

## ZIENKIEWICZ–ZHU ERROR ESTIMATORS ON ANISOTROPIC TETRAHEDRAL AND TRIANGULAR FINITE ELEMENT MESHES

GERD KUNERT<sup>1</sup> AND SERGE NICAISE<sup>2</sup>

**Abstract.** We consider *a posteriori* error estimators that can be applied to *anisotropic* tetrahedral finite element meshes, *i.e.* meshes where the aspect ratio of the elements can be arbitrarily large. Two kinds of Zienkiewicz–Zhu (ZZ) type error estimators are derived which originate from different backgrounds. In the course of the analysis, the first estimator turns out to be a special case of the second one, and both estimators can be expressed using some recovered gradient. The advantage of keeping two different analyses of the estimators is that they allow different and partially novel investigations and results. Both rigorous analytical approaches yield the equivalence of each ZZ error estimator to a known residual error estimator. Thus reliability and efficiency of the ZZ error estimation is obtained. The anisotropic discretizations require analytical tools beyond the standard isotropic methods. Particular attention is paid to the requirements on the anisotropic mesh. The analysis is complemented and confirmed by extensive numerical examples. They show that good results can be obtained for a large class of problems, demonstrated exemplary for the Poisson problem and a singularly perturbed reaction diffusion problem.

**Mathematics Subject Classification.** 65N15, 65N30, 65N50.

Received: April 17, 2002.

### 1. INTRODUCTION

Several classes of boundary value problems intrinsically give rise to solutions that exhibit lower dimensional, *anisotropic* behaviour. Such anisotropic solutions show little variation in certain space directions but much variation otherwise. For example, singularly perturbed problems often result in solutions with boundary layers. Even the solution of the Poisson problem in three space dimensions is generically anisotropic along some concave edge, see the numerical experiments of Section 5. Within the finite element method, such anisotropic solutions can be favorably resolved with *anisotropic meshes*. By this we understand meshes with stretched elements which are characterized by an unbounded aspect ratio, *i.e.* the ratio of the diameters of the circumscribed and inscribed sphere can be arbitrarily large.

---

*Keywords and phrases.* Anisotropic mesh, error estimator, Zienkiewicz–Zhu estimator, recovered gradient.

<sup>1</sup> Fakultät für Mathematik, TU Chemnitz, 09107 Chemnitz, Germany. e-mail: gkunert@mathematik.tu-chemnitz.de

<sup>2</sup> Université de Valenciennes et du Hainaut Cambrésis, MACS, B.P. 311, 59304 Valenciennes Cedex, France.  
e-mail: snicaise@univ-valenciennes.fr

Our emphasis is on *error estimators* which form a basis of any adaptive solution algorithm. The theory of error estimation is nowadays well established for conventional, isotropic meshes (*i.e.* where the aspect ratio of the elements is bounded). The books [1, 24] provide a comprehensive overview of the state of the art.

On *anisotropic meshes* the error estimation theory is much less developed. Recently the intensive research has led to several estimators that can be applied to different boundary value problems as well as norms, see [10, 12–15, 17, 22]. Exemplary we mention residual error estimators and local problem error estimators for the Poisson problem or a singularly perturbed problem; the error can be estimated in the energy norm or in the  $L_2$ -norm.

There is one popular estimator for isotropic meshes that did not have yet a counterpart for anisotropic meshes. This so-called Zienkiewicz–Zhu (ZZ) estimator has been invented in [26] and later been improved in [27]; many more variants have been developed since. The basic idea consists in computing an improvement of the gradient of the numerical solution by some post-processing procedure. The difference between this so-called *recovered gradient* and the original gradient serves as error estimator. This idea of ZZ error estimation has been very appealing and popular in the finite element community since

- the estimator is comparatively cheap because a recovered gradient is often computed anyway;
- the estimator is astonishingly robust (in numerical experiments) for a wide range of problems, see *e.g.* [3, 4].

Our work here is devoted to the extension of the ZZ estimator to *anisotropic meshes*. We start with a recapitulation of the existing isotropic analyses and discuss their suitability for anisotropic meshes. The theoretical approaches to ZZ error estimators (on isotropic meshes) can be divided roughly into three classes:

- proving equivalence to residual error estimators;
- utilizing superconvergence properties;
- minimization approach.

Each of these approaches will now be discussed briefly.

*Equivalence to residual error estimator:* Here the ZZ error estimator is proven to be equivalent to a residual error estimator, thus transferring reliability and efficiency to the first estimator. This approach goes back to [20] and is repeated in [24, Sect. 1.5]. In our paper these ideas will be generalized to *anisotropic meshes*. Of course several modifications and extensions are necessary:

- Although some recovered gradient is still applied, it is now scaled with weights that depend on the stretching directions (*i.e.* the alignment) of the anisotropic elements.
- The anisotropic meshes have to meet additional requirements which are due to the anisotropy. These requirements roughly mean that the anisotropic meshes should not be totally unstructured but instead obey some “sensible” geometrical principles. These demands also seem reasonable in the light of superconvergence properties discussed below.

*Superconvergence approach:* It forms the basis of most proofs for ZZ estimators. Exemplary we refer to [1] and the citations therein. In suitable, specialized settings even asymptotic exactness of the (global) ZZ estimator can be shown. This requires:

- consistence, localization, boundedness and linearity of the recovery operator;
- and a superconvergence property of the finite element scheme.

Unfortunately the superconvergence approach inherits two drawbacks. Firstly, the theoretical analysis requires very specialized meshes which are rarely found in practice (*e.g.* in adaptive refinement procedures). Secondly local equivalences cannot be proven.

The application of such a superconvergence analysis to *anisotropic meshes* is unclear up to now. Superconvergence results are not known for general meshes but only for special Shishkin (type) meshes, see [21, 25]. For example, the authors of [21] prove a certain kind of superconvergence for 2D Shishkin-type meshes consisting of axiparallel rectangles, bilinear finite elements, and a singularly perturbed reaction-convection-diffusion problem in the unit square. Most likely the results can be employed to define a ZZ estimator, even if this is not presented in the aforementioned work.

Summarizing, we do not pursue the superconvergence approach because of the high demands on the meshes which are hardly consistent with anisotropic solutions.

*Minimization approach:* A third kind of analysis utilizes a close relation between the ZZ estimator and a minimum formulation, cf. [5,8]. It allows to investigate general averaging operators which define the estimator, and it avoids superconvergence assumptions. The resulting error bounds involve so-called “higher order terms” that contain the unknown solution. Hence these bounds can only be interpreted in an asymptotic sense. Moreover the constants in the reliability result depend on the shape of the finite elements.

After presenting different techniques to analyse ZZ estimators, we will consider from now on exclusively the first approach, namely the equivalence to a residual error estimator. As it has been explained, this analysis seems most promising for anisotropic meshes.

The outline of this paper is as follows. The model problem, some notation as well as the assumptions on the mesh are introduced in Section 2. In Section 3 we first recall a known residual error estimator that is required for the subsequent analysis. Afterwards two kinds of ZZ error estimators are derived and rigorously analysed. The first estimator is based on a global projection property which corresponds to a *particular choice* of the underlying recovered gradient. The second ZZ estimator is an improvement because the recovered gradient can be defined with *arbitrary weights*. Our novel analysis additionally yields *local elementwise estimates*. Although the first estimator is a special case of the second one, we present them both because of the different background and the different and partially new analytical approaches and results. Section 4 is devoted to a detailed examination of the mesh assumptions. These investigations indicate that the anisotropic discretization is the main technicality when deriving the ZZ estimator. The numerical examples of Section 5 complement and confirm the theoretical analysis. There we also consider a (more realistic) reaction-diffusion problem (for which boundary layers appear naturally), state first theoretical results and present some numerical tests for that model problem.

## 2. MODEL PROBLEM AND NOTATION

### 2.1. Model problem

We consider a Poisson model problem with homogeneous Dirichlet boundary conditions in a polyhedral domain  $\Omega \subset \mathbb{R}^d$ ,  $d = 2, 3$ :

$$-\Delta u = f \quad \text{in } \Omega, \quad u = 0 \quad \text{on } \Gamma_D := \partial\Omega. \tag{2.1}$$

Our analysis is presented for three dimensional domains ( $d = 3$ ); the application to two dimensional domains ( $d = 2$ ) is readily possible. The corresponding variational formulation reads

$$\text{Find } u \in H_o^1(\Omega) : \quad \int_{\Omega} \nabla u \nabla v = \int_{\Omega} f v \quad \forall v \in H_o^1(\Omega), \tag{2.2}$$

where  $H_o^1(\Omega)$  denotes the usual Sobolev space of functions of  $H^1(\Omega)$  that vanish on  $\Gamma_D$ . For  $f \in L_2(\Omega)$  problem (2.2) admits a unique solution.

In order to discretize (2.2), let  $\mathcal{F} = \{\mathcal{T}_h\}$  be a family of triangulations  $\mathcal{T}_h$  of  $\Omega$ . We assume a conforming triangulation (cf. [9, Chap. 2]) that consists of tetrahedra ( $d = 3$ ) or triangles ( $d = 2$ ). Let  $V_h \subset H_o^1(\Omega)$  be the finite element space of piecewise affine linear functions on  $\mathcal{T}_h$  that vanish on  $\Gamma_D$ . The finite element solution  $u_h$  is uniquely obtained *via*

$$\text{Find } u_h \in V_h : \quad \int_{\Omega} \nabla u_h \nabla v_h = \int_{\Omega} f v_h \quad \forall v_h \in V_h. \tag{2.3}$$

### 2.2. Notation

The following paragraphs now introduce most of the notation required. For some domain  $\omega \subset \mathbb{R}^2$  or  $\omega \subset \mathbb{R}^3$  let  $\|\cdot\|_{\omega} := \|\cdot\|_{L_2(\omega)}$  be the usual  $L_2(\omega)$  norm. The space of polynomials of order at most  $k$  is denoted by  $\mathbb{P}^k(\omega)$ . For some (column) vectors  $\underline{v}, \underline{w}$  let  $(\underline{v}, \underline{w})$  be the Euclidean scalar product and  $|\underline{v}| := (\underline{v}, \underline{v})^{1/2}$  be the Euclidean

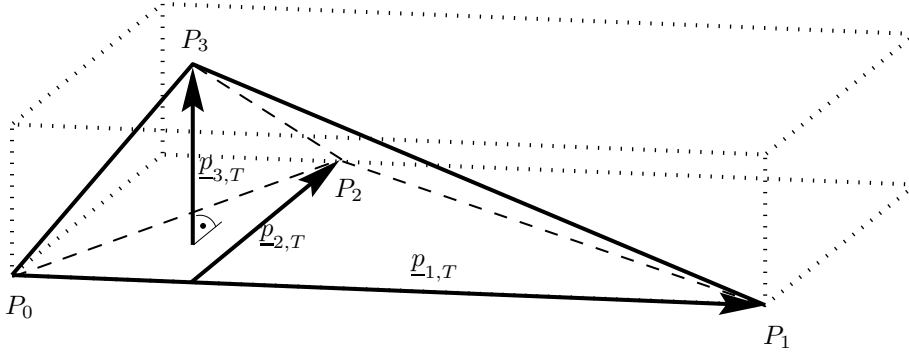


FIGURE 1. Notation of tetrahedron  $T$ .

length. Instead of  $x \leq c \cdot y$  or  $c_1 y \leq x \leq c_2 y$  (with positive constants independent of  $x, y$  or  $\mathcal{T}_h$ ) we use the shorthand notation  $x \lesssim y$  or  $x \sim y$ , respectively.

The next paragraph presents notation that is related to the triangulation  $\mathcal{T}_h$  and its elements. Tetrahedra are denoted by  $T, T'$  or  $T_i$ , faces are denoted by  $E$ , and nodes of  $\mathcal{T}_h$  are denoted by  $x$ . Next, define nodal sets  $\mathcal{N}_T, \mathcal{N}_E, \mathcal{N}_{\bar{\Omega}}$  that contain all nodes of a tetrahedron  $T$ , a face  $E$ , or of  $\bar{\Omega}$  (*i.e.* including boundary nodes), respectively. Let  $\mathcal{E}_{\bar{\Omega}}$  be the set of all interior edges (2D) or faces (3D) of  $\bar{\Omega}$ . For a node  $x$  we introduce a local neighbourhood patch  $\omega_x := \bigcup_{T: x \in \mathcal{N}_T} T \subset \mathbb{R}^3$  which is the union of all tetrahedra having this node. Similarly for some face  $E$  let  $\omega_E \subset \mathbb{R}^3$  be the union of both tetrahedra having this face (with obvious boundary modifications). For a tetrahedron  $T$ , a face  $E$  or a patch  $\omega_x$  set  $|T| = \text{meas}_3(T)$ ,  $|E| = \text{meas}_2(E)$  or  $|\omega_x| = \text{meas}_3(\omega_x)$ , respectively (the distinction from the Euclidean vector length is obvious from the context).

The four vertices of an arbitrary but fixed tetrahedron  $T \in \mathcal{T}_h$  are temporarily denoted by  $P_0, \dots, P_3$  such that  $P_0P_1$  is the longest edge of  $T$ ,  $\text{meas}_2(\triangle P_0P_1P_2) \geq \text{meas}_2(\triangle P_0P_1P_3)$ , and  $\text{meas}_1(P_1P_2) \geq \text{meas}_1(P_0P_2)$ . Additionally define three pairwise orthogonal vectors  $\underline{p}_{i,T}$  with lengths  $h_{i,T} := |\underline{p}_{i,T}|$ , see Figure 1. Observe  $h_{1,T} > h_{2,T} \geq h_{3,T}$  and set  $h_{\min,T} := \min_{i=1\dots 3} h_{i,T} = h_{3,T}$ . The matrix  $C_T \in \mathbb{R}^{3 \times 3}$  is defined as

$$C_T := \left( \underline{p}_{1,T}, \underline{p}_{2,T}, \underline{p}_{3,T} \right),$$

and describes (roughly speaking) the anisotropic orientations of the tetrahedron  $T$ .

For a face  $E$  of a tetrahedron  $T$  let  $h_{E,T} := 3|T|/|E|$  be the length of the *height* over  $E$  in  $T$ . Note that  $h_{E,T}$  is *not* the diameter of  $E$ , as in the usual convention.

The quantities  $h_{\min,T}$  and  $h_{E,T}$  are associated with a *tetrahedron*  $T$ . Often it is more convenient to utilize equivalent data that are related to a *face*  $E$  or *node*  $x$ . To this end define averaged terms by

$$\begin{aligned} h_E &:= (h_{E,T_1} + h_{E,T_2})/2 && \text{for } E = T_1 \cap T_2 \\ h_{\min,E} &:= (h_{\min,T_1} + h_{\min,T_2})/2 && \text{for } E = T_1 \cap T_2 \\ h_{i,x} &:= \frac{1}{n} \sum_{T \subset \omega_x} h_{i,T} && h_{\min,x} := \frac{1}{n} \sum_{T \subset \omega_x} h_{\min,T}, \end{aligned}$$

where  $n$  is the number of elements  $T$  in  $\omega_x$ . Note that  $h_{\min,E}$  is not the minimal dimension of the face  $E$ . For boundary faces  $E \subset \partial\Omega$  modify  $h_E := h_{E,T}$  and  $h_{\min,E} := h_{\min,T}$ , where  $\partial T \supset E$ .

