

Quasi-Optimal Convergence Rate for an Adaptive Finite Element Method

J. Manuel Cascon, Christian Kreuzer, Ricardo H. Nochetto, Kunibert G. Siebert

Angaben zur Veröffentlichung / Publication details:

Cascon, J. Manuel, Christian Kreuzer, Ricardo H. Nochetto, and Kunibert G. Siebert. 2007.
"Quasi-Optimal Convergence Rate for an Adaptive Finite Element Method." Augsburg:
Universität Augsburg.

Nutzungsbedingungen / Terms of use:

licgercopyright

Dieses Dokument wird unter folgenden Bedingungen zur Verfügung gestellt: / This document is made available under these conditions:

Deutsches Urheberrecht

Weitere Informationen finden Sie unter: / For more information see:

<https://www.uni-augsburg.de/de/organisation/bibliothek/publizieren-zitieren-archivieren/publiz/>





Universität Augsburg

Institut für
Mathematik

J. Manuel Cascon, Christian Kreuzer, Ricardo H. Nochetto, Kunibert G. Siebert

Quasi-Optimal Convergence Rate for an Adaptive Finite Element Method

Preprint Nr. 009/2007 — 22. Mai 2007

Institut für Mathematik, Universitätsstraße, D-86135 Augsburg

<http://www.math.uni-augsburg.de/>

Impressum:

Herausgeber:

Institut für Mathematik

Universität Augsburg

86135 Augsburg

<http://www.math.uni-augsburg.de/forschung/preprint/>

ViSdP:

Kunibert G. Siebert

Institut für Mathematik

Universität Augsburg

86135 Augsburg

Preprint: Sämtliche Rechte verbleiben den Autoren © 2007

QUASI-OPTIMAL CONVERGENCE RATE FOR AN ADAPTIVE FINITE ELEMENT METHOD

J. MANUEL CASCON, CHRISTIAN KREUZER, RICARDO H. NOCHETTO,
AND KUNIBERT G. SIEBERT

ABSTRACT. We analyze the simplest and most standard adaptive finite element method (AFEM), with any polynomial degree, for general second order linear, symmetric elliptic operators. As it is customary in practice, AFEM marks exclusively according to the error estimator and performs a minimal element refinement without the interior node property. We prove that AFEM is a contraction for the sum of energy error and scaled error estimator, between two consecutive adaptive loops. This geometric decay is instrumental to derive optimal cardinality of AFEM. We show that AFEM yields a decay rate of energy error plus oscillation in terms of number of degrees of freedom as dictated by the best approximation for this combined nonlinear quantity.

1. INTRODUCTION

Let Ω be a bounded, polyhedral domain in \mathbb{R}^d , $d \geq 2$. We consider a homogeneous Dirichlet boundary value problem for a selfadjoint second order elliptic partial differential equation (PDE)

$$(1.1) \quad \begin{aligned} \mathcal{L}u &:= -\operatorname{div}(\mathbf{A}\nabla u) + cu = f && \text{in } \Omega, \\ u &= 0 && \text{on } \partial\Omega. \end{aligned}$$

The choice of boundary condition is made for ease of presentation, since similar results are valid for other boundary conditions. Precise conditions on given data $\mathbf{D} := (\mathbf{A}, c)$ and f of \mathcal{L} are stated in § 2.1. We refer to [4] for more general operators.

We analyze here a standard adaptive finite element method (AFEM) of the form

$$(1.2) \quad \text{SOLVE} \rightarrow \text{ESTIMATE} \rightarrow \text{MARK} \rightarrow \text{REFINE}.$$

Even though adaptivity has been a fundamental tool of engineering and scientific computing for about three decades, the convergence analysis is rather recent. It started with Dörfler [7], who introduced a crucial marking, from now on called Dörfler's marking, and proved strict energy error reduction for the Laplacian provided the initial mesh \mathcal{T}_0 satisfies a fineness assumption. Morin, Nochetto, and Siebert [16, 17] showed that such strict energy error reduction cannot be expected in general. Introducing the concept of data oscillation and the interior node property, they proved convergence of AFEM without restrictions on \mathcal{T}_0 . The latter result, however, is only valid for \mathbf{A} in (1.1) piecewise constant on \mathcal{T}_0 and vanishing c . Inspired by the work by Chen and Feng [5], Mekchay and Nochetto [15] extended

Date: Version: May 3, 2007— 17:36.

2000 *Mathematics Subject Classification.* 65N30, 65N50, 65N15, 65N12, 41A25.

Key words and phrases. Error reduction, convergence, optimal cardinality, adaptive algorithm.

this result to general second order elliptic operators upon dealing with the new concept of *total error*, namely the sum of energy error plus oscillation. They proved that AFEM is a contraction for the total error, a property that will turn out to be essential in this paper as well. Recently, Diening and Kreuzer [6] proved a similar property for the p -Laplacian but avoiding marking for oscillation.

In this paper we go back to the basis of AFEM philosophy and consider the simplest possible approach to adaptivity for (1.1). In §2 we first state precisely the problem and describe the modules of AFEM: **SOLVE** computes the Ritz-Galerkin approximation, i. e. we assume exact solution and integration; **ESTIMATE** computes the standard residual estimator; **MARK** resorts to Dörfler marking solely based on the estimator; and **REFINE** utilizes bisectioning of elements with the minimal refinement condition that marked elements are bisected at least once. In this respect AFEM is a really standard algorithm in that it avoids marking for oscillation and circumvents the interior node property of [15, 16, 17] for marked elements.

When marking with respect to two quantities such as estimator and oscillation, the role of marking becomes critical for proving optimality. To shed light on this issue, we discuss in §3 the simultaneous adaptive approximation of two functions with distinct asymptotic error decays and compare separate and collective marking. The discussion reveals that separate marking might be sub-optimal and should thus be avoided, whereas collective marking leads to optimal convergence rates.

To summarize the main results, let $\{\mathcal{T}_k, \mathbb{V}_k, U_k, \eta_k, \text{osc}_k\}_{k \geq 0}$ be the sequence of meshes, finite element spaces, discrete solutions, estimators, and oscillations produced by AFEM in the k th step. Even though the energy error is monotone, strict error reduction fails when $U_{k+1} = U_k$, i. e. $\|u - U_{k+1}\|_\Omega = \|u - U_k\|_\Omega$; see [16, 17] for further details. On the other hand, the residual estimator $\eta_k = \eta_k(U_k)$ exhibits a strict reduction when $U_{k+1} = U_k$ but no monotone behavior in general; this is shown in §4. The new insight of this paper is that the sum of energy error and scaled estimator, the so-called *quasi-error* $(\|u - U_k\|_\Omega^2 + \gamma \eta_k^2)^{1/2}$, is strictly reduced by AFEM even though each term may not be. In fact, we prove in §5 that AFEM is a contraction for this new error notion:

Main Result 1. *There exist constants $\gamma > 0$, and $0 < \alpha < 1$, such that*

$$\|u - U_{k+1}\|_\Omega^2 + \gamma \eta_{k+1}^2 \leq \alpha^2 \left(\|u - U_k\|_\Omega^2 + \gamma \eta_k^2 \right).$$

Quasi-optimal convergence rates for AFEM, expressing energy error decay in terms of number of degrees of freedom (DOFs) as dictated by nonlinear approximation theory, were first proved by Binev, Dahmen and DeVore [3]. They resorted to a crucial, but somewhat artificial, coarsening step. Coarsening was later removed by Stevenson [22], who developed an optimality theory for a much more realistic AFEM but still including an inner loop to deal with oscillation; see §3.3 for basic details. Both papers [3, 22] are restricted to the Laplace operator and rely on suitable marking by oscillation and the interior node property.

To derive convergence rates we need to seek a suitable error quantity and associated approximation class \mathbb{A}_s . First, we observe that the estimator dominates oscillation which in conjunction with the *global lower bound* gives

$$\text{osc}_k \leq \eta_k \preceq \|u - U_k\|_\Omega + \text{osc}_k,$$

whence

$$\|u - U_k\|_\Omega + \eta_k \approx \|u - U_k\|_\Omega + \text{osc}_k.$$

Hence, the sum of energy error plus oscillation, the so-called *total error*, is equivalent to the quantity being reduced by AFEM. This motivates the definition of the approximation class \mathbb{A}_s in §6 which roughly states that (u, f, \mathbf{D}) belongs to \mathbb{A}_s if the total error can be approximated within any tolerance $\epsilon > 0$ with $O(\epsilon^{-s})$ DOFs. Note that for a linear PDE with variable coefficients, oscillation and solution couple in a nonlinear fashion. As an outcome, an important pending issue is a complete characterization of \mathbb{A}_s ; we refer to §6.1 for further discussion.

We conclude the article by proving a *quasi-optimal* convergence rate for AFEM in §6. Assuming certain restrictions on the initial triangulation and the marking parameter $\theta \in (0, \theta_*)$, we show that AFEM achieves the same asymptotic decay rate s in terms of DOFs as for elements in \mathbb{A}_s :

Main Result 2. *If $(u, f, \mathbf{D}) \in \mathbb{A}_s$, then there exists a constant C such that*

$$\|u - U_k\|_{\Omega} + \text{osc}_k \leq C (\#\mathcal{T}_k - \#\mathcal{T}_0)^{-s}.$$

The contraction property of AFEM, written as Main Result 1, plays an essential role. In contrast to former optimality proofs [3, 22], the present analysis stays within the class of conforming meshes. This is a necessary framework for the analysis when dealing with oscillation in the jump residual. The proof of Main Result 2 relies on a quasi-monotonicity property of oscillation as well as a localized upper bound, both proved in §4. The latter refines a similar bound by Stevenson [22].

In this paper we start with a standard AFEM as it is used in practice and provide a theory for convergence and optimality. The first convergence result for such a standard AFEM is due to Morin, Siebert, and Veiser [19]. They proved recently, for a larger problem class and more general marking strategies, plain convergence of AFEM but without an error reduction property. Relying on Dörfler marking we are able to prove a stronger result for selfadjoint elliptic operators of the form (1.1), namely contraction of the quasi-error and quasi-optimal cardinality of AFEM.

In all other convergence and optimality results, the standard form of AFEM is first altered and then the modified algorithms are analyzed. Because of theoretical needs additional ingredients, such as the interior node property and marking for oscillation [15, 16, 17, 18], a coarsening step [3], or an additional inner loop to decrease oscillation relative to the estimator [22], are added to AFEM.

We would like to stress that removing these ingredients is very important from a practical point of view. The interior node property enforces six bisections of marked elements in three space dimensions and thus increases the number of DOFs between two iterations drastically. In fact, computational resources can be used more efficiently with fewer element refinements. Additionally, since oscillation is not used by AFEM, it does not have to be computed. Although computing time for oscillation may be negligible, this strongly improves on implementation requirements for AFEM. Computing oscillation gets inevitably more involved for higher order discretizations, since one has to solve small linear systems on elements and sides to find local projections. Last but not least, in contrast to modified versions of AFEM, the standard form only needs *one single* parameter, namely the parameter θ of Dörfler marking. Hence, we do not need to fit several parameters, which in turn makes the resulting algorithm more robust.

2. PROBLEM AND ADAPTIVE FINITE ELEMENT METHOD

We first introduce the underlying problem and state assumptions on given data. We then describe the refinement framework and AFEM along with its modules.

2.1. Weak Formulation. Let Ω be a bounded, polyhedral domain in \mathbb{R}^d , $d \geq 2$, that is triangulated by a conforming triangulation \mathcal{T}_0 . We assume that data of (1.1) has the following properties:

- (a) $\mathbf{A}: \Omega \mapsto \mathbb{R}^{d \times d}$ is piecewise Lipschitz over \mathcal{T}_0 and is symmetric positive definite with eigenvalues in $0 < a_* \leq a^* < \infty$, i.e.,

$$a_*(x) |\xi|^2 \leq \mathbf{A}(x) \xi \cdot \xi \leq a^*(x) |\xi|^2, \quad \forall \xi \in \mathbb{R}^d, x \in \Omega;$$

- (b) $c \in L^\infty(\Omega)$ is nonnegative, i.e. $c \geq 0$ in Ω ;

- (c) $f \in L^2(\Omega)$.

Now we turn to the weak formulation of (1.1). For any set $\omega \subset \mathbb{R}^d$ with non-empty interior we denote by $H^1(\omega)$ the usual Sobolev space of functions in $L^2(\omega)$ whose first derivatives are also in $L^2(\omega)$, endowed with the norm

$$\|u\|_{H^1(\omega)} := \left(\|u\|_{L^2(\omega)}^2 + \|\nabla u\|_{L^2(\omega)}^2 \right)^{1/2}.$$

Moreover, we denote $\langle \cdot, \cdot \rangle_\omega$ the $L^2(\omega)$ scalar product. Finally we let $\mathbb{V} := H_0^1(\Omega)$ be the space of functions in $H^1(\Omega)$ with vanishing trace on $\partial\Omega$. A weak solution of (1.1) is a function u satisfying

$$(2.1) \quad u \in \mathbb{V}: \quad \mathcal{B}[u, v] = \langle f, v \rangle_\Omega \quad \forall v \in \mathbb{V},$$

where the bilinear form is defined to be

$$\mathcal{B}[u, v] := \langle \mathbf{A} \nabla u, \nabla v \rangle_\Omega + \langle c u, v \rangle_\Omega \quad \forall u, v \in \mathbb{V}.$$

In view of Poincaré-Friedrichs inequality, one has *coercivity* in \mathbb{V}

$$\mathcal{B}[v, v] = \langle \mathbf{A} \nabla v, \nabla v \rangle_\Omega + \langle c v, v \rangle_\Omega \geq \int_\Omega a_* |\nabla v|^2 + c v^2 \geq c_B \|v\|_{H^1(\Omega)}^2,$$

and c_B depends only on data and Ω . The bilinear form \mathcal{B} induces the so-called *energy norm*:

$$\|v\|_\omega := \mathcal{B}[v, v]^{1/2} \quad \forall v \in H^1(\omega).$$

Note that the bilinear form also fulfills the *local continuity*

$$\mathcal{B}[v, w] \leq \sqrt{C_B} \|v\|_\omega \|w\|_{H^1(\omega)} \quad \forall v, w \in H^1(\omega), \text{supp}(w) \subset \omega \subset \Omega.$$

This local continuity is essential in deriving *local lower bounds* in the a posteriori error analysis. Furthermore it implies continuity of $\mathcal{B}[\cdot, \cdot]$ on $H^1(\Omega)$ at once. Thanks to coercivity and continuity of \mathcal{B} , the norm $\|\cdot\|_\Omega$ is equivalent to $\|\cdot\|_{H^1(\Omega)}$ on $H_0^1(\Omega)$. Existence and uniqueness of (2.1) thus follows from Lax-Milgram theorem [8]. The restriction to a symmetric bilinear form and $c \geq 0$ can be relaxed provided (2.1) admits a unique solution. This extension will be studied in [4].

2.2. Refinement Framework. Refinement is based on shape-regular bisection of single elements. Any given simplex is subdivided into two subsimplices of same size such that the minimal angle is uniformly bounded from below; we refer to [2, 12, 13, 14, 23, 24] or the monograph [20] and the references therein. In 2d, this is the newest vertex bisection. Bisectioning creates a unique *master forest* \mathbb{F} of binary trees with infinite depth, where each node is a simplex, its two successors are the two children created by bisection, and the roots of the binary trees are the elements of the initial conforming triangulation \mathcal{T}_0 . The master forest \mathbb{F} contains full information of all possible subdivisions created from \mathcal{T}_0 by bisection, i.e. information about vertices, neighbors, refinement edges, etc. for any simplex that can be generated.

