

Why Adaptive Finite Element Methods Outperform Classical Ones

Ricardo H. Nochetto*

Abstract

Adaptive finite element methods (AFEM) are a fundamental numerical tool in science and engineering. They are known to outperform classical FEM in practice and deliver optimal convergence rates when the latter cannot. This paper surveys recent progress in the theory of AFEM which explains their success and provides a solid mathematical framework for further developments.

Mathematics Subject Classification (2010). Primary 65N30, 65N50, 65N15; Secondary 41A25.

Keywords. finite element methods, a posteriori error estimates, adaptivity, contraction, approximation class, nonlinear approximation, convergence rates.

1. Introduction

Mathematically sound adaptive finite element methods (AFEM) have been the subject of intense research since the late 70's, starting with the pioneering work of Babuška [4, 3]. It is known to practitioners that AFEM can achieve optimal performance, measured as error vs degrees of freedom, in situations when classical FEM cannot. However, it took about 30 years to develop a theory for the energy norm that explains this behavior and provides solid mathematical foundations for further development. This paper presents this theory [10, 33], and its connection to nonlinear approximation [17], for the *model elliptic PDE*

$$-\operatorname{div}(\mathbf{A}\nabla u) = f \quad \text{in } \Omega, \quad (1)$$

with Ω a polyhedral domain of \mathbb{R}^d ($d \geq 2$), homogeneous Dirichlet boundary condition on $\partial\Omega$, and \mathbf{A} symmetric, bounded, and uniformly positive definite.

*Partially supported by NSF grant DMS-0807811.

Department of Mathematics and Institute of Physical Science and Technology, University of Maryland, College Park, MD 20742. E-mail: rhn@math.umd.edu.

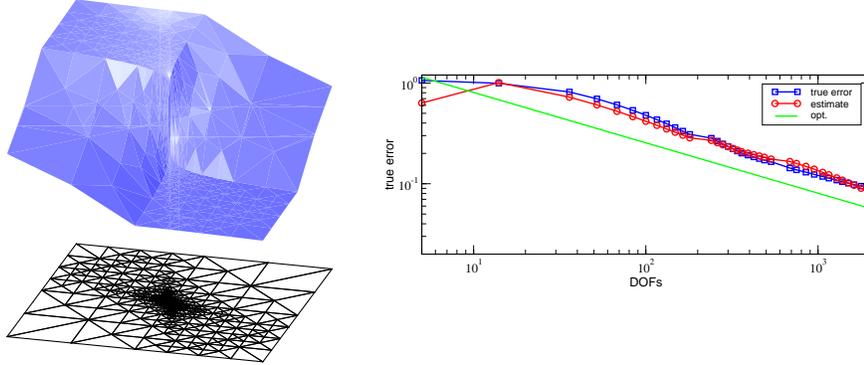


Figure 1. Discontinuous coefficients in checkerboard pattern: (a) graph of the discrete solution u , which is $u \approx r^{0.1}$, and underlying strongly graded grid \mathcal{T} towards the origin (notice the steep gradient of u at the origin); (b) estimate and true error in terms of $\#\mathcal{T}$ (the optimal decay for piecewise linear elements in 2d is indicated by the straight line with slope $-1/2$).

We start with a simple yet quite demanding example with discontinuous coefficients for $d = 2$ due to Kellogg [20], and used by Morin, Nochetto, and Siebert [25, 26] as a benchmark for AFEM. We consider $\Omega = (-1, 1)^2$, $\mathbf{A} = a_1 \mathbf{I}$ in the first and third quadrants, and $\mathbf{A} = a_2 \mathbf{I}$ in the second and fourth quadrants. This checkerboard pattern is the worst for the regularity of the solution u at the origin. For $f = 0$, a function of the form $u(r, \theta) = r^\gamma \mu(\theta)$ in polar coordinates solves (1) with nonvanishing Dirichlet condition for suitable $0 < \gamma < 1$ and μ [25, 26, 28]. We choose $\gamma = 0.1$, which leads to $u \in H^s(\Omega)$ for $s < 1.1$ and piecewise in W_p^1 for some $p > 1$. This corresponds to diffusion coefficients $a_1 \cong 161.44$ and $a_2 = 1$, which can be computed via Newton's method; the closer γ is to 0, the larger is the ratio a_1/a_2 . The solution u and a sample mesh are depicted in Figure 1(a).

Figure 1(b) documents the optimal performance of AFEM: both the energy error and estimator exhibit optimal decay $(\#\mathcal{T})^{-1/2}$ in terms of the cardinality $\#\mathcal{T}$ of the underlying mesh \mathcal{T} for piecewise linear finite elements. On the other hand, Figure 2 displays a strongly graded mesh \mathcal{T} towards the origin generated by AFEM using bisection, and three zooms which reveal a self-similar structure. It is worth stressing that the meshsize is of order 10^{-10} at the origin and $\#\mathcal{T} \approx 2 \times 10^3$, whereas to reach a similar resolution with a uniform mesh \mathcal{T} we would need $\#\mathcal{T} \approx 10^{20}$. This example clearly reveals that adaptivity can restore optimal performance even with modest computational resources.

Classical FEM with quasi-uniform meshes \mathcal{T} require regularity $u \in H^2(\Omega)$ to deliver an optimal convergence rate $(\#\mathcal{T})^{-1/2}$. Since $u \notin H^s(\Omega)$ for any $s > 1.1$, this is not possible for the example above. However, the problem is not quite the lack of second derivatives, but rather the fact that they are not

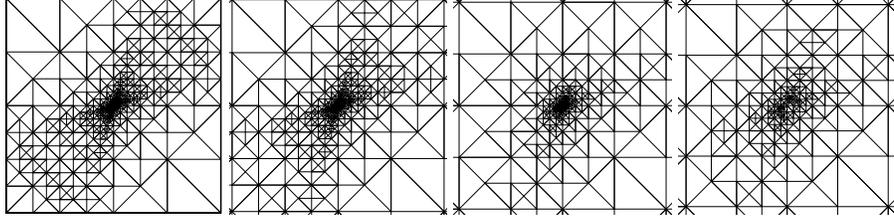


Figure 2. Discontinuous coefficients in checkerboard pattern: (a) final grid \mathcal{T} highly graded towards the origin with $\#\mathcal{T} \approx 2000$; (b) zoom to $(-10^{-3}, 10^{-3})^2$; (c) zoom to $(-10^{-6}, 10^{-6})^2$; (d) zoom to $(-10^{-9}, 10^{-9})^2$. For a similar resolution, a uniform grid \mathcal{T} would require $\#\mathcal{T} \approx 10^{20}$.

square integrable. In fact, the function u is in W_p^2 for $p > 1$ in each quadrant, and so over the initial mesh \mathcal{T}_0 , namely $u \in W_p^2(\Omega; \mathcal{T}_0)$.

To measure the performance of AFEM we introduce an approximation class \mathcal{A}_s for $s > 0$. Given an initial grid \mathcal{T}_0 , and the set \mathbb{T}_N of all conforming refinements \mathcal{T}_0 by *bisection* with at most N elements more than \mathcal{T}_0 , we consider the best error

$$\sigma_N(u) := \inf_{\mathcal{T} \in \mathbb{T}_N} \inf_{V \in \mathbb{V}(\mathcal{T})} \|u - V\|_{\Omega} \quad (2)$$

in the *energy norm* $\|\cdot\|_{\Omega} = \|\mathbf{A}^{1/2} \nabla \cdot\|_{L^2(\Omega)}$, where $\mathbb{V}(\mathcal{T}) \subset H_0^1(\Omega)$ is the conforming finite element space of piecewise polynomials of degree $\leq n$ with $n \geq 1$ over \mathcal{T} . We say that $u \in \mathcal{A}_s$ if

$$\sigma_N(u) \lesssim N^{-s}. \quad (3)$$

We wonder whether or not AFEM is able to deliver this asymptotic error decay. If we have access to the local energy error, we give a constructive proof in §3 of the fact that for $d = 2$

$$u \in W_p^2(\Omega; \mathcal{T}_0) \cap H_0^1(\Omega) \quad \Rightarrow \quad u \in \mathcal{A}_{1/2}. \quad (4)$$

This shows that piecewise linear finite element approximations can deliver optimal error decay. However, we only have indirect access to the solution u of (1) via the error estimators, so it is highly nontrivial whether a similar result holds for the Galerkin solution given by AFEM. The answer to this question requires two steps:

- *Contraction property*: we show in §5.1 that the energy error contracts provided the data is piecewise constant (so that the oscillation vanishes) and the interior node property holds. Otherwise, we identify in §5.2 a novel contractive quantity for general data, the so-called *quasi-error*, and prove that AFEM contracts it.
- *Convergence rate*: we show in §6.3 that the class \mathcal{A}_s is adequate provided the oscillation vanishes. However, the concept of approximation class for AFEM

is generally more involved than just \mathcal{A}_s because it entails dealing with the *total error*, namely the sum of energy error and oscillation. We discuss this issue in §6.1 and §6.2, and next prove that AFEM delivers a convergence rate similar to (3) up a multiplicative constant in §6.4.