Next consider an arbitrary interior face  $E$ . Let  $\underline{n}_E$  be any of the two unit normal vectors for  $E$ , and keep it fixed from here on. For a piecewise continuous (scalar or vector valued) function  $v$  denote by  $[[v]]_E$  the jump of  $v$  across  $E$  in the direction  $\underline{n}_E$ . Let  $\partial_{n_E} v := \underline{n}_E \cdot \nabla v$  be the (unitary) directional derivative. Note that the orientation of  $\underline{n}_E$  influences terms like  $[[v]]_E$  but not  $[[\partial_{n_E} v]]_E$ .

### 2.3. Mesh requirements

In addition to the usual conformity conditions of the mesh (see Ciarlet [9, Chap. 2]) we demand the following assumptions. They are explained in more detail in Remark 2.1 below and in Section 4.

- (A1) The number of tetrahedra containing a node  $x$  is bounded uniformly.  
 (A2) The dimensions of adjacent tetrahedra must not change rapidly, *i.e.*

$$h_{i,T'} \sim h_{i,T} \quad \forall T, T' \quad \text{with} \quad T \cap T' \neq \emptyset, i = 1 \dots d.$$

- (A3) For each node  $x$  there exists a matrix  $C_x \in \mathbb{R}^{d \times d}$  such that

$$|C_x^{-1} \underline{v}| \sim |C_T^{-1} \underline{v}| \quad \forall \underline{v} \in \mathbb{R}^d, \forall T \subset \omega_x.$$

- (A4) An assumption on the shape of each element:

$$|C_T^{-1} \underline{n}_E| \sim h_{E,T}^{-1} \quad \forall E \subset \partial T.$$

- (A5) The  $L_2$  projection is stable in the sense of [17, Sect. 4]. For self-containment we repeat the definition given there. Start with two (distinct) elements  $T_1, T_2 \in \mathcal{T}_h$  and define their (topological) edge distance by  $l(T_1, T_2) := 1 +$  minimal edge number of all edge paths connecting  $T_1$  and  $T_2$ .

Set  $l(T, T) := 0$ . Note that in both the 2D and 3D case the *edges* count. Next, for a given element  $T$  introduce neighbourhood (ring) patches by

$$R_k(T) := \{T' : l(T', T) = k\}, \quad k \in \mathbb{N}.$$

Then assumption (A5) is satisfied if there exist positive constants  $c_1, c_2, \alpha, \beta, r$  such that

$$\left\{ \begin{array}{l} h_{\min, T_1} / h_{\min, T_2} \leq c_1 \cdot \alpha^{l(T_1, T_2)} \quad \forall T_1, T_2 \in \mathcal{T}_h \\ \text{card}(R_k(T)) \leq c_2 \cdot k^r \beta^k \quad \forall T \in \mathcal{T}_h, \forall k \in \mathbb{N}_+ \\ \alpha \cdot \beta < \begin{cases} \sqrt{2} + \sqrt{3} \approx 3.146 & \text{if } d = 2 \\ (3 + \sqrt{5})/2 \approx 2.618 & \text{if } d = 3. \end{cases} \end{array} \right. \quad (2.4)$$

The mesh assumptions (A1) and (A2) imply several convenient equivalences.

$$\left\{ \begin{array}{l} \text{(A2)} \Rightarrow \quad h_{i,x} \sim h_{i,T} \quad \forall T \subset \omega_x \\ \text{(A2)} \Rightarrow \quad h_{\min,x} \sim h_{\min,T} \quad \forall T \subset \omega_x \\ \text{(A2)} \Rightarrow \quad h_{\min,x} \sim h_{\min,E} \quad \forall E : x \in \mathcal{N}_E \\ \text{(A2)} \Rightarrow \quad h_E \sim h_{E,T} \quad \forall E \subset \partial T \\ \text{(A1)+(A2)} \Rightarrow \quad |T| \sim |\omega_x| \quad \forall T \subset \omega_x \end{array} \right. \quad (2.5)$$

Furthermore, with the help of (A2) and (A3) we can rewrite assumption (A4) as

$$|C_x^{-1} \underline{n}_E| \sim h_E^{-1} \quad \forall E : x \in \mathcal{N}_E. \quad (2.6)$$

**Remark 2.1.** The mesh assumptions are scrutinized in detail in Section 4. Here some remarks may facilitate the understanding.

Assumption (A3) roughly means that there exists a transformation  $C_x^{-1}$  which maps the patch  $\omega_x$  onto an isotropic patch of size  $\mathcal{O}(1)$ .

Assumption (A4) roughly demands for an anisotropic tetrahedron that small faces are almost perpendicular to long edges, depending on the aspect ratio.

Finally the stability assumption (A5) is only a sufficient condition to derive the residual error estimation. Recent research [6,7,23] suggests that the restrictions of (A5) can be weakened; some results of the aforementioned work already apply to our setting here.

## 2.4. Matching function

*Reliability* and *efficiency* are highly desirable properties in a *posteriori* error estimation. They basically mean that the error  $\|u - u_h\|_*$  (in some suitable norm) can be bounded from above and below, respectively, with constants independent of  $u$ ,  $u_h$  or  $\mathcal{T}_h$ .

Most standard error estimators on *isotropic* finite element meshes are reliable and efficient at the same time, cf. [1,24]. Unfortunately the situation is much less obvious on *anisotropic* meshes. The analysis as well as numerical experiments strongly suggest that reliability and efficiency cannot be achieved simultaneously on *arbitrary* anisotropic meshes. However if the anisotropy of the solution is sufficiently well aligned with the anisotropy of the mesh then one can expect both properties at the same time. Intuitively all applications of anisotropic finite elements follow this concept: an element should be stretched in that direction where the function (or more precisely, its derivative) exhibits little variation.

In order to measure the alignment of an anisotropic mesh  $\mathcal{T}_h$  with an anisotropic function  $v$ , a so-called *matching function* has been proposed by Kunert [11,12].

**Definition 2.2** (Matching function). Let  $v \in H^1(\Omega)$ , and  $\mathcal{T}_h \in \mathcal{F}$  be a triangulation of  $\Omega$ . Define the *matching function*  $m_1 : H^1(\Omega) \times \mathcal{F} \mapsto \mathbb{R}$  by

$$m_1(v, \mathcal{T}_h) := \left( \sum_{T \in \mathcal{T}_h} h_{\min, T}^{-2} \cdot \|C_T^\top \nabla v\|_T^2 \right)^{1/2} / \|\nabla v\|_\Omega. \quad (2.7)$$

The vital importance of the matching function for anisotropic error estimation can be seen in the error bounds (3.2, 3.3) and (3.7, 3.8) below.

The matching function is not central to our analysis here. Hence we refer to [12] for a comprehensive discussion, and restrict ourselves to a brief explanation of basic features. Firstly the definition immediately implies  $m_1(v, \mathcal{T}_h) \geq 1$ . On *isotropic* meshes one obtains easily  $m_1(v, \mathcal{T}_h) \sim 1$ ; then the matching function merges with other constants and becomes invisible. In contrast to this more care is necessary for *anisotropic* meshes. If the mesh is suitably aligned with the anisotropic solution one still achieves  $m_1 \sim 1$  and thus reliable and efficient error estimation. If however the anisotropic mesh is not aligned with the solution then  $m_1$  can be arbitrarily large (cf. [13, Numerical experiment 2] or [11, Rem. 3.3]). Hence upper and lower error bounds may differ by an arbitrarily large factor; thus the error estimator is useless for error control and adaptive refinement.

## 3. ERROR ESTIMATORS

We start by presenting the residual error estimator of [17] which forms the basis of the subsequent analysis. Afterwards two kinds of ZZ error estimators are presented whose *element-based* definitions are analogous. The main difference of both estimators is the transformation to *nodal-based* terms which are required for the analysis. The first ZZ estimator follows the lines of [20] (also described in [24] (Sect. 1.5)). It is based on a recovered gradient  $\nabla^{R_1}$  which satisfies a global projection property. In contrast, the second ZZ estimator utilizes a *new* technique to derive *nodal-based* terms. Consequently a much more flexible recovered gradient  $\nabla^{R_2}$  can be employed to define this ZZ estimator. Accordingly a novel analysis is required (cf. Lem. 3.12 and Th. 3.13) which is based on different techniques than for the first estimator.

Note that all estimators are given in several forms. The first representation is the one used in practice, and is related either to a face  $E$  or an element  $T$ . The other, equivalent representation is related to a node  $x$ , and is required for analytical purposes.

3.1. Residual error estimator

In [17] a face residual based error estimator is introduced for interior faces by

$$\eta_{R,E} := h_{\min,E} h_E^{-1/2} \cdot \| [\![\partial_{n_E} u_h]\!] \|_E, \quad E \in \mathcal{E}_\Omega. \tag{3.1}$$

The corresponding lower and upper error bounds are given in [17, Th. 5.1]. Provided that mesh assumptions (A1), (A2) and (A5) are satisfied, one has

$$\eta_{R,E} \lesssim \|\nabla(u - u_h)\|_{\omega_E} + \inf_{f_h \in V_h} h_{\min,E} \|f - f_h\|_{\omega_E} \quad \forall E \in \mathcal{E}_\Omega \tag{3.2}$$

$$\|\nabla(u - u_h)\|_\Omega \lesssim m_1(u - u_h, \mathcal{T}_h) \cdot \left( \sum_{E \in \mathcal{E}_\Omega} \eta_{R,E}^2 + \inf_{f_h \in V_h} \sum_{T \in \mathcal{T}_h} h_{\min,T}^2 \|f - f_h\|_T^2 \right)^{1/2}. \tag{3.3}$$

Clearly  $\eta_{R,E}$  is associated with a face  $E$ . For our purposes, however, node related quantities are much better suited. Therefore we fix a patch  $\omega_x$  and combine all its (interior) faces. The first expression below introduces the local, node related error estimator. The second definition introduces the global error estimator whereas the remaining definition facilitates our exposition later on.

**Definition 3.1** (Residual error estimators). The local and global residual error estimators are given by

$$\eta_{R,x}^2 := h_{\min,x}^2 |\omega_x| \sum_{E: x \in \mathcal{N}_E} h_E^{-2} [\![\partial_{n_E} u_h]\!]_E^2 \tag{3.4}$$

$$\eta_R^2 := \sum_{x \in \mathcal{N}_\Omega} \eta_{R,x}^2 \tag{3.5}$$

$$\eta_{\tilde{R},x}^2 := h_{\min,x}^2 |\omega_x| \sum_{E: x \in \mathcal{N}_E} |C_x^{-1} [\![\nabla u_h]\!]_E|^2. \tag{3.6}$$

**Lemma 3.2.** *Let the mesh assumptions (A1), (A2) be satisfied. Then*

$$\eta_{R,x}^2 \sim \sum_{E: x \in \mathcal{N}_E} \eta_{R,E}^2.$$

*Proof.* The assertion follows immediately from the fact that the dimensions of neighbouring elements must not change rapidly, cf. (2.5). □

The error estimation by means of the node related error estimator  $\eta_{R,x}$  can now be derived easily.

**Lemma 3.3** (Residual error estimation). *Assume that the mesh assumptions (A1), (A2) and (A5) are satisfied. The error is bounded locally from below and globally from above.*

$$\eta_{R,x} \lesssim \|\nabla(u - u_h)\|_{\omega_x} + \inf_{f_h \in V_h} h_{\min,x} \|f - f_h\|_{\omega_x} \quad \forall x \in \mathcal{N}_\Omega \tag{3.7}$$

$$\|\nabla(u - u_h)\|_\Omega \lesssim m_1(u - u_h, \mathcal{T}_h) \left( \eta_R^2 + \inf_{f_h \in V_h} \sum_{T \in \mathcal{T}_h} h_{\min,T}^2 \|f - f_h\|_T^2 \right)^{1/2}. \tag{3.8}$$

*Proof.* The inequalities follow immediately from (3.2, 3.3) and Lemma 3.2. □

The next lemma presents a sufficient condition for the equivalence of  $\eta_{R,x}$  and  $\eta_{\tilde{R},x}$ . This lemma will be essential for further analysis.

**Lemma 3.4.** *Let the mesh assumption (A2)–(A4) be satisfied, then it holds:*

$$\eta_{R,x} \sim \eta_{\tilde{R},x}. \tag{3.9}$$

*Proof.* Consider an arbitrary face  $E$ ,  $x \in \mathcal{N}_E$ , and any one of its two unit normal vectors  $\underline{n}_E$ . Then there exist two further unit vectors  $\underline{\tau}_1, \underline{\tau}_2$  such that  $(\underline{n}_E, \underline{\tau}_1, \underline{\tau}_2)$  forms an orthonormal vector system. Hence

$$\begin{aligned} & \underline{n}_E \underline{n}_E^\top + \underline{\tau}_1 \underline{\tau}_1^\top + \underline{\tau}_2 \underline{\tau}_2^\top = I_{3 \times 3} \\ \text{giving} \quad & \underline{n}_E \cdot \partial_{n_E} u_h + \underline{\tau}_1 \cdot \partial_{\underline{\tau}_1} u_h + \underline{\tau}_2 \cdot \partial_{\underline{\tau}_2} u_h = \nabla u_h. \end{aligned}$$

Both terms  $\partial_{\underline{\tau}_i} u_h$  are continuous across  $E$ ; only  $\partial_{n_E} u_h$  jumps. Thus we conclude

$$[[\partial_{n_E} u_h]]_E \cdot \underline{n}_E = [[\nabla u_h]]_E.$$

Together with assumptions (A2)–(A4) which imply (2.6) one obtains

$$\begin{aligned} \sum_{E:x \in \mathcal{N}_E} |C_x^{-1} [[\nabla u_h]]_E|^2 &= \sum_{E:x \in \mathcal{N}_E} |C_x^{-1} [[\partial_{n_E} u_h]]_E \cdot \underline{n}_E|^2 \\ &= \sum_{E:x \in \mathcal{N}_E} |C_x^{-1} \underline{n}_E|^2 \cdot [[\partial_{n_E} u_h]]_E^2 \stackrel{(2.6)}{\sim} \sum_{E:x \in \mathcal{N}_E} h_E^{-2} \cdot [[\partial_{n_E} u_h]]_E^2 \end{aligned}$$

which proves the assertion. □

### 3.2. First ZZ error estimator

Let us first define the recovered gradient  $\nabla^{R_1}$  by means of a projection with respect to a particular scalar product. For a precise definition of this inner product, let  $W_h$  be the space of piecewise linear vector fields on the triangulation, and set  $V_h := W_h \cap C(\Omega, \mathbb{R}^d)$ , cf. also [24]. In order to shorten the notation we temporarily introduce the matrices

$$B_x := h_{\min,x} C_x^{-1} \quad \text{and} \quad B_T := h_{\min,T} C_T^{-1}.$$

On  $W_h$ , we introduce the mesh dependent inner product  $(\cdot, \cdot)_h$  by

$$(\underline{v}, \underline{w})_h := \sum_{T \in \mathcal{T}_h} |T| \sum_{x \in \mathcal{N}_T} (B_x \underline{v}|_T(x), B_x \underline{w}|_T(x)), \tag{3.10}$$

where  $\underline{v}|_T(x) = \lim_{y \rightarrow x, y \in T} \underline{v}(y)$ .