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

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

Any node of \mathcal{F} is thus uniquely connected with a node of the initial triangulation \mathcal{T}_0 . Furthermore, any forest may have interior nodes, i.e. nodes with successors, and does have leaf nodes, i.e. nodes without successors. Finally, each forest \mathcal{F} corresponds one-to-one to a subdivision $\mathcal{T}(\mathcal{F})$ of Ω into simplices by defining $\mathcal{T}(\mathcal{F})$ as the set of leaf nodes of \mathcal{F} . Note, that such a subdivision may be non-conforming.

If $\mathcal{F} \subset \mathcal{F}_*$ are two forests, we call $\mathcal{T}_* = \mathcal{T}(\mathcal{F}_*)$ a refinement of $\mathcal{T} = \mathcal{T}(\mathcal{F})$ and denote this by $\mathcal{T} \leq \mathcal{T}_*$. We define

$$\mathcal{R}_{\mathcal{T} \rightarrow \mathcal{T}_*} := \mathcal{T} \setminus (\mathcal{T}_* \cap \mathcal{T})$$

as the set of *refined elements*, i.e. $\mathcal{R}_{\mathcal{T} \rightarrow \mathcal{T}_*}$ is the set of leaf nodes of \mathcal{F} that are interior nodes of \mathcal{F}_* . The class of all *conforming* refinements by bisection of \mathcal{T}_0 is

$$\mathbb{T} := \{\mathcal{T}(\mathcal{F}) \mid \mathcal{F} \subset \mathbb{F} \text{ is a forest and } \mathcal{T}(\mathcal{F}) \text{ is a conforming triangulation of } \Omega\}.$$

Given two forests $\mathcal{F}, \mathcal{F}_* \subset \mathbb{F}$ corresponding to subdivisions \mathcal{T} and \mathcal{T}_* of Ω , we define $\mathcal{F} \cup \mathcal{F}_*$ to be the union of the nodes of \mathcal{F} and \mathcal{F}_* . Obviously, $\mathcal{F} \cup \mathcal{F}_* \subset \mathbb{F}$ and $\mathcal{T}_0 \subset \mathcal{F} \cup \mathcal{F}_*$. By construction, all nodes of $(\mathcal{F} \cup \mathcal{F}_*) \setminus \mathcal{T}_0$ have a predecessor, and all nodes of the union have either two successors or none: this implies that $\mathcal{F} \cup \mathcal{F}_*$ is a forest. We call the resulting unique subdivision of Ω the *overlay of \mathcal{T} and \mathcal{T}_** :

$$\mathcal{T} \oplus \mathcal{T}_* := \mathcal{T}(\mathcal{F} \cup \mathcal{F}_*).$$

For two conforming triangulations $\mathcal{T}, \mathcal{T}_* \in \mathbb{T}$ we prove in Lemma 4.7 below that $\mathcal{T} \oplus \mathcal{T}_* \in \mathbb{T}$ is the smallest conforming refinement of $\mathcal{T}, \mathcal{T}_*$ and satisfies

$$(2.2) \quad \#(\mathcal{T} \oplus \mathcal{T}_*) \leq \#\mathcal{T} + \#\mathcal{T}_* - \#\mathcal{T}_0.$$

We finally introduce some notations related to triangulations. For $T \in \mathcal{T}$ we denote by $h_T := |T|^{\frac{1}{d}}$ the local meshsize, and by ω_T the union of all elements in \mathcal{T} sharing one side with T . We further denote by \mathcal{S} the skeleton of \mathcal{T} , i.e. the union of the inter-element sides and for an interior side $\sigma \in \mathcal{S}$ we let ω_σ be the union of the two adjacent elements sharing σ .

Thanks to properties of bisection, all constants depending on shape regularity of $\mathcal{T} \in \mathbb{T}$ are uniformly bounded by a constant solely depending on \mathcal{T}_0 [13, 24].

2.3. The Module SOLVE. Given any conforming triangulation $\mathcal{T} \in \mathbb{T}$ we define the finite element space

$$\mathbb{V}(\mathcal{T}) := \{V \in \mathbb{V} \mid V|_T \in \mathbb{P}_n(T), T \in \mathcal{T}\},$$

where $n \in \mathbb{N}$ is a fixed polynomial degree and \mathbb{P}_n denotes the space of all polynomials of degree $\leq n$. Since continuity and coercivity of \mathcal{B} are inherited by any subspace of \mathbb{V} , the Lax-Milgram theorem implies existence and uniqueness of the Ritz-Galerkin approximation in $\mathbb{V}(\mathcal{T})$ defined by

$$(2.3) \quad U \in \mathbb{V}(\mathcal{T}) : \quad \mathcal{B}[U, V] = \langle f, V \rangle_\Omega \quad \forall V \in \mathbb{V}(\mathcal{T}).$$

We suppose that the module **SOLVE** outputs the exact Ritz-Galerkin solution on \mathcal{T}

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

i. e. U is computed via exact linear algebra and exact integration. Optimal multi-level solvers can be incorporated as in [22], but quadrature is more delicate.

We observe that for any pair $\mathcal{T}, \mathcal{T}_* \in \mathbb{T}$ with $\mathcal{T} \leq \mathcal{T}_*$ there holds $\mathbb{V}(\mathcal{T}) \subset \mathbb{V}(\mathcal{T}_*)$, i. e. the discrete spaces are *nested*. This turns out to be a crucial property in the subsequent analysis in that the following *orthogonality* relation holds:

$$(2.4) \quad \|u - U\|_\Omega^2 = \|u - U_*\|_\Omega^2 + \|U_* - U\|_\Omega^2.$$

This is a consequence of $\mathcal{B}[u - U_*, U_* - U] = 0$, or equivalently that $u - U_*$ and $U_* - U$ are orthogonal with respect to the scalar product induced by $\mathcal{B}[\cdot, \cdot]$.

2.4. The Module ESTIMATE. For $\mathcal{T} \in \mathbb{T}$ and $V \in \mathbb{V}(\mathcal{T})$ we define the *element residual* and *jump residual* for V by

$$R(V)|_T := (f + \mathcal{L}V)|_T, \quad T \in \mathcal{T} \quad \text{and} \quad J(V)|_\sigma := (\llbracket \mathbf{A} \nabla V \rrbracket \cdot \boldsymbol{\nu})|_\sigma \quad \sigma \in \mathcal{S}.$$

Hereafter, $\llbracket q \rrbracket$ is the jump of q across an interior side σ , and $\boldsymbol{\nu}$ denotes a unit normal vector associate to side σ . The error indicator for V on $T \in \mathcal{T}$ is given by

$$(2.5) \quad \eta_T^2(V, T) := h_T^2 \|R(V)\|_{L^2(T)}^2 + h_T \|J(V)\|_{L^2(\partial T \cap \Omega)}^2,$$

where we recall $h_T = |T|^{1/d}$. We assume that, given a triangulation \mathcal{T} and the Ritz-Galerkin solution $U \in \mathbb{V}(\mathcal{T})$, the module **ESTIMATE** outputs the indicators

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

The lower bound involves oscillation, which we define next. For $m \in \mathbb{N}_0$, we denote by Π_m^p the L^p -best approximation operator onto the set of discontinuous polynomials of degree $\leq m$ over either $T \in \mathcal{T}$ or $\sigma \in \mathcal{S}$, depending on the context. If $P_m^p := \text{id} - \Pi_m^p$, then we define oscillation of $V \in \mathbb{V}(\mathcal{T})$ to be

$$\text{osc}_T^2(V, T) := h_T^2 \|P_{2n-2}^2 R(V)\|_{L^2(T)}^2 + h_T \|P_{2n-1}^2 J(V)\|_{L^2(\partial T \cap \Omega)}^2.$$

Finally, for any subset $\mathcal{T}' \subset \mathcal{T}$ we set

$$\eta_T^2(V, \mathcal{T}') := \sum_{T \in \mathcal{T}'} \eta_T^2(V, T) \quad \text{and} \quad \text{osc}_T^2(V, \mathcal{T}') := \sum_{T \in \mathcal{T}'} \text{osc}_T^2(V, T).$$

Remark 2.1. Let $\mathcal{T} \in \mathbb{T}$ and $V \in \mathbb{V}(\mathcal{T})$ be given. We first observe that the indicator $\eta_T(V, T)$ *dominates* oscillation $\text{osc}_T(V, T)$, i. e. $\text{osc}_T(V, T) \leq \eta_T(V, T)$ for all $T \in \mathcal{T}$. In addition, the definition of error indicator and oscillation are fully localized to T , i. e. for any triangulation $\mathcal{T}_* \in \mathbb{T}$ with $T \in \mathcal{T}_*$ there holds $\eta_T(V, T) = \eta_{\mathcal{T}_*}(V, T)$ and $\text{osc}_T(V, T) = \text{osc}_{\mathcal{T}_*}(V, T)$. Moreover, if $\mathcal{T}_* \geq \mathcal{T}$ is any refinement of \mathcal{T} , then combining the monotonicity of local meshsizes and properties of the local L^2 projection we deduce the *monotonicity* properties

$$\eta_{\mathcal{T}_*}(V, \mathcal{T}_*) \leq \eta_{\mathcal{T}}(V, \mathcal{T}) \quad \text{and} \quad \text{osc}_{\mathcal{T}_*}(V, \mathcal{T}_*) \leq \text{osc}_{\mathcal{T}}(V, \mathcal{T}) \quad \forall V \in \mathbb{V}(\mathcal{T}).$$

We now recall the well-known upper and lower bounds for the energy error in terms of the residual-type estimator [1, 15, 16, 17, 22, 26].

Lemma 2.2 (Global A Posteriori Upper and Lower Bounds). *Let $u \in \mathbb{V}$ be the solution of (2.1), $\mathcal{T} \in \mathbb{T}$, and $U \in \mathbb{V}(\mathcal{T})$ be the Ritz-Galerkin solution (2.3).*

Then there exist a constant C_1 , solely depending on \mathcal{T}_0 and $1/c_B$, such that

$$(2.6) \quad \|u - U\|_{\Omega}^2 \leq C_1 \eta_{\mathcal{T}}^2(U, \mathcal{T}),$$

and a constant C_2 , solely depending on \mathcal{T}_0 and $1/C_B$, such that

$$(2.7) \quad C_2 \eta_{\mathcal{T}}^2(U, \mathcal{T}) \leq \|u - U\|_{\Omega}^2 + \text{osc}_{\mathcal{T}}^2(U, \mathcal{T}).$$

2.5. The Module MARK. In the selection of elements we rely on Dörfler marking. Given a grid \mathcal{T} , the set of indicators $\{\eta_{\mathcal{T}}(U, T)\}_{T \in \mathcal{T}}$, and marking parameter $\theta \in (0, 1]$, we suppose that MARK outputs a subset of *marked elements* $\mathcal{M} \subset \mathcal{T}$, i. e.

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

such that \mathcal{M} satisfies Dörfler property

$$(2.8) \quad \eta_{\mathcal{T}}(U, \mathcal{M}) \geq \theta \eta_{\mathcal{T}}(U, \mathcal{T}).$$

2.6. The Module REFINE. We suppose that a function REFINE is at our disposal that implements iterative or recursive bisection, see [2, 12, 13, 14, 23, 24]. When relying on recursive bisection, the distribution of refinement edges has to fulfill some compatibility conditions on \mathcal{T}_0 . Given a fixed number $b \geq 1$, for any $\mathcal{T} \in \mathbb{T}$ and a subset $\mathcal{M} \subset \mathcal{T}$ of marked elements

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

outputs a conforming triangulation $\mathcal{T}_* \in \mathbb{T}$, where at least all elements of \mathcal{M} are bisected b times. In particular, this implies $\mathcal{M} \subset \mathcal{R}_{\mathcal{T} \rightarrow \mathcal{T}_*}$.

To ensure conformity of \mathcal{T}_* one usually has to refine additional elements in $\mathcal{T} \setminus \mathcal{M}$ once or several times. In general, the number of these additionally refined elements is not controlled by $\#\mathcal{M}$, that is $\#\mathcal{T}_* - \#\mathcal{T}$ cannot be bounded by $C \#\mathcal{M}$ with a constant C independent of \mathcal{T} ; C may depend on the refinement level. On the other hand, arguing with the entire sequence $\{\mathcal{T}_k\}_{k \geq 0}$ of refinements, Binev, Dahmen, and DeVore showed in 2d that the *cumulative* number of elements added by conformity does not inflate the total number of marked elements [3, Theorem 2.4]. Stevenson generalized this result to higher dimensions [23, Theorem 6.1] and also weakened the assumptions on \mathcal{T}_0 in 2d.

Lemma 2.3 (Complexity of REFINE). *Assume that \mathcal{T}_0 verifies condition (b) of §4 in [23]. For $k \geq 0$ let $\{\mathcal{T}_k\}_{k \geq 0}$ be any sequence of refinements of \mathcal{T}_0 where \mathcal{T}_{k+1} is generated from \mathcal{T}_k by $\mathcal{T}_{k+1} = \text{REFINE}(\mathcal{T}_k, \mathcal{M}_k)$ with a subset $\mathcal{M}_k \subset \mathcal{T}_k$.*

Then, there exists a constant C_0 solely depending on \mathcal{T}_0 and b such that

$$\#\mathcal{T}_k - \#\mathcal{T}_0 \leq C_0 \sum_{j=0}^{k-1} \#\mathcal{M}_j \quad \forall k \geq 1.$$

The original result is proved for a sequence of meshes generated by a loop

$$\text{Mark} \rightarrow \text{Subdivision} \rightarrow \text{Complete},$$

where **Mark** selects certain elements, **Subdivision** bisects marked elements only once ($b = 1$), and **Complete** adds some additional subdivisions to generate a conforming

triangulation. To apply the above result 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 Lemma 2.3 is a direct consequence of [23, Theorem 6.1] because

$$\sum_{j=1}^b \#\mathcal{M}_k(j) \leq \sum_{j=1}^b 2^{j-1} \#\mathcal{M}_k = (2^b - 1) \#\mathcal{M}_k.$$

2.7. Adaptive Algorithm. In this section we now collect the modules described in the previous sections into the adaptive finite element method. In doing this, we replace the dependence on the actual triangulation \mathcal{T} by the iteration counter $k \geq 0$, for instance $\mathbb{V}_k = \mathbb{V}(\mathcal{T}_k)$. The basic loop of AFEM is then given by the following iteration:

AFEM

Given the initial grid \mathcal{T}_0 and marking parameter $0 < \theta \leq 1$
 set $k := 0$ and iterate

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

We want to stress that AFEM is really a *standard algorithm* in that it only employs the error indicators $\{\eta_k(U_k, T)\}_{T \in \mathcal{T}_k}$ and does not use the oscillation indicators $\{\text{osc}_k(U_k, T)\}_{T \in \mathcal{T}_k}$. In addition, AFEM only relies on a minimal refinement, i. e. an interior node property is not enforced for marked elements.