It is worth stressing that AFEM learns about the decay rate $s > 0$ via the estimator. In fact, this exponent is never used in the design of AFEM. We discuss the basic modules of AFEM along with their key properties in §4, and the properties of bisection in §2.

2. The Bisection Method

We briefly discuss the *bisection* method, the most elegant and successful technique for subdividing Ω in any dimension into a conforming mesh made of simplices. We mention the recursive algorithms by Mitchell [24] for $d = 2$ and Kossaczky [21] for $d = 3$. We focus on the special case $d = 2$, and follow Binev, Dahmen, and DeVore [5] and Nochetto and Veiser [29], but the key Theorem 2 holds for any $d \geq 2$ as shown by Stevenson [34]. We refer to Nochetto, Siebert, and Veiser [28] for a rather complete discussion for $d \geq 2$.

2.1. Definition and Properties of Bisection. 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 3.

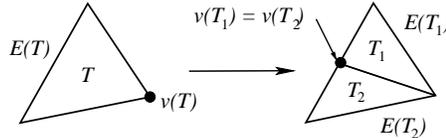


Figure 3. Triangle $T \in \mathcal{T}$ with vertex $v(T)$ and opposite refinement edge $E(T)$. The bisection rule for $d = 2$ consists of connecting $v(T)$ with the midpoint of $E(T)$, thereby giving rise to children T_1, T_2 with common vertex $v(T_1) = v(T_2)$, the newly created vertex, and opposite 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 follow Mitchell [24] and consider the *newest vertex bisection* as depicted in Figure 3. For $d > 2$ the situation is more complicated and one needs the concepts of type and vertex order [21, 28, 34].

Let \mathbb{T} be the set of all conforming bisection refinements of \mathcal{T}_0 . If $\mathcal{T}_* \in \mathbb{T}$ is a conforming refinement of $\mathcal{T} \in \mathbb{T}$, we write $\mathcal{T}_* \geq \mathcal{T}$. For instance, Figure 4 displays a sequence $\{\mathcal{T}_k\}_{k=0}^2$ with $\mathcal{T}_0 = \{T_i\}_{i=1}^4$ and $\mathcal{T}_k \geq \mathcal{T}_{k-1}$ obtained by bisecting the longest edge.

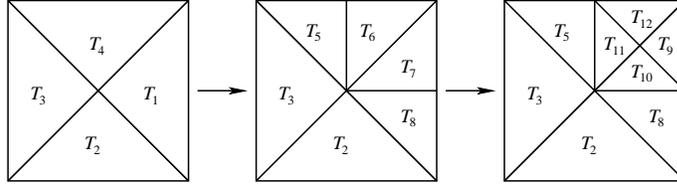


Figure 4. Sequence of bisection meshes $\{\mathcal{T}_k\}_{k=0}^2$ starting from the initial mesh $\mathcal{T}_0 = \{T_i\}_{i=1}^4$ with longest edges labeled for bisection. Mesh \mathcal{T}_1 is created from \mathcal{T}_0 upon bisecting T_1 and T_4 , whereas mesh \mathcal{T}_2 arises from \mathcal{T}_1 upon refining T_6 and T_7 . The bisection rule is described in Figure 3.

The following assertion about element shape is valid for $d \geq 2$ but we state it for $d = 2$.

Lemma 1 (Shape regularity). *The partitions $\mathcal{T} \in \mathbb{T}$ generated by newest vertex bisection satisfy a uniform minimal angle condition, or equivalently the maximal ratio of element diameter over diameter of largest inscribed ball for all $T \in \mathcal{T}$ is uniformly bounded, only depending on the initial partition \mathcal{T}_0 .*

We define the *generation* $g(T)$ 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

$$h_T = |T|^{1/2} = 2^{-g(T)/2} h_{T_0} \quad \text{for all } T \in \mathcal{T}. \quad (5)$$

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 [5], but for $d > 2$ we refer to Condition (b) of Section 4 of [34]. 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 [5]:

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

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 4. Condition (6) can be enforced for $d = 2$ upon bisecting twice each element of \mathcal{T}_0 and labeling 0 the two newest edges [5]. For $d > 2$ the construction is much trickier [34].

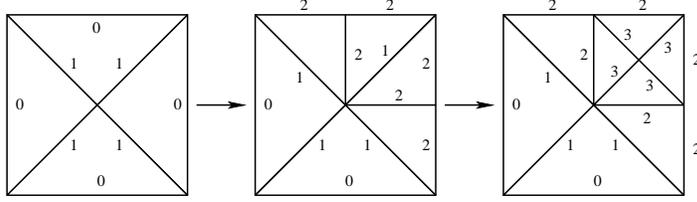


Figure 5. Initial labeling and its evolution for the sequence of conforming refinements of Figure 4.

2.2. Complexity of Bisection. Given $\mathcal{T} \in \mathbb{T}$ and a subset $\mathcal{M} \subset \mathcal{T}$ of marked elements to be refined, 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.

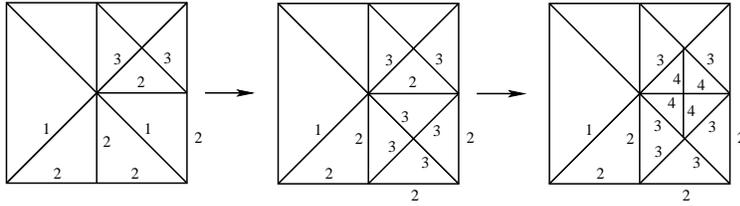


Figure 6. Recursive refinement of $T_{10} \in \mathcal{T}$ in Figures 4 and 5. This entails refining the chain $\{T_{10}, T_8, T_2\}$, starting from the last element $T_2 \in \mathcal{T}$, which form alone a compatible bisection patch because its refinement edge is on the boundary, and continuing with $T_8 \in \mathcal{T}$ and finally $T_{10} \in \mathcal{T}$. Note that the successive meshes are always conforming, that each element in the chain is bisected twice before getting back to T_{10} , and that $\#\{T_{10}, T_8, T_2\} = g(T_{10}) = 3$.

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. Figure 6 shows 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 because the constant may be as large as $g(T)$ with $T \in \mathcal{M}$.

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 [5] for $d = 2$ and Stevenson [34] for $d > 2$. We refer to Nochetto, Siebert and Veiser [28] for a complete discussion for $d \geq 2$.

Theorem 2 (Complexity of REFINE). *If \mathcal{T}_0 satisfies the initial labeling (6) for $d = 2$, or that in [34, 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 2 remains valid with Λ_0 also depending on b .

3. Piecewise Polynomial Interpolation

3.1. Quasi-interpolation. 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 those 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

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

which satisfies $\max_{T \in \mathcal{T}} \#N_{\mathcal{T}}(T) \leq C(\mathcal{T}_0)$ and $\max_{T' \in N_{\mathcal{T}}(T)} \frac{|T|}{|T'|} \leq C(\mathcal{T}_0)$ with $C(\mathcal{T}_0)$ depending only on the shape coefficient of \mathcal{T}_0 . We consider now the *quasi-interpolation operator* $I_{\mathcal{T}} : W_1^1(\Omega) \rightarrow \mathbb{V}(\mathcal{T})$ due to Scott and Zhang [9, 30]. For $n = 1$ it reads

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

where $\{\phi_z^*\}_{z \in \mathcal{N}(\mathcal{T})}$ is a suitable set of dual functions for each node z so that $I_{\mathcal{T}}v = 0$ on $\partial\Omega$ provided $v = 0$ on $\partial\Omega$. We recall the notion of *Sobolev number*: $\text{sob}(W_p^s) = s - d/p$.