From mesh assumptions (A2), (A3) we have concluded (2.5), i.e.  $h_{\min,x} \sim h_{\min,T}$ ,  $|C_x^{-1} \underline{v}| \sim |C_T^{-1} \underline{v}|$  for all  $T \subset \omega_x$  and thus also  $|B_x \underline{v}| \sim |B_T \underline{v}|$  for all  $T \subset \omega_x$  and all vectors  $\underline{v} \in \mathbb{R}^d$ . For an arbitrary but fixed tetrahedron  $T$  and for any piecewise linear function  $\underline{v} \in W_h$  we can further conclude

$$|T| \sum_{x \in \mathcal{N}_T} |B_x \underline{v}|_T(x)|^2 \stackrel{(A3)}{\sim} |T| \sum_{x \in \mathcal{N}_T} |B_T \underline{v}|_T(x)|^2 \sim \|B_T \underline{v}\|_T^2.$$

Therefore the mesh assumptions (A2), (A3) imply

$$(\underline{v}, \underline{v})_h \sim \sum_{T \in \mathcal{T}_h} \|B_T \underline{v}\|_T^2. \tag{3.11}$$

This last result also shows that  $(\cdot, \cdot)_h$  is a scalar product indeed since all  $B_T$  are regular matrices. Now the recovered gradient can be defined.



**Definition 3.5** (First recovered gradient). The *recovered gradient*  $\nabla^{R_1} : W_h \rightarrow V_h$  is defined as the projection of  $\nabla u_h$  onto  $V_h$  with respect to the inner product  $(\cdot, \cdot)_h$ , i.e.  $\nabla^{R_1} u_h \in V_h$  is uniquely determined by the condition

$$(\nabla^{R_1} u_h - \nabla u_h, \underline{v}_h)_h = 0 \quad \forall \underline{v}_h \in V_h. \tag{3.12}$$

The recovered gradient  $\nabla^{R_1} u_h$  is piecewise linear and continuous. Its nodal values can be computed locally and coincide with the usual recovered gradient as presented, for example, in [24, equality (1.80)]. Details are given in the next lemma.

**Lemma 3.6.** *The value of the recovered gradient at a node  $x$  can be determined locally by*

$$(\nabla^{R_1} u_h)(x) = \sum_{T \subset \omega_x} \mu_T \nabla u_h|_T \quad \text{with weight} \quad \mu_T := \frac{|T|}{|\omega_x|} \in \mathbb{R}, T \subset \omega_x. \tag{3.13}$$

*Proof.* The proof utilizes standard ideas as presented e.g. in [24]. Fix the node  $x$  and apply the definition of the recovered gradient with  $\underline{v}_h := \varphi_x \cdot \underline{e}_i$ , where  $\varphi_x$  is the standard (piecewise linear) basis function of  $V_h$  for node  $x$ , and  $\underline{e}_i \in \mathbb{R}^d$  is the  $i$ th unit vector. Then

$$\begin{aligned} 0 &= (\nabla^{R_1} u_h - \nabla u_h, \varphi_x \cdot \underline{e}_i)_h \\ &= \sum_{T \in \mathcal{T}_h} |T| \sum_{x' \in \mathcal{N}_T} (B_{x'}^\top B_{x'} (\nabla^{R_1} u_h(x') - \nabla u_h|_T(x')), \varphi_x|_T(x') \underline{e}_i) \\ &= \sum_{T \subset \omega_x} |T| \cdot (B_x^\top B_x (\nabla^{R_1} u_h(x) - \nabla u_h|_T(x)), \underline{e}_i) \end{aligned}$$

holds for  $i = 1 \dots d$ . Furthermore  $B_x^\top B_x$  is regular, and hence

$$\underline{0} = \sum_{T \subset \omega_x} |T| \cdot (\nabla^{R_1} u_h(x) - \nabla u_h|_T(x)) = |\omega_x| \nabla^{R_1} u_h(x) - \sum_{T \subset \omega_x} |T| \cdot \nabla u_h|_T(x)$$

which proves the assertion.

Note that the choice of the regular matrix  $B_x$  in the definition of the scalar product  $(\cdot, \cdot)_h$  has no influence on the nodal value of the recovered gradient. □

Now we are ready to define our anisotropic version of the first ZZ estimator. Again, the first two terms are given in a form that can be used in practice. The third quantity is a node related term which can be utilized in further analysis.

**Definition 3.7** (First anisotropic ZZ estimators). The local and global ZZ estimators are given by

$$\eta_{Z_1, T} := h_{\min, T} \|C_T^{-1} (\nabla^{R_1} u_h - \nabla u_h)\|_T \tag{3.14}$$

$$\eta_{Z_1}^2 := \sum_{T \in \mathcal{T}_h} \eta_{Z_1, T}^2 \tag{3.15}$$

$$\eta_{Z_1, x}^2 := h_{\min, x}^2 |\omega_x| \left( \sum_{T \subset \omega_x} \frac{|T|}{|\omega_x|} |C_x^{-1} \nabla u_h|_T|^2 - \left| \sum_{T \subset \omega_x} \frac{|T|}{|\omega_x|} C_x^{-1} \nabla u_h|_T \right|^2 \right). \tag{3.16}$$

Similar to the residual error estimator we first establish a relation between the global estimator  $\eta_{Z_1}$  and the node related quantity  $\eta_{Z_1, x}$ . To achieve this, assume that mesh assumptions (A2) and (A3) hold which

imply (3.11). Furthermore utilize the projection property (3.12), recall the definition of the matrices  $B_x, B_T$  and of the scalar product to obtain

$$\begin{aligned}
 \eta_{Z_1}^2 &= \sum_{T \in \mathcal{T}_h} h_{\min, T}^2 \|C_T^{-1} (\nabla^{R_1} u_h - \nabla u_h)\|_T^2 \\
 &\stackrel{(3.11)}{\sim} (\nabla^{R_1} u_h - \nabla u_h, \nabla^{R_1} u_h - \nabla u_h)_h \\
 &\stackrel{(3.12)}{=} (\nabla u_h, \nabla u_h)_h - (\nabla^{R_1} u_h, \nabla^{R_1} u_h)_h \\
 &\stackrel{(3.10)}{=} \sum_{T \in \mathcal{T}_h} |T| \sum_{x \in \mathcal{N}_T} h_{\min, x}^2 \left( |C_x^{-1} \nabla u_h|_T|^2 - |C_x^{-1} \nabla^{R_1} u_h(x)|^2 \right).
 \end{aligned}$$

Insert now the nodal value of  $\nabla^{R_1} u_h$  and change the summation order from  $\sum_{T \in \mathcal{T}_h} \sum_{x \in \mathcal{N}_T}$  to  $\sum_{x \in \mathcal{N}_\Omega} \sum_{T \subset \omega_x}$  to conclude

$$\begin{aligned}
 \eta_{Z_1}^2 &\sim \sum_{x \in \mathcal{N}_\Omega} h_{\min, x}^2 \sum_{T \subset \omega_x} |T| \cdot \left( |C_x^{-1} \nabla u_h|_T|^2 - |C_x^{-1} \nabla^{R_1} u_h(x)|^2 \right) \\
 &\stackrel{(3.13)}{=} \sum_{x \in \mathcal{N}_\Omega} h_{\min, x}^2 |\omega_x| \left( \sum_{T \subset \omega_x} \frac{|T|}{|\omega_x|} |C_x^{-1} \nabla u_h|_T|^2 - \left| \sum_{T \subset \omega_x} \frac{|T|}{|\omega_x|} C_x^{-1} \nabla u_h|_T \right|^2 \right).
 \end{aligned}$$

Hence the following relation between the global estimator  $\eta_{Z_1}$  and the node related estimator  $\eta_{Z_1, x}$  is obtained provided that the mesh assumptions (A2) and (A3) hold:

$$\eta_{Z_1}^2 \sim \sum_{x \in \mathcal{N}_\Omega} \eta_{Z_1, x}^2. \tag{3.17}$$

Let us start the analysis of the estimator with a general equivalence lemma which is already known from isotropic investigations.

**Lemma 3.8.** *Let mesh assumption (A1) be satisfied, and consider an arbitrary node  $x$  and the associated patch  $\omega_x$ . Let  $\underline{v}$  be a (scalar or vector valued) function defined on  $\omega_x$  such that  $\underline{v}|_T \in \mathbb{P}^0(T)$ , i.e.  $\underline{v}$  is piecewise constant. Let further  $\mu_T, T \subset \omega_x$  be arbitrary positive weights such that all  $\mu_T$  are uniformly bounded away from zero,  $\mu_T \geq c > 0$ , and that satisfy  $\sum_{T \subset \omega_x} \mu_T = 1$ . Define the ZZ averaged value by  $\underline{v}_{ZZ} := \sum_{T \subset \omega_x} \mu_T \underline{v}|_T$ . Then*

$$\sum_{E: x \in \mathcal{N}_E} |[[\underline{v}]]_E|^2 \sim \sum_{T \subset \omega_x} \mu_T |\underline{v}|_T|^2 - |\underline{v}_{ZZ}|^2. \tag{3.18}$$

*Proof.* For two dimensional domains ( $d = 2$ ) this lemma has been proven in [20]; the proof is also repeated in [24] (Sect. 1.5). An extension to three dimensional domains ( $d = 3$ ) is readily possible with the ideas from the proof of Lemma 3.12. □

The main result follows now.

**Theorem 3.9** (Equivalences with first ZZ estimator). *Let the mesh assumptions (A1)–(A4) be satisfied. Then the residual error estimator and the first ZZ error estimator are equivalent:*

$$\eta_{R, x} \sim \eta_{Z_1, x} \tag{3.19}$$

$$\eta_R \sim \eta_{Z_1}. \tag{3.20}$$

*Proof.* We apply the previous lemma 3.8 with  $\underline{v} := C_x^{-1} \nabla u_h$  and  $\mu_T = |T|/|\omega_x|$  as well as lemma 3.4 to derive

$$\begin{aligned} \eta_{R,x}^2 &\stackrel{(3.9)}{\sim} \eta_{R,x}^2 = h_{\min,x}^2 |\omega_x| \sum_{E:x \in \mathcal{N}_E} |C_x^{-1} [[\nabla u_h]]_E|^2 \\ &\stackrel{(3.18)}{\sim} h_{\min,x}^2 |\omega_x| \sum_{T \subset \omega_x} \frac{|T|}{|\omega_x|} |C_x^{-1} \nabla u_h|_T|^2 - \left| \sum_{T \subset \omega_x} \frac{|T|}{|\omega_x|} C_x^{-1} \nabla u_h|_T \right|^2 \\ &\stackrel{(3.16)}{=} \eta_{Z_1,x}^2. \end{aligned}$$

This yields (3.19); the equivalence (3.20) follows thanks to (3.17). □

Note that this is only an equivalence between the *global* estimators. An equivalence involving the *local* estimator  $\eta_{Z_1,T}$  cannot be proven in this way since the projection property (3.12) is given globally. The procedure of the second ZZ error estimator avoids this drawback.

### 3.3. Second ZZ error estimator

A different approach to describe a ZZ error estimator is given now. It avoids the global projection property (3.12) at the cost of a refined analysis. As a consequence local elementwise relations can be derived. We start with the definition of an arbitrary recovered gradient.

**Definition 3.10** (Arbitrary recovered gradient).

The *arbitrary recovered gradient*  $\nabla^{R_2} : W_h \rightarrow V_h$  is defined by the nodal values

$$(\nabla^{R_2} u_h)(x) := \sum_{T \subset \omega_x} \mu_{T,x} \nabla u_h|_T \tag{3.21}$$

where the weights  $\mu_{T,x} \geq 0$  can be chosen arbitrarily such that  $\sum_{T \subset \omega_x} \mu_{T,x} = 1$ .

The corresponding second ZZ estimator is given next. Again the first two definitions describe the local (element related) estimator and its global counterpart. The third term is a node related quantity required for the subsequent analysis.

**Definition 3.11** (Second anisotropic ZZ estimator). The local and global ZZ estimators are given by

$$\eta_{Z_2,T} := h_{\min,T} \|C_T^{-1} (\nabla^{R_2} u_h - \nabla u_h)\|_T \tag{3.22}$$

$$\eta_{Z_2}^2 := \sum_{T \in \mathcal{T}_h} \eta_{Z_2,T}^2 \tag{3.23}$$

$$\eta_{Z_2,x}^2 := h_{\min,x}^2 |\omega_x| \sum_{T \subset \omega_x} |C_x^{-1} (\nabla^{R_2} u_h(x) - \nabla u_h|_T(x))|^2. \tag{3.24}$$

In order to establish a relation between the node related term  $\eta_{Z_2,x}$  and the element related estimator  $\eta_{Z_2,T}$ , recall that  $\nabla^{R_2} u_h - \nabla u_h$  is linear on  $T$ . Together with mesh assumptions (A1)–(A3) we conclude

$$\begin{aligned} \eta_{Z_2,T}^2 &\sim h_{\min,T}^2 |T| \sum_{x \in \mathcal{N}_T} |C_T^{-1} (\nabla^{R_2} u_h(x) - \nabla u_h|_T(x))|^2 \\ &\stackrel{(2.5),(A3)}{\sim} \sum_{x \in \mathcal{N}_T} h_{\min,x}^2 |\omega_x| \cdot |C_x^{-1} (\nabla^{R_2} u_h(x) - \nabla u_h|_T(x))|^2. \end{aligned}$$

Note that equivalences (2.5), ASSc have been applied to switch from element related data  $h_{\min,T}, C_T^{-1}$  to node related data  $h_{\min,x}, C_x^{-1}$ . This yields immediately the desired inequalities

$$\eta_{Z_2,x}^2 \lesssim \sum_{T \subset \omega_x} \eta_{Z_2,T}^2 \tag{3.25}$$

$$\eta_{Z_2,T}^2 \lesssim \sum_{x \in \mathcal{N}_T} \eta_{Z_2,x}^2 \tag{3.26}$$

provided that the mesh assumptions (A1)–(A3) are satisfied. Note that the sums on the right-hand side of (3.25, 3.26) are necessary because  $\eta_{Z_2,x}$  depends on  $u_h|_{\omega_x}$  whereas  $\eta_{Z_2,T}$  depends on  $u_h$  on  $\bigcup_{x \in \mathcal{N}_T} \omega_x$ .

The next lemma states a novel equivalence that is similar to the one of Lemma 3.8. The main difference is that now the weights  $\mu_T$  do not have to be bounded away from 0. The technique to prove this lemma seems to be partially new: the transformation to a matrix eigenvalue problem is standard (at least in 2D, cf. [24, Sect. 1.5]) while the subsequent eigenvalue analysis is novel.

**Lemma 3.12.** *Let mesh assumption (A1) be satisfied, and consider an arbitrary node  $x$  and the associated patch  $\omega_x$ . Let  $\underline{v}$  be a (scalar or vector valued) function defined on  $\omega_x$  such that  $\underline{v}|_T \in \mathbb{P}^0(T)$ , i.e.  $\underline{v}$  is piecewise constant. Let further  $\mu_T, T \subset \omega_x$ , be arbitrary non-negative weights such that  $\sum_{T \subset \omega_x} \mu_T = 1$ . Define  $\underline{v}_{ZZ}$  as in*

*Lemma 3.8. Then*

$$\sum_{E:x \in \mathcal{N}_E} |[\underline{v}]_E|^2 \sim \sum_{T \subset \omega_x} |\underline{v}_{ZZ} - \underline{v}|_T|^2. \tag{3.27}$$

*Proof.* Note first that it suffices to prove (3.27) component wise, i.e. assume that  $\underline{v} \equiv v$  is a scalar, piecewise constant function on  $\omega_x$ . For simplicity of notation denote the elements of  $\omega_x$  temporarily by  $T_1 \dots T_n$ . Accordingly set  $\mu_i := \mu_{T_i}$  and  $v^i := v|_{T_i}$ . The mesh assumption (A1) states that  $n$  is bounded uniformly on  $\mathcal{T}_h$ .