Several remarks are now in order:

- Recently, Morin, Siebert and Veiser proved convergence of the above iteration for general marking strategies, including maximum and equidistribution besides Dörfler strategy [19]. The main result is a plain convergence result, i. e.

$$\lim_{k \rightarrow \infty} \|u - U_k\|_{L^2(\Omega)} = \lim_{k \rightarrow \infty} \eta_k(U_k, \mathcal{T}_k) = 0.$$

This does not provide a strict total error reduction between two successive iterations, which is one key ingredient in the optimality proof of §6.

- Using similar techniques as Diening and Kreuzer for the p -Laplacian [6], we prove in §5 a contraction property for $\|u - U_k\|_{L^2(\Omega)}^2 + \gamma \eta_k^2(U_k, \mathcal{T}_k)$ between two successive iterations, with $\gamma > 0$ a suitable scaling factor; we observe that [6] still requires the interior node property. This in turn implies *linear convergence*. Crucial ingredients are the a posteriori *global upper bound* and *Dörfler marking*.
- We prove in §6 that \mathcal{T}_k exhibits *quasi-optimal cardinality* provided that
 - (a) the marking parameter θ satisfies $\theta < \theta_*$ where θ_* , defined explicitly in Assumption 6.6, depends on the ratio $\sqrt{C_2/C_1}$ of constants in (2.6) and (2.7);
 - (b) \mathcal{M}_k satisfies (2.8) with minimal cardinality;
 - (c) \mathcal{T}_{k+1} is the smallest refinement of \mathcal{T}_k such that $\mathcal{M}_k \subset \mathcal{R}_{\mathcal{T}_k \rightarrow \mathcal{T}_{k+1}}$.

The proof hinges on the a posteriori *lower bound* and a *localized upper bound*.

3. THE CRITICAL ROLE OF MARKING

Former convergence proofs of AFEM are based on the AFEM designed by Morin, Nochetto, and Siebert [16, 17], from now on called MNS. This is an adaptive loop of the form (1.2), which marks separately for both estimator and oscillation in each iteration. It turns out that separate marking appears to be problematic for proving optimality and thus variants of MNS are studied instead [3, 22] (see also §3.3). The key issue is the delicate choice of marking parameters for estimator and oscillation. To shed light on this intrinsic difficulty we first present some numerical experiments, then discuss the effect of separate and collective marking, and finally comment on marking for an optimal AFEM.

3.1. Separate Marking and MNS. The procedure **ESTIMATE** of MNS calculates both error and oscillation indicators $\{\eta_k(U_k, T), \text{osc}_k(U_k, T)\}_{T \in \mathcal{T}_k}$ and the procedure **MARK** uses Dörfler marking for both estimator and oscillation. More precisely, the routine **MARK** of MNS is of the form: *Given parameters* $0 < \theta_{\text{est}}, \theta_{\text{osc}} < 1$

(3.1a) *Mark any subset* $\mathcal{M}_k \subset \mathcal{T}_k$ *such that* $\eta_k(U_k, \mathcal{M}_k) \geq \theta_{\text{est}} \eta_k(U_k, \mathcal{T}_k)$;

(3.1b) *If necessary enlarge* \mathcal{M}_k *to satisfy* $\text{osc}_k(U_k, \mathcal{M}_k) \geq \theta_{\text{osc}} \text{osc}_k(U_k, \mathcal{T}_k)$.

Since oscillation is generically of higher order than the estimator, the issue at stake is whether elements added by oscillation, even though immaterial relative to the error, could ruin the optimality of MNS observed in experiments [16, 17, 15]. If $\eta_k(U_k, \mathcal{T}_k)$ has large indicators in a small area, then Dörfler marking for the estimator (3.1a) could select a set \mathcal{M}_k with a small number of elements relative to \mathcal{T}_k . However, if $\text{osc}_k(U_k, \mathcal{T}_k)$ were globally distributed in \mathcal{T}_k , then MNS would require additional marking of a large percentage of all elements to satisfy (3.1b), i.e. $\#\mathcal{M}_k$ could be large relative to $\#\mathcal{T}_k$.

To explore this idea computationally, we consider a simple modification of the 2D Example 5.3 in [11, 16]

$$-\text{div } \mathbf{A} \nabla u = f \quad \text{in } \Omega, \quad u = g \quad \text{on } \Omega,$$

where $\Omega = (-1, 1)^2$, and $\mathbf{A} := a_i \mathbf{I}$ is a piecewise constant (checkerboard pattern) with $a_1 = a_3$ in the first and third quadrants, and $a_2 = a_4$ in the second and fourth quadrants. The exact solution u is given as a sum: $u = u_K + u_S$. Function $u_K \in H^1(\Omega)$ is the weak solution of $\text{div}(\mathbf{A} \nabla u_K) = 0$ for parameters [11, 16, 17],

$$a_1 = a_3 = 161.4476387975881, \quad a_2 = a_4 = 1;$$

the singularity $u_K \approx r^{0.1}$, at the origin is so extreme that typically leads to marking of a handful of elements per step to satisfy (3.1a) (see [16, 17] for details). Function u_S is given in each quadrant by

$$u_S(x, y) = 10^{-2} a_i^{-1} (x^2 + y^2) \sin^2(4\pi x) \sin^2(4\pi y) \quad 1 \leq i \leq 4,$$

is smooth and of comparable magnitude with u_K , while the corresponding $f = -\text{div } \mathbf{A} \nabla u_S$ exhibits an increasing amount of data oscillation away from the origin. Procedure **MARK** takes the usual value of $\theta_{\text{est}} = 0.5$ [7, 16, 17, 20], and procedure **REFINE** subdivides all elements in \mathcal{M}_k using two bisections.

The behavior of MNS for several values of θ_{osc} is depicted in Figure 3.1. We can visualize the sensitivity of MNS with respect to parameter θ_{osc} . For values of

$\theta_{\text{osc}} \leq 0.4$ the rate of convergence appears to be quasi-optimal. However, beyond this threshold the curves for both error and estimator flatten out thereby showing lack of optimality. The threshold value $\theta_{\text{osc}} = 0.4$, even though consistent with practice of MNS, is tricky to find in general since it is problem dependent.

Lemma 3.1 below provides theoretical insight on this matter. Altogether this indicates that an optimality proof for MNS is out of reach without restrictions on θ_{osc} , and makes marking by oscillation (3.1b) questionable.

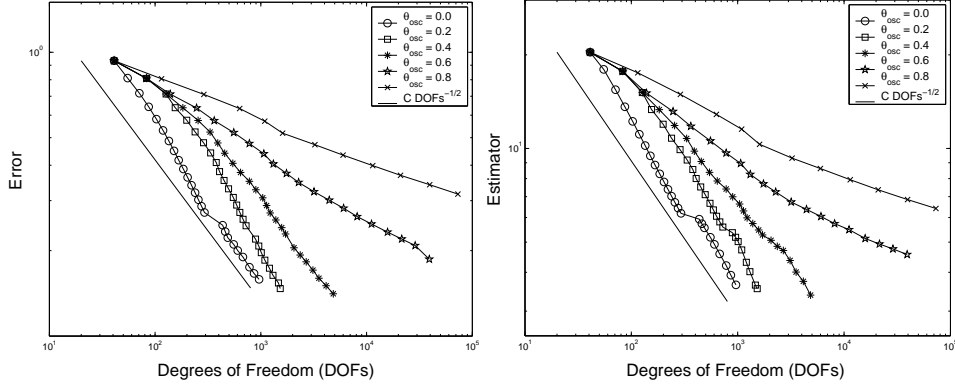


FIGURE 3.1. Decay of error (left) and estimator (right) vs. DOFs for $\theta_{\text{est}} = 0.5$ and values $\theta_{\text{osc}} = 0.0, 0.2, 0.4, 0.6, 0.8$. For values of $\theta_{\text{osc}} \leq 0.4$ the rate of convergence is quasi-optimal, but for $\theta_{\text{osc}} > 0.4$ the curves flatten out thereby showing lack of optimality.

3.2. Separate vs. Collective Marking. In this section we focus on the adaptive approximation of two given functions in an idealized scenario. We compare the effect of separate marking, similar to (3.1), with collective marking, i.e. one single marking for the combined error. The discussion reveals that collective marking leads to optimal convergence rates whereas separate marking might be sub-optimal. However, a suitable choice of marking parameters may restore optimality for separate marking. The experiments of §3.1 confirm this theoretical insight in the non-idealized case of MNS.

Throughout this section we use the notation $a \lesssim b$ to indicate $a \leq Cb$, with generic constants C not depending on the iteration counter. We denote $a \lesssim b \lesssim a$ by $a \approx b$.

For the discussion, we assume that we have two functions u_i , for $i = 1, 2$, and have access to their local approximation error $e_{\mathcal{T}}(u_i; T)$

$$e_{\mathcal{T}}(u_i; T) = |u_i - I_{\mathcal{T}} u_i|_{i; T} \quad \forall T \in \mathcal{T},$$

and global error $e_{\mathcal{T}}^2(u_i) = \sum_{T \in \mathcal{T}} e_{\mathcal{T}}^2(u_i; T)$; hereafter $|\cdot|_i$ are unspecified norms and $I_{\mathcal{T}}$ is a local interpolation operator over $T \in \mathbb{T}$. We define the *total error* to be

$$e_{\mathcal{T}}^2 := e_{\mathcal{T}}^2(u_1) + e_{\mathcal{T}}^2(u_2),$$

and are interested in its asymptotic decay. If $\mathcal{T} = \mathcal{T}_k$, then we denote $e_k = e_{\mathcal{T}_k}$.

To explore the use of (3.1), we examine the effect of separate marking for $e_k(u_i)$ on a sequence of meshes \mathcal{T}_k^i for $i = 1, 2$. We put ourselves in an idealized, but

plausible, situation governed by the following three simplifying assumptions:

- (3.2a) *Independence:* \mathcal{T}_k^1 and \mathcal{T}_k^2 are generated from \mathcal{T}_0 and are independent of each other;
- (3.2b) *Marking:* Separate Dörfler marking with parameters $\theta_i \in (0, 1)$ imply $e_k(u_i) \approx \alpha_i^k$ on \mathcal{T}_k^i with $\alpha_i \in (0, 1)$;
- (3.2c) *Approximability:* $e_k(u_i) \approx (\#\mathcal{T}_k^i - \#\mathcal{T}_0)^{-s_i}$ with $s_1 \leq s_2$ maximal.

We are interested in the decay of the total error e_k on the overlay $\mathcal{T}_k := \mathcal{T}_k^1 \oplus \mathcal{T}_k^2$. This scenario is a simplification of the more realistic approximation of u_1 and u_2 with separate Dörfler marking on the same sequence of grids \mathcal{T}_k but avoids the complicated interaction of the two marking procedures.

Lemma 3.1 (Separate Marking). *Let assumptions (3.2) be satisfied.*

Then the decay of the total error e_k on the overlay $\mathcal{T}_k = \mathcal{T}_k^1 \oplus \mathcal{T}_k^2$ for separate marking is always suboptimal except when α_1 and α_2 satisfy

$$\alpha_2 \leq \alpha_1 \leq \alpha_2^{s_1/s_2}.$$

Proof. Assumption (3.2b) on the average reduction rate implies for the total error

$$(3.3) \quad e_k \approx e_k(u_1) + e_k(u_2) \approx \max\{e_k(u_1), e_k(u_2)\} \approx \max\{\alpha_1^k, \alpha_2^k\}.$$

Combining (3.2b) and (3.2c) yields $\alpha_i^k \approx (\#\mathcal{T}_k^i - \#\mathcal{T}_0)^{-s_i}$, whence

$$(3.4) \quad \#\mathcal{T}_k^1 - \#\mathcal{T}_0 \approx \alpha_1^{-k/s_1} = \beta^k \alpha_2^{-k/s_2} \approx \beta^k (\#\mathcal{T}_k^2 - \#\mathcal{T}_0)$$

with $\beta = \alpha_1^{-1/s_1} \alpha_2^{1/s_2}$. In view of (2.2), this gives for the overlay $\mathcal{T}_k = \mathcal{T}_k^1 \oplus \mathcal{T}_k^2$

$$(3.5) \quad \#\mathcal{T}_k - \#\mathcal{T}_0 \approx \begin{cases} \#\mathcal{T}_k^1 - \#\mathcal{T}_0, & \beta \geq 1, \\ \#\mathcal{T}_k^2 - \#\mathcal{T}_0, & \beta < 1. \end{cases}$$

The optimal decay of total error e_k corresponds to $e_k \approx (\#\mathcal{T}_k - \#\mathcal{T}_0)^{-s_1}$ because $s_1 \leq s_2$. In analyzing the relation of e_k to the number of elements $\#\mathcal{T}_k$ in the overlay \mathcal{T}_k we distinguish three cases and employ (3.3), (3.4), and (3.5).

Case 1: $\alpha_1 < \alpha_2$. We note that $\alpha_1 < \alpha_2$ and $s_1 \leq s_2$ yields $\beta \geq 1$. We thus deduce

$$\begin{aligned} e_k &\approx \max\{\alpha_1^k, \alpha_2^k\} = \alpha_2^k = (\alpha_2/\alpha_1)^k \alpha_1^k \\ &\approx (\alpha_2/\alpha_1)^k (\#\mathcal{T}_k^1 - \#\mathcal{T}_0)^{-s_1} \approx (\alpha_2/\alpha_1)^k (\#\mathcal{T}_k - \#\mathcal{T}_0)^{-s_1}. \end{aligned}$$

Since $\alpha_2/\alpha_1 > 1$, the approximation of e_k on \mathcal{T}_k is suboptimal.

Case 2: $\alpha_1 \geq \alpha_2$ and $\beta < 1$. We obtain

$$\begin{aligned} e_k &\approx \max\{\alpha_1^k, \alpha_2^k\} = \alpha_1^k \approx (\#\mathcal{T}_k^1 - \#\mathcal{T}_0)^{-s_1} \\ &\approx \beta^{-ks_1} (\#\mathcal{T}_k^2 - \#\mathcal{T}_0)^{-s_1} \approx \beta^{-ks_1} (\#\mathcal{T}_k - \#\mathcal{T}_0)^{-s_1}, \end{aligned}$$

whence the approximation of the total error on \mathcal{T}_k is again suboptimal.