Proposition 3 (Local interpolation error). *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 $\text{sob}(W_p^s) > \text{sob}(W_q^t)$. The quasi-interpolation operator $I_{\mathcal{T}}$ is invariant in $\mathbb{V}(\mathcal{T})$ and satisfies for $s \geq 1$*

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

provided \mathcal{T} is shape regular. Moreover, if $\text{sob}(W_p^2) > 0$, then v is continuous and (7) remains valid with $I_{\mathcal{T}}$ replaced by the Lagrange interpolation operator and $N_{\mathcal{T}}(T)$ by T .

3.2. Principle of Error Equidistribution. We investigate the relation between local meshsize and regularity for the design of graded meshes adapted to a given function $v \in H^1(\Omega)$ 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 [4]. Let

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

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

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

be dictated by (7) 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

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

with Lagrange multiplier $\lambda \in \mathbb{R}$. The optimality condition reads $h(x)^{2(p-1)+2} |D^2v(x)|^p = \Lambda$, where $\Lambda > 0$ is a constant. To interpret this expression, we compute the interpolation error E_T incurred in element $T \in \mathcal{T}$. According to Proposition 3, E_T is given by

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

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

$$\text{A graded mesh is quasi-optimal if the local error is equidistributed.} \quad (8)$$

Meshes satisfying (8) have been constructed by Babuška et al [2] for *corner singularities* and $d = 2$; see also [19]. If $0 < \gamma < 1$ and the function v behaves like $v(x) \approx r(x)^\gamma$, where $r(x)$ is the distance from $x \in \Omega$ to a reentrant corner of Ω , then

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

is the optimal mesh grading. This in turn implies

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

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,

$$\|\nabla(v - I_{\mathcal{T}}v)\|_{L^2(\Omega)}^2 = \sum_{T \in \mathcal{T}} E_T^2 = \Lambda^{\frac{2}{p}}(\#\mathcal{T}) \approx (\#\mathcal{T})^{-1} \quad (9)$$

gives the optimal decay rate for $d = 2, n = 1$. 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 $\mathcal{T} \in \mathbb{T}$.

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

$$v \in W_p^2(\Omega; \mathcal{T}_0) \cap H_0^1(\Omega), \quad p > 1. \quad (10)$$

Following Binev, Dahmen, DeVore and Petrushev [6], we use a thresholding algorithm that is based on the knowledge of the element errors and on bisection. The algorithm hinges on (8): 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 , the procedure 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( $\mathcal{T}$ )
    