We start the proof for an interior node  $x$  and follow [24, Sect. 1.5]. Consider first the left hand side of (3.27) which now reads

$$\sum_{E:x \in \mathcal{N}_E} |[\underline{v}]_E|^2 = \sum_{\substack{i,j \\ x \in \mathcal{N}_E, E=T_i \cap T_j}} |v^i - v^j|^2$$

i.e. we sum over all elements  $T_i$  and  $T_j$  that share a common face  $E$  (in 3D) or a common edge (in 2D). The last sum can be written in matrix notation as

$$0 \leq \sum_{\substack{i,j \\ x \in \mathcal{N}_E, E=T_i \cap T_j}} |v^i - v^j|^2 = (A\underline{w}, \underline{w}) \tag{3.28}$$

with  $\underline{w} := (v^1, v^2, \dots, v^n)^\top$  and

$$A = (a_{i,j})_{i,j=1}^n \in \mathbb{R}^{n \times n}, \quad a_{i,j} = \begin{cases} d & \text{if } i = j \\ -1 & \text{if } T_i \text{ and } T_j \text{ share a common face (3D) or edge (2D)} \\ 0 & \text{otherwise.} \end{cases}$$

Obviously  $A = A^\top$  is positively semidefinite and weakly diagonally dominant. From (3.28) we further conclude that  $A$  has exactly one eigenvalue 0 corresponding to the eigenvector  $\underline{w} = \underline{1} := (1, 1, \dots, 1)^\top \in \mathbb{R}^n$ ; all other eigenvalues are positive. The matrix  $A$  depends solely on the topology of the patch  $\omega_x$  but not on its geometry. Since the number of such topologies is finite ( $n$  is bounded because of mesh assumption (A1)), there is only a finite number of possibilities for the corresponding matrices  $A$ . Hence all positive eigenvalues of  $A$  are bounded from above and below (and away from 0). Note that in 2D the matrix  $A$  simplifies to a circulant tridiagonal matrix consisting of  $(-1, 2, -1)$ .

Consider next the right hand side of (3.27) which can be rewritten as

$$\sum_{T \subset \omega_x} \left| \underline{v}_{ZZ} - \underline{v}_T \right|^2 = \sum_{i=1}^n \left| \left( \sum_{j=1}^n \mu_j v^j \right) - v^i \right|^2 = (B \underline{w}, \underline{w}),$$

with  $B = (b_{i,j})_{i,j=1}^n \in \mathbb{R}^{n \times n}$  ,  $b_{i,j} = \delta_{ij} + n\mu_i \mu_j - \mu_i - \mu_j$ .

Introducing  $\underline{\mu} := (\mu_1, \dots, \mu_n)^\top \in \mathbb{R}^n$  one derives

$$B = I + n\underline{\mu} \underline{\mu}^\top - \underline{\mu} \underline{1}^\top - \underline{1} \underline{\mu}^\top$$

$$B - I = \underline{\nu} \underline{\mu}^\top + \underline{\mu} \underline{\nu}^\top \quad \text{with} \quad \underline{\nu} := \left( \frac{n}{2} \underline{\mu} - \underline{1} \right).$$

Since  $B - I$  is symmetric, it has a full system of eigenvectors. Because  $B - I$  is of rank 2, it has  $n - 2$  eigenvalues 0. For every other eigenvalue  $\lambda$  of  $B - I$  the corresponding eigenvector is a linear combination of  $\underline{\mu}$  and  $\underline{\nu}$ . A simple calculation reveals that then  $\lambda$  is also an eigenvalue of the matrix

$$\begin{bmatrix} \underline{\mu}^\top \underline{\nu} & \underline{\mu}^\top \underline{\mu} \\ \underline{\nu}^\top \underline{\nu} & \underline{\nu}^\top \underline{\mu} \end{bmatrix} = \begin{bmatrix} \frac{n}{2} \underline{\mu}^\top \underline{\mu} - 1 & \underline{\mu}^\top \underline{\mu} \\ \frac{n^2}{4} \underline{\mu}^\top \underline{\mu} & \frac{n}{2} \underline{\mu}^\top \underline{\mu} - 1 \end{bmatrix} \in \mathbb{R}^{2 \times 2},$$

*i.e.*  $\lambda_1 = -1$  and  $\lambda_2 = n \underline{\mu}^\top \underline{\mu} - 1$ . Hence the eigenvalues of  $B$  are

$$\lambda(B) = \begin{cases} 0 & \text{single eigenvalue} \\ n \underline{\mu}^\top \underline{\mu} & \text{single eigenvalue} \\ 1, \dots, 1 & n - 2 \text{ times.} \end{cases}$$

The arithmetic quadratic mean inequality gives

$$1 \leq n \underline{\mu}^\top \underline{\mu} = n \sum_{i=1}^n \mu_i^2 \leq n,$$

hence all positive eigenvalues of  $B$  lie in the range  $[1, n]$ . The eigenvalue 0 is associated with the eigenvector  $\underline{1}$ .

Summarizing,  $A$  and  $B$  both have a single eigenvalue 0 corresponding to the same eigenvector  $\underline{1}$ . All other eigenvalues are positive and bounded from above and below. This implies

$$A \sim B \quad \text{and} \quad (A \underline{w}, \underline{w}) \sim (B \underline{w}, \underline{w}) \quad \forall \underline{w} \in \mathbb{R}^n$$

which proves the lemma for an interior node  $x$ .

For a boundary node  $x$  we can proceed in almost the same way. The only difference consists in a slight modification of the matrix  $A$ , namely,  $a_{i,i} = d - k$  where  $k$  is the number of boundary faces (of the element  $T_i$ ) that contain the node  $x$ . The properties of  $A$  and the remainder of the proof stay exactly the same as before.  $\square$

Now we are able to prove equivalences with the second ZZ estimator (involving the arbitrary recovered gradient  $\nabla^{R_2}$ ).

**Theorem 3.13** (Equivalences with second ZZ estimator). *Let the mesh assumptions (A1)–(A4) be satisfied. Then the following local and global relations hold (for all  $x \in \mathcal{N}_{\tilde{\Omega}}$  or  $T \in \mathcal{T}_h$ ).*

$$\eta_{R,x} \sim \eta_{Z_2,x} \tag{3.29}$$

$$\eta_R \sim \eta_{Z_2} \tag{3.30}$$

$$\eta_{R,x}^2 \lesssim \sum_{T \subset \omega_x} \eta_{Z_2,T}^2 \tag{3.31}$$

$$\eta_{Z_2,T}^2 \lesssim \sum_{x \in \mathcal{N}_T} \eta_{R,x}^2. \tag{3.32}$$

*Proof.* To prove (3.29), fix an arbitrary node  $x \in \mathcal{N}_{\tilde{\Omega}}$  and consider  $\underline{v} := C_x^{-1} \nabla u_h$  on the patch  $\omega_x$ . Then the ZZ averaged value at the node  $x$  becomes

$$\underline{v}_{ZZ} = \sum_{T \subset \omega_x} \mu_{T,x} C_x^{-1} \nabla u_h|_T = C_x^{-1} \nabla^{R_2} u_h(x).$$

Since  $\underline{v}$  is piecewise constant on  $\omega_x$ , Lemma 3.12 can be applied. In conjunction with lemma 3.4 this yields

$$\begin{aligned} \eta_{R,x}^2 &\stackrel{(3.9)}{\sim} \eta_{\tilde{R},x}^2 = h_{\min,x}^2 |\omega_x| \sum_{E:x \in \mathcal{N}_E} |[[C_x^{-1} \nabla u_h]]_E|^2 \\ &\stackrel{(3.27)}{\sim} h_{\min,x}^2 |\omega_x| \sum_{T \subset \omega_x} |C_x^{-1} \nabla^{R_2} u_h(x) - C_x^{-1} \nabla u_h|_T|^2 = \eta_{Z_2,x}^2. \end{aligned}$$

Next, (3.31) is a direct consequence of (3.29) and (3.25). The converse relation (3.32) can be concluded similarly from (3.26) and (3.29).

Finally the global equivalence (3.30) can be proven *via* (3.31, 3.32).

$$\begin{aligned} \eta_R^2 &\stackrel{(3.31)}{\lesssim} \sum_{x \in \mathcal{N}_{\tilde{\Omega}}} \sum_{T \subset \omega_x} \eta_{Z_2,T}^2 = (d+1) \sum_{T \in \mathcal{T}_h} \eta_{Z_2,T}^2 = (d+1) \eta_{Z_2}^2 \\ &\stackrel{(3.32)}{\lesssim} \sum_{T \in \mathcal{T}_h} \sum_{x \in \mathcal{N}_T} \eta_{R,x}^2 \lesssim \sum_{x \in \mathcal{N}_{\tilde{\Omega}}} \eta_{R,x}^2 = \eta_R^2. \end{aligned} \tag{3.32} \quad \square$$

Note that the sums in (3.31, 3.32) appear because  $\eta_{R,x}$  is a node related term whereas  $\eta_{Z_2,T}$  is an element related quantity.

**Theorem 3.14** (ZZ error estimation). *Assume that mesh assumptions (A1)–(A5) are satisfied. Then the error is bounded locally from below and globally from above.*

$$\eta_{Z_2,x} \lesssim \|\nabla(u - u_h)\|_{\omega_x} + \inf_{f_h \in V_h} h_{\min,x} \|f - f_h\|_{\omega_x} \quad \forall x \in \mathcal{N}_{\tilde{\Omega}} \tag{3.33}$$

$$\|\nabla(u - u_h)\|_{\Omega} \lesssim m_1(u - u_h, \mathcal{T}_h) \left( \eta_{Z_2}^2 + \inf_{f_h \in V_h} \sum_{T \in \mathcal{T}_h} h_{\min,T}^2 \|f - f_h\|_T^2 \right)^{1/2}. \tag{3.34}$$

*Proof.* These are immediate consequences of Lemma 3.3 and Theorem 3.13. □

**Corollary 3.15** (ZZ error estimation on isotropic meshes). *Assume that an isotropic mesh satisfies mesh assumption (A5). Then the ZZ error estimator  $\eta_{Z_2,x}$  is reliable and efficient.*

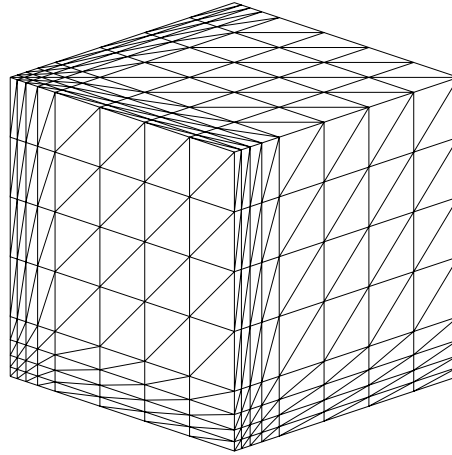


FIGURE 2. Anisotropic tensor product type mesh.

This holds even when the corresponding recovered gradient  $\nabla^{R_2}$  is defined with arbitrary weights (non-negative with sum 1).

This result seems to be new even for isotropic meshes (at least the authors have not found a proof anywhere else). So far special weights had to be chosen for the recovered gradient in order to prove equivalence with the residual error estimator and, in turn, reliability and local efficiency of the ZZ error estimator, cf. [24, Sect. 1.5]. Now there is the freedom to choose arbitrary weights.

Note that reliability alone for an arbitrary recovered gradient has been shown in [8]. Global efficiency (up to higher order terms) is obtained in the sequel [5].

#### 4. THE MESH ASSUMPTIONS REVISITED

As we have seen, the analysis of the ZZ error estimators required several mesh assumptions that were introduced in Section 2.3. These assumptions are now discussed in more detail.

In Section 4.1 it is shown that there exist meshes which satisfy all assumptions. Sections 4.2 and 4.3 are devoted to mesh assumption (A3) while Section 4.4 investigates mesh assumption (A4). With that help we can prove in Section 4.5 that the mesh assumptions are satisfied for another class of meshes. In Section 4.6 the role of the mesh assumptions is examined by showing that assumption (A4) is a necessary condition for error estimation.

##### 4.1. Rectangular tensor product type meshes satisfy the mesh assumptions

In this section we prove that the mesh assumptions (A1)–(A5) can be satisfied. To this end we consider rectangular tensor product type tetrahedral meshes. By this we understand that the tetrahedra of  $\mathcal{T}_h$  can be grouped such that a set of six of them forms a rectangular hexahedron, cf. also Figure 2. Assumption (A1) then clearly holds.

At this stage of generality, assumption (A2) obviously can be satisfied, so we assume that it holds. It states that the dimension of the tetrahedra (in each of the three anisotropic directions) must not change rapidly across neighbouring elements. This assumption is quite weak. It allows, for example, meshes that resolve boundary layers, see e.g. Figure 2.

For rectangular tensor product type meshes we now prove that (A1) and (A2) imply (A3)–(A5). Let us start with (A3), i.e. we construct a matrix  $C_x$  and show the corresponding properties. Our exposition describes the 3D case; the 2D analogies are straightforward.

*Proof of assumption (A3).* Start with a node  $x$  of  $\mathcal{T}_h$  and an arbitrary tetrahedron  $T \subset \omega_x$ . Since we consider tensor product type meshes there exists a circumscribing rectangular brick (*i.e.* hexahedron)  $B \supset T$ . The three edge lengths of this brick  $B$  are denoted by  $h_{1,B} \geq h_{2,B} \geq h_{3,B}$ . Choose corresponding edge vectors  $\underline{p}_{i,B}$ ,  $i = 1, 2, 3$ , *i.e.* such that  $|\underline{p}_{i,B}| = h_{i,B}$ . The orientation of these orthogonal vectors does not matter. Define next the matrix  $C_B \in \mathbb{R}^{3 \times 3}$  whose columns are formed by the vectors  $\underline{p}_{i,B}$ ,

$$C_B := \left( \underline{p}_{1,B}, \underline{p}_{2,B}, \underline{p}_{3,B} \right),$$

as well as three orthogonal vectors

$$\underline{p}_{i,x} := \frac{h_{i,x}}{h_{i,B}} \cdot \underline{p}_{i,B} = h_{i,x} \cdot \frac{\underline{p}_{i,B}}{|\underline{p}_{i,B}|} \quad i = 1, 2, 3$$

which are oriented along the edges vectors  $\underline{p}_{i,B}$  of  $B$  but which have a different length  $|\underline{p}_{i,x}| = h_{i,x}$ . Define the matrix  $C_x \in \mathbb{R}^{3 \times 3}$  by

$$C_x := \left( \underline{p}_{1,x}, \underline{p}_{2,x}, \underline{p}_{3,x} \right).$$

This immediately implies  $C_x^\top C_x = \text{diag}(h_{1,x}^2, h_{2,x}^2, h_{3,x}^2)$ . Furthermore the geometric properties as well as relations (2.5) yield the equivalences

$$h_{i,T} \sim h_{i,B} \sim h_{i,x} \quad i = 1, 2, 3. \tag{4.1}$$

Now we are ready to prove assumption (A3). Let us start with investigations of the linear transformations associated with  $C_B$  and  $C_T^{-1}$ . Recall first that  $\underline{e}_i \in \mathbb{R}^3$  is the  $i$ th unit vector. Because of  $C_B \underline{e}_i = \underline{p}_{i,B} \in \mathbb{R}^3$ , the transformation *via*  $C_B$  maps the unit cube  $[0, 1]^3$  onto the brick  $B$  (or more precisely onto the corresponding brick at the origin of the coordinate system). Since the four vertices of  $T \subset B$  are also vertices of  $B$ , the transformation *via*  $C_B$  thus maps  $\tilde{T} \rightarrow T$ , where  $\tilde{T} \subset [0, 1]^3$  is a tetrahedron whose four vertices are also vertices of the unit cube  $[0, 1]^3$ . Therefore the diameter  $\varrho(\tilde{T})$  of the inscribed sphere of  $\tilde{T}$  is of order  $\mathcal{O}(1)$ , *i.e.*,  $\varrho(\tilde{T}) \sim 1$ . Similarly the second transformation *via*  $C_T^{-1}$  is examined. It maps  $T \rightarrow \hat{T}$ , where the tetrahedron  $\hat{T}$  has vertices  $(0, 0, 0)^\top$ ,  $(1, 0, 0)^\top$ ,  $(x_2, 1, 0)^\top$  and  $(x_3, y_3, 1)^\top$ , with  $0 \leq x_2, x_3 \leq 1$  and  $|y_3| \leq 1$ , *cf.* the definition of  $C_T$  or [12, Sect. 1.2]. Thus the diameter  $\text{diam}(\hat{T})$  of the tetrahedron  $\hat{T}$  satisfies  $1 < \text{diam}(\hat{T}) \leq \sqrt{6}$ .