Case 3: $\alpha_1 \geq \alpha_2$ and $\beta \geq 1$. We infer that

$$e_k \approx \max\{\alpha_1^k, \alpha_2^k\} = \alpha_1^k \approx (\#\mathcal{T}_k^1 - \#\mathcal{T}_0)^{-s_1} \approx (\#\mathcal{T}_k - \#\mathcal{T}_0)^{-s_1}$$

and that \mathcal{T}_k exhibits optimal cardinality. This exceptional case corresponds to the assertion and concludes the proof. \square

We next investigate the effect of collective marking. We assume the following properties for $\mathcal{T}, \mathcal{T}_* \in \mathbb{T}$ with $\mathcal{T} \leq \mathcal{T}_*$:

$$(3.6a) \quad \text{Monotonicity: } e_{\mathcal{T}_*}(u_i; \mathcal{T}) \leq e_{\mathcal{T}}(u_i; \mathcal{T}) \text{ for all } \mathcal{T} \in \mathcal{T}_*;$$

$$(3.6b) \quad \text{Error reduction: there exist constants } \lambda_i < 1, \text{ such that } e_{\mathcal{T}_*}^2(u_i; \mathcal{T}) \leq \lambda_i e_{\mathcal{T}}^2(u_i; \mathcal{T}) \text{ for all } \mathcal{T} \in \mathcal{R}_{\mathcal{T} \rightarrow \mathcal{T}_*};$$

$$(3.6c) \quad \text{Localized upper bound: } |I_{\mathcal{T}_*} u_i - I_{\mathcal{T}} u_i|_{i; \Omega} \leq e_{\mathcal{T}}(u_i; \mathcal{R}_{\mathcal{T} \rightarrow \mathcal{T}_*});$$

$$(3.6d) \quad \text{Approximability: there exists } s_i > 0 \text{ such that for any } \epsilon > 0 \text{ there is a } \mathcal{T}_\epsilon^i \in \mathbb{T} \text{ with } \#\mathcal{T}_\epsilon^i - \#\mathcal{T}_0 \preccurlyeq \epsilon^{-\frac{1}{s_i}} \text{ and } e_{\mathcal{T}_\epsilon^i} \leq \epsilon.$$

Variants of the assumptions on monotonicity (3.6a), error reduction (3.6b) and localized upper bound (3.6c) are given for AFEM in §4. The proof of the following lemma can be seen as a simplification of the convergence and optimality proofs in §5 and §6. It thus serves as a road-map for the rest of the paper.

Lemma 3.2 (Collective Marking). *Let assumptions (3.6) hold.*

Then Dörfler marking with minimal cardinality for the total error yields a sequence of meshes $\mathcal{T}_k \in \mathbb{T}$ with

$$e_k \preccurlyeq (\#\mathcal{T}_k - \#\mathcal{T}_0)^{-s}$$

for all marking parameters $0 < \theta < 1$. The convergence rate $s := \min\{s_1, s_2\}$ is quasi-optimal.

Proof. Let $\mathcal{T} \in \mathbb{T}$ and $\mathcal{M} \subset \mathcal{T}$ satisfy Dörfler property with parameter θ

$$(3.7) \quad e_{\mathcal{T}}(\mathcal{M}) \geq \theta e_{\mathcal{T}}(\mathcal{T})$$

for the total error $e_{\mathcal{T}}(\mathcal{T})$. Let $\mathcal{T}_* \geq \mathcal{T}$ satisfy $\mathcal{M} \subset \mathcal{R} := \mathcal{R}_{\mathcal{T}_k \rightarrow \mathcal{T}_*}$. We prove now the same key properties that are needed to derive the main results in §5 and §6:

1. *Contraction Property.* Let $\lambda := \min\{\lambda_1, \lambda_2\}$. In view of (3.6b) and (3.7), we have

$$e_{\mathcal{T}_*}^2(\mathcal{T}_*) \leq \lambda e_{\mathcal{T}}^2(\mathcal{R}) + e_{\mathcal{T}}^2(\mathcal{T} \setminus \mathcal{R}) = e_{\mathcal{T}}^2(\mathcal{T}) - (1 - \lambda) e_{\mathcal{T}}^2(\mathcal{R}) \leq \alpha^2 e_{\mathcal{T}}^2(\mathcal{T}),$$

with $\alpha^2 = (1 - \theta^2(1 - \lambda)) < 1$. This mimics the estimate of Theorem 5.1.

2. *Optimal Marking.* If $\mu := 1 - \theta$ and \mathcal{T}_* satisfies $e_{\mathcal{T}_*} \leq \mu e_{\mathcal{T}}$, then (3.7) is valid with \mathcal{R} instead of \mathcal{M} ; this mimics Lemma 6.7. In fact, combining the triangle inequality with (3.6c), we have

$$(1 - \mu) e_{\mathcal{T}} \leq e_{\mathcal{T}} - e_{\mathcal{T}_*} \leq e_{\mathcal{T}}(\mathcal{R}).$$

3. *Quasi-Optimal Decay.* Given $\epsilon > 0$, let $\mathcal{T}_\epsilon^i \in \mathbb{T}$ satisfy (3.6d). Then (2.2) gives

$$\#\mathcal{T}_\epsilon - \#\mathcal{T}_0 \preccurlyeq \epsilon^{-\frac{1}{s}}, \quad e_{\mathcal{T}_\epsilon} \leq 2\epsilon,$$

for the overlay $\mathcal{T}_\epsilon := \mathcal{T}_\epsilon^1 \oplus \mathcal{T}_\epsilon^2$. Equivalently, we infer that $e_{\mathcal{T}_\epsilon} \preccurlyeq (\#\mathcal{T}_\epsilon - \#\mathcal{T}_0)^{-s}$, whence the quasi-optimal decay of the total error is dictated by s .

4. *Cardinality of \mathcal{M} .* Let $\epsilon := \frac{1}{2}\mu e_{\mathcal{T}}$, $\mu = 1 - \theta$, let \mathcal{T}_ϵ be as in step 3, and let $\mathcal{T}_* := \mathcal{T}_\epsilon \oplus \mathcal{T}$. In view of (3.6a), we deduce that $e_{\mathcal{T}_*} \leq e_{\mathcal{T}_\epsilon} \leq 2\epsilon = \mu e_{\mathcal{T}}$ and, in view of step 2, that \mathcal{T}_* satisfies (3.7) with \mathcal{R} instead of \mathcal{M} . This, in conjunction with (2.2) and the minimality of $\#\mathcal{M}$, yields a result similar to Lemma 6.8:

$$\#\mathcal{M} \leq \#\mathcal{R} \preccurlyeq \#\mathcal{T}_* - \#\mathcal{T} \leq \#\mathcal{T}_\epsilon - \#\mathcal{T}_0 \preccurlyeq \epsilon^{-\frac{1}{s}} \preccurlyeq \mu^{-\frac{1}{s}} e_{\mathcal{T}}^{-\frac{1}{s}}.$$

5. *Counting DOFs.* Since $e_k \leq \alpha^{k-j} e_j$, according to step 1, Lemma 2.3 gives

$$\#\mathcal{T}_k - \#\mathcal{T}_0 \preceq \sum_{j=0}^{k-1} \#\mathcal{M}_j \preceq \mu^{-\frac{1}{s}} \sum_{j=0}^{k-1} e_j^{-\frac{1}{s}} \leq \mu^{-\frac{1}{s}} e_k^{-\frac{1}{s}} \sum_{j=1}^k \alpha^{\frac{j}{s}} \preceq \mu^{-\frac{1}{s}} \frac{\alpha^{\frac{1}{s}}}{1 - \alpha^{\frac{1}{s}}} e_k^{-\frac{1}{s}};$$

this mimics Theorem 6.9. This is the asserted quasi-optimal estimate in disguise, and completes the proof. \square

3.3. Marking for Optimal Cardinality. We conclude from §3.2 that collective marking is preferable to separate marking in computing adaptive approximations of functions with different asymptotic error decay. According to Lemma 3.1, separate marking requires a critical choice of parameters θ_i to obtain optimal cardinality of grids with respect to the total error e_k . Revisiting MNS in light of Lemma 3.1, we could identify the estimator η_k with the error $e_k(u_1)$ and oscillation osc_k with the error $e_k(u_2)$. We observe that $\text{osc}_k \leq \eta_k$ combined with (3.2b) implies $\alpha_2 \leq \alpha_1$ and that osc_k is generically of higher order than η_k , thereby yielding $s_1 < s_2$.

We wonder whether or not the optimality condition $\alpha_1 \leq \alpha_2^{s_1/s_2}$ is valid. Note that $\alpha_2^{s_1/s_2}$ increases as the gap between s_1 and s_2 increases. Since the oscillation reduction estimate of [16] reveals that α_2 increases as θ_{osc} decreases, we see that separate marking may be optimal for a wide range of marking parameters $\theta_{\text{est}}, \theta_{\text{osc}}$; this is confirmed by the numerical experiments in §3.1 even though it is unclear whether η_k and osc_k satisfy (3.2). However, choosing marking parameters $\theta_{\text{est}}, \theta_{\text{osc}}$ is rather tricky in practice because neither the explicit dependence of average reduction rates α_1, α_2 on $\theta_{\text{est}}, \theta_{\text{osc}}$ is known nor the optimal exponents s_1, s_2 .

In contrast, Lemma 3.2 shows that collective marking is always optimal. Using the crucial observation that estimator dominates oscillation [6], we obtain

$$\eta_k^2(U_k, \mathcal{T}_k) + \text{osc}_k^2(U_k, \mathcal{T}_k) \approx \eta_k^2(U_k, \mathcal{T}_k).$$

Hence, collective marking for estimator and oscillation simply reduces to just marking for the estimator, as proposed for AFEM in §2.7. We stress that oscillation does not have to be computed, which turns out to be quite advantageous for its practical realization.

We conclude by reviewing how alternative optimal AFEM compensate for the lack of optimality of separate marking. Binev, Dahmen, and DeVore [3] added a coarsening step to MNS to prove optimality. Veiser [25] was the first to mark oscillation relative to the error estimator to prove convergence of AFEM for the p -Laplacian. More recently, Stevenson [22] resorted to a similar marking to prove optimality of AFEM. This algorithm is discussed next in more detail.

For the Poisson problem, Stevenson [22] replaces the separate marking for both estimator and oscillation by the following conditional inner loop:

```

 $\delta_k = 2 C_1 \eta_k(U_k, \mathcal{T}_k);$ 
do  $\delta_k = \delta_k/2$ 
   $[\mathcal{T}_k, f_k] = \text{RHS}(f, \mathcal{T}_k, \delta_k);$ 
   $U_k = \text{SOLVE}(\mathcal{T}_k);$ 
   $\{\eta_k(U_k, T)\}_{T \in \mathcal{T}_k} = \text{ESTIMATE}(U_k, \mathcal{T}_k);$ 
  if  $C_1 \eta_k(U_k, \mathcal{T}_k) < \text{tol}$ , STOP;
until  $\delta_k < \omega \eta_k(U_k, \mathcal{T}_k).$ 

```

The purpose of this iteration is to render data oscillation $\text{osc}_k = \|f - f_k\|_{H^{-1}(\Omega)}$ small relative to the error estimator $\eta_k(U_k, \mathcal{T}_k)$, namely $\text{osc}_k \leq \delta_k < \omega \eta_k(U_k, \mathcal{T}_k)$ with $\omega > 0$ sufficiently small. The output of RHS is a conforming mesh $\mathcal{T}_k \in \mathbb{T}$ and a

piecewise constant function f_k over \mathcal{T}_k such that $\text{osc}_k \leq \delta_k$. Stevenson assumes that oscillation can at least be approximated with the same rate as the energy error. If oscillation is small, it may happen that RHS does not change the actual grid. However, if oscillation is large, the inner loop may be traversed several times. This loop enables Stevenson to prove optimality for non-zero oscillation at the price that SOLVE and ESTIMATE may be called repeatedly in one adaptive iteration. This is inconsistent with practice; in particular, SOLVE is typically the most expensive procedure of AFEM and should be called only once per adaptive loop k .

4. AUXILIARY RESULTS

In this section we prove a quasi-reduction property of the estimator, which is instrumental for convergence. We also show a perturbation result of oscillation and a localized upper bound, both between two discrete functions, that play an essential role in the optimality proof. Finally we discuss the overlay of two meshes and thereby prove (2.2). In the remainder of the paper the constants hidden in ‘ \lesssim ’ solely depend on shape regularity, thus on \mathcal{T}_0 , the number of bisections b of marked elements, and the polynomial degree n . Dependence on data $\mathbf{D} = (\mathbf{A}, c)$ and f of (1.1) is traced explicitly.

4.1. Reduction of Error Estimator and Oscillation. We relate the error indicators and oscillation of two nested triangulations to each other. The link involves weighted maximum-norms of the coefficient functions \mathbf{D} , or their oscillation. These results are the basis for our analyses in §5 and §6.

We start by defining the weighted maximum-norm of the coefficients \mathbf{D} and their oscillation. For $m \in \mathbb{N}_0$, $\mathcal{T} \in \mathbb{T}$, and $v \in L^\infty(\Omega)$, we recall that $\Pi_m^\infty v$ is the best $L^\infty(\Omega)$ -approximation in the space of discontinuous polynomials of degree $\leq m$. We further set $\Pi_{-1}^\infty v = 0$, $P_m^\infty v = v - \Pi_m^\infty v$ and

$$\begin{aligned} \eta_{\mathcal{T}}^2(\mathbf{D}, T) &:= h_T^2 \left(\|\text{div } \mathbf{A}\|_{L^\infty(T)}^2 + h_T^{-2} \|\mathbf{A}\|_{L^\infty(\omega_T)}^2 + \|c\|_{L^\infty(T)}^2 \right), \\ \text{osc}_{\mathcal{T}}^2(\mathbf{D}, T) &:= h_T^2 \left(\|P_{n-1}^\infty \text{div } \mathbf{A}\|_{L^\infty(T)}^2 + h_T^{-2} \|P_n^\infty \mathbf{A}\|_{L^\infty(\omega_T)}^2 \right. \\ &\quad \left. + h_T^2 \|P_{n-2}^\infty c\|_{L^\infty(T)}^2 + \|P_{n-1}^\infty c\|_{L^\infty(T)}^2 \right); \end{aligned}$$

note that P_m^∞ is defined elementwise. For any subset $\mathcal{T}' \subset \mathcal{T}$ we finally set

$$\eta_{\mathcal{T}}(\mathbf{D}, \mathcal{T}') := \max_{T \in \mathcal{T}'} \eta_{\mathcal{T}}(\mathbf{D}, T) \quad \text{and} \quad \text{osc}_{\mathcal{T}}(\mathbf{D}, \mathcal{T}') := \max_{T \in \mathcal{T}'} \text{osc}_{\mathcal{T}}(\mathbf{D}, T).$$

Remark 4.1 (Monotonicity). The use of best approximation in L^∞ in the definitions of $\eta_{\mathcal{T}}(\mathbf{D}, T)$ and $\text{osc}_{\mathcal{T}}(\mathbf{D}, T)$ implies the following monotonicity property: For all $\mathcal{T}, \mathcal{T}_* \in \mathbb{T}$ with $\mathcal{T} \leq \mathcal{T}_*$, there holds

$$\eta_{\mathcal{T}_*}(\mathbf{D}, \mathcal{T}_*) \leq \eta_{\mathcal{T}}(\mathbf{D}, \mathcal{T}) \quad \text{and} \quad \text{osc}_{\mathcal{T}_*}(\mathbf{D}, \mathcal{T}_*) \leq \text{osc}_{\mathcal{T}}(\mathbf{D}, \mathcal{T}).$$

The following variant of the Bramble-Hilbert Lemma allows us to avoid any smoothness assumption on the coefficients \mathbf{D} of the PDE.