```

Since $W_p^2(\Omega; \mathcal{T}_0) \cap H_0^1(\Omega) \subset C^0(\bar{\Omega})$, because $p > 1$, we can use the Lagrange interpolant and local estimate (7) with $r = \text{sob}(W_p^2) - \text{sob}(H^1) = 2 - 2/p > 0$ and $N_{\mathcal{T}}(T) = T$:

$$E_T \lesssim h_T^r \|D^2 v\|_{L^p(T)}. \quad (11)$$

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

Theorem 4 (Thresholding). *If v verifies (10), then the output $\mathcal{T} \in \mathbb{T}$ of THRESHOLD satisfies*

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

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

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

We proceed in several steps.

[1] 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

$$2^{-(j+1)} \#\mathcal{P}_j \leq |\Omega| \quad \Rightarrow \quad \#\mathcal{P}_j \leq |\Omega| 2^{j+1}. \quad (12)$$

[2] In light of (11), we have for $T \in \mathcal{P}_j$

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

Therefore

$$\delta^{2p} \#\mathcal{P}_j \lesssim 2^{-(j/2)rp} \sum_{T \in \mathcal{P}_j} \|D^2v\|_{L^p(T)}^p \leq 2^{-(j/2)rp} \|D^2v\|_{L^p(\Omega; \mathcal{T}_0)}^p$$

whence

$$\#\mathcal{P}_j \lesssim \delta^{-2p} 2^{-(j/2)rp} \|D^2v\|_{L^p(\Omega; \mathcal{T}_0)}^p. \quad (13)$$

[3] The two bounds for $\#\mathcal{P}$ in (12) and (13) 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

$$2^{j_0+1} |\Omega| = \delta^{-2p} 2^{-j_0(rp/2)} \|D^2v\|_{L^p(\Omega; \mathcal{T}_0)}^p \quad \Rightarrow \quad 2^{j_0} \approx \delta^{-2} \frac{\|D^2v\|_{L^p(\Omega; \mathcal{T}_0)}}{|\Omega|^{1/p}}.$$

[4] We now compute

$$\#\mathcal{M} = \sum_j \#\mathcal{P}_j \lesssim \sum_{j \leq j_0} 2^j |\Omega| + \delta^{-2p} \|D^2v\|_{L^p(\Omega; \mathcal{T}_0)}^p \sum_{j > j_0} (2^{-rp/2})^j.$$

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

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

We finally apply Theorem 2 to arrive at

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

□ 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. □

Upon relating the threshold δ and the number of elements N , we obtain a convergence rate. In particular, this implies (4): $\sigma_N(v) \lesssim \|D^2v\|_{L^p(\Omega;\mathcal{T}_0)} N^{-1/2}$ for all $N \geq \#\mathcal{T}_0$.

Corollary 5 (Convergence rate). *Let v satisfy (10). Then for $N > \#\mathcal{T}_0$ integer there exists $\mathcal{T} \in \mathbb{T}$ such that*

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

Proof. Choose $\delta^2 = |\Omega|^{1-1/p} \|D^2v\|_{L^p(\Omega;\mathcal{T}_0)} N^{-1}$ in Theorem 4. Then, there exists $\mathcal{T} \in \mathbb{T}$ such that $\#\mathcal{T} - \#\mathcal{T}_0 \lesssim N$ and

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

because $\#\mathcal{T} \lesssim N$. This finishes the Corollary. □

Remark 6 (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 $B_{p,p}^{n+1}(\Omega)$ [17]. The argument of Theorem 4 works provided we replace (11) by a modulus of regularity [6].

Remark 7 (Isotropic elements). Corollary 5 shows that isotropic graded meshes can always deal with geometric singularities for $d = 2$. This is no longer true for $d > 2$ due to *edge singularities*: if $d = 3$ and $v(x) \approx r(x)^\gamma$ near an edge, then $n = 1$ requires $\gamma > \frac{1}{3}$ whereas $n = 2$ needs $\gamma > \frac{2}{3}$. The latter corresponds to a dihedral angle $\omega < \frac{3\pi}{2}$.

4. Adaptive Finite Element Methods (AFEM)

We now present the four basic modules of AFEM for (1) and discuss their main properties.

4.1. Modules of AFEM. They are SOLVE, ESTIMATE, MARK, and REFINE.

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,

$$U \in \mathbb{V} : \int_{\Omega} \mathbf{A} \nabla U \cdot \nabla V = \int_{\Omega} f V \quad \text{for all } V \in \mathbb{V}. \quad (14)$$

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

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

are the element indicators defined as follows: for any $V \in \mathbb{V}$

$$\mathcal{E}_{\mathcal{T}}^2(V, T) = h_T^2 \|r(V)\|_T^2 + h_T \|j(V)\|_{\partial T}^2 \quad \text{for all } T \in \mathcal{T}, \quad (15)$$

where the *interior* and *jump residuals* are given by

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

and $j(V)|_S = 0$ for boundary sides $S \in \mathcal{S}$. We denote $\mathcal{E}_{\mathcal{T}}^2(V, \mathcal{P}) = \sum_{T \in \mathcal{P}} \mathcal{E}_{\mathcal{T}}^2(V, T)$ for any subset \mathcal{P} of \mathcal{T} and $\mathcal{E}_{\mathcal{T}}(V) = \mathcal{E}_{\mathcal{T}}(V, \mathcal{T})$.

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) [18], i. e., using a fixed parameter $\theta \in (0, 1]$ the output \mathcal{M} of

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

satisfies

$$\mathcal{E}_{\mathcal{T}}(U, \mathcal{M}) \geq \theta \mathcal{E}_{\mathcal{T}}(U, \mathcal{T}). \quad (16)$$

This marking guarantees that \mathcal{M} contains a substantial part of the total (or bulk), thus its name. 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.

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, the piecewise constant meshsize functions satisfy $h_{\mathcal{T}_*} \leq h_{\mathcal{T}}$ and the strict reduction property

$$h_{\mathcal{T}_*}|_T \leq 2^{-b/d} h_{\mathcal{T}}|_T \quad \text{for all } T \in \mathcal{M}. \quad (17)$$

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. Given an initial grid \mathcal{T}_0 , set $k = 0$ and iterate

$$\begin{aligned} U_k &= \text{SOLVE}(\mathcal{T}_k); \\ \{\mathcal{E}_k(U_k, T)\}_{T \in \mathcal{T}_k} &= \text{ESTIMATE}(U_k, \mathcal{T}_k); \\ \mathcal{M}_k &= \text{MARK}(\{\mathcal{E}_k(U_k, T)\}_{T \in \mathcal{T}_k, \mathcal{T}_k}); \\ \mathcal{T}_{k+1} &= \text{REFINE}(\mathcal{T}_k, \mathcal{M}_k); k \leftarrow k + 1. \end{aligned}$$

4.2. Basic Properties of AFEM. We next follow Cascón, Kreuzer, Nochetto, and Siebert [10] and 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 (1).

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

Lemma 8 (Pythagoras). *Let $\mathcal{T}, \mathcal{T}_* \in \mathbb{T}$ be 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*

$$\|u - U\|_{\Omega}^2 = \|u - U_*\|_{\Omega}^2 + \|U_* - U\|_{\Omega}^2. \quad (18)$$

Property (18) is valid for (1) 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 continue with the concept of *oscillation*. We denote by $\text{osc}_{\mathcal{T}}(V, T)$ the *element oscillation* for any $V \in \mathbb{V}$

$$\text{osc}_{\mathcal{T}}(V, T) = \|h(r(V) - \overline{r(V)})\|_{L^2(T)} + \|h^{1/2}(j(V) - \overline{j(V)})\|_{L^2(\partial T \cap \Omega)}, \quad (19)$$

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, $\text{osc}_{\mathcal{T}}(V, T) = \|h(f - \bar{f})\|_{L^2(T)}$ is called *data oscillation* [25, 26].

Proposition 9 (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*

$$\|u - U\|_{\Omega}^2 \leq C_1 \mathcal{E}_{\mathcal{T}}^2(U) \quad (20a)$$

$$C_2 \mathcal{E}_{\mathcal{T}}^2(U) \leq \|u - U\|_{\Omega}^2 + \text{osc}_{\mathcal{T}}^2(U). \quad (20b)$$

This Proposition is essentially due to Babuška and Miller [3]; see also [1, 8, 28, 35]. The constants C_1 and C_2 depend on the smallest and largest global eigenvalues of \mathbf{A} as well as interpolation estimates. The definitions of $\overline{r(V)}$ and $\overline{j(V)}$, as well as the lower bound (20b), are immaterial for deriving a contraction property of §5 but are important for proving convergence rates in §6; we refer to [28] for a discussion of oscillation.

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. This is tackled in the next two lemmas [10, 27].