The combined transformation *via*  $C_T^{-1} C_B$  now maps  $C_T^{-1} C_B : \tilde{T} \rightarrow \hat{T}$ . Hence the spectral norm of this matrix can be bounded from above by

$$\|C_T^{-1} C_B\| \leq \frac{\text{diam}(\hat{T})}{\varrho(\tilde{T})} \lesssim 1.$$

This inequality can be used to derive the matrix bound

$$\begin{aligned} \|C_T^{-1} C_x C_x^\top C_T^{-\top}\| &= \|C_T^{-1} C_B \cdot C_B^{-1} C_x C_x^\top C_B^{-\top} \cdot C_B^\top C_T^{-\top}\| \\ &\leq \|C_B^{-1} C_x\|^2 \cdot \|C_T^{-1} C_B\|^2 \lesssim \max_{i=1,2,3} \frac{h_{i,x}^2}{h_{i,B}^2} \cdot 1 \lesssim 1 \end{aligned}$$

since  $C_B^{-1} C_x = \text{diag}(h_{1,x}/h_{1,B}, h_{2,x}/h_{2,B}, h_{3,x}/h_{3,B})$ , and because of (4.1). The first matrix  $M := C_T^{-1} C_x C_x^\top C_T^{-\top}$  is symmetric and positive definite. For such matrices the largest eigenvalue is  $\lambda_{\max}(M) = \|M\|$  and hence

$$\lambda_{\max}(C_T^{-1} C_x C_x^\top C_T^{-\top}) \lesssim 1.$$

In a completely analogous fashion one treats  $M^{-1} = C_T^\top C_x^{-\top} C_x^{-1} C_T$  to obtain

$$\lambda_{\max}(M^{-1}) = \|C_T^\top C_x^{-\top} C_x^{-1} C_T\| \lesssim 1.$$



This implies  $\lambda_{\min}(M) = (\lambda_{\max}(M^{-1}))^{-1} \gtrsim 1$ , *i.e.* all eigenvalues of  $M$  are of order  $\mathcal{O}(1)$ . Since the eigenvalues of  $M = C_T^{-1}C_x C_x^\top C_T^{-\top}$  and of  $(C_x^{-\top}C_x^{-1})^{-1}C_T^{-\top}C_T^{-1}$  are the same, one further concludes

$$\underline{v}^\top C_x^{-\top} C_x^{-1} \underline{v} \sim \underline{v}^\top C_T^{-\top} C_T^{-1} \underline{v} \quad \forall \underline{v} \in \mathbb{R}^3.$$

This finally gives  $|C_x^{-1} \underline{v}| \sim |C_T^{-1} \underline{v}|$  for all  $\underline{v} \in \mathbb{R}^3$  which proves (A3). □

*Proof of assumption (A4).* We now prove that (A2) also yields (A4). Thus let  $T$  be an arbitrary tetrahedron and  $E$  be any face thereof. Employ the notation of the previous paragraphs and consider the brick  $B$  that circumscribes  $T$ . Then  $C_B^{-1}$  maps  $T$  onto  $\tilde{T}$  (see above). Next we consider the vector  $h_{E,T} \underline{n}_E$  in a geometric way. If the unit vector  $\underline{n}_E$  points inward (with respect to  $T$ ) then  $h_{E,T} \underline{n}_E$  points from the face  $E$  of  $T$  (or its plane) to the opposite vertex of  $T$ . If  $\underline{n}_E$  is the outward vector then consider  $-h_{E,T} \underline{n}_E$  instead.

Therefore  $C_B^{-1} h_{E,T} \underline{n}_E$  is a vector that points from the face  $\tilde{E} := C_B^{-1} E$  of  $\tilde{T}$  to the opposite vertex of  $\tilde{T}$ . This results in

$$1 \sim \varrho(\tilde{T}) < |C_B^{-1} h_{E,T} \underline{n}_E| < \sqrt{3} \quad , \textit{i.e.} \quad |C_B^{-1} \underline{n}_E| \sim h_{E,T}^{-1}.$$

Next recall that  $C_x^{-1} C_B$  is a diagonal matrix. Apply the equivalence  $h_{i,B} \sim h_{i,x}$  from above to conclude

$$\min_{i=1,2,3} \frac{h_{i,B}}{h_{i,x}} \cdot |C_B^{-1} \underline{n}_E| \leq |C_x^{-1} \underline{n}_E| = |C_x^{-1} C_B \cdot C_B^{-1} \underline{n}_E| \leq \max_{i=1,2,3} \frac{h_{i,B}}{h_{i,x}} \cdot |C_B^{-1} \underline{n}_E|.$$

In conjunction with (A3) one finally arrives at the desired equivalence

$$h_{E,T}^{-1} \sim |C_B^{-1} \underline{n}_E| \sim |C_x^{-1} \underline{n}_E| \stackrel{(A3)}{\sim} |C_T^{-1} \underline{n}_E|. \quad \square$$

*Proof of assumption (A5).* For (A5) to hold we have to specify assumption (A2) slightly more precisely, namely we demand

$$\frac{h_{\min,T_1}}{h_{\min,T_2}} < \alpha_d := \begin{cases} \sqrt{2} + \sqrt{3} \approx 3.146 & \text{if } d = 2 \\ (3 + \sqrt{5})/2 \approx 2.618 & \text{if } d = 3 \end{cases} \quad \forall T_1 \cap T_2 \neq \emptyset.$$

This slightly more restrictive assumption on the change of  $h_{\min,T}$  across neighbouring elements immediately implies the first inequality of (2.4) in (A5).

In order to investigate the neighbourhood patches  $R_k(T)$  observe first that  $\bigcup_{l=0}^k R_l(T)$  contains  $\mathcal{O}(k^d)$  elements. Hence  $R_k(T)$  contains  $\mathcal{O}(k^{d-1})$  elements, and the second inequality of (2.4) in (A5) holds with  $r = d - 1, \beta = 1$ . With these values of  $\alpha_d$  and  $\beta$  the third inequality of (2.4) in (A5) is satisfied as well. □

### 4.2. Assumption (A3) implies (A1) and (A2)

In this section we state that assumptions (A1) and (A2) are already consequences of assumption (A3). The proof uses some standard arguments (as in the above section) and is therefore omitted for the sake of shortness (we refer to [16] for the details).

**Theorem 4.1.** *Assumption (A3) implies (A1) and (A2).*

**Remark 4.2.** The converse implication does not hold as a comparatively simple counterexample can show. Thus (A3) is a stronger assumption.

**4.3. Necessary and sufficient condition for mesh assumption (A3)**

Here we state a geometrical condition which is necessary and sufficient for assumption (A3) on *unstructured* tetrahedral meshes. We start with some technical equivalences.

**Lemma 4.3.** *The assumption (A3) is equivalent to the condition*

$$\|C_x^{-1}C_T\| \sim 1 \quad \text{and} \quad \|C_T^{-1}C_x\| \sim 1 \quad \forall T \subset \omega_x \text{ and all nodes } x. \tag{4.2}$$

*Proof.*  $\Rightarrow$ : Starting from (A3) and taking  $\underline{v} := C_T \underline{w}$ , we get

$$|C_x^{-1}C_T \underline{w}| = |C_x^{-1} \underline{v}| \stackrel{(A3)}{\sim} |C_T^{-1} \underline{v}| = |\underline{w}| \quad \forall \underline{w} \in \mathbb{R}^d, \forall T \subset \omega_x.$$

This yields

$$\|C_x^{-1}C_T\| = \max_{|\underline{w}|=1} |C_x^{-1}C_T \underline{w}| \sim 1.$$

We obtain similarly the second bound by taking  $\underline{v} := C_x \underline{w}$ .

$\Leftarrow$ : Define the symmetric, positive definite matrix  $M := C_T^{-1}C_x C_x^T C_T^{-T}$ . Completely analogous to Section 4.1 one concludes

$$\begin{aligned} \lambda_{max}(M) &= \|C_T^{-1}C_x C_x^T C_T^{-T}\| \leq \|C_T^{-1}C_x\|^2 \sim 1 \\ \lambda_{min}(M) &= (\lambda_{max}(M^{-1}))^{-1} = \|C_T^T C_x^{-T} C_x^{-1} C_T\|^{-1} \geq \|C_x^{-1}C_T\|^{-2} \sim 1. \end{aligned}$$

Hence all eigenvalues of  $M$  are of order  $\mathcal{O}(1)$ . Following once more the arguments of Section 4.1 yields

$$|C_T^{-1} \underline{v}| \sim |C_x^{-1} \underline{v}| \quad \forall \underline{v} \in \mathbb{R}^d,$$

which is nothing else than (A3). □

**Corollary 4.4.** The assumption (A3) is equivalent to the condition

$$\|C_{T_1}^{-1}C_{T_2}\| \lesssim 1 \quad \forall T_1, T_2 \subset \omega_x \text{ and all nodes } x. \tag{4.3}$$

*Proof.* For the necessity of the condition (4.3) apply Lemma 4.3 and write

$$\|C_{T_1}^{-1}C_{T_2}\| = \|C_{T_1}^{-1}C_x C_x^{-1}C_{T_2}\| \leq \|C_{T_1}^{-1}C_x\| \cdot \|C_x^{-1}C_{T_2}\| \sim 1 \quad \forall T_1, T_2 \subset \omega_x.$$

The sufficiency of (4.3) follows directly by the choice  $C_x := C_{T'}$  for an arbitrary element  $T' \subset \omega_x$ . □

**Theorem 4.5** (Equivalent formulation of (A3)). *Assume that for all patches  $\omega_x$  and any two elements  $T_1, T_2 \subset \omega_x$  the inequality*

$$\left| \cos \angle \left[ \underline{p}_{i,T_1}, \underline{p}_{j,T_2} \right] \right| \lesssim \frac{h_{i,T_1}}{h_{j,T_2}} \quad \forall 1 \leq i, j \leq d \tag{4.4}$$

*is satisfied. Then we can fix an arbitrary element  $T' \subset \omega_x$  and set  $C_x := C_{T'}$ . This choice implies assumption (A3), i.e.*

$$|C_x^{-1} \underline{v}| \sim |C_{T'}^{-1} \underline{v}| \quad \forall \underline{v} \in \mathbb{R}^d, \forall T \subset \omega_x.$$

*Conversely the assumption (A3) implies (4.4) for all  $T_1, T_2 \subset \omega_x$  and all nodes  $x$ .*

*Proof.* Let us first derive an equivalent formulation of inequality (4.4). Fix an arbitrary patch  $\omega_x$  and two arbitrary elements  $T_1, T_2 \subset \omega_x$ . Since the vectors  $\underline{p}_{i,T_1}$  are mutually orthogonal, there exists a unique decomposition

$$\underline{p}_{j,T_2} = \sum_{i=1}^d \alpha_{ij} \cdot \underline{p}_{i,T_1} \quad \forall j = 1 \dots d.$$

The real coefficients  $\alpha_{ij}$  satisfy  $(\underline{p}_{j,T_2}, \underline{p}_{i,T_1}) = \alpha_{ij} \cdot (\underline{p}_{i,T_1}, \underline{p}_{i,T_1}) = \alpha_{ij} \cdot h_{i,T_1}^2$ . Utilizing the definition of  $h_{i,T_k}$  one obtains

$$\alpha_{ij} = \frac{h_{j,T_2} h_{i,T_1} \cdot \cos \angle [\underline{p}_{i,T_1}, \underline{p}_{j,T_2}]}{h_{i,T_1}^2} = \frac{h_{j,T_2}}{h_{i,T_1}} \cdot \cos \angle [\underline{p}_{i,T_1}, \underline{p}_{j,T_2}].$$

Condition (4.4) of the theorem is thus equivalent to

$$|\alpha_{ij}| \lesssim 1 \quad \forall 1 \leq i, j \leq d.$$

Recall next that the matrices  $C_{T_1}, C_{T_2}$  are formed by  $C_{T_k} := (\underline{p}_{1,T_k}, \underline{p}_{2,T_k}, \underline{p}_{3,T_k})$ , cf. Section 2.2, which results in

$$\begin{aligned} C_{T_1}^{-1} \underline{p}_{j,T_2} &= C_{T_1}^{-1} \sum_{i=1}^d \alpha_{ij} \underline{p}_{i,T_1} = \sum_{i=1}^d \alpha_{ij} \underline{e}_i \\ C_{T_1}^{-1} C_{T_2} &= (\alpha_{ij})_{i,j=1}^d \\ \|C_{T_1}^{-1} C_{T_2}\| &\sim \max_{i,j=1 \dots d} |\alpha_{ij}|. \end{aligned}$$

Hence  $\|C_{T_1}^{-1} C_{T_2}\| \lesssim 1$  is equivalent to  $|\alpha_{ij}| \lesssim 1 \forall i, j$  and to (4.4). From here we conclude the desired result thanks to Corollary 4.4.  $\square$

**Remark 4.6.** The previous theorem provides the means for practical tests whether assumption (A3) is satisfied on a real mesh. For neighbouring elements one has to compute the angle between the main anisotropic direction vectors  $\underline{p}_{i,T_1}$  and  $\underline{p}_{j,T_2}$  and compare its cosine with the stretching ratio  $h_{i,T_1}/h_{j,T_2}$ .

#### 4.4. Necessary and sufficient condition for mesh assumption (A4)

In this section we give equivalent formulations of mesh assumption (A4), both of which are geometrically characterized.