Lemma 4.2 (Implicit Interpolation). *Let ω be either a d or a $(d-1)$ simplex. For $\ell \in \mathbb{N}$ denote by $\Pi_m^p : L^p(\omega, \mathbb{R}^\ell) \rightarrow \mathbb{P}_m(\omega, \mathbb{R}^\ell)$ the operator of best L^p -approximation in ω and set $P_m^p = \text{id} - \Pi_m^p$. Then, for all $v \in L^\infty(\omega, \mathbb{R}^\ell)$, $V \in \mathbb{P}_n(\omega, \mathbb{R}^\ell)$ and $m \geq n$, there holds*

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

Proof. Since $V \in \mathbb{P}_n$ we obtain $\Pi_{m-n}^\infty(v) \cdot V \in \mathbb{P}_m$ and thus the orthogonality of the L^2 -projection yields

$$\begin{aligned} \int_{\omega} (vV - \Pi_m^2(vV))(vV - \Pi_m^2(vV)) &= \int_{\omega} (vV - \Pi_m^2(vV))(v - \Pi_{m-n}^\infty(v)) \cdot V \\ &\leq \|P_m^2(vV)\|_{L^2(\omega)} \|P_{m-n}^\infty v\|_{L^\infty(\omega)} \|V\|_{L^2(\omega)} \end{aligned}$$

which finishes the proof. \square

Proposition 4.3 (Local Perturbation). *Let $\mathcal{T} \in \mathbb{T}$. For all $T \in \mathcal{T}$ and for any pair of discrete functions $V, W \in \mathbb{V}(\mathcal{T})$, we have*

$$(4.1) \quad \eta_{\mathcal{T}}(V, T) \leq \eta_{\mathcal{T}}(W, T) + \bar{\Lambda}_1 \eta_{\mathcal{T}}(\mathbf{D}, T) \|V - W\|_{H^1(\omega_T)},$$

$$(4.2) \quad \text{osc}_{\mathcal{T}}(V, T) \leq \text{osc}_{\mathcal{T}}(W, T) + \bar{\Lambda}_1 \text{osc}_{\mathcal{T}}(\mathbf{D}, T) \|V - W\|_{H^1(\omega_T)},$$

where ω_T is the union of elements in \mathcal{T} sharing a side with T . The constant $\bar{\Lambda}_1 > 0$ only depends on the shape-regularity of \mathcal{T}_0 and the polynomial degree n .

Proof. We only prove (4.2), estimate (4.1) is somewhat simpler and can be derived similarly. Recall the definition of the element residual $R(V) = f + \mathcal{L}(V)$ and the notation $P_m^p = \text{id} - \Pi_m^p$. Since the L^2 projection Π_m^2 is linear, by adding and subtracting W and using the triangle inequality, we obtain

$$\text{osc}_{\mathcal{T}}(V, T) \leq \text{osc}_{\mathcal{T}}(W, T) + h_T \|P_{2n-2}^2 \mathcal{L}(E)\|_{L^2(T)} + h_T^{1/2} \|P_{2n-1}^2 J(E)\|_{L^2(\partial T)}$$

with $E := V - W$. It remains to show that the second and third terms are bounded by $\text{osc}_{\mathcal{T}}(\mathbf{D}, T)$ times the local H^1 -norm of E . For $\mathcal{L}(E) = \text{div } \mathbf{A} \nabla E - cE$ we have

$$\|P_{2n-2}^2 \mathcal{L}(E)\|_{L^2(T)} \leq \|P_{2n-2}^2 (\text{div } \mathbf{A} \nabla E)\|_{L^2(T)} + \|P_{2n-2}^2 (cE)\|_{L^2(T)}.$$

We split the divergence term as

$$\text{div}(\mathbf{A} \nabla E) = \text{div } \mathbf{A} \cdot \nabla E + \mathbf{A} : D^2 E,$$

where $D^2 E$ is the Hessian of E . Invoking Lemma 4.2 with $\omega = T$, and observing that the polynomial degree of ∇E is $n-1$, we infer for the first term that

$$\|P_{2n-2}^2 (\text{div } \mathbf{A} \cdot \nabla E)\|_{L^2(T)} \leq \|P_{n-1}^\infty \text{div } \mathbf{A}\|_{L^\infty(T)} \|\nabla E\|_{L^2(T)}.$$

Since $D^2 E$ is a polynomial of degree $\leq n-2$, applying Lemma 4.2 again in conjunction with an inverse inequality, we have for the second term

$$\|P_{2n-2}^2 (\mathbf{A} : D^2 E)\|_{L^2(T)} \preccurlyeq h_T^{-1} \|P_n^\infty \mathbf{A}\|_{L^\infty(T)} \|\nabla E\|_{L^2(T)}.$$

To analyse the reaction term we write

$$\|P_{2n-2}^2 (cE)\|_{L^2(T)} \leq \|P_{2n-2}^2 (c\Pi_0^2 E)\|_{L^2(T)} + \|P_{2n-2}^2 (cP_0^2 E)\|_{L^2(T)}.$$

Applying again Lemma 4.2 we have for the first term

$$\|P_{2n-2}^2 (c\Pi_0^2 E)\|_{L^2(T)} \leq \|P_{2n-2}^\infty c\|_{L^\infty(T)} \|E\|_{L^2(T)},$$

and for the last one

$$\|P_{2n-2}^2 (cP_0^2 E)\|_{L^2(T)} \preccurlyeq h_T \|P_{n-2}^\infty c\|_{L^\infty(T)} \|\nabla E\|_{L^2(T)}.$$

We now deal with the jump residual. Let $T' \in \mathcal{T}$ share an interior side σ with T . We write $J(E) = ((\mathbf{A}\nabla E)|_T - (\mathbf{A}\nabla E)|_{T'}) \cdot \boldsymbol{\nu}$, use linearity of Π_{2n-1}^2 , Lemma 4.2 with $\omega = \sigma$ and the inverse inequality $\|\nabla E\|_{L^2(\sigma)} \lesssim h_T^{-1/2} \|\nabla E\|_{L^2(T)}$, to deduce

$$\|P_{2n-1}^2(\mathbf{A}\nabla E)|_T \cdot \boldsymbol{\nu}\|_{L^2(\sigma)} \lesssim h_T^{-\frac{1}{2}} \|P_n^\infty \mathbf{A}\|_{L^\infty(T)} \|\nabla E\|_{L^2(T)}.$$

The same argument holds for T' and, since \mathcal{T} is shape regular, we can replace $h_{T'}$ by h_T . Finally, collecting the above estimates for T and all its neighbors yields the assertion (4.2). \square

The following two Corollaries are global forms of the above Lemma.

Corollary 4.4 (Estimator Reduction). *For $\mathcal{T} \in \mathbb{T}$ and $\mathcal{M} \subset \mathcal{T}$ let $\mathcal{T}_* \in \mathbb{T}$ be given by $\mathcal{T}_* := \text{REFINE}(\mathcal{T}, \mathcal{M})$. If $\Lambda_1 := (d+1)\bar{\Lambda}_1^2/c_B$ with $\bar{\Lambda}_1$ from Proposition 4.3 and $\lambda := 1 - 2^{-\frac{b}{d}} > 0$, then there holds for all $V \in \mathbb{V}(\mathcal{T})$, $V_* \in \mathbb{V}(\mathcal{T}_*)$ and any $\delta > 0$*

$$\begin{aligned} \eta_{\mathcal{T}_*}^2(V_*, \mathcal{T}_*) &\leq (1 + \delta) \{ \eta_{\mathcal{T}}^2(V, \mathcal{T}) - \lambda \eta_{\mathcal{T}}^2(V, \mathcal{M}) \} \\ &\quad + (1 + \delta^{-1}) \Lambda_1 \eta_{\mathcal{T}_0}^2(\mathbf{D}, \mathcal{T}_0) \|V_* - V\|_\Omega^2. \end{aligned}$$

Proof. Applying Proposition 4.3 with $V_*, V \in \mathbb{V}(\mathcal{T}_*)$ over $T \in \mathcal{T}_*$, and using Young's inequality with parameter δ , we derive

$$\eta_{\mathcal{T}_*}^2(V_*, T) \leq (1 + \delta) \eta_{\mathcal{T}_*}^2(V, T) + (1 + \delta^{-1}) \bar{\Lambda}_1^2 \eta_{\mathcal{T}_*}^2(\mathbf{D}, T) \|V_* - V\|_{H^1(\omega_T)}^2.$$

Summing over all elements $T \in \mathcal{T}_*$, using the finite overlap property of patches ω_T , and the equivalence of the H^1 -norm and the energy-norm in Ω , we have the following direct by-product:

$$\eta_{\mathcal{T}_*}^2(V_*, \mathcal{T}_*) \leq (1 + \delta) \eta_{\mathcal{T}_*}^2(V, \mathcal{T}_*) + (1 + \delta^{-1}) \Lambda_1 \eta_{\mathcal{T}_*}^2(\mathbf{D}, \mathcal{T}_*) \|V_* - V\|_\Omega^2.$$

For a marked element $T \in \mathcal{M}$, we set $\mathcal{T}_{*,T} := \{T' \in \mathcal{T}_* \mid T' \subset T\}$. Since $V \in \mathbb{V}(\mathcal{T})$ and \mathbf{A} jumps only across sides of \mathcal{T}_0 we see that $J(V) = 0$ on sides of $\mathcal{T}_{*,T}$ in the interior of T . We then obtain

$$\sum_{T' \in \mathcal{T}_{*,T}} \eta_{\mathcal{T}_*}^2(V, T') \leq 2^{-\frac{b}{d}} \eta_{\mathcal{T}}^2(V, T),$$

because refinement by bisection implies $h_{T'} = |T'|^{\frac{1}{d}} \leq (2^{-b}|T|)^{\frac{1}{d}} \leq 2^{-\frac{b}{d}} h_T$ for all $T' \in \mathcal{T}_{*,T}$. For an element $T \in \mathcal{T} \setminus \mathcal{M}$, instead, Remark 2.1 yields $\eta_{\mathcal{T}_*}(V, T) \leq \eta_{\mathcal{T}}(V, T)$. Hence, summing over all $T \in \mathcal{T}_*$, we arrive at

$$\eta_{\mathcal{T}_*}^2(V, \mathcal{T}_*) \leq \eta_{\mathcal{T}}^2(V, \mathcal{T} \setminus \mathcal{M}) + 2^{-\frac{b}{d}} \eta_{\mathcal{T}}^2(V, \mathcal{M}) = \eta_{\mathcal{T}}^2(V, \mathcal{T}) - \lambda \eta_{\mathcal{T}}^2(V, \mathcal{M}).$$

The assertion finally follows from the monotonicity $\eta_{\mathcal{T}_*}(\mathbf{D}, \mathcal{T}_*) \leq \eta_{\mathcal{T}_0}(\mathbf{D}, \mathcal{T}_0)$ stated in Remark 4.1. \square

Corollary 4.5 (Perturbation of Oscillation). *Let $\mathcal{T}, \mathcal{T}_* \in \mathbb{T}$ with $\mathcal{T} \leq \mathcal{T}_*$ and let $\Lambda_1 = (d+1)\bar{\Lambda}_1^2/c_B$ be as in Corollary 4.4. For all $V \in \mathbb{V}(\mathcal{T})$, $V_* \in \mathbb{V}(\mathcal{T}_*)$, we have*

$$\text{osc}_{\mathcal{T}}^2(V, \mathcal{T} \cap \mathcal{T}_*) \leq 2 \text{osc}_{\mathcal{T}_*}^2(V_*, \mathcal{T} \cap \mathcal{T}_*) + 2\Lambda_1 \text{osc}_{\mathcal{T}_0}^2(\mathbf{D}, \mathcal{T}_0) \|V - V_*\|_\Omega^2.$$

Proof. Remark 2.1 yields $\text{osc}_{\mathcal{T}}(V, T) = \text{osc}_{\mathcal{T}_*}(V, T)$ for all $T \in \mathcal{T} \cap \mathcal{T}_*$, whence

$$\text{osc}_{\mathcal{T}}^2(V, T) \leq 2 \text{osc}_{\mathcal{T}_*}^2(V_*, T) + 2\bar{\Lambda}_1^2 \text{osc}_{\mathcal{T}_*}^2(\mathbf{D}, T) \|V - V_*\|_{H^1(\omega_T)}^2,$$

by (4.2) and Young's inequality. We now sum over $T \in \mathcal{T} \cap \mathcal{T}_*$, use the equivalence of H^1 and energy norms in Ω , as well as the monotonicity property $\text{osc}_{\mathcal{T}_*}(\mathbf{D}, \mathcal{T}_*) \leq \text{osc}_{\mathcal{T}_0}(\mathbf{D}, \mathcal{T}_0)$ stated in Remark 4.1, to prove the assertion. \square

4.2. Localized Upper Bound. To prove the optimality of AFEM, we need a localized upper bound for the distance between two nested solutions. This slightly improves a similar result by Stevenson [22] in the sense that the error can be estimated here only using the indicators of refined elements, without a buffer layer.

Lemma 4.6 (Localized Upper Bound). *For $\mathcal{T}, \mathcal{T}_* \in \mathbb{T}$ with $\mathcal{T} \leq \mathcal{T}_*$ let $\mathcal{R} = \mathcal{R}_{\mathcal{T} \rightarrow \mathcal{T}_*}$ be the set of refined elements. Let $U \in \mathbb{V}(\mathcal{T})$ and $U_* \in \mathbb{V}(\mathcal{T}_*)$ be the discrete solutions of (2.1). Then the following localized upper bound is valid*

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

Proof. We start with an observation. Let $\Omega_* = \bigcup\{T : T \in \mathcal{R}\}$ be the union of refined elements and let $\mathbb{V}(\mathcal{R})$ be the restriction of $\mathbb{V}(\mathcal{T})$ to Ω_* . Denote by $P_{\mathcal{R}} : H^1(\Omega_*) \rightarrow \mathbb{V}(\mathcal{R})$ the Scott-Zhang interpolation operator over the triangulation \mathcal{R} [21]. The following interpolation estimate holds for all $v \in H^1(\Omega_*)$

$$(4.3) \quad \sum_{T \in \mathcal{R}} h_T^{-2} \|v - P_{\mathcal{R}} v\|_{L_2(T)}^2 + h_T \|v - P_{\mathcal{R}} v\|_{L_2(\partial T)}^2 \preccurlyeq \|\nabla v\|_{L^2(\Omega_*)}^2,$$

where the constant hidden in ' \preccurlyeq ' does not depend on Ω_* but only on the shape-regularity of the underlying triangulation \mathcal{R} , and thus on \mathcal{T}_0 . The operator $P_{\mathcal{R}}$ is a projection, i.e. $P_{\mathcal{R}} V = V$ for all $V \in \mathbb{V}(\mathcal{R})$, and it preserves homogeneous boundary values. Hence, it also preserves conforming boundary values, i.e. $P_{\mathcal{R}} v = v$ on $\partial\Omega_*$ whenever $v = V$ on $\partial\Omega_*$ for some $V \in \mathbb{V}(\mathcal{R})$. For the error $E_* := U_* - U \in \mathbb{V}(\mathcal{T}_*)$ we construct an approximation $V \in \mathbb{V}(\mathcal{T})$ by

$$V := \begin{cases} E_* & \text{in } \Omega \setminus \Omega_* \\ P_{\mathcal{R}} E_* & \text{in } \Omega_*. \end{cases}$$

Since E_* has conforming boundary values on $\partial\Omega_*$ in $\mathbb{V}(\mathcal{R})$, we conclude that V is continuous in Ω implying $V \in \mathbb{V}(\mathcal{T})$ and V is an H^1 -stable approximation to E_* .