Lemma 10 (Reduction of $\mathcal{E}_{\mathcal{T}}(V)$ with respect to \mathcal{T}). *If $\lambda = 1 - 2^{-b/d}$, then*

$$\mathcal{E}_{\mathcal{T}_*}^2(V, \mathcal{T}_*) \leq \mathcal{E}_{\mathcal{T}}^2(V, \mathcal{T}) - \lambda \mathcal{E}_{\mathcal{T}}^2(V, \mathcal{M}) \quad \text{for all } V \in \mathbb{V}(\mathcal{T}). \quad (21)$$

Lemma 11 (Lipschitz property of $\mathcal{E}_{\mathcal{T}}(V)$ with respect to V). *Let $\text{div } \mathbf{A}$ be the divergence of \mathbf{A} computed by rows, and $\eta_{\mathcal{T}}(\mathbf{A}) := \max_{\mathcal{T}} (h_T \|\text{div } \mathbf{A}\|_{L^\infty(T)} + \|\mathbf{A}\|_{L^\infty(T)})$. Then the following estimate is valid*

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

Upon combining Lemmas 10 and 11 we obtain the following crucial property.

Proposition 12 (Estimator reduction). *Given $\mathcal{T} \in \mathbb{T}$ and a subset $\mathcal{M} \subset \mathcal{T}$ of marked elements, let $\mathcal{T}_* = \text{REFINE}(\mathcal{T}, \mathcal{M})$. 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*

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

5. Contraction Property of AFEM

A key question to ask is what is (are) the quantity(ies) that AFEM may contract. In light of (18), an obvious candidate is the energy error $\|u - U_k\|_\Omega$; see Dörfler [18]. We first show in §5.1, 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. However, the energy error may not contract in general unless REFINE enforces several levels of refinement. We discuss this in §5.2, and present an approach that eliminates the interior node property at the expense of a more complicated contractive quantity, the quasi-error; see Theorem 16.

5.1. 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 [25].

Definition 13 (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 [25, 26], is the following *discrete* lower a posteriori bound:

$$C_2 \mathcal{E}_k^2(U_k, \mathcal{M}_k) \leq \|U_k - U_{k+1}\|_\Omega^2 + \text{osc}_k^2(U_k). \quad (22)$$

Lemma 14 (Contraction property for piecewise constant data). *If \mathcal{T}_{k+1} satisfies an interior node property with respect to \mathcal{T}_k and $\text{osc}_k(U_k) = 0$, then for $\alpha := (1 - \theta^2 \frac{C_2}{C_1})^{1/2} < 1$*

$$\|u - U_{k+1}\|_\Omega \leq \alpha \|u - U_k\|_\Omega, \quad (23)$$

where $0 < \theta < 1$ is the parameter in (16) and $C_1 \geq C_2$ are the constants in (20).

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 (18), namely $e_{k+1}^2 = e_k^2 - E_k^2$, and show that E_k is a significant portion of e_k . Since (22) together with $\text{osc}_k(U_k) = 0$ imply $C_2 \mathcal{E}_k^2(\mathcal{M}_k) \leq E_k^2$, applying Dörfler marking (16) and the upper bound (20a), 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 (23). \square

We wonder whether or not the interior node property is necessary for (23). We present an example, introduced in [25, 26] to justify such a property for constant data and $n = 1$.

Example 15 (Lack of 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 7. If ϕ_0 denotes the basis function associated with the only interior node of the initial mesh \mathcal{T}_0 , then $U_0 = U_1 = \frac{1}{12} \phi_0$ and $U_2 \neq U_1$.

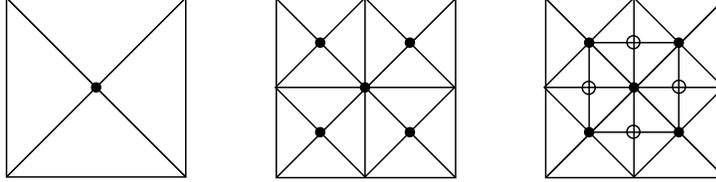


Figure 7. Grids \mathcal{T}_0 , \mathcal{T}_1 , and \mathcal{T}_2 of Example 15. The mesh \mathcal{T}_1 has nodes in the middle of sides of \mathcal{T}_0 , but only \mathcal{T}_2 has nodes in the interior of elements of \mathcal{T}_0 . Hence, \mathcal{T}_2 satisfies the interior node property of Definition 13 with respect to \mathcal{T}_0 whereas \mathcal{T}_1 does not.

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 does 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, because then Lemma 14 holds.

5.2. 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 (18), 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 (21) 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*

$$\|u - U_k\|_{\Omega}^2 + \gamma \mathcal{E}_k^2(U_k) \quad (24)$$

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 16 (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 (1). 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 $k \geq 0$*

$$\|u - U_{k+1}\|_{\Omega}^2 + \gamma \mathcal{E}_{k+1}^2(U_{k+1}) \leq \alpha^2 \left(\|u - U_k\|_{\Omega}^2 + \gamma \mathcal{E}_k^2(U_k) \right). \quad (25)$$

Proof. We split the proof into four steps and use the notation in Lemma 14.

□ The error orthogonality (18) reads

$$e_{k+1}^2 = e_k^2 - E_k^2. \quad (26)$$

Employing Proposition 12 with $\mathcal{T} = \mathcal{T}_k$, $\mathcal{T}_* = \mathcal{T}_{k+1}$, $V = U_k$ and $V_* = U_{k+1}$ gives

$$\mathcal{E}_{k+1}^2 \leq (1 + \delta)(\mathcal{E}_k^2 - \lambda \mathcal{E}_k^2(\mathcal{M}_k)) + (1 + \delta^{-1}) \Lambda_0 E_k^2, \quad (27)$$

where $\Lambda_0 = \Lambda \eta_{T_0}^2(\mathbf{A}) \geq \Lambda \eta_{T_k}^2(\mathbf{A})$. After multiplying (27) by $\gamma > 0$, to be determined later, we add (26) and (27) to obtain

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

□ 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

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

□ 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 (20a), 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 $\alpha^2 := \max\left\{1 - \frac{\gamma\lambda\theta^2}{4C_1}, 1 - \frac{\lambda\theta^2}{4}\right\} < 1$, and thus concludes the theorem. □

Remark 17 (Basic ingredients). This proof solely uses Dörfler marking, Pythagoras identity (18), the a posteriori upper bound (20a), and the estimator reduction property (Proposition 12). The proof does not use the lower bound (20b).

Remark 18 (Separate marking). MARK is driven by \mathcal{E}_k exclusively, as it happens in all practical AFEM. Previous proofs in [14, 23, 25, 26] require separate marking by estimator and oscillation. It is shown in [10] that separate marking may lead to suboptimal convergence rates. On the other hand, we will prove in §6 that the present AFEM yields quasi-optimal convergence rates.

6. Convergence Rates of AFEM

A crucial insight for the simplest scenario, the Laplacian and piecewise constant forcing f , is due to Stevenson [33]:

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. (28)

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 §6.3. However, this is not enough to handle the model problem (1) with variable data \mathbf{A} and f .

The objective of this section is to study (1) 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 [33] to variable data. In doing so, we follow closely Cascón, Kreuzer, Nochetto, and Siebert [10]. 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 (1) 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 $\mathcal{T} \in \mathbb{T}$.

6.1. The Total Error. We first present the concept of *total error* for the Galerkin function $U \in \mathbb{V}(\mathcal{T})$, introduced by Mekchay and Nochetto [23],

$$\|u - U\|_{\Omega}^2 + \text{osc}_{\mathcal{T}}^2(U), \quad (29)$$

and next assert its equivalence to the quasi error (24). In fact, in view of the upper and lower a posteriori error bounds (20), and $\text{osc}_{\mathcal{T}}^2(U) \leq \mathcal{E}_{\mathcal{T}}^2(U)$, we have

$$C_2 \mathcal{E}_{\mathcal{T}}^2(U) \leq \|u - U\|_{\Omega}^2 + \text{osc}_{\mathcal{T}}^2(U) \leq \|u - U\|_{\Omega}^2 + \mathcal{E}_{\mathcal{T}}^2(U) \leq (1 + C_1) \mathcal{E}_{\mathcal{T}}^2(U),$$

whence

$$\mathcal{E}_{\mathcal{T}}^2(U) \approx \|u - U\|_{\Omega}^2 + \text{osc}_{\mathcal{T}}^2(U). \quad (30)$$