**Theorem 4.7** (Equivalent formulation of (A4)). *The assumption (A4) holds if and only if for all elements  $T$  and all faces  $E \subset \partial T$  one has*

$$\max_{i=1, \dots, d} h_{i,T}^{-1} \cdot \left| \cos \angle [\underline{p}_{i,T}, \underline{n}_E] \right| \lesssim h_{E,T}^{-1}. \tag{4.5}$$

*Proof.* Fix an element  $T$  and a face  $E \subset \partial T$ . As before we may write

$$\underline{n}_E = \sum_{i=1}^d \alpha_i \cdot \underline{p}_{i,T},$$

with  $(\underline{n}_E, \underline{p}_{i,T}) = \alpha_i \cdot h_{i,T}^2$  and  $\alpha_i = h_{i,T}^{-1} \cos \angle [\underline{p}_{i,T}, \underline{n}_E]$ . Since  $C_T^{-1} \underline{p}_{i,T} = \underline{e}_i$  we obtain  $C_T^{-1} \underline{n}_E = (\alpha_1, \alpha_2, \alpha_3)^\top$ . From the equivalence of norms in  $\mathbb{R}^d$  we conclude

$$\|C_T^{-1} \underline{n}_E\| \sim \max_{i=1, \dots, d} h_{i,T}^{-1} \cdot \left| \cos \angle [\underline{p}_{i,T}, \underline{n}_E] \right|$$

which finishes the proof. □

Next we derive a purely geometrical formulation of (A4). This assumption states

$$|C_T^{-1}h_{E,T\underline{n}_E}| \sim 1 \quad \forall E \subset \partial T.$$

Thus fix an arbitrary element  $T$ . Given a face  $E \subset \partial T$ , denote temporarily its opposite vertex by  $V_E$ . Let  $U_E$  be the orthogonal projection of  $V_E$  onto  $E$  (or the plane that contains  $E$ ). Hence  $\overrightarrow{U_E V_E}$  is the height of  $V_E$  onto the plane of  $E$ .

Next we have to define an appropriate neighbourhood of  $E$ . To this end denote by  $E^\alpha$  the face  $E$  scaled by the real factor  $\alpha$  with respect to the midpoint  $M_E$  of  $E$ . In vector notation this can be written as  $E^\alpha := \{\vec{M}_E + \alpha \cdot (\underline{y} - \vec{M}_E) : \underline{y} \in E\}$ . In other words,  $E^\alpha$  is contained in the plane of  $E$ , and  $E^1 \equiv E$ . With that definition we can reformulate (A4) as an equivalent geometrical condition.

**Theorem 4.8** (Equivalent formulation of (A4)). *(A4) holds if and only if  $U_E \in E^\alpha$  is satisfied for all  $E \subset \partial T$  with some  $\alpha \lesssim 1$ .*

*Proof.* Obviously the vector  $\overrightarrow{U_E V_E}$  equals  $U_E \overrightarrow{V_E} = \pm h_{E,T\underline{n}_E}$ . Since  $C_T^{-1}$  maps  $T$  onto  $\hat{T}$ , the vector  $\overrightarrow{U_E V_E}$  is mapped onto a vector from the point  $\hat{U}_E := C_T^{-1}(U_E)$  of the face  $\hat{E} := C_T^{-1}(E)$  of  $\hat{T}$  to the opposite vertex  $\hat{V}_E := C_T^{-1}(V_E)$ . Utilizing  $\pm C_T^{-1}h_{E,T\underline{n}_E} = C_T^{-1}(U_E \overrightarrow{V_E}) = \hat{U}_E \overrightarrow{\hat{V}_E}$ , assumption (A4) can be rewritten as

$$\left| \hat{U}_E \overrightarrow{\hat{V}_E} \right| \sim 1.$$

Because  $\hat{T}$  is an isotropic tetrahedron of size  $\mathcal{O}(1)$  and  $\hat{V}_E$  is a vertex thereof, this is equivalent to  $\hat{U}_E \in \hat{E}^\alpha$  and  $U_E \in E^\alpha$ , with  $\alpha \lesssim 1$ . □

#### 4.5. Prismatic tensor product meshes satisfy the mesh assumptions

In Section 4.1 we have shown that the mesh assumptions are satisfied for tetrahedral meshes which are the tensor product of *three 1D meshes*. In this section we state that the assumptions (A3)–(A5) hold also for anisotropic tensor product meshes of a *prismatic domain*  $\Omega = G \times (a, b)$  with  $a < b$ , obtained using a 2D refined isotropic mesh of Raugel’s type in  $G$  and a uniform mesh in the third direction. Examples of such meshes are given in the right part of Figure 5 and in [2].

We define families of meshes  $\mathcal{T}_h$  of  $\Omega$  by introducing in  $G$  the standard mesh grading for two-dimensional corner problems, see for example [18, 19]. Let  $\mathcal{T}_G = \{K\}$  be a regular isotropic triangulation of  $G$ ; the elements  $K$  are triangles. Let  $r_K$  be the distance of  $K$  to the corner,

$$r_K := \inf_{(x_1, x_2) \in K} (x_1^2 + x_2^2)^{1/2},$$

(note that  $\Omega$  is scaled such that  $r_K < 1$ ). With  $h$  being a global mesh parameter and  $\mu \in (0, 1]$  being a grading parameter, we assume that the element size  $h_K := \text{diam } K$  satisfies

$$h_K \sim \begin{cases} h^{1/\mu} & \text{for } r_K = 0, \\ hr_T^{1-\mu} & \text{for } r_K > 0. \end{cases}$$

This graded two-dimensional mesh is now extended in the third dimension using the uniform mesh size  $h$ . In this way we obtain a pentahedral (*i.e.* prismatic) triangulation and, by dividing each pentahedron into three tetrahedra, we further get a tetrahedral triangulation  $\mathcal{T}_h$  of  $\Omega$ , see the right part of Figure 5 for an illustration.

Using Lemma 4.3 and Theorem 4.7 one can show the following theorem (see [16] for the details).

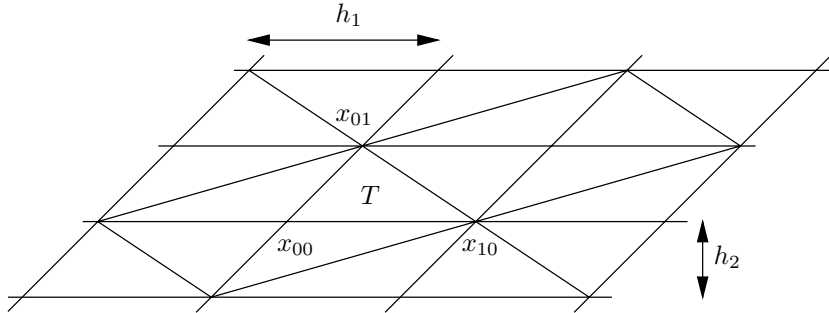


FIGURE 3. Mesh for the counterexample.

**Theorem 4.9.** *The above family of meshes satisfies assumptions (A3) and (A4).*

Note that the assumption (A5) holds under exactly the same conditions as described for the rectangular tensor product type meshes of Section 4.1.

**4.6. Mesh assumption (A4) is necessary for error estimation**

In the previous sections we investigated what meshes satisfy the mesh assumptions. In contrast, this section sheds light on the *role* that the mesh assumptions play in error estimation.

Our main Theorem 3.13 states that mesh assumptions (A1)–(A4) are *sufficient* to prove equivalences between the residual error estimator and the ZZ error estimator. Here we prove that mesh assumption (A4) is also a *necessary* condition. To this end we present a 2D counterexample where (A4) is violated and consequently the desired equivalences no longer hold.

Consider a criss-cross type mesh with nodal points located at

$$x_{ik} = \begin{pmatrix} i \cdot h_1 + k \cdot h_2 \\ k \cdot h_2 \end{pmatrix} = i \cdot h_1 \begin{pmatrix} 1 \\ 0 \end{pmatrix} + k \cdot h_2 \begin{pmatrix} 1 \\ 1 \end{pmatrix} \quad i, k \in \mathbf{Z},$$

where  $0 < h_2 \ll h_1$  are fixed parameters, cf. Figure 3. Elementary calculations yield

**Theorem 4.10.** *The above mesh (cf. Fig. 3) satisfies assumptions (A1), (A2) and (A3) but not assumption (A4). Moreover there exists a finite element solution  $u_h := \max\{0, y - x, x - y - h_1\}$  which has in particular the nodal values*

$$u_h(x_{ik}) = \begin{cases} 0 & \text{for } i = 0 \text{ or } i = 1 \\ h_1 & \text{for } i = 2 \text{ or } i = -1, \end{cases}$$

see Figure 4, for which relation (3.32) of Theorem 3.13 does not hold. Consequently the local equivalence (3.29) is violated at certain nodes  $x$  of  $\mathcal{T}_h$ .

**5. NUMERICAL EXPERIMENTS**

The aims of the numerical experiments are threefold. Firstly we investigate the mesh assumptions. Secondly the main theoretical predictions are to be verified. Lastly the constants that are involved in most inequalities/equivalences are examined numerically, and the asymptotic behaviour is observed.

To this end we present five experiments. The first one features an *isotropic* solution on an *isotropic* mesh, and thus tells what can reasonably be expected. The second experiment exhibits an *anisotropic* solution on tensor product type, rectangular anisotropic mesh. We believe such structured meshes to be best suited for ZZ error estimation. Finally the third experiment involves an anisotropic solution on a more irregular anisotropic mesh (which is unstructured in the  $xy$  directions, cf. also Sect. 4.5).

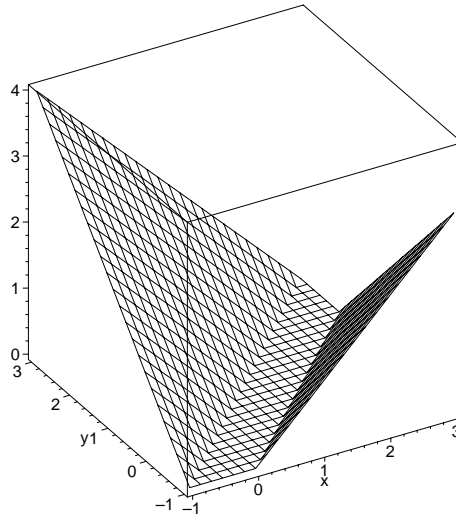


FIGURE 4. Finite element solution  $u_h$  for the counterexample.

In Section 5.1 we present the details of each experiment. Section 5.2 is devoted to the mesh assumptions (A3) and (A4). Finally in Section 5.3 the main theoretical results are tested numerically. We restrict ourselves to the second ZZ error estimator  $\eta_{Z_2}$  because it is more general than the first ZZ estimator  $\eta_{Z_1}$ , and since the second estimator allows *local* equivalences/estimates.

### 5.1. Description of the experiments

*Experiment 1: Isotropic solution + uniform mesh*

This experiment utilizes the most favorite settings; thus one can observe which results reasonably can be expected. Here we solve the Poisson problem

$$-\Delta u = f \quad \text{in } \Omega := (0, 1)^3, \quad u = u_D \quad \text{on } \partial\Omega.$$

The exact *isotropic* solution  $u$  is prescribed to be

$$u = e^{-x} + e^{-y} + e^{-z},$$

and the data  $f, u_D$  are chosen accordingly. We employ *isotropic*, uniform tetrahedral meshes  $\mathcal{T}_l, l = 1 \dots 5$ , which are the tensor product of three uniform 1D meshes of mesh size  $h = 2^{-l}$ . The table below displays some interesting information about mesh and solution.

Level $l$	Elements	$\ \nabla(u - u_h)\ _\Omega$	$\max_{T \in \mathcal{T}_l} h_{1,T}/h_{3,T}$	$m_1(u - u_h, \mathcal{T}_l)$
1	48	$1.61E - 1$	2.45	1.71
2	384	$8.16E - 2$	2.45	1.71
3	3 072	$4.10E - 2$	2.45	1.71
4	24 576	$2.05E - 2$	2.45	1.71
5	196 508	$1.03E - 2$	2.45	1.71

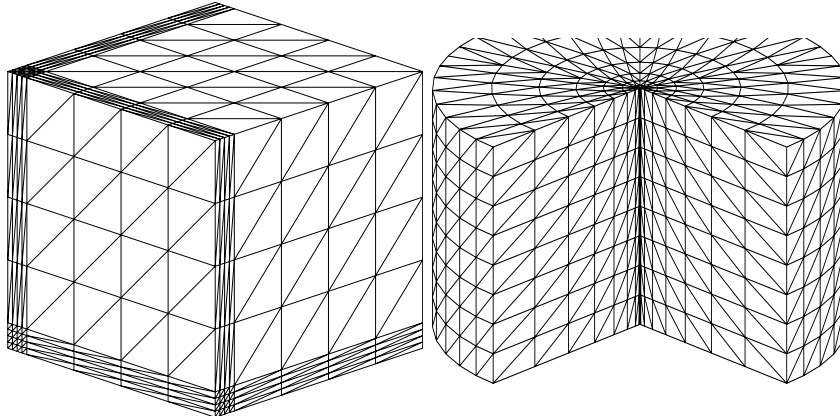


FIGURE 5. Meshes  $\mathcal{T}_3$  of experiment 2 (left) and 3 (right).

*Experiment 2: Anisotropic solution + structured anisotropic mesh*

Here again the Poisson problem with inhomogeneous Dirichlet boundary conditions is solved in  $\Omega := (0, 1)^3$ . The exact *anisotropic* solution  $u$  is here prescribed to be

$$u = e^{-x/\varepsilon} + e^{-y/\varepsilon} + e^{-z/\varepsilon}, \quad \varepsilon := 10^{-2},$$

and thus exhibits sharp boundary layers along the planes  $x = 0, y = 0$  and  $z = 0$ . The data  $f, u_D$  are chosen accordingly. We employ structured *anisotropic* meshes  $\mathcal{T}_l, l = 1 \dots 5$ , cf. left part of Figure 5. These meshes are formed by the tensor product of three 1D Shishkin type meshes with transition point at  $\tau = 2\varepsilon |\ln \varepsilon|$ .

In a similar fashion as before we present details of mesh and solution.

Level $l$	Elements	$\ \nabla(u - u_h)\ _\Omega$	$\max_{T \in \mathcal{T}_l} h_{1,T}/h_{3,T}$	$m_1(u - u_h, \mathcal{T}_l)$
1	48	$9.91E + 0$	14.1	1.61
2	384	$8.82E + 0$	14.3	1.71
3	3 072	$6.28E + 0$	14.4	1.70
4	24 576	$3.67E + 0$	14.5	1.67
5	196 508	$1.94E + 0$	14.5	1.62

Note first that the problem is comparatively poorly resolved. This is mainly due to the right hand side  $f = -\Delta u \equiv \varepsilon^{-2}u$  which has large and steep boundary layers (although still  $f \in L_2(\Omega)$ ). Secondly, the maximum aspect ratio of the anisotropic meshes is about 1:15. These meshes are well suited to the anisotropic solution, as the small matching number  $m_1(u - u_h, \mathcal{T}_l) \approx 1.7$  confirms (cf. also exp. 1).

*Experiment 3: Anisotropic solution + semi-structured anisotropic mesh*

The domain  $\Omega$  here consists of 3/4 of a cylinder of height and radius 1, cf. the right part of Figure 5. The exact *anisotropic* solution  $u$  is prescribed to be

$$u(r, \varphi, z) = r^\lambda \cdot \sin(\lambda\varphi) \cdot \begin{cases} 1 + 2z(2z - 1) & \text{for } z \in (0, 1/2) \\ 1 + (3 - 4z)(2z - 1) & \text{for } z \in (1/2, 1) \end{cases}, \quad \lambda = 2/3.$$

This function behaves anisotropically along the concave edge, and is piecewise quadratic in the  $z$  direction. The data  $f$  and  $u_D$  are chosen accordingly.