Since $\mathbb{V}(\mathcal{T}) \subset \mathbb{V}(\mathcal{T}_*)$ are nested subspaces of \mathbb{V} , by Galerkin orthogonality $\mathcal{B}[E_*, E_*] = \mathcal{B}[E_*, E_* - V]$, we obtain by standard arguments

$$\mathcal{B}[E_*, E_*] = \sum_{T \in \mathcal{R}} \langle R(U), E_* - V \rangle_T + \frac{1}{2} \langle J(U), E_* - V \rangle_{\partial T} \preccurlyeq \eta_{\mathcal{T}}(U, \mathcal{R}) \|\nabla E_*\|_{L^2(\Omega_*)},$$

where we have used (4.3) in the last step. This, in conjunction with the coercivity of \mathcal{B} , proves the proposition. \square

4.3. Overlay of Meshes. We finish the auxiliary results with a counting argument for the overlay $\mathcal{T}_1 \oplus \mathcal{T}_2$ of two triangulations $\mathcal{T}_1, \mathcal{T}_2$. As a consequence of the following lemma we see that for two conforming triangulations $\mathcal{T}_1, \mathcal{T}_2 \in \mathbb{T}$ the overlay is the smallest conforming triangulation $\mathcal{T} \in \mathbb{T}$ with $\mathcal{T}_1, \mathcal{T}_2 \leq \mathcal{T}$.

Lemma 4.7 (Overlay of Meshes). *For $\mathcal{T}_1, \mathcal{T}_2 \in \mathbb{T}$ the overlay $\mathcal{T} := \mathcal{T}_1 \oplus \mathcal{T}_2$ is conforming, i.e. $\mathcal{T} \in \mathbb{T}$, and satisfies*

$$\#\mathcal{T} \leq \#\mathcal{T}_1 + \#\mathcal{T}_2 - \#\mathcal{T}_0.$$

Proof. Assume, that \mathcal{T} contains a non-conforming vertex z . Then there exist $T_1, T_2 \in \mathcal{T}$ with a common edge such that z is a vertex of T_1 and $z \in T_2$ but z is no vertex of T_2 . Without loss of generality let $T_1 \in \mathcal{T}_1$. Since \mathcal{T}_1 is conforming, there exists a $T' \in \mathcal{T}_1$, $T' \subset T_2$ such that z is a vertex of T' . Hence, T_2 cannot be a leaf node of $\mathcal{F}(\mathcal{T})$, i.e. $T_2 \notin \mathcal{T}$, a contradiction.

For $T \in \mathcal{T}_0$ and $i = 1, 2$ we denote by $\mathcal{F}_i(T) \subset \mathcal{F}(T)$ the binary trees with root T corresponding to \mathcal{T}_i and let $\mathcal{T}_i(T)$ be the triangulation given by the leaf nodes of $\mathcal{F}_i(T)$. Since $\mathcal{T}(T) \subset \mathcal{T}_1(T) \cup \mathcal{T}_2(T)$, we infer that $\#\mathcal{T}(T) \leq \#\mathcal{T}_1(T) + \#\mathcal{T}_2(T)$. We now show $\#\mathcal{T}(T) \leq \#\mathcal{T}_1(T) + \#\mathcal{T}_2(T) - 1$ by distinguishing two cases:
Case 1: $\mathcal{T}_1(T) \cap \mathcal{T}_2(T) \neq \emptyset$. Then there exists $T' \in \mathcal{T}_1(T) \cap \mathcal{T}_2(T)$, and so $T' \in \mathcal{T}(T)$. Counting T' only once in $\#(\mathcal{T}_1(T) \cup \mathcal{T}_2(T))$ we get $\#\mathcal{T}(T) \leq \#\mathcal{T}_1(T) + \#\mathcal{T}_2(T) - 1$.
Case 2: $\mathcal{T}_1(T) \cap \mathcal{T}_2(T) = \emptyset$. Then there exists $T' \in \mathcal{T}_1(T)$ (resp. $T' \in \mathcal{T}_2(T)$) so that $T' \notin \mathcal{T}(T)$, for otherwise $T' \in \mathcal{T}_2(T)$ (resp. $T' \in \mathcal{T}_1(T)$) thereby contradicting the assumption. We obtain again $\#\mathcal{T}(T) \leq \#\mathcal{T}_1(T) + \#\mathcal{T}_2(T) - 1$.

Finally, since $\mathcal{T}_i = \bigcup_{T \in \mathcal{T}_0} \mathcal{T}_i(T)$, the assertion follows by adding over $T \in \mathcal{T}_0$. \square

5. CONTRACTION PROPERTY OF AFEM

We now prove that AFEM is a contraction with respect to the sum of energy error plus scaled error estimator, the so-called *quasi-error*. Consequently, the quasi-error is reduced by a fixed rate at every step. This can be motivated heuristically as follows: in light of (2.4) the energy error $\|u - U_k\|_\Omega$ decreases strictly, unless $U_{k+1} = U_k$, in which case the estimator $\eta_k(U_k, \mathcal{T}_k)$ does according to Corollary 4.4. The theorem below makes this observation quantitative.

Theorem 5.1 (Contraction Property). *Let $\theta \in (0, 1]$ and let $\{\mathcal{T}_k, \mathbb{V}_k, U_k\}_{k \geq 0}$ be the sequence of meshes, finite element spaces, and discrete solutions produced by AFEM.*

Then, there exist constants $\gamma > 0$, and $0 < \alpha < 1$, depending solely on the shape-regularity of \mathcal{T}_0 , b , and the marking-parameter $0 < \theta \leq 1$, such that

$$\|u - U_{k+1}\|_\Omega^2 + \gamma \eta_{k+1}^2(U_{k+1}, \mathcal{T}_{k+1}) \leq \alpha^2 \left(\|u - U_k\|_\Omega^2 + \gamma \eta_k^2(U_k, \mathcal{T}_k) \right).$$

Proof. For convenience, we use the notation

$$e_k := \|u - U_k\|_\Omega, \quad E_k := \|U_{k+1} - U_k\|_\Omega, \\ \eta_k := \eta_k(U_k, \mathcal{T}_k), \quad \eta_k(\mathcal{M}_k) := \eta_k(U_k, \mathcal{M}_k), \quad \eta_0(\mathbf{D}) := \eta_0(\mathbf{D}, \mathcal{T}_0).$$

We combine the orthogonality (2.4) with Corollary 4.4 to write

$$e_{k+1}^2 + \gamma \eta_{k+1}^2 \leq e_k^2 - E_k^2 + (1 + \delta) \gamma (\eta_k^2 - \lambda \eta_k^2(\mathcal{M}_k)) + (1 + \delta^{-1}) \gamma \Lambda_1 \eta_0^2(\mathbf{D}) E_k^2.$$

We choose γ dependent on δ to be

$$(5.1) \quad \gamma := \frac{1}{(1 + \delta^{-1}) \Lambda_1 \eta_0^2(\mathbf{D})} \quad \Leftrightarrow \quad \gamma (1 + \delta) = \frac{\delta}{\Lambda_1 \eta_0^2(\mathbf{D})}$$

to obtain

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

Invoking Dörfler marking (2.8), we deduce

$$e_{k+1}^2 + \gamma \eta_{k+1}^2 \leq e_k^2 + (1 + \delta) \gamma \eta_k^2 - (1 + \delta) \lambda \theta^2 \gamma \eta_k^2.$$

We rewrite this inequality as follows with any $\beta \in (0, 1)$

$$e_{k+1}^2 + \gamma \eta_{k+1}^2 \leq e_k^2 + (1 + \delta) \gamma \eta_k^2 - \beta (1 + \delta) \lambda \theta^2 \gamma \eta_k^2 - (1 - \beta) (1 + \delta) \lambda \theta^2 \gamma \eta_k^2,$$

apply the upper bound (2.6) and replace γ according to (5.1) to obtain

$$e_{k+1}^2 + \gamma \eta_{k+1}^2 \leq \alpha_1^2(\delta, \beta) e_k^2 + \gamma \alpha_2^2(\delta, \beta) \eta_k^2$$

with

$$\alpha_1^2(\delta, \beta) := 1 - \beta \frac{\lambda \theta^2}{C_1 \Lambda_1 \eta_0^2(\mathbf{D})} \delta, \quad \alpha_2^2(\delta, \beta) := (1 + \delta)(1 - (1 - \beta) \lambda \theta^2).$$

Now choosing $\delta > 0$ small enough yields

$$\alpha^2 := \max\{\alpha_1^2, \alpha_2^2\} < 1,$$

which is the desired result. \square

Remark 5.2 (Ingredients for Convergence). We stress that this new proof of linear convergence relies exclusively on the upper bound (2.6), the orthogonality relation (2.4), the error estimator reduction property of Corollary 4.4 and the Dörfler marking (2.8) for the estimator. It does not need any marking due to oscillation which turns out to be problematic for optimality in light of the discussion in §3. Equality (2.4) is only used to cancel, via (5.1), the contribution involving $\eta_0(\mathbf{D}, \mathcal{T}_0)$ in Corollary 4.4. Its role is much less prominent than in [7, 15, 16, 17, 22].

Moreover, we neither use the lower bound (2.7) nor a discrete lower bound for proving convergence. The latter hinges on the rather demanding *interior node property*: every element of \mathcal{M}_k , as well as its adjacent elements, contains a node of \mathcal{T}_{k+1} in their interior as well as in the interior of their common sides. However, the global lower bound (2.7) will be instrumental to prove optimality in §6.

The treatment of oscillation and the interior node property is essential in [15, 16, 17], and so in [3, 22]. Our new approach simplifies the analysis and directly applies to any polynomial degree $n \geq 1$.

Remark 5.3 (Optimal Contraction Factor α). Consider

$$D := \{(\delta, \beta) \in \mathbb{R}_+ \times [0, 1] : 0 \leq \alpha_1^2(\delta, \beta), \alpha_2^2(\delta, \beta) \leq 1\}.$$

This set is nonempty, according to the proof of Theorem 5.1, is closed and bounded; thus D is compact. Since α_1^2, α_2^2 are continuous functions in D , $\alpha^2 = \max\{\alpha_1^2, \alpha_2^2\}$ attains its minimum in D . It turns out that $\alpha^2|_{\partial D} = 1$ by definition of D and $\alpha^2 < 1$ in the interior of D from the proof of Theorem 5.1. Consequently, α^2 attains an absolute minimum smaller than 1 and satisfies

$$\alpha^2 = \alpha_1^2 = \alpha_2^2.$$

To see this assume $\alpha_1^2 < \alpha_2^2$, and decrease the value of δ slightly. Since α_1^2 increases whereas α_2^2 decreases, this yields a contradiction to the minimality of α^2 . The case $\alpha_1^2 > \alpha_2^2$ is similar. In principle, this optimal value of α^2 can be computed explicitly.

Remark 5.4 (Range of γ and α). We see from (5.1) that $\gamma \approx \eta_0^{-2}(\mathbf{D}, \mathcal{T}_0)$ provided that $\eta_0(\mathbf{D}, \mathcal{T}_0)$ is large; this provides a lower bound for γ . An upper bound results from the condition $\alpha_1^2 > 0$ and $\delta \leq 1$, namely,

$$\alpha_1^2 = 1 - \frac{(1 + \delta)\beta\lambda\theta^2}{C_1} \gamma \quad \Rightarrow \quad \gamma < \frac{C_1}{2\beta\lambda\theta^2} \approx C_1.$$

On the other hand, it is clear from the definitions of α_1 and α_2 that α deteriorates if either $\eta_0(\mathbf{D}, \mathcal{T}_0)$ increases or θ decreases. In fact, take $\beta = \frac{1}{2}$ and $\delta = \frac{1}{2}\lambda\theta^2$ to get

$$\alpha_1^2 = 1 - \frac{1}{C_1 \Lambda_1 \eta_0^2(\mathbf{D}, \mathcal{T}_0)} \left(\frac{\lambda\theta^2}{2}\right)^2, \quad \alpha_2^2 = 1 - \left(\frac{\lambda\theta^2}{2}\right)^2.$$

Thus by the definition of α there exist constants $c, C > 0$ not depending on θ with

$$1 - C\theta^4 \leq \alpha^2 \leq 1 - c\theta^4.$$

6. QUASI-OPTIMAL CARDINALITY OF AFEM

Building on fundamental ideas of Binev, Dahmen and DeVore [3] and Stevenson [22], used in the first optimality proofs for variants of AFEM for the Poisson equation, we prove now optimal cardinality of a *standard* AFEM for the general symmetric elliptic problem (1.1). We thus improve and extend the results of [3, 22].

6.1. Approximation Class. To prove optimality in the present context, we need to seek a suitable error quantity being controlled by AFEM and its associated approximation class \mathbb{A}_s . On the one hand, oscillation is dominated by the estimator according to Remark 2.1, thereby yielding

$$\|u - U_k\|_\Omega^2 + \text{osc}_k^2(U_k, \mathcal{T}_k) \leq \|u - U_k\|_\Omega^2 + \eta_k^2(U_k, \mathcal{T}_k).$$

On the other hand, the lower bound (2.7) implies

$$\|u - U_k\|_\Omega^2 + \eta_k^2(U_k, \mathcal{T}_k) \leq (1 + C_2^{-1}) (\|u - U_k\|_\Omega^2 + \text{osc}_k^2(U_k, \mathcal{T}_k)).$$

We thus realize that

$$\|u - U_k\|_\Omega^2 + \eta_k^2(U_k, \mathcal{T}_k) \approx \|u - U_k\|_\Omega^2 + \text{osc}_k^2(U_k, \mathcal{T}_k),$$

and call the square root of the right-hand side the *total error*. This is equivalent to the quantity being reduced by AFEM, the quasi-error, and satisfies a Cea's Lemma.