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 (20a) 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 16). Finally, the total error satisfies the following Cea's type-lemma, or equivalently AFEM is quasi-optimal regarding the total error [10].

Lemma 19 (Quasi-optimality of total error). *Let $\Lambda_1 = 2\Lambda$ with Λ the constant in Proposition 12, and let $C_3 := \Lambda_1 \text{osc}_{\mathcal{T}_0}^2(\mathbf{A})$ and $\Lambda_2 := \max\{2, 1 + C_3\}$. Then, for any $\mathcal{T} \in \mathbb{T}$ and the corresponding Galerkin solution $U \in \mathbb{V}(\mathcal{T})$, there holds*

$$\|u - U\|_{\Omega}^2 + \text{osc}_{\mathcal{T}}^2(U) \leq \Lambda_2 \inf_{V \in \mathbb{V}(\mathcal{T})} \left(\|u - V\|_{\Omega}^2 + \text{osc}_{\mathcal{T}}^2(V) \right).$$

6.2. Approximation Classes. In view of (30) and Lemma 19, the definition of approximation class \mathbb{A}_s depends on the triple (u, f, \mathbf{A}) , not just u , and hinges on the concept of best total error for meshes \mathcal{T} with N elements more than \mathcal{T}_0 , namely $\mathcal{T} \in \mathbb{T}_N$:

$$\sigma_N(u, f, \mathbf{A}) := \inf_{\mathcal{T} \in \mathbb{T}_N} \inf_{V \in \mathbb{V}(\mathcal{T})} \left(\|u - V\|_{\Omega}^2 + \text{osc}_{\mathcal{T}}^2(V) \right)^{1/2}.$$

We say that $(u, f, \mathbf{A}) \in \mathbb{A}_s$ for $s > 0$ if and only if $\sigma_N(u, f, \mathbf{A}) \lesssim N^{-s}$, and denote $|u, f, \mathbf{A}|_s := \sup_{N > 0} (N^s \sigma_N(u, f, \mathbf{A}))$. We point out the upper bound $s \leq n/d$ for polynomial degree $n \geq 1$; this can be seen with full regularity $H^{n+1}(\Omega)$ and uniform refinement. Note that if $(u, f, \mathbf{A}) \in \mathbb{A}_s$ 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

$$\|v - V_{\varepsilon}\|_{\Omega}^2 + \text{osc}_{\mathcal{T}_{\varepsilon}}^2(V_{\varepsilon}) \leq \varepsilon^2 \quad \text{and} \quad \#\mathcal{T}_{\varepsilon} - \#\mathcal{T}_0 \leq |v, f, \mathbf{A}|_s^{1/s} \varepsilon^{-1/s}. \quad (31)$$

Mesh Overlay. For the subsequent discussion it will be convenient to merge (or superpose) two conforming meshes $\mathcal{T}_1, \mathcal{T}_2 \in \mathbb{T}$, thereby giving rise to the so-called *overlay* $\mathcal{T}_1 \oplus \mathcal{T}_2$. This operation corresponds to the union in the sense of trees [10, 33]. 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 [10, 33].

Lemma 20 (Overlay). *The overlay $\mathcal{T} = \mathcal{T}_1 \oplus \mathcal{T}_2$ is conforming and*

$$\#\mathcal{T} \leq \#\mathcal{T}_1 + \#\mathcal{T}_2 - \#\mathcal{T}_0. \quad (32)$$

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 \mathbf{A} piecewise polynomial of degree $\leq n$ over \mathcal{T}_0 . In this simplified scenario, the oscillation $\text{osc}_{\mathcal{T}}(U)$ of (19) reduces to *data oscillation* $\text{osc}_{\mathcal{T}}(f) :=$

$\|h(f - P_{2n-2}f)\|_{L^2(\Omega)}$. We then have the following characterization of \mathbb{A}_s in terms of the approximation class \mathcal{A}_s and [5, 6, 33]:

$$\mathcal{B}_s := \left\{ g \in L^2(\Omega) \mid |g|_{\mathcal{B}_s} := \sup_{N>0} (N^s \inf_{\mathcal{T} \in \mathbb{T}_N} \text{osc}_{\mathcal{T}}(g)) < \infty \right\}.$$

Lemma 21 (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 $(u, f) \in \mathcal{A}_s \times \mathcal{B}_s$ and*

$$|u, f, \mathbf{A}|_s \approx |u|_{\mathcal{A}_s} + |f|_{\mathcal{B}_s}. \quad (33)$$

Corollary 22 (Membership in $\mathbb{A}_{1/2}$ with piecewise constant \mathbf{A}). *Let $d = 2$, $n = 1$, $p > 1$. If $f \in L^2(\Omega)$, \mathbf{A} is piecewise constant over \mathcal{T}_0 , and the solution $u \in W_p^2(\Omega; \mathcal{T}_0) \cap H_0^1(\Omega)$ of (1) is piecewise W_p^2 over the initial grid \mathcal{T}_0 , then $(u, f, \mathbf{A}) \in \mathbb{A}_{1/2}$ and*

$$|u, f, \mathbf{A}|_{1/2} \lesssim \|D^2u\|_{L^p(\Omega; \mathcal{T}_0)} + \|f\|_{L^2(\Omega)}.$$

Proof. Since $f \in L^2(\Omega)$, we realize that for all quasi-uniform refinements $\mathcal{T} \in \mathbb{T}$

$$\text{osc}_{\mathcal{T}}(f) = \|h(f - P_0f)\|_{L^2(\Omega)} \leq h_{\max}(\mathcal{T})\|f\|_{L^2(\Omega)} \lesssim (\#\mathcal{T})^{-1/2}\|f\|_{L^2(\Omega)}.$$

This implies $f \in \mathcal{B}_{1/2}$ with $|f|_{\mathcal{B}_{1/2}} \lesssim \|f\|_{L^2(\Omega)}$. On the other hand, for $u \in W_p^2(\Omega; \mathcal{T}_0)$ we learn from Corollary 5 that $u \in \mathcal{A}_{1/2}$ and $|u|_{\mathcal{A}_{1/2}} \lesssim \|D^2u\|_{L^2(\Omega; \mathcal{T}_0)}$. The assertion then follows from Lemma 21. \square

Corollary 23 (Membership in $\mathbb{A}_{1/2}$ with variable \mathbf{A}). *Let $d = 2$, $n = 1$, $p > 1$. If $f \in L^2(\Omega)$, $\mathbf{A} \in W_{\infty}^1(\Omega, \mathcal{T}_0)$ is piecewise Lipschitz over \mathcal{T}_0 , and $u \in W_p^2(\Omega; \mathcal{T}_0) \cap H_0^1(\Omega)$ is piecewise W_p^2 over \mathcal{T}_0 , then $(u, f, \mathbf{A}) \in \mathbb{A}_{1/2}$ and*

$$|u, f, \mathbf{A}|_{1/2} \lesssim \|D^2u\|_{L^p(\Omega; \mathcal{T}_0)} + \|f\|_{L^2(\Omega)} + \|\mathbf{A}\|_{W_{\infty}^1(\Omega; \mathcal{T}_0)}.$$

6.3. Quasi-Optimal Cardinality: Vanishing Oscillation. In this section we follow the ideas of Stevenson [33] for the simplest scenario with vanishing oscillation $\text{osc}_{\mathcal{T}}(U) = 0$, and thereby explore the insight (28). We recall that in this case the a posteriori error estimates (20) become

$$C_2 \mathcal{E}_{\mathcal{T}}^2(U) \leq \|u - U\|_{\Omega}^2 \leq C_1 \mathcal{E}_{\mathcal{T}}^2(U). \quad (34)$$

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 for which this theory extends [12, 22].

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 24 (Marking parameter: vanishing oscillation). *The parameter θ of Dörfler marking satisfies $\theta \in (0, \theta_*)$ with $\theta_* := \sqrt{C_2/C_1}$.*

Assumption 25 (Cardinality of \mathcal{M}). *MARK selects a set \mathcal{M} with minimal cardinality.*

Assumption 26 (Initial labeling). *The labeling of the initial mesh \mathcal{T}_0 satisfies (6) for $d = 2$ [24, 5] or its multidimensional counterpart for $d > 2$ [33, 28].*

A few comments about these assumptions are now in order.

Remark 27 (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 $\theta_* < 1$ in Assumption 24.

Remark 28 (Minimal \mathcal{M}). According to the equidistribution principle (8) and the local lower bound $C_2 \mathcal{E}_{\mathcal{T}}(U, T) \leq \|u - U\|_{N_{\mathcal{T}}(T)}$ without oscillation, it is natural to mark elements with largest error indicators. This leads to a minimal set \mathcal{M} , as stated in Assumption 25, and turns out to be crucial to link AFEM with optimal meshes.

Remark 29 (Initial triangulation). Assumption 26 guarantees the complexity estimate of module REFINE stated in Theorem 2: $\#\mathcal{T}_k - \#\mathcal{T}_0 \leq \Lambda_0 \sum_{j=0}^{k-1} \#\mathcal{M}_j$.

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}_* [33, 10, 28].

Lemma 30 (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*

$$\|U_* - U\|_{\Omega}^2 \leq C_1 \mathcal{E}_{\mathcal{T}}^2(U, \mathcal{R}). \quad (35)$$

We are now ready to explore Stevenson's insight (28) for the simplest scenario with vanishing oscillation $\text{osc}_{\mathcal{T}}(U) = 0$.

Lemma 31 (Dörfler marking: vanishing oscillation). *Let θ satisfy Assumption 24 and set $\mu := 1 - \theta^2/\theta_*^2 > 0$. Let $\mathcal{T}_* \geq \mathcal{T}$ and $U_* \in \mathbb{V}(\mathcal{T}_*)$ satisfy*

$$\|u - U_*\|_{\Omega}^2 \leq \mu \|u - U\|_{\Omega}^2. \quad (36)$$

Then the refined set $\mathcal{R} = \mathcal{R}_{\mathcal{T} \rightarrow \mathcal{T}_}$ satisfies the Dörfler property*

$$\mathcal{E}_{\mathcal{T}}(U, \mathcal{R}) \geq \theta \mathcal{E}_{\mathcal{T}}(U, \mathcal{T}). \quad (37)$$