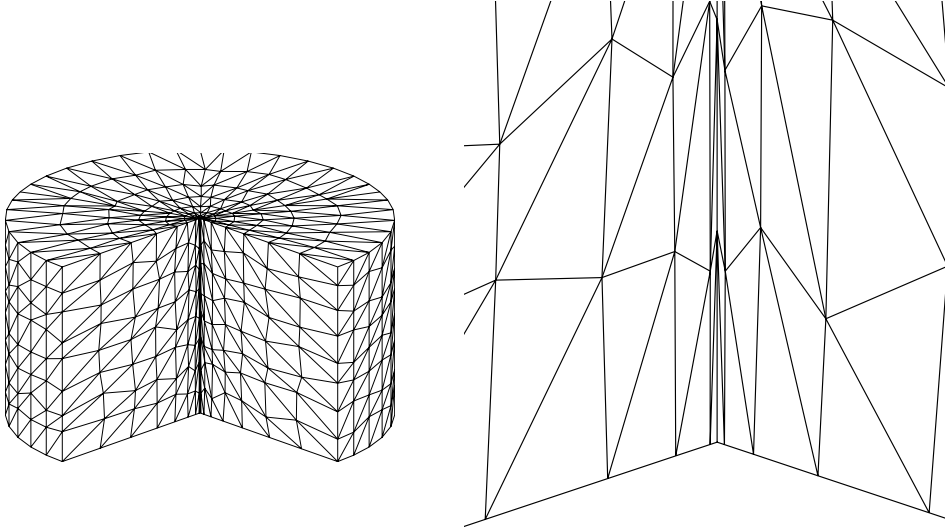


FIGURE 6. Mesh  $\mathcal{T}_3$  of experiment 4 (left) and zoom at bottom corner (right).

The sequence of meshes  $\mathcal{T}_l$ ,  $l = 1 \dots 5$ , is constructed by first generating an isotropic, uniform mesh in the domain  $\Omega$ . The subsequent nodal coordinate transformation

$$(x, y, z)^\top := (\rho \cdot \hat{x}, \rho \cdot \hat{y}, \hat{z})^\top \quad \text{with } \rho := \{\hat{x}^2 + \hat{y}^2\}^{(1-\mu)/2\mu}, \quad \mu = 0.4,$$

yields the final, *anisotropic* mesh, see right part of Figure 5. Hence the semi-structured meshes  $\mathcal{T}_l$  are the tensor product of an unstructured, graded 2D mesh in the  $xy$  plane, and a uniform 1D mesh in the  $z$  direction, as in Section 4.5.

The details of the meshes and of the solution are displayed below. The problem is well resolved, and all anisotropic meshes are well adapted to the solution, *i.e.*  $m_1 < 2$ .

Level $l$	Elements	$\ \nabla(u - u_h)\ _\Omega$	$\max_{T \in \mathcal{T}_l} h_{1,T}/h_{3,T}$	$m_1(u - u_h, \mathcal{T}_l)$
1	96	$1.59E + 0$	5.4	1.91
2	768	$8.60E - 1$	9.7	1.86
3	6 144	$4.50E - 1$	27.3	1.83
4	49 152	$2.33E - 1$	77.0	1.83
5	393 216	$1.21E - 1$	217.7	1.83

*Experiment 4: Anisotropic solution + randomly perturbed mesh*

The domain  $\Omega$  and the solution  $u$  are the same as in experiment 3. The meshes  $\mathcal{T}_l$ ,  $l = 1 \dots 5$ , are here a random perturbation of those of experiment 3. More precisely we start with an anisotropic mesh as in experiment 3, perturb the nodal coordinates randomly (uniform distribution,  $\max = \pm 0.25h_{\text{isotrop}} \approx 0.25 \cdot 2^{-l}$ ) and then use the same nodal coordinate transformation as above. See the left part of Figure 6.

Although the mesh still looks quite structured, this is not really the case. Near the concave edge the random perturbations of the nodal coordinates have quite a dramatic effect on the shape of the elements, which can be clearly seen in the zoom at the right part of Figure 6. This will also be observed in the mesh assumptions (A3) and (A4), *cf.* Section 5.2 below.

The next table presents some informations on the meshes and the solution. From this, we may conclude a good decrease of the error (similar to exp. 3). As already pointed out, the random perturbations induce strongly



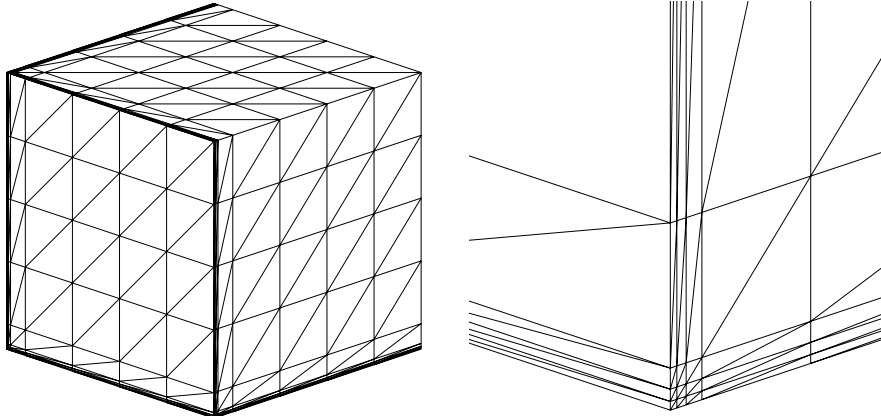


FIGURE 7. Mesh  $\mathcal{T}_3$  of experiment 5 (left) and zoom at bottom corner (right).

anisotropic meshes (the aspect ratio becoming larger and larger). Despite of this fact, the matching functions remain small and thus reliable error estimations may be expected.

Level $l$	Elements	$\ \nabla(u - u_h)\ _\Omega$	$\max_{T \in \mathcal{T}_l} h_{1,T}/h_{3,T}$	$m_1(u - u_h, \mathcal{T}_l)$
1	96	$1.66E + 0$	8.3	1.94
2	768	$9.11E - 1$	10.2	1.94
3	6 144	$5.02E - 1$	11.6	1.89
4	49 152	$2.62E - 1$	288.5	1.87
5	393 216	$1.35E - 1$	732.3	1.87

*Experiment 5: Singularly perturbed reaction diffusion problem*

According to the results from [17] and the theory developed in the previous sections, we may extend our results to reaction-diffusion equations (see below). These are problems where boundary layers appear naturally and for which the use of anisotropic meshes is recommended and particularly advantageous. This example illustrates that our theory is not restricted to the comparatively basic Poisson problem but also applies to real-world problems.

As a model problem we consider the singularly perturbed reaction-diffusion equation:

$$-\varepsilon \Delta u + \gamma u = f \quad \text{in } \Omega := (0, 1)^3,$$

with  $\varepsilon := 10^{-4}, \gamma := 1, f := 0$  and inhomogeneous Dirichlet boundary conditions such that the prescribed exact solution is

$$u = e^{-x/\sqrt{\varepsilon}} + e^{-y/\sqrt{\varepsilon}} + e^{-z/\sqrt{\varepsilon}}.$$

In contrast to experiment 2, we now employ a tensor product type mesh which consists of three 1D *Bakhvalov* type meshes, cf. left part of Figure 7. The main difference with a Shishkin type mesh is that here the mesh is exponentially graded inside the layer region. This can be observed in the right part of Figure 7 which depicts a zoom of the mesh. For a comprehensive description we refer to [15].

The energy norm over some domain  $\omega$  becomes

$$\|v\|_\omega^2 := \varepsilon \|\nabla v\|_\omega^2 + \gamma \|v\|_\omega^2.$$

For the reaction diffusion problem we have to modify the error estimators. To this end define

$$\mu_T := \min \left\{ \frac{h_{\min,T}}{\sqrt{\varepsilon}}, \frac{1}{\sqrt{\gamma}} \right\}, \quad \mu_x := \min \left\{ \frac{h_{\min,x}}{\sqrt{\varepsilon}}, \frac{1}{\sqrt{\gamma}} \right\}, \quad \mu_E := \min \left\{ \frac{h_{\min,E}}{\sqrt{\varepsilon}}, \frac{1}{\sqrt{\gamma}} \right\}.$$

On the basis of [17] and the theory of the previous sections one concludes fairly easily the definitions and results below.

**Definition 5.1** (Error estimators for a reaction diffusion problem).

$$\eta_{R,E}^2 := \mu_E \varepsilon^{3/2} h_{\min,E} h_E^{-1} \left\| \left[ \left[ \partial_{n_E} u_h \right] \right]_E \right\|_E^2 \tag{5.1}$$

$$\eta_{R,x}^2 := \mu_x \varepsilon^{3/2} h_{\min,x} |\omega_x| \sum_{E:x \in \mathcal{N}_E} h_E^{-2} \left\| \left[ \left[ \partial_{n_E} u_h \right] \right]_E \right\|_E^2 \tag{5.2}$$

$$\eta_{Z_2,x}^2 := \mu_x \varepsilon^{3/2} h_{\min,x} |\omega_x| \sum_{T \subset \omega_x} \left| C_x^{-1} (\nabla^{R_2} u_h(x) - \nabla u_h|_T(x)) \right|^2 \tag{5.3}$$

$$\eta_{Z_2,T}^2 := \mu_T \varepsilon^{3/2} h_{\min,T} \left\| C_T^{-1} (\nabla^{R_2} u_h - \nabla u_h) \right\|_T^2 \tag{5.4}$$

$$\eta_R^2 := \sum_{x \in \mathcal{N}_\Omega} \eta_{R,x}^2 \quad \eta_{Z_2}^2 := \sum_{T \in \mathcal{T}_h} \eta_{Z_2,T}^2 \tag{5.5}$$

Note that we use the same notation as for the Poisson problem. In the numerical experiments below it will be clear from the context which formula is to be used.

**Theorem 5.2** (ZZ error estimation for a reaction diffusion problem). *With the definitions from above, the following error equivalences and error bounds hold.*

$$\eta_{R,x} \sim \eta_{Z_2,x} \tag{5.6}$$

$$\eta_R \sim \eta_{Z_2} \tag{5.7}$$

$$\eta_{Z_2,x} \lesssim \|u - u_h\|_{\omega_x} + \inf_{f_h \in V_h} \mu_x \|f - f_h\|_{\omega_x} \quad \forall x \in \mathcal{N}_\Omega \tag{5.8}$$

$$\|u - u_h\|_\Omega \lesssim m_1(u - u_h, \mathcal{T}_h) \left( \eta_{Z_2}^2 + \inf_{f_h \in V_h} \sum_{T \in \mathcal{T}_h} \mu_T^2 \|f - f_h\|_T^2 \right)^{1/2}. \tag{5.9}$$

As before we give below some details on the meshes and the solution.

Level $l$	Elements	$\ u - u_h\ _\Omega$	$\max_{T \in \mathcal{T}_l} h_{1,T}/h_{3,T}$	$m_1(u - u_h, \mathcal{T}_l)$
1	48	$2.62E + 1$	14.1	1.60
2	384	$7.57E + 0$	62.0	1.85
3	3 072	$2.87E + 0$	74.6	2.19
4	24 576	$1.26E + 0$	80.3	2.46
5	196 508	$5.93E - 1$	83.1	2.63

Again, the error decreases at a quasi optimal rate. Although the meshes are relatively anisotropic, they are well adapted to the solution. This is reflected by small values of the matching function,  $m_1(u - u_h, \mathcal{T}_l) \approx 1.6 \dots 2.6$ . Thus reliable error estimation is to be expected.

**5.2. Mesh Assumptions (A3) and (A4)**

Here the mesh assumptions are investigated numerically. An additional, graphical presentation of some results is given in [16].

*Mesh Assumption (A3)*

This assumption can be reformulated as

$$c_1 \cdot |C_x^{-1} \underline{v}| \leq |C_T^{-1} \underline{v}| \leq c_2 \cdot |C_x^{-1} \underline{v}| \quad \forall \underline{v} \in \mathbb{R}^d, \forall T \subset \omega_x.$$

TABLE 1. Values of  $c_1, c_2$  for assumption (A3); all experiments.

Level	Experiment 1		Experiment 2		Experiment 3		Experiment 4		Experiment 5	
	$c_1$	$c_2$	$c_1$	$c_2$	$c_1$	$c_2$	$c_1$	$c_2$	$c_1$	$c_2$
1	0.500	1.856	0.082	11.908	0.451	6.169	0.345	10.184	0.082	11.908
2	0.500	1.856	0.078	11.730	0.323	7.328	0.290	13.393	0.092	12.553
3	0.500	1.856	0.078	11.640	0.379	8.178	0.234	18.147	0.103	10.997
4	0.500	1.856	0.078	13.444	0.372	8.336	0.238	21.070	0.109	9.092
5	0.500	1.856	0.078	13.419	0.353	8.357	0.227	39.618	0.127	7.517

TABLE 2. Values of  $c_3, c_4$  for assumption (A4); all experiments.

Level	Experiment 1		Experiment 2		Experiment 3		Experiment 4		Experiment 5	
	$c_3$	$c_4$	$c_3$	$c_4$	$c_3$	$c_4$	$c_3$	$c_4$	$c_3$	$c_4$
1	0.754	1.202	0.901	1.492	0.760	1.415	0.723	2.476	0.901	1.492
2	0.754	1.202	0.754	1.497	0.723	1.564	0.638	6.828	0.754	1.705
3	0.754	1.202	0.754	1.500	0.714	1.690	0.609	13.912	0.754	1.679
4	0.754	1.202	0.754	1.501	0.714	1.714	0.566	42.960	0.754	1.706
5	0.754	1.202	0.754	1.502	0.712	1.717	0.552	70.767	0.754	1.711

In order to investigate this condition numerically we have to specify the matrix  $C_x$  for a given node  $x$ . In view of Theorem 4.5 choose that element  $T \subset \omega_x$  that has the smallest aspect ratio  $h_{1,T}/h_{3,T}$ , and set  $C_x := C_T$ . Table 1 gives the corresponding values of  $c_1, c_2$  for all five experiments, and all meshes  $\mathcal{T}_l$ .

On isotropic meshes (experiment 1) one always has  $c_1 \sim c_2 \sim 1$  which is confirmed by the moderate values. For the anisotropic meshes of experiments 2 and 5, the theoretical considerations of Section 4.1 reveal that (A3) holds as well. The values of  $c_1, c_2$  are mainly of the same size over the different levels, although they are less favorable than in the isotropic case. This mainly seems to be due to relatively large changes of the element sizes  $h_{i,T}$  across neighbouring elements. This observation is strengthened by the results of experiments 3 and 4. Indeed experiment 3 presents a gradual change of the element sizes and the values of  $c_1, c_2$  are more moderate. In experiment 4, the random perturbations imply a large change of the element sizes which in turn seems to induce larger  $c_2$  so that (A3) tends to be violated.

Summarizing, well structured meshes are more advantageous for (A3) to hold.

*Mesh Assumption (A4)*

The assumption (A4) on the shape of the elements can be rewritten as

$$c_3 \leq |C_T^{-1} \underline{n}_E| \cdot h_{E,T} \leq c_4 \quad \forall T \in \mathcal{T}_l, \forall E \subset \partial T.$$

Utilizing the theory of Section 4.4, we can apply Theorem 4.8 to all experiments except exp. 4, which yields  $c_3 \sim c_4 \sim 1$  (alternatively employ the results of Section 4.1 for experiments 1, 2 and 5, as well as the results of Section 4.5 for experiment 3). This is verified impressively by the numerical results presented in Table 2. This table also reveals a very large value of  $c_4$  for experiment 4. This is readily explained by the very unstructured (*i.e.* perturbed) meshes of that example.

