplacian equation*, SIAM J. Numer. Anal., 46 (2008), pp. 614–638.

24. W. DÖRFLER, *A convergent adaptive algorithm for Poisson's equation*, SIAM J. Numer. Anal., 33 (1996), pp. 1106–1124.
25. W. DÖRFLER AND M. RUMPF, *An adaptive strategy for elliptic problems including a posteriori controlled boundary approximation*, Math. Comp., 67 (1998), pp. 1361–1382.
26. T. DUPONT AND L.R. SCOTT, *Polynomial approximation of functions in Sobolev spaces*, Math. Comp., 34 (1980), pp. 441–463.
27. F. FIERRO AND A. VEESER, *A posteriori error estimators, gradient recovery by averaging, and superconvergence*, Numer. Math., 103 (2006), pp. 267–298.
28. E.M. GARAU, P. MORIN, C. ZUPPA, *Convergence of adaptive finite element methods for eigenvalue problems*, Math. Models Methods Appl. Sci. 19 (2009), pp. 721–747.
29. F. GASPOZ AND P. MORIN, *Approximation classes for adaptive higher order finite element approximation*, (in preparation), (2010).
30. P. GRISVARD, *Elliptic Problems in Nonsmooth Domains, Monographs and Studies in Mathematics*, vol. 24. Pitman (Advanced Publishing Program), Boston, MA (1985).
31. M. HOLST, G. TSOGTGEREL, AND Y. ZHU, *Local convergence of adaptive methods for nonlinear partial differential equations*, preprint arXiv:math.NA/1001.1382 (2009).
32. R. H. W. HOPPE, G. KANSCHAT, AND T. WARBURTON, *Convergence analysis of an adaptive interior penalty discontinuous Galerkin method*, SIAM J. Numer. Anal., 47 (2009), pp. 534–550.
33. O.A. KARAKASHIAN AND F. PASCAL, *Convergence of adaptive discontinuous Galerkin approximations of second-order elliptic problems*, SIAM J. Numer. Anal. 45 (2007), 641–665.
34. R.B. KELLOGG, *On the Poisson equation with intersecting interfaces*, Applicable Anal., 4 (1974/75), 101–129.
35. CH. KREUZER AND K.G. SIEBERT, *Decay rates of adaptive finite elements with Dörfler marking*, (submitted).
36. J.M. MAUBACH, *Local bisection refinement for n -simplicial grids generated by reflection*, SIAM J. Sci. Comput. 16 (1995), 210–227.
37. K. MEKCHAY AND R. H. NOCHETTO, *Convergence of adaptive finite element methods for general second order linear elliptic PDEs*, SIAM J. Numer. Anal., 43 (2005), pp. 1803–1827 (electronic).
38. K. MEKCHAY, P. MORIN, AND R. H. NOCHETTO, *AFEM for the Laplace-Beltrami operator on graphs: design and conditional contraction property*, Math. Comp. (to appear).
39. W.F. MITCHELL, *Unified multilevel adaptive finite element methods for elliptic problems*, Ph.D. thesis, Department of Computer Science, University of Illinois, Urbana (1988)
40. P. MORIN, R. H. NOCHETTO, AND K. G. SIEBERT, *Data oscillation and convergence of adaptive FEM*, SIAM J. Numer. Anal., 38 (2000), pp. 466–488.
41. P. MORIN, R. H. NOCHETTO, AND K. G. SIEBERT, *Convergence of adaptive finite element methods*, SIAM Review, 44 (2002), pp. 631–658.
42. P. MORIN, R. H. NOCHETTO, AND K. G. SIEBERT, *Local problems on stars: a posteriori error estimators, convergence, and performance*, Math. Comp., 72 (2003), pp. 1067–1097 (electronic).
43. P. MORIN, K. G. SIEBERT, AND A. VEESER, *Convergence of finite elements adapted for weak norms*, in Applied and Industrial Mathematics in Italy II, Vinvenzo Cutello, Giorgio Fotia, Luigia Puccio eds, Series on Advances in Mathematics for Applied Sciences, 75 (2007), pp. 468–479.
44. P. MORIN, K. G. SIEBERT, AND A. VEESER, *A basic convergence result for conforming adaptive finite elements*, Math. Mod. Meth. Appl. Sci., 5 (2008), pp. 707–737.
45. R.H. NOCHETTO, K. G. SIEBERT, AND A. VEESER, *Theory of adaptive finite element methods: an introduction*, in *Multiscale, Nonlinear and Adaptive Approximation*, Springer, 2009, pp. 409–542.
46. R. H. NOCHETTO, A. SCHMIDT, K. G. SIEBERT, AND A. VEESER, *Pointwise a posteriori error estimates for monotone semi-linear equations*, Numer. Math., 104 (2006), pp. 515–538.
47. R. SACCHI AND A. VEESER, *Locally efficient and reliable a posteriori error estimators for Dirichlet problems*, Math. Models Methods Appl., 16 (2006), pp. 319–346.

48. L.R. SCOTT AND S. ZHANG, *Finite element interpolation of nonsmooth functions satisfying boundary conditions*, Math. Comp., 54 (1990), pp. 483–493.
49. K.G. SIEBERT, *A convergence proof for adaptive finite elements without lower bound*, Preprint Universität Duisburg-Essen and Universität Freiburg No. 1/2009.
50. K.G. SIEBERT, *Mathematically founded design of adaptive finite element software*, in Multiscale and Adaptivity: Modeling, Numerics and Applications, CIME-EMS Summer School in Applied Mathematics, G. Naldi and G. Russo eds., Springer (2010).
51. K.G. SIEBERT AND A. VEESER, *A unilaterally constrained quadratic minimization with adaptive finite elements*, SIAM J. Optim., 18 (2007), pp. 260–289.
52. R. STEVENSON, *Optimality of a standard adaptive finite element method*, Found. Comput. Math., 7 (2007), pp. 245–269.
53. R. STEVENSON, *The completion of locally refined simplicial partitions created by bisection*, Math. Comput., 77 (2008), pp. 227–241.
54. C.T. TRAXLER, *An algorithm for adaptive mesh refinement in n dimensions*, Computing 59 (1997), 115–137.
55. A. VEESER, *Convergent adaptive finite elements for the nonlinear Laplacian*, Numer. Math. 92 (2002), pp. 743–770.
56. A. VEESER AND R. VERFÜRTH, *Explicit upper bounds for dual norms of residuals*, SIAM J. Numer. Anal., 47 (2009), pp. 2387–2405.
57. R. VERFÜRTH, *A posteriori error estimators for the Stokes equations*, Numer. Math., 55 (1989), pp. 309–325.
58. R. VERFÜRTH, *A Review of A Posteriori Error Estimation and Adaptive Mesh-Refinement Techniques*, Adv. Numer. Math. John Wiley, Chichester, UK (1996).
59. H. WU AND Z. CHEN, *Uniform convergence of multigrid V-cycle on adaptively refined finite element meshes for second order elliptic problems*, Science in China: Series A Mathematics, 49 (2006), pp. 1–28.
60. J. XU, L. CHEN, AND R.H. NOCHETTO, *Adaptive multilevel methods on graded bisection grids*, in *Multiscale, Nonlinear and Adaptive Approximation*, Springer, 2009, pp. 599–659.