Proof. Since $\mu < 1$ we use the lower bound in (34), in conjunction with (36) and Pythagoras equality (18), to derive

$$(1 - \mu)C_2\mathcal{E}_{\mathcal{T}}^2(U, \mathcal{T}) \leq (1 - \mu)\|u - U\|_{\Omega}^2 \leq \|u - U\|_{\Omega}^2 - \|u - U_*\|_{\Omega}^2 = \|U - U_*\|_{\Omega}^2.$$

In view of Lemma 30, 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 31. We show this now.

Lemma 32 (Cardinality of \mathcal{M}_k). *Let Assumptions 24 and 25 hold. If $u \in \mathcal{A}_s$ then*

$$\#\mathcal{M}_k \lesssim |u|_s^{1/s} \|u - U_k\|_{\Omega}^{-1/s} \quad \text{for all } k \geq 0. \quad (38)$$

Proof. We invoke that $u \in \mathcal{A}_s$ and (31) 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

$$\|u - U_{\varepsilon}\|_{\Omega}^2 \leq \varepsilon^2, \quad \#\mathcal{T}_{\varepsilon} - \#\mathcal{T}_0 \lesssim |u|_s^{\frac{1}{s}} \varepsilon^{-\frac{1}{s}}.$$

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

$$\|u - U_*\|_{\Omega}^2 \leq \|u - U_{\varepsilon}\|_{\Omega}^2 \leq \varepsilon^2 = \mu \|u - U\|_{\Omega}^2.$$

Therefore, Lemma 31 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 25, whence

$$\#\mathcal{M}_k \leq \#\mathcal{R} \leq \#\mathcal{T}_* - \#\mathcal{T}_k \leq \#\mathcal{T}_{\varepsilon} - \#\mathcal{T}_0 \lesssim |u|_s^{\frac{1}{s}} \varepsilon^{-\frac{1}{s}},$$

where we use Lemma 20 to account for the overlay. The proof is complete. \square

Proposition 33 (Quasi-optimality: vanishing oscillation). *Let Assumptions 24-26 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*

$$\|u - U_k\|_{\Omega} \lesssim |u|_s (\#\mathcal{T}_k - \#\mathcal{T}_0)^{-s} \quad \text{for all } k \geq 1.$$

Proof. We make use of Assumption 26, along with Theorem 2, to infer that

$$\#\mathcal{T}_k - \#\mathcal{T}_0 \leq \Lambda_0 \sum_{j=0}^{k-1} \#\mathcal{M}_j \lesssim |u|_s^{\frac{1}{s}} \sum_{j=0}^{k-1} \|u - U_j\|_{\Omega}^{-\frac{1}{s}}.$$

We now use the contraction property $\|u - U_k\|_\Omega \leq \alpha^{k-j} \|u - U_j\|_\Omega$ of Lemma 14 to replace the sum above by

$$\sum_{j=0}^{k-1} \|u - U_j\|_\Omega^{-\frac{1}{s}} \leq \|u - U_k\|_\Omega^{-\frac{1}{s}} \sum_{j=0}^{k-1} \alpha^{\frac{k-j}{s}} < \frac{\alpha^{\frac{1}{s}}}{1 - \alpha^{\frac{1}{s}}} \|u - U_k\|_\Omega^{-\frac{1}{s}},$$

because $\alpha < 1$ and the series is summable. This completes the proof. \square

6.4. 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 §6.1 and §6.2. 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}) ; see (31) for the definition of \mathbb{A}_s . To account for the presence of general data f and \mathbf{A} , we need to make an even more stringent assumption on the threshold θ_* .

Assumption 34 (Marking parameter: general data). *Let $C_3 = \Lambda_1 \text{osc}_{\mathcal{T}_0}^2(\mathbf{A})$ be the constant in Lemma 19. The marking parameter θ satisfies $\theta \in (0, \theta_*)$ with*

$$\theta_* = \sqrt{\frac{C_2}{1 + C_1(1 + C_3)}}.$$

We now proceed along the same lines as those of §6.3.

Lemma 35 (Dörfler marking: general data). *Let Assumption 34 hold and set $\mu := \frac{1}{2}(1 - \frac{\theta^2}{\theta_*^2}) > 0$. If $\mathcal{T}_* \geq \mathcal{T}$ and $U_* \in \mathbb{V}(\mathcal{T}_*)$ satisfy*

$$\|u - U_*\|_\Omega^2 + \text{osc}_{\mathcal{T}_*}^2(U_*) \leq \mu(\|u - U\|_\Omega^2 + \text{osc}_{\mathcal{T}}^2(U)), \quad (39)$$

then the refined set $\mathcal{R} = \mathcal{R}_{\mathcal{T} \rightarrow \mathcal{T}_*}$ satisfies the Dörfler property

$$\mathcal{E}_{\mathcal{T}}(U, \mathcal{R}) \geq \theta \mathcal{E}_{\mathcal{T}}(U, \mathcal{T}). \quad (40)$$

Proof. We split the proof into four steps.

\square In view of the global lower bound (20b) and (39), we can write

$$\begin{aligned} (1 - 2\mu) C_2 \mathcal{E}_{\mathcal{T}}^2(U) &\leq (1 - 2\mu)(\|u - U\|_\Omega^2 + \text{osc}_{\mathcal{T}}^2(U)) \\ &\leq (\|u - U\|_\Omega^2 - 2\|u - U_*\|_\Omega^2) + (\text{osc}_{\mathcal{T}}^2(U) - 2\text{osc}_{\mathcal{T}_*}^2(U_*)). \end{aligned}$$

\square Combining the Pythagoras orthogonality relation (18)

$$\|u - U\|_\Omega^2 - \|u - U_*\|_\Omega^2 = \|U - U_*\|_\Omega^2.$$

with the localized upper bound (35) yields

$$\|u - U\|_\Omega^2 - 2\|u - U_*\|_\Omega^2 \leq \|U - U_*\|_\Omega^2 \leq C_1 \mathcal{E}_{\mathcal{T}}^2(U, \mathcal{R}).$$

[3] 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 a variant of Lemma 11 for $\text{osc}_{\mathcal{T}}(U)$ together with Lemma 30, to get

$$\text{osc}_{\mathcal{T}}^2(U, \mathcal{T} \setminus \mathcal{R}) - 2 \text{osc}_{\mathcal{T}_*}^2(U_*, \mathcal{T} \setminus \mathcal{R}) \leq C_3 \|U - U_*\|_{\Omega}^2 \leq C_1 C_3 \mathcal{E}_{\mathcal{T}}^2(U, \mathcal{R}).$$

Adding these two estimates gives

$$\text{osc}_{\mathcal{T}}^2(U) - 2 \text{osc}_{\mathcal{T}_*}^2(U_*) \leq (1 + C_1 C_3) \mathcal{E}_{\mathcal{T}}^2(U, \mathcal{R}).$$

[4] Returning to [1] 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 (40) in disguise. \square

Lemma 36 (Cardinality of \mathcal{M}_k : general data). *Let Assumptions 25 and 34 hold. If the triple $(u, f, \mathbf{A}) \in \mathbb{A}_s$, then*

$$\#\mathcal{M}_k \lesssim |u, f, \mathbf{A}|_s^{1/s} (\|u - U_k\|_{\Omega} + \text{osc}_k(U_k))^{-1/s} \quad \text{for all } k \geq 0. \quad (41)$$

Proof. We split the proof into three steps.

[1] 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 35 and Λ_2 given Lemma 19. Since $(u, f, \mathbf{A}) \in \mathbb{A}_s$, in view of (31) 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 $\mathcal{T}_* = \mathcal{T}_k \oplus \mathcal{T}_{\varepsilon}$.

[2] 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 19 to obtain

$$\begin{aligned} \|u - U_*\|_{\Omega}^2 + \text{osc}_{\mathcal{T}_*}^2(U_*) &\leq \Lambda_2 \left(\|u - U_{\varepsilon}\|_{\Omega}^2 + \text{osc}_{\varepsilon}^2(U_{\varepsilon}) \right) \\ &\leq \Lambda_2 \varepsilon^2 = \mu (\|u - U_k\|_{\Omega}^2 + \text{osc}_k^2(U_k)). \end{aligned}$$

Upon applying Lemma 35 we conclude that the set $\mathcal{R} = \mathcal{R}_{\mathcal{T}_k \rightarrow \mathcal{T}_*}$ of refined elements satisfies a Dörfler marking (40) with parameter $\theta < \theta_*$.

[3] According to Assumption 25, MARK selects a minimal set \mathcal{M}_k satisfying this property. Therefore, employing Lemma 20 to account for the cardinality of the overlay, we deduce

$$\#\mathcal{M}_k \leq \#\mathcal{R} \leq \#\mathcal{T}_* - \#\mathcal{T}_k \leq \#\mathcal{T}_\varepsilon - \#\mathcal{T}_0 \lesssim |u, f, \mathbf{A}|_s^{1/s} \varepsilon^{-1/s}.$$

Finally, recalling the definition of ε we end up with the asserted estimate (41). \square

We are ready to prove the main result of this section, which combines Theorem 16 and Lemma 36.

Theorem 37 (Quasi-optimality: general data). *Let Assumptions 25, 26 and 34 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*

$$\|u - U_k\|_\Omega + \text{osc}_k(U_k) \lesssim |u, f, \mathbf{A}|_s (\#\mathcal{T}_k - \#\mathcal{T}_0)^{-s} \quad \text{for all } k \geq 1.$$

Proof. Since no confusion arises, we use the notation $\text{osc}_j = \text{osc}_j(U_j)$ and $\mathcal{E}_j = \mathcal{E}_j(U_j)$.

[1] In light of Assumption 26, which yields Theorem 2, and (41) we have

$$\#\mathcal{T}_k - \#\mathcal{T}_0 \lesssim \sum_{j=0}^{k-1} \#\mathcal{M}_j \lesssim |u, f, \mathbf{A}|_s^{1/s} \sum_{j=0}^{k-1} (\|u - U_j\|_\Omega^2 + \text{osc}_j^2)^{-1/(2s)}.$$

[2] Let $\gamma > 0$ be the scaling factor in the (contraction) Theorem 16. The lower bound (20b) along with $\text{osc}_j \leq \mathcal{E}_j$ implies

$$\|u - U_j\|_\Omega^2 + \gamma \text{osc}_j^2 \leq \|u - U_j\|_\Omega^2 + \gamma \mathcal{E}_j^2 \leq \left(1 + \frac{\gamma}{C_2}\right) (\|u - U_j\|_\Omega^2 + \text{osc}_j^2).$$