Summarizing, (A4) does not cause problems for well shaped elements.

TABLE 3. Values of  $c_5, c_6$  for *local* estimator equivalence (3.29); all experiments.

Level	Experiment 1		Experiment 2		Experiment 3		Experiment 4		Experiment 5	
	$c_5$	$c_6$	$c_5$	$c_6$	$c_5$	$c_6$	$c_5$	$c_6$	$c_5$	$c_6$
1	0.855	1.309	0.197	3.420	0.844	4.354	1.082	6.070	0.185	2.082
2	0.826	1.309	0.843	14.854	0.848	5.512	0.682	6.031	0.370	4.482
3	0.817	1.309	0.562	15.576	0.859	5.607	0.806	9.599	0.398	3.173
4	0.815	1.309	0.541	14.546	0.797	5.500	0.674	8.694	0.468	2.321
5	0.815	1.309	0.598	13.833	0.725	5.440	0.620	8.530	0.289	2.055

TABLE 4. *Global* estimator equivalence (3.30); all experiments.

Level	$\eta_R/\eta_{Z_2}$				
	Experiment 1	Experiment 2	Experiment 3	Experiment 4	Experiment 5
1	25.1	39.9	43.7	40.8	19.3
2	28.7	35.1	36.7	39.0	31.9
3	31.2	28.2	39.5	41.1	31.8
4	32.8	30.5	43.0	46.0	36.0
5	33.8	33.6	45.8	45.8	38.7

### 5.3. Main numerical results

In this section the main theoretical results for the second ZZ error estimator are tested numerically. The corresponding recovered gradient  $\nabla^{R_2}$  is defined with weights  $\mu_{T,x} := |T|/|\omega_x|$ , cf. definition 3.10.

First we investigate relations (3.29, 3.30) of Theorem 3.13 which state a local and global equivalence between the residual error estimator and the ZZ error estimator, respectively. Afterwards the results of the actual ZZ error estimation of Theorem 3.14 are presented.

#### Results for Theorem 3.13

The *local* equivalence (3.29) can be rewritten as

$$c_5 \cdot \eta_{Z_2,x} \leq \eta_{R,x} \leq c_6 \cdot \eta_{Z_2,x} \quad \forall x \in \mathcal{N}_{\bar{\Omega}}.$$

The values of  $c_5, c_6$  are given in Table 3. One observes that the equivalence between both error estimators diminishes for anisotropic meshes but is still acceptable even for the unstructured meshes of experiment 4 (note that  $c_5, c_6$  describe only the worst cases over all  $x \in \mathcal{N}_{\bar{\Omega}}$ ). The comparatively large values of  $c_6$  in experiments 2 and 4 seem to be caused by the sharp change of the element sizes. This observation is supported by experiment 5 which features a smoother change of the element sizes and reveals smaller values of  $c_6$ .

In view of experiment 4, the mesh assumption (A3) and (A4) seem to be of less influence on  $c_5$  and  $c_6$ .

The *global* equivalence (3.30) between the residual estimator and the ZZ estimator reads  $\eta_R \sim \eta_{Z_2}$ . Thus we present  $\eta_R/\eta_{Z_2}$  for all meshes and experiments. The results of Table 4 impressively confirm the theoretically proven equivalence and underline the weak influence of assumptions (A3) and (A4) on these results. Note that the comparatively large values of  $\eta_R/\eta_{Z_2}$  are mainly due to the different range of the sums, cf. (3.5) and (3.23). Furthermore the summand  $\eta_{R,x}$  contains the factor  $|\omega_x|$  while  $\eta_{Z_2,T}$  is related to  $|T|$ .

TABLE 5. Lower ZZ error bounds: (3.33) of Theorem 3.14 and (5.8) of Theorem 5.2; all experiments.

Level	Lower error bound (3.33)				Lower error bound (5.8)
	$\max_{x \in \mathcal{N}_{\bar{\Omega}}} \frac{\eta_{Z_2,x}}{\ \nabla(u - u_h)\ _{\omega_x} + \zeta_x}$				$\max_{x \in \mathcal{N}_{\bar{\Omega}}} \frac{\eta_{Z_2,x}}{\ u - u_h\ _{\omega_x} + \zeta_x}$
	Exp. 1	Exp. 2	Exp. 3	Exp. 4	Exp. 5
1	7.094	0.379	2.990	2.617	0.169
2	7.968	4.180	6.567	6.415	4.909
3	8.235	5.462	7.866	8.741	7.452
4	8.302	10.744	8.265	10.819	9.605
5	8.319	8.201	8.382	8.600	9.644

Results for Theorem 3.14

In order to present the results of the ZZ error estimation clearly, let us denote the data approximation terms of Theorems 3.14 and 5.2 by

Poisson problem (Th. 3.14):  $\zeta_x := h_{\min,x} \|f - L_h f\|_{\omega_x} \quad \zeta^2 := \sum_{T \in \mathcal{T}_h} h_{\min,T}^2 \|f - L_h f\|_T^2 \sim \sum_{x \in \mathcal{N}_{\bar{\Omega}}} \zeta_x^2$

Reaction diffusion problem (Th. 5.2):  $\zeta_x := \mu_x \|f - L_h f\|_{\omega_x} \quad \zeta^2 := \sum_{T \in \mathcal{T}_h} \mu_T \|f - L_h f\|_T^2 \sim \sum_{x \in \mathcal{N}_{\bar{\Omega}}} \zeta_x^2$

with  $L_h$  being the linear Lagrange interpolation operator. Note that for the reaction diffusion problem of experiment 5 we have  $f \equiv 0$  and thus  $\zeta_x = \zeta = 0$ .

Next, inequalities (3.33, 3.34) of Theorem 3.14 and (5.8, 5.9) of Theorem 5.2 can be reformulated as

	Lower error bound (3.33) or (5.8)	Upper error bound (3.34) or (5.9)
Poisson problem:	$\frac{\eta_{Z_2,x}}{\ \nabla(u - u_h)\ _{\omega_x} + \zeta_x} \lesssim 1 \quad \forall x \in \mathcal{N}_{\bar{\Omega}}$	$\frac{\ \nabla(u - u_h)\ _{\Omega}}{m_1(\eta_{Z_2}^2 + \zeta^2)^{1/2}} \lesssim 1$
Reaction diffusion problem:	$\frac{\eta_{Z_2,x}}{\ u - u_h\ _{\omega_x} + \zeta_x} \lesssim 1 \quad \forall x \in \mathcal{N}_{\bar{\Omega}}$	$\frac{\ u - u_h\ _{\Omega}}{m_1(\eta_{Z_2}^2 + \zeta^2)^{1/2}} \lesssim 1$

*i.e.* all ratios have to be bounded from above.

The numerical results for the lower error bounds (3.33) and (5.8) are given in Table 5. Clearly, all values are bounded from above, as predicted by the theory. Furthermore, the actual size is even similar to the values of the isotropic experiment 1. This is quite a positive surprise, in particular in view of the anisotropic elements or the unstructured meshes of experiment 4.

The exceptionally small values on the coarsest level 1 for experiments 2 and 5 seem to be due to an insufficient resolution of the boundary layer. Lastly, the influence of the mesh assumptions is again barely noticeable.

The numerical results for the upper error bounds (3.34) and (5.9) are presented in Table 6. All values but one are bounded from above by the value of the isotropic experiment 1 (*i.e.* bounded by 1.9). Thus the theory is confirmed again. Nevertheless the situation gives rise to several interpretations.

Firstly, the large value for experiment 5 level 1 seems to be due to unresolved boundary layers (*cf.* also the corresponding entry of Table 5 which shows an exceptionally small value). Secondly, on the coarse levels 1–3 of experiment 2 one observes small values. They are caused by a dominating approximation term  $\zeta$ , *i.e.* the discretization is not yet fine enough.

Summarizing, the numerical results are in good agreement with the theory. Again, the influence of the mesh assumptions seems to be rather weak.

TABLE 6. Upper ZZ error bounds: (3.34) of Theorem 3.14 and (5.9) of Theorem 5.2; all experiments.

Level	Upper error bound (3.34)				Upper error bound (5.9)
	Exp. 1	Exp. 2	Exp. 3	Exp. 4	Exp. 5
1	1.819	0.015	0.110	0.155	42.140
2	1.935	0.048	0.167	0.207	1.574
3	1.864	0.180	0.246	0.324	1.270
4	1.834	0.678	0.356	0.454	1.060
5	1.823	1.479	0.505	0.635	0.965

## 6. SUMMARY

Zienkiewicz–Zhu error estimators are popular because of their cheap implementation and their astonishing robustness. We have proposed and rigorously analysed two kinds of ZZ error estimators that can be applied to *anisotropic* tetrahedral finite element meshes. Both estimators have been defined by scaling the components of the original gradient  $\nabla u_h$  and some recovered gradient  $\nabla^R u_h$ . Although the first estimator has turned out to be a special case of the second one, both have been presented because of their different background, analysis, and results.

While our first ZZ estimator is related to a particular choice of the recovered gradient, our second ZZ estimator is much more flexible because *arbitrary weights* can be employed to define the recovered gradient. Hence our novel analysis proves that each averaging technique yields *reliable* and *efficient* error control.

The basis of examination is formed by standard isotropic (2D) arguments. However, our analysis requires a refinement and improvement of these ideas. Moreover, the main technicality of the analysis stems from the anisotropic nature of the discretization. Particular emphasis has been given to the requirements on the anisotropic mesh.

The analysis has been complemented and confirmed by extensive numerical examples. The experiments show that ZZ error estimation is possible on anisotropic meshes. The estimators can be defined for a large class of problems; here the Poisson problem and a singularly perturbed reaction diffusion problem have been treated. Judging from the numerical experiments, structured meshes are advantageous but unstructured ones also give satisfactory results. The examples further show that the mesh assumptions are not too dominant in practice.

## REFERENCES

- [1] M. Ainsworth and J.T. Oden, *A posteriori error estimation in finite element analysis*. Wiley (2000).
- [2] Th. Apel, *Anisotropic finite elements: Local estimates and applications*, Advances in Numerical Mathematics. Teubner, Stuttgart (1999).
- [3] I. Babuška, T. Strouboulis and C.S. Upadhyay, A model study of the quality of *a posteriori* error estimators for linear elliptic problems. Error estimation in the interior of patchwise uniform grids of triangles. *Comput. Methods Appl. Mech. Engrg.* **114** (1994) 307–378.
- [4] I. Babuška, T. Strouboulis, C.S. Upadhyay, S.K. Gangaraj and K. Copps, Validation of *a posteriori* error estimators by numerical approach. *Int. J. Numer. Methods Eng.* **37** (1994) 1073–1123.
- [5] S. Bartels and C. Carstensen, Each averaging technique yields reliable *a posteriori* error control in FEM on unstructured grids. Part II: High order FEM. *Math. Comp.* **71** (2002) 971–994.
- [6] J.H. Bramble, J.E. Pasciak and O. Steinbach, On the stability of the  $L_2$ -projection in  $H^1(\omega)$ . *Math. Comp.* **71** (2002) 147–156.
- [7] C. Carstensen, Merging the Bramble-Pasciak-Steinbach and the Crouzeix-Thomée criterion for  $H^1$ -stability of the  $L^2$ -projection onto finite element spaces. *Math. Comp.* **71** (2002) 157–163.
- [8] C. Carstensen and S. Bartels, Each averaging technique yields reliable *a posteriori* error control in FEM on unstructured grids. Part I: Low order conforming, nonconforming, and mixed FEM. *Math. Comp.* **71** (2002) 945–969.
- [9] P.G. Ciarlet, *The finite element method for elliptic problems*. North-Holland, Amsterdam (1978).

- [10] M. Dobrowolski, S. Gräf and C. Pflaum, On *a posteriori* error estimators in the finite element method on anisotropic meshes. *Electron. Trans. Numer. Anal.* **8** (1999) 36–45.
- [11] G. Kunert, *A posteriori* error estimation for anisotropic tetrahedral and triangular finite element meshes. Logos Verlag, Berlin (1999). Also Ph.D. thesis, TU Chemnitz, <http://archiv.tu-chemnitz.de/pub/1999/0012/index.html>
- [12] G. Kunert, An *a posteriori* residual error estimator for the finite element method on anisotropic tetrahedral meshes. *Numer. Math.* **86** (2000) 471–490, DOI 10.1007/s002110000170.
- [13] G. Kunert, A local problem error estimator for anisotropic tetrahedral finite element meshes. *SIAM J. Numer. Anal.* **39** (2001) 668–689.
- [14] G. Kunert, *A posteriori*  $L_2$  error estimation on anisotropic tetrahedral finite element meshes. *IMA J. Numer. Anal.* **21** (2001) 503–523.
- [15] G. Kunert, Robust *a posteriori* error estimation for a singularly perturbed reaction–diffusion equation on anisotropic tetrahedral meshes. *Adv. Comput. Math.* **15** (2001) 237–259.
- [16] G. Kunert and S. Nicaise, *Zienkiewicz–Zhu error estimators on anisotropic tetrahedral and triangular finite element meshes*, preprint SFB393/01–20, TU Chemnitz, July 2001. Also <http://archiv.tu-chemnitz.de/pub/2001/0059/index.html>
- [17] G. Kunert and R. Verfürth, Edge residuals dominate *a posteriori* error estimates for linear finite element methods on anisotropic triangular and tetrahedral meshes. *Numer. Math.* **86** (2000) 283–303, DOI 10.1007/s002110000152.
- [18] L.A. Oganessian and L.A. Rukhovets, *Variational-difference methods for the solution of elliptic equations*. Izd. Akad. Nauk Armyanskoi SSR, Jerevan (1979), in Russian.
- [19] G. Raugel, Résolution numérique par une méthode d’éléments finis du problème de Dirichlet pour le Laplacien dans un polygone. *C. R. Acad. Sci. Paris, Sér. I Math* **286** (1978) A791–A794.
- [20] R. Rodriguez, Some remarks on the Zienkiewicz–Zhu estimator. *Numer. Meth. PDE* **10** (1994) 625–635.
- [21] H.G. Roos and T. Linß, Gradient recovery for singularly perturbed boundary value problems II: Two-dimensional convection-diffusion. *Math. Models Methods Appl. Sci.* **11** (2001) 1169–1179.
- [22] K.G. Siebert, An *a posteriori* error estimator for anisotropic refinement. *Numer. Math.* **73** (1996) 373–398.
- [23] O. Steinbach, On the stability of the  $L_2$ -projection in fractional Sobolev spaces. *Numer. Math.* **88** (2001) 367–379.
- [24] R. Verfürth, *A review of a posteriori error estimation and adaptive mesh–refinement techniques*. Wiley-Teubner, Chichester, Stuttgart (1996).
- [25] Zh. Zhang, *Superconvergent finite element method on a Shishkin mesh for convection-diffusion problems*. Report 98-006, Texas Tech University (1998).
- [26] O.C. Zienkiewicz and J.Z. Zhu, A simple error estimator and adaptive procedure for practical engineering analysis. *Internat. J. Numer. Methods Engrg.* **24** (1987) 337–357.
- [27] O.C. Zienkiewicz and J.Z. Zhu, The superconvergent patch recovery (SPR) and adaptive finite element refinement. *Comput. Methods Appl. Mech. Engrg.* **101** (1992) 207–224.