Lemma 6.1 (Quasi-Optimality of the Total Error). *Let u be the solution of (2.1) and for $\mathcal{T} \in \mathbb{T}$ let $U \in \mathbb{V}(\mathcal{T})$ be the Ritz-Galerkin approximation of (2.3).*

Then, there exists a constant $C_{\mathbf{D}}$ only depending on data \mathbf{D} and shape-regularity of \mathcal{T}_0 such that

$$\|u - U\|_\Omega^2 + \text{osc}_{\mathcal{T}}^2(U, \mathcal{T}) \leq C_{\mathbf{D}} \inf_{V \in \mathbb{V}(\mathcal{T})} (\|u - V\|_\Omega^2 + \text{osc}_{\mathcal{T}}^2(V, \mathcal{T})).$$

Proof. For $\epsilon > 0$ choose $V_\epsilon \in \mathbb{V}(\mathcal{T})$ with

$$\|u - V_\epsilon\|_\Omega^2 + \text{osc}_{\mathcal{T}}^2(V_\epsilon, \mathcal{T}) \leq (1 + \epsilon) \inf_{V \in \mathbb{V}(\mathcal{T})} (\|u - V\|_\Omega^2 + \text{osc}_{\mathcal{T}}^2(V, \mathcal{T})).$$

Corollary 4.5 with $\mathcal{T}_* = \mathcal{T}$, $V = U$, and $V_* = V_\epsilon$ yields

$$\text{osc}_{\mathcal{T}}^2(U, \mathcal{T}) \leq 2 \text{osc}_{\mathcal{T}}^2(V_\epsilon, \mathcal{T}) + 2 \Lambda_1 \text{osc}_{\mathcal{T}_0}^2(\mathbf{D}, \mathcal{T}_0) \|U - V_\epsilon\|_\Omega^2.$$

Since U is the Galerkin solution, $\|u - U\|_\Omega^2 + \|U - V_\epsilon\|_\Omega^2 = \|u - V_\epsilon\|_\Omega^2$, whence

$$\begin{aligned} \|u - U\|_\Omega^2 + \text{osc}_{\mathcal{T}}^2(U, \mathcal{T}) &\leq (1 + 2 \Lambda_1 \text{osc}_{\mathcal{T}_0}^2(\mathbf{D}, \mathcal{T}_0)) \|u - V_\epsilon\|_\Omega^2 + 2 \text{osc}_{\mathcal{T}}^2(V_\epsilon, \mathcal{T}) \\ &\leq (1 + \epsilon) C_{\mathbf{D}}^2 \inf_{V \in \mathbb{V}(\mathcal{T})} (\|u - V\|_\Omega^2 + \text{osc}_{\mathcal{T}}^2(V, \mathcal{T})) \end{aligned}$$

with $C_{\mathbf{D}} = \max \{2, 1 + 2 \Lambda_1 \text{osc}_{\mathcal{T}_0}^2(\mathbf{D}, \mathcal{T}_0)\}$. The assertion follows from $\epsilon \rightarrow 0$. \square

This motivates the following definition of \mathbb{A}_s . Let $\mathbb{T}_N \subset \mathbb{T}$ be the set of all possible conforming triangulations generated from \mathcal{T}_0 with at most N elements more than \mathcal{T}_0 :

$$\mathbb{T}_N := \{\mathcal{T} \in \mathbb{T} \mid \#\mathcal{T} - \#\mathcal{T}_0 \leq N\}.$$

The quality of the best approximation to the total error in the set \mathbb{T}_N is given by

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

Notice that the solution u and coefficients \mathbf{D} interact in a nonlinear fashion through the oscillation term. We now define the nonlinear approximation class \mathbb{A}_s to be

$$\mathbb{A}_s := \left\{ (u, f, \mathbf{D}) \mid |u, f, \mathbf{D}|_s := \sup_{N > 0} (N^s \sigma(N; u, f, \mathbf{D})) < \infty \right\}.$$

An important pending issue is the characterization of \mathbb{A}_s . This is beyond the scope of this paper as well as unnecessary to examine optimality of AFEM. However, a few remarks are in order to clarify the nature of \mathbb{A}_s .

Remark 6.2 (Regularity H^{1+r}). If $u \in H^{1+r}(\Omega)$, which is always true for some $0 < r \leq 1$ in any dimension d [10, Theorem 3], then quasi-uniform refinement yields a decay rate $\|u - U_k\|_{\Omega} \leq C (\#\mathcal{T}_k)^{-r/d}$ for the energy error and $\text{osc}_k(U_k, \mathcal{T}_k) \leq C (\#\mathcal{T}_k)^{-1/d}$ for oscillation. In fact, applying (4.2) with $V = U_k$ and $W = 0$, we obtain

$$\begin{aligned} \text{osc}_k(U_k, T) &\leq h_T \|P_{2n-2}^2 f\|_T + \bar{\Lambda}_1 \text{osc}_k(\mathbf{D}, T) \|U_k\|_{H^1(\omega_T)} \\ &\leq C h_T (\|f\|_T + \|U_k\|_{H^1(\omega_T)}) \end{aligned}$$

by definition of oscillation $\text{osc}_k(\mathbf{D}, T)$ and regularity of \mathbf{A} . Thus $(u, f, \mathbf{D}) \in \mathbb{A}_{r/d}$ for any polynomial degree $n \geq 1$ and dimension d .

Remark 6.3 (Regularity W_1^2). Let $d = 2$, $n = 1$, and $u \in W_1^2(\Omega)$; this is true provided that in addition to the assumptions on data of (1.1) stated in §2.1 \mathbf{A} and Ω are Lipschitz [9, Theorem 5.2.2]. Assume now that a mesh \mathcal{T} equidistributes the quantity $\Lambda = \|D^2 u\|_{L^1(T)}$. Then

$$\|\nabla(u - I_T u)\|_T \preccurlyeq \|D^2 u\|_{L^1(\omega_T)} \approx \Lambda,$$

and

$$\|\nabla(u - I_T u)\|_{\Omega}^2 = \sum_{T \in \mathcal{T}} \|\nabla(u - I_T u)\|_T^2 \preccurlyeq \Lambda^2 \#\mathcal{T}.$$

Since $\|D^2 u\|_{L^1(\Omega)} \approx \Lambda \#\mathcal{T}$, we deduce the error decay

$$\|\nabla(u - I_T u)\|_{\Omega} \preccurlyeq \|D^2 u\|_{L^1(\Omega)} (\#\mathcal{T})^{-1/2}.$$

On the other hand, as in Remark 6.2 we have an oscillation decay $\text{osc}_{\mathcal{T}}(U, \mathcal{T}) \leq C (\#\mathcal{T})^{-1/2}$ for any quasi-uniform mesh \mathcal{T} . Given N DOFs, we let \mathcal{T}_1 be an optimal graded mesh for the energy error and \mathcal{T}_2 be a quasi-uniform mesh for oscillation, such that $\#\mathcal{T}_1 - \#\mathcal{T}_0, \#\mathcal{T}_2 - \#\mathcal{T}_0 \leq N$. According to Lemma 4.7, the overlay $\mathcal{T} = \mathcal{T}_1 \oplus \mathcal{T}_2$ is a mesh with $\#\mathcal{T} - \#\mathcal{T}_0 \leq 2N$ that yields a total error decay $N^{-1/2}$ and we infer that $(u, f, \mathbf{D}) \in \mathbb{A}_{1/2}$.

Remark 6.4 (Preasymptotics). Let $\mathbf{D} = (\mathbf{I}, 0)$, let f be the oscillating function with checkerboard pattern of Example 3.6 in [16], and let $n = 1$ be the polynomial degree. In this case, we have

$$\sigma(N; u, f, \mathbf{D}) := \inf_{\mathcal{T} \in \mathbb{T}_N} \inf_{V \in \mathbb{V}(\mathcal{T})} \left(\|u - V\|_{\Omega}^2 + \|h(f - \Pi_0^2 f)\|_{\Omega}^2 \right)^{1/2},$$

the discrete solution $U_k = 0$ and $\|u - U_k\|_{\Omega}$ is constant for as many steps $k \leq K$ as desired. In contrast $\eta_k(U_k, \mathcal{T}_k) = \|hf\|_{\Omega}$ reduces strictly for $k \leq K$ but overestimates the energy error because $\|u\|_{\Omega}$ can be made arbitrarily small by increasing

K . On the other hand, for $k > K$ we have $\|h(f - \Pi_0^2 f)\|_\Omega = 0$ and the total error asymptotics is dictated by the energy error alone; thus $(u, f, \mathbf{D}) \in \mathbb{A}_{1/2}$. The fact that the preasymptotic regime $k \leq K$ can be made arbitrarily long is crucial for adaptivity, but is not described by membership in \mathbb{A}_s .

In practice, this effect is typically less dramatic because f is not orthogonal to $\mathbb{V}(\mathcal{T}_k)$. Figure 6.1 displays the behavior of AFEM for the smooth solution u_S of §3.1 with frequencies 5, 10, 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 with optimal rate. The class \mathbb{A}_s misses to describe this behavior.

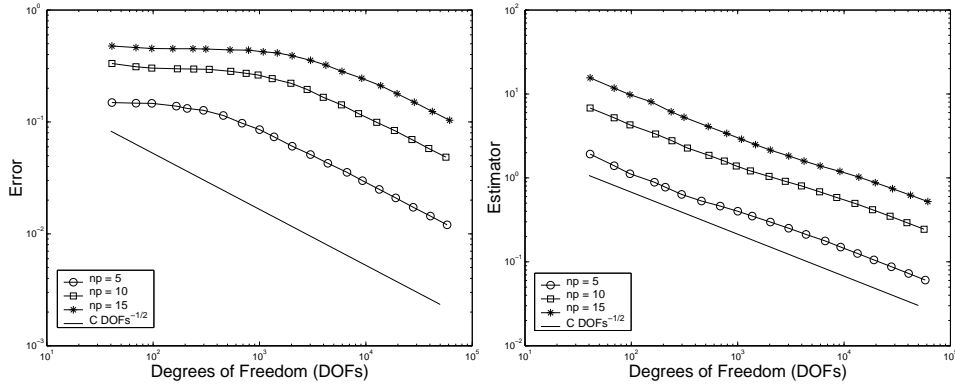


FIGURE 6.1. Decay of error (left) and estimator (right) for the smooth solution u_S of §3.1 with frequencies 5, 10, 15. The error exhibits a frequency-dependent plateau in the preasymptotic regime and later an optimal decay. This behavior is not described by \mathbb{A}_s .

Remark 6.5 (Conforming vs Non-conforming Meshes). In contrast to [3, 22], our approach relies exclusively on conforming triangulations. When $\mathbf{D} = (\mathbf{I}, 0)$ and $n = 1$, the approximation class \mathbb{A}_s is the same regardless of conformity [3]. The situation is quite different, however, in dealing with oscillation of jump residual unless the depth of nonconforming refinement is restricted beforehand.

We now assume that $(u, f, \mathbf{D}) \in \mathbb{A}_s$ for some $0 < s \leq n/d$, and prove that the approximation U_k generated by AFEM converges to u with the same rate $(\#\mathcal{T}_k - \#\mathcal{T}_0)^{-s}$ as the best approximation up to a multiplicative constant. We need to count elements marked by the estimator (the cardinality of \mathcal{M}_k) as well as those added to keep mesh conformity (see Lemma 2.3). To this end, we impose more stringent requirements than for convergence of AFEM.

Assumption 6.6 (Optimality). *We assume the following properties of AFEM:*

(a) *The marking parameter θ satisfies $\theta \in (0, \theta_*)$ with*

$$\theta_*^2 = \frac{C_2}{1 + C_1(1 + 2\Lambda_1 \text{osc}_{\mathcal{T}_0}^2(\mathbf{D}, \mathcal{T}_0))}.$$

(b) *Procedure MARK selects a set \mathcal{M}_k of marked elements with minimal cardinality.*

(c) *The distribution of refinement edges on \mathcal{T}_0 satisfies condition (b) of §4 in [23].*

The limit value θ_* depends of the ratio $\sqrt{C_2/C_1} \leq 1$, which quantifies the quality of estimator $\eta_{\mathcal{T}_k}(U_k, \mathcal{T}_k)$, as well as the oscillation $\text{osc}_{\mathcal{T}_0}(\mathbf{D}, \mathcal{T}_0)$ of coefficients of the PDE on \mathcal{T}_0 .

6.2. Cardinality of \mathcal{M}_k . The following Lemma establishes a link between non-linear approximation theory and AFEM through the Dörfler marking strategy. Roughly speaking we prove that if an approximation satisfies a suitable total error reduction from \mathcal{T} to $\mathcal{T}_* \geq \mathcal{T}$, then the error indicators of the coarser solutions must satisfy a Dörfler property on the set $\mathcal{R}_{\mathcal{T} \rightarrow \mathcal{T}_*}$. In other words, Dörfler marking and total error reduction are intimately connected.

Lemma 6.7 (Optimal Marking). *Assume that the marking parameter θ verifies (a) of Assumption 6.6. Let $\mathcal{T} \in \mathbb{T}$ and $U \in \mathbb{V}(\mathcal{T})$ be the discrete solution of (2.3). Set $\mu := \frac{1}{2}(1 - \frac{\theta^2}{\theta_*^2}) > 0$ and let $\mathcal{T}_* \in \mathbb{T}$ be any refinement of \mathcal{T} , i. e. $\mathcal{T} \leq \mathcal{T}_*$, such that the discrete solution $U_* \in \mathbb{V}(\mathcal{T}_*)$ satisfies*

$$(6.1) \quad \|u - U_*\|_{\Omega}^2 + \text{osc}_{\mathcal{T}_*}^2(U_*, \mathcal{T}_*) \leq \mu \left\{ \|u - U\|_{\Omega}^2 + \text{osc}_{\mathcal{T}}^2(U, \mathcal{T}) \right\}.$$

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

$$(6.2) \quad \eta_{\mathcal{T}}(U, \mathcal{R}) \geq \theta \eta_{\mathcal{T}}(U, \mathcal{T}).$$

Proof. We first combine the lower bound (2.7) with (6.1) to obtain

$$(6.3) \quad \begin{aligned} (1 - 2\mu) C_2 \eta_{\mathcal{T}}^2(U, \mathcal{T}) &\leq (1 - 2\mu) \left(\|u - U\|_{\Omega}^2 + \text{osc}_{\mathcal{T}}^2(U, \mathcal{T}) \right) \\ &\leq \|u - U\|_{\Omega}^2 - \|u - U_*\|_{\Omega}^2 + \text{osc}_{\mathcal{T}}^2(U, \mathcal{T}) - 2 \text{osc}_{\mathcal{T}_*}^2(U_*, \mathcal{T}_*). \end{aligned}$$

We estimate separately error and oscillation terms. We invoke the orthogonality (2.4) and the localized upper bound of Lemma 4.6 to arrive at

$$(6.4) \quad \|u - U\|_{\Omega}^2 - \|u - U_*\|_{\Omega}^2 = \|U_* - U\|_{\Omega}^2 \leq C_1 \eta_{\mathcal{T}}^2(U, \mathcal{R}).$$

For the oscillation terms we argue according to whether an element $T \in \mathcal{T}$ belongs to the set of refined elements \mathcal{R} or not. For $T \in \mathcal{R}$ we use the dominance

$$\text{osc}_{\mathcal{T}}^2(U, T) \leq \eta_{\mathcal{T}}^2(U, T)$$

of Remark 2.1. For $T \in \mathcal{T} \cap \mathcal{T}_*$, Corollary 4.5 with $V = U$ and $V_* = U_*$ yields

$$\text{osc}_{\mathcal{T}}^2(U, \mathcal{T} \cap \mathcal{T}_*) - 2 \text{osc}_{\mathcal{T}_*}^2(U_*, \mathcal{T} \cap \mathcal{T}_*) \leq 2\Lambda_1 \text{osc}_{\mathcal{T}_0}^2(\mathbf{D}, \mathcal{T} \cap \mathcal{T}_*) \|U_* - U\|_{\Omega}^2.$$

Combining these two estimates with (6.4) we infer that

$$(6.5) \quad \text{osc}_{\mathcal{T}}^2(U, \mathcal{T}) - 2 \text{osc}_{\mathcal{T}_*}^2(U_*, \mathcal{T}_*) \leq (1 + 2 C_1 \Lambda_1 \text{osc}_{\mathcal{T}_0}^2(\mathbf{D}, \mathcal{T}_0)) \eta_{\mathcal{T}}^2(U, \mathcal{R}).$$

Collecting (6.3), (6.4) and (6.5), we finally deduce

$$\eta_{\mathcal{T}}^2(U, \mathcal{R}) \geq \frac{(1 - 2\mu) C_2}{1 + C_1(1 + 2 \Lambda_1 \text{osc}_{\mathcal{T}_0}^2(\mathbf{D}, \mathcal{T}_0))} \eta_{\mathcal{T}}^2(U, \mathcal{T}) = \theta^2 \eta_{\mathcal{T}}^2(U, \mathcal{T}),$$

in light of the definitions of θ_* , θ and μ . This concludes the proof. \square

The key to relate the best mesh with AFEM triangulations is the fact that procedure MARK selects the marked set \mathcal{M}_k with *minimal* cardinality. This crucial idea, due to Stevenson [23], is used next.