[3] Theorem 16 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 $\sum_{j=0}^{k-1} \alpha^{(k-j)/s} < \sum_{j=1}^\infty \alpha^{j/s} < \infty$ because $\alpha < 1$, the assertion follows easily. \square

We conclude this section with a couple of applications of Theorem 37. The first one is valid for the example of §1.

Corollary 38 (W_p^2 -regularity with piecewise constant \mathbf{A}). *Let $d = 2$, the polynomial degree be $n = 1$, $f \in L^2(\Omega)$, and let \mathbf{A} be piecewise constant 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 22 with Theorem 37. \square

Corollary 39 (W_p^2 -regularity with variable \mathbf{A}). *Besides the assumptions of Corollary 38, let \mathbf{A} be piecewise Lipschitz over the initial grid \mathcal{T}_0 . Then AFEM gives rise to a sequence $\{\mathcal{T}_k, \mathbb{V}_k, U_k\}_{k=0}^\infty$ 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 23 with Theorem 37. \square

7. Extensions and Limitations

Nonconforming Meshes. Bonito and Nochetto [7] have shown that Theorem 2 extends to *admissible* nonconforming meshes for $d \geq 2$ (those with a fixed level of nonconformity), along with the theory of §5 and §6.

Discontinuous Galerkin Methods (dG). Bonito and Nochetto [7] have also shown that such theory extends to the interior penalty dG method for the model problem (1) and for $d \geq 2$. This relies on a result of independent interest:

the approximation classes for discontinuous and continuous elements of any degree $n \geq 1$ coincide.

Non-residual Estimators. Cascón and Nochetto [12] and Kreuzer and Siebert [22] have extended the above theory to non-residual estimators (hierarchical estimators, Zienkiewicz-Zhu and Braess-Schoerbel estimators, and those based on the solution of local problems).

Other Norms. The above theory is just for the energy norm. We refer to Demlow [15] for local energy norms and Demlow and Stevenson [16] for the L^2 -norm. The theory for more practical norms, such as L^∞ or W_∞^1 , is open.

Other Problems and Markings. The theory above relies strongly on the Pythagoras equality (18) and Dörfler marking (16), and extends to symmetric

problems in $H(\text{div})$ [11] and $H(\text{curl})$ [37] as well as to non-symmetric coercive problems [12]. For non-coercive problems, for which we just have an inf-sup condition, as well as markings other than Dörfler, the theory is mostly lacking except for mixed AFEM for (1) [13]. We refer to Morin, Siebert, and Vesser [27] and Siebert [31] for convergence results without rates.

Multilevel Methods on Graded Meshes. We refer to Xu, Chen, and Nochetto [36] for a theory of multilevel methods on graded meshes created by bisection. The analysis uses several geometric properties of bisection, discussed in [36], and is valid for any d and n .

References

- [1] M. AINSWORTH AND J.T. ODEN, *A Posteriori Error Estimation in Finite Element Analysis*, Pure and Applied Mathematics (New York), Wiley-Interscience, 2000.
- [2] 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.
- [3] 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.
- [4] I. BABUŠKA AND W. RHEINBOLDT, *Error estimates for adaptive finite element computations* SIAM J. Numer. Anal., 15 (1978), pp. 736–754.
- [5] P. BINEV, W. DAHMEN, AND R. DEVORE, *Adaptive finite element methods with convergence rates*, Numer. Math., 97 (2004), pp. 219–268.
- [6] 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.
- [7] A. BONITO AND R.H. NOCHETTO, *Quasi-optimal convergence rate for an adaptive discontinuous Galerkin method*, SIAM J. Numer. Anal. (to appear).
- [8] D. BRAESS, *Finite Elements. Theory, Fast Solvers, and Applications in Solid Mechanics*, 2nd edition edn. Cambridge University Press (2001).
- [9] S. BRENNER AND L.R. SCOTT, *The Mathematical Theory of Finite Element Methods*, Springer Texts in Applied Mathematics 15 (2008).
- [10] 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.
- [11] J. M. CASCÓN, L. CHEN, R. H. NOCHETTO, AND K. G. SIEBERT, *Quasi-optimal convergence rates for AFEM in $H(\text{div})$* , (in preparation).
- [12] J. M. CASCÓN AND R. H. NOCHETTO, *Convergence and quasi-optimality for AFEM based on non-residual a posteriori error estimators*, (in preparation).

-
- [13] L. CHEN, M. HOLST, AND J. XU, *Convergence and optimality of adaptive mixed finite element methods*, Math. Comp., 78 (2009), pp. 35–53.
- [14] 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.
- [15] A. DEMLOW, *Convergence of an adaptive finite element method for controlling local energy errors*, (submitted).
- [16] A. DEMLOW AND R. STEVENSON, *Convergence and quasi-optimality of an adaptive finite element method for controlling L_2 errors*, (submitted).
- [17] R.A. DEVORE, *Nonlinear approximation*, Acta Numerica, 7, (1998), pp. 51–150.
- [18] W. DÖRFLER, *A convergent adaptive algorithm for Poisson’s equation*, SIAM J. Numer. Anal., 33 (1996), pp. 1106–1124.
- [19] P. GRISVARD, *Elliptic Problems in Nonsmooth Domains, Monographs and Studies in Mathematics*, vol. 24. Pitman (Advanced Publishing Program), Boston, MA (1985).
- [20] R.B. KELLOGG, *On the Poisson equation with intersecting interfaces*, Applicable Anal., 4 (1974/75), 101–129.
- [21] I. KOSSACZKÝ, *A recursive approach to local mesh refinement in two and three dimensions*, J. Comput. Appl. Math., 55 (1994), pp. 275–288.
- [22] CH. KREUZER AND K.G. SIEBERT, *Decay rates of adaptive finite elements with Dörfler marking*, (submitted).
- [23] 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).
- [24] W. F. MITCHELL, *A Comparison of adaptive refinement techniques for elliptic problems*, ACM Trans. Math. Softw., 15 (1989), pp. 326 - 347.
- [25] 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.
- [26] P. MORIN, R. H. NOCHETTO, AND K. G. SIEBERT, *Convergence of adaptive finite element methods*, SIAM Review, 44 (2002), pp. 631–658.
- [27] 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.
- [28] R.H. NOCHETTO, K. G. SIEBERT, AND A. VEESER, *Theory of adaptive finite element methods: an introduction*, in *Multiscale, Nonlinear and Adaptive Approximation*, A. Kunoth and R. DeVore eds, Springer, 2009, pp. 409–542.
- [29] R. H. NOCHETTO AND A. VEESER, *Primer of adaptive finite element methods*, in *Multiscale and Adaptivity: Modeling, Numerics and Applications*, CIME-EMS Summer School in Applied Mathematics, G. Naldi and G. Russo eds., Springer (2010).
- [30] L. R. SCOTT AND S. ZHANG, *Finite element interpolation of nonsmooth functions satisfying boundary conditions*, Math. Comp., 54 (1990), pp. 483–493.

- [31] K. G. SIEBERT, *A convergence proof for adaptive finite elements without lower bound*, Preprint Universität Duisburg-Essen and Universität Freiburg No. 1/2009.
- [32] 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).
- [33] R. STEVENSON, *Optimality of a standard adaptive finite element method*, *Found. Comput. Math.*, 7 (2007), pp. 245–269.
- [34] R. STEVENSON, *The completion of locally refined simplicial partitions created by bisection*, *Math. Comput.*, 77 (2008), pp. 227–241.
- [35] R. VERFÜRTH, *A Review of A Posteriori Error Estimation and Adaptive Mesh-Refinement Techniques*, *Adv. Numer. Math.* John Wiley, Chichester, UK (1996).
- [36] J. XU, L. CHEN, AND R. H. NOCHETTO, *Adaptive multilevel methods on graded bisection grids*, in *Multiscale, Nonlinear and Adaptive Approximation*, A. Kunoth and R. DeVore eds, Springer, 2009, pp. 599–659.
- [37] L. ZHONG, L. CHEN, S. SHU, G. WITTUM, AND J. XU, *Quasi-optimal convergence of adaptive edge finite element methods for three dimensional indefinite time-harmonic Maxwell's equations*, (submitted).