Lemma 6.8 (Cardinality of \mathcal{M}_k). *Assume that the marking parameter θ verifies (a) and procedure MARK satisfies (b) of Assumption 6.6. Let u be the solution of (1.1), and let $\{\mathcal{T}_k, \mathbb{V}_k, U_k\}_{k \geq 0}$ be the sequence of meshes, finite element spaces, and discrete solutions produced by AFEM.*

If $(u, f, \mathbf{D}) \in \mathbb{A}_s$, then the following estimate is valid

$$\#\mathcal{M}_k \leq \left(1 - \frac{\theta^2}{\theta_*^2}\right)^{-\frac{1}{2s}} |u, f, \mathbf{D}|_s^{\frac{1}{s}} C_{\mathbf{D}}^{\frac{1}{2s}} \left\{ \|u - U_k\|_{\Omega}^2 + \text{osc}_{\mathcal{T}_k}^2(U_k, \mathcal{T}_k) \right\}^{-\frac{1}{2s}}.$$

Proof. We set $\epsilon^2 := \mu C_{\mathbf{D}}^{-1} (\|u - U_k\|_{\Omega}^2 + \text{osc}_{\mathcal{T}_k}^2(U_k, \mathcal{T}_k))$, where $\mu = \frac{1}{2}(1 - \frac{\theta^2}{\theta_*^2}) > 0$ and $C_{\mathbf{D}}$ is the constant in Lemma 6.1. Since $(u, f, \mathbf{D}) \in \mathbb{A}_s$, there exists a $\mathcal{T}_{\epsilon} \in \mathbb{T}$ and $V_{\epsilon} \in \mathbb{V}(\mathcal{T}_{\epsilon})$ such that

$$(6.6) \quad \#\mathcal{T}_{\epsilon} - \#\mathcal{T}_0 \leq |u, f, \mathbf{D}|_s^{1/s} \epsilon^{-1/s},$$

$$(6.7) \quad \|u - V_{\epsilon}\|_{\Omega}^2 + \text{osc}_{\mathcal{T}_{\epsilon}}^2(V_{\epsilon}, \mathcal{T}_{\epsilon}) \leq \epsilon^2.$$

Let $\mathcal{T}_* := \mathcal{T}_{\epsilon} \oplus \mathcal{T}_k$ be the overlay of \mathcal{T}_{ϵ} and \mathcal{T}_k , and let $U_* \in \mathbb{V}(\mathcal{T}_*)$ be the discrete solution of (2.3) on \mathcal{T}_* . To show that there is a reduction by a factor μ of the total error between U_* and U_k , since $\mathcal{T}_* \geq \mathcal{T}_{\epsilon}$, we argue as in Lemma 6.1 and obtain

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

Hence, we deduce from Lemma 6.7 that the subset $\mathcal{R} := \mathcal{R}_{\mathcal{T}_k \rightarrow \mathcal{T}_*} \subset \mathcal{T}_k$ verifies the Dörfler property (6.2) for $\theta < \theta_*$. The fact that procedure MARK selects a subset $\mathcal{M}_k \subset \mathcal{T}_k$ with minimal cardinality satisfying the same property translates into

$$(6.8) \quad \#\mathcal{M}_k \leq \#\mathcal{R} \leq \#\mathcal{T}_* - \#\mathcal{T}_k \leq \#\mathcal{T}_{\epsilon} - \#\mathcal{T}_0,$$

where Lemma 4.7 have been employed in the last step. Finally, combining (6.8), (6.6), and the definition of ϵ , we end up with

$$\#\mathcal{M}_k \leq \#\mathcal{T}_{\epsilon} - \#\mathcal{T}_0 \leq \mu^{-\frac{1}{2s}} |u, f, \mathbf{D}|_s^{\frac{1}{s}} C_{\mathbf{D}}^{\frac{1}{2s}} \left\{ \|u - U_k\|_{\Omega}^2 + \text{osc}_{\mathcal{T}_k}^2(U_k, \mathcal{T}_k) \right\}^{-\frac{1}{2s}},$$

which is the asserted estimate. \square

6.3. Quasi-Optimality. The following result is a consequence of the previous estimates and the fact that AFEM is a contraction for the quasi-error, namely the sum of energy error and scaled error estimator.

Theorem 6.9 (Quasi-Optimality). *Let Assumption 6.6 be satisfied by AFEM. Let u be the solution of (1.1), and let $\{\mathcal{T}_k, \mathbb{V}_k, U_k\}_{k \geq 0}$ be the sequence of meshes, finite element spaces, and discrete solutions produced by AFEM.*

Let $(u, f, \mathbf{D}) \in \mathbb{A}_s$ and $\Theta(\theta, s) := \theta^{-4s} \left(1 - \frac{\theta^2}{\theta_^2}\right)^{-1/2}$ describe the asymptotics of AFEM as $\theta \rightarrow \theta_*$, 0 or $s \rightarrow 0$. Then there exists a constant C , depending on data, the refinement depth b , and \mathcal{T}_0 , but independent of s , such that*

$$\left\{ \|u - U_k\|_{\Omega}^2 + \gamma \text{osc}_{\mathcal{T}_k}^2(U_k, \mathcal{T}_k) \right\}^{1/2} \leq C \Theta(\theta, s) |u, f, \mathbf{D}|_s (\#\mathcal{T}_k - \#\mathcal{T}_0)^{-s}.$$

Proof. Combining Lemmas 2.3 and 6.8 we deduce

$$(6.9) \quad \#\mathcal{T}_k - \#\mathcal{T}_0 \leq \sum_{j=0}^{k-1} \#\mathcal{M}_j \leq M \sum_{j=0}^{k-1} \left\{ \|u - U_j\|_{\Omega}^2 + \text{osc}_{\mathcal{T}_j}^2(U_j, \mathcal{T}_j) \right\}^{-\frac{1}{2s}},$$

with $M := (1 - \frac{\theta^2}{\theta_*^2})^{-\frac{1}{2s}} |u, f, \mathbf{D}|_s C_{\mathbf{D}}^{\frac{1}{2s}}$. We infer from the lower bound (2.7)

$$(6.10) \quad \begin{aligned} \|u - U_j\|_{\Omega}^2 + \gamma \operatorname{osc}_j^2(U_j, \mathcal{T}_j) &\leq \|u - U_j\|_{\Omega}^2 + \gamma \eta_j^2(U_j, \mathcal{T}_j) \\ &\leq \left(1 + \frac{\gamma}{C_2}\right) \left\{ \|u - U_j\|_{\Omega}^2 + \operatorname{osc}_j^2(U_j, \mathcal{T}_j) \right\}. \end{aligned}$$

On the other hand, the linear rate $\alpha = \alpha(\theta) < 1$ of convergence of Theorem 5.1 for the sum of energy error and scaled error estimator implies for $0 \leq j \leq k-1$

$$(6.11) \quad \|u - U_k\|_{\Omega}^2 + \gamma \eta_k^2(U_k, \mathcal{T}_k) \leq \alpha^{2(k-j)} \left\{ \|u - U_j\|_{\Omega}^2 + \gamma \eta_j^2(U_j, \mathcal{T}_j) \right\}.$$

We combine (6.9), (6.10) and (6.11) to obtain

$$\#\mathcal{T}_k - \#\mathcal{T}_0 \leq M \left(1 + \frac{\gamma}{C_2}\right)^{\frac{1}{2s}} \left\{ \|u - U_k\|_{\Omega}^2 + \gamma \eta_k^2(U_k, \mathcal{T}_k) \right\}^{-\frac{1}{2s}} \sum_{j=1}^k \alpha^{\frac{j}{s}}.$$

Since $\alpha < 1$ the geometric series is bounded by the constant $S_{\theta} = \alpha^{1/s} (1 - \alpha^{1/s})^{-1}$. Recalling that the element residual dominates the oscillation, we end up with

$$(6.12) \quad \#\mathcal{T}_k - \#\mathcal{T}_0 \leq S_{\theta} M \left(1 + \frac{\gamma}{C_2}\right)^{\frac{1}{2s}} \left\{ \|u - U_k\|_{\Omega}^2 + \gamma \operatorname{osc}_k^2(U_k, \mathcal{T}_k) \right\}^{-\frac{1}{2s}}.$$

To examine the asymptotics as $\theta, s \rightarrow 0$, hidden in S_{θ} , we use Remark 5.4 to get $\alpha \leq (1-t)^{\frac{1}{2}}$ with $t = c\theta^4$ and observe that $\lim_{t \rightarrow 0} \frac{t}{1-(1-t)^{\frac{1}{2s}}} = 2s$. Therefore

$$S_{\theta}^s \leq \frac{\alpha}{t^s} \left(\frac{t}{1-(1-t)^{1/2s}} \right)^s \approx \alpha s^s \theta^{-4s},$$

whence $S_{\theta}^s M^s \approx \Theta(\theta, s) |u, f, \mathbf{D}|_s C_{\mathbf{D}}^{1/2}$. Raising (6.12) to the s -th power and re-ordering, we finally obtain the desired estimate with $C > 0$ independent of s . \square

REFERENCES

- [1] M. AINSWORTH AND J. T. ODEN, *A Posteriori Error Estimation in Finite Element Analysis*, John Wiley & Sons, Inc., 2000.
- [2] E. BÄNSCH, *Local mesh refinement in 2 and 3 dimensions*, IMPACT Comput. Sci. Engrg. 3 (1991), 181–191.
- [3] P. BINEV, W. DAHMEN, AND R. DEVORE, *Adaptive finite element methods with convergence rate*, Numer. Math. 97 (2004), 219–268.
- [4] J. M. CASCÓN, C. KREUZER, R. H. NOCHETTO, AND K. G. SIEBERT, *Quasi-optimal AFEM for general elliptic operators*, in preparation.
- [5] Z. CHEN, J. FENG, *An adaptive finite element algorithm with reliable and efficient error control for linear parabolic problems*, Math. Comp. 73 (2004), 1167–1193.
- [6] L. DIENING AND C. KREUZER, *Linear Convergence of an adaptive finite element method for the p -Laplacian equation*, Preprint No. 03/2007, Institut für Mathematik, Universität Freiburg.
- [7] W. DÖRFLER, *A convergent adaptive algorithm for Poisson's equation*, SIAM J. Numer. Anal. 33 (1996), 1106–1124.
- [8] D. GILBARG AND N. S. TRUDINGER, *Elliptic Partial Differential Equations of Second Order*, Springer-Verlag, Germany, 1983.
- [9] P. GRISVARD, *Elliptic Problems in Nonsmooth Domains*, Pitman, Boston 1985.
- [10] F. JOCHMANN, *An H^s regularity result for the gradient of solutions to elliptic equations with mixed boundary conditions*, J. Math. Anal. Appl. 238, (1999), 429–450.
- [11] R. B. KELLOGG, *On the Poisson equation with intersecting interfaces*, Appl. Anal. 4 (1975), 101–129.
- [12] I. KOSSACZKÝ, *A recursive approach to local mesh refinement in two and three dimensions*, J. Comput. Appl. Math. 55 (1994), 275–288.

- [13] J. MAUBACH, *Local bisection refinement for n -simplicial grids generated by reflection*, SIAM J. Sci. Comput. 16, (1995), 210–227.
- [14] W. F. MITCHELL, *A comparison of adaptive refinement techniques for elliptic problems*, ACM Trans. Math. Softw. 15 (1989), 326–347.
- [15] 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), 1043–1068.
- [16] P. MORIN, R. H. NOCHETTO, AND K. G. SIEBERT, *Data oscillation and convergence of adaptive FEM*, SIAM J. Numer. Anal. 38 (2000), 466–488.
- [17] P. MORIN, R. H. NOCHETTO, AND K. G. SIEBERT, *Convergence of adaptive finite element methods*, SIAM Rev. 44, (2002), 631–658.
- [18] P. MORIN, R. H. NOCHETTO, AND K. G. SIEBERT, *Local problems on stars: a posteriori error estimation, convergence, and performance*, Math. Comp. 72 (2003), 1067–1097.
- [19] P. MORIN, K.G. SIEBERT, A. VEESER, *A basic convergence result for conforming adaptive finite elements*, preprint no. 1/2007, Dipartimento di Matematica “F. Enriques”, Via C. Saldini 50, 20133 Milano, Italy.
- [20] A. SCHMIDT AND K. G. SIEBERT. *Design of Adaptive Finite Element Software: The Finite Element Toolbox ALBERTA*, LNCSE 42, Springer-Verlag Berlin Heidelberg 2005.
- [21] L. R. SCOTT AND S. ZHANG, *Finite element interpolation of nonsmooth functions satisfying boundary conditions*, Math. Comp. 54 (1990), 483–493.
- [22] R. STEVENSON, *Optimality of a standard adaptive finite element method*, Found. Comput. Math. Online. DOI 10.1007/s10208-005-0183-0, (2006).
- [23] R. STEVENSON, *The completion of locally refined simplicial partitions created by bisection*, Preprint No. 1336, 13 pages, Department of Mathematics, Utrecht University, September 2005. Corrected version, October 2006.
- [24] C. T. TRAXLER, *An algorithm for adaptive mesh refinement in n dimensions*, Computing 59, (1997), 115–137.
- [25] A. VEESER, *Convergent adaptive finite elements for the nonlinear Laplacian*, Numer. Math. 92, (2002), 743–770.
- [26] R. VERFÜRTH, *A Review of A Posteriori Error Estimation and Adaptive Mesh-Refinement Technique*, Wiley-Teubner, Chichester, 1996.

J. MANUEL CASCON, DEPARTAMENTO DE MATEMÁTICAS, UNIVERSIDAD DE SALAMANCA, 37008, SALAMANCA, SPAIN.

Partially supported by Grants CGL2004-06171-C03-03/CLI, Ministerio de Ciencia y Tecnología (Spain) and SA078A05, Junta de Castilla y León (Spain), FEDER funds (European Union), and NSF Grant DMS-0505454.

URL: <http://matematicas.fis.usal.es/~casbar>

E-mail address: casbar@usal.es

CHRISTIAN KREUZER, INSTITUT FÜR MATHEMATIK, UNIVERSITÄT AUGSBURG, UNIVERSITÄTS-STRASSE 14, 86159 AUGSBURG, GERMANY.

Partially supported by DAAD Grant “Efficient Finite Element Methods for Solid and Fluid Mechanics Computations”

URL: http://www.math.uni-augsburg.de/de/prof/lam/mitarbeiter/christian_kreuzer

E-mail address: kreuzer@math.uni-augsburg.de

RICARDO H. NOCHETTO, DEPARTMENT OF MATHEMATICS AND INSTITUTE OF PHYSICAL SCIENCE AND TECHNOLOGY, UNIVERSITY OF MARYLAND, COLLEGE PARK, MD 20742, USA.

Partially supported by NSF Grant DMS-0505454.

URL: <http://www.math.umd.edu/~rhn>

E-mail address: rhn@math.umd.edu

KUNIBERT G. SIEBERT, INSTITUT FÜR MATHEMATIK, UNIVERSITÄT AUGSBURG, UNIVERSITÄTS-STRASSE 14, 86159 AUGSBURG, GERMANY.

Partially supported by DAAD Grant “Efficient Finite Element Methods for Solid and Fluid Mechanics Computations”

URL: <http://scicomp.math.uni-augsburg.de/siebert/>

E-mail address: siebert@math.uni-augsburg.de