



TECHNISCHE UNIVERSITÄT CHEMNITZ

## Sonderforschungsbereich 393

Parallele Numerische Simulation für Physik und Kontinuumsmechanik

Gerd Kunert

Zoubida Mghazli

Serge Nicaise

### A posteriori error estimation for a finite volume discretization on anisotropic meshes

Preprint SFB393/03-16

#### Abstract

A singularly perturbed reaction diffusion problem is considered. The small diffusion coefficient generically leads to solutions with boundary layers.

The problem is discretized by a vertex-centered finite volume method. The anisotropy of the solution is reflected by using *anisotropic meshes* which can improve the accuracy of the discretization considerably.

The main focus is on *a posteriori* error estimation. A residual type error estimator is proposed and rigorously analysed. It is shown to be robust with respect to the small perturbation parameter. The estimator is also robust with respect to the mesh anisotropy as long as the anisotropic mesh sufficiently reflects the anisotropy of the solution (which is almost always the case for sensible discretizations).

Altogether, reliable and efficient *a posteriori* error estimation is achieved for the finite volume method on anisotropic meshes.

**Keywords:** error estimator, anisotropic solution, finite volume method, singular perturbations

**AMS:** 65N15, 74S10, 35B25, 65N30

Preprintreihe des Chemnitzer SFB 393

ISSN 1619-7178 (Print)

ISSN 1619-7186 (Internet)

SFB393/03-16

September 2003

# Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
<b>2</b>	<b>Singularly perturbed model problem</b>	<b>2</b>
<b>3</b>	<b>Finite Volume discretization</b>	<b>3</b>
3.1	Notation . . . . .	3
3.2	Discretization . . . . .	4
3.3	Existence and uniqueness of a discrete solution . . . . .	5
<b>4</b>	<b>Anisotropic meshes</b>	<b>6</b>
<b>5</b>	<b>Anisotropic interpolation error estimates</b>	<b>9</b>
<b>6</b>	<b>Residual error estimation</b>	<b>11</b>
<b>7</b>	<b>Summary</b>	<b>15</b>

Author's addresses:

Gerd Kunert, TU Chemnitz, Fakultät für Mathematik, 09107 Chemnitz, Germany  
<http://www.tu-chemnitz.de/sfb393/>

Zoubida Mghazli, Université Ibn Tofail, Faculté des Sciences, Laboratoire SIANO,  
B.P. 133 Kenitra, Maroc  
[mghazli\\_zoubida@yahoo.com](mailto:mghazli_zoubida@yahoo.com)

Serge Nicaise, Université de Valenciennes et du Hainaut Cambrésis, MACS, B.P. 311, 59304  
Valenciennes Cedex, France  
[snicaise@univ-valenciennes.fr](mailto:snicaise@univ-valenciennes.fr)

# 1 Introduction

Certain classes of partial differential equations (PDEs) generically lead to solutions with strong directional features. These solutions are characterized by much variation in one spatial direction but little variation otherwise. Such functions are frequently termed *anisotropic*. Examples of PDEs with anisotropic solutions include, for example, the Poisson problem in 3D domains with concave edges [Ape99], singularly perturbed diffusion convection reaction problems [RST96, MOS96, Mor96], or the (linearized) Navier Stokes problem.

When solving such PDEs numerically, it is natural (and in some cases even necessary) to reflect the anisotropy of the solution by a proper, *anisotropic discretization*. In the context of the finite element method, this leads to anisotropic meshes, cf. the extensive discussion in [Ape99]. The anisotropic elements are characterized by a (vary) large aspect ration, i.e. the ratio of the diameters of the circumscribed sphere and the inscribed sphere is large (ranging e.g. from 10 to more than  $10^6$ ). Figure 1 visualizes this notion. To give an illustrative example, the well-known Shishkin meshes are anisotropic (see e.g. [RST96, Ape99]).

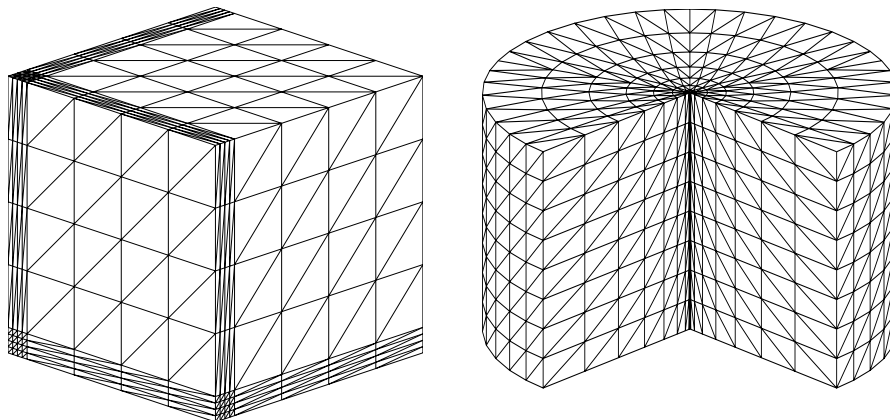


Figure 1: Examples of anisotropic meshes: Shishkin type mesh (left), graded mesh (right)

In order to employ anisotropic discretizations efficiently and comfortably, *adaptivity* becomes a major issue. While this topic is quite mature by now for standard, isotropic discretizations, the matter is much more complicated when anisotropic meshes occur. For example, the mesh construction/refinement/coarsening is quite more technical. Furthermore, *a posteriori error estimation* has to be reinvestigated completely since standard, isotropic estimators usually fail when applied on anisotropic discretizations. This failure is caused by the large aspect ratio which is no longer bounded, as the standard theory requires.

From now on, our focus is on *a posteriori* error estimation on anisotropic discretizations. This topic is well investigated and mature for the finite element method (FEM),

see the textbooks [Ver96, AO00]. For *finite volume method (FVM)*, the literature about *a posteriori* error estimation is less voluminous; exemplarily we refer to [MMM93, Ang95, KO00, Ohl01a, Ohl01b, LT02, BM00, BMV03] and the citations therein.

In this paper we show that *a posteriori* error estimation for the finite volume method can also be achieved for anisotropic discretizations. Up to now, no comparable result is known to the authors. The analysis here employs several tools that already have proved useful in the context of the finite element method; this facilitates our investigations to some extent.<sup>1</sup>

The remainder of this work is organised as follows. Section 2 introduces the model problem. Some notation and the finite volume discretization are presented in Section 3. The question of existence and uniqueness of the discrete solution is answered there as well. Section 4 is devoted to the anisotropic elements and introduces appropriate notation for them. The new anisotropic interpolation error estimates of Section 5 are similar in structure to estimates that are used to analyse the finite element method. Correspondingly, some known techniques can be employed with comparatively little modifications. Finally, Section 6 presents the main, novel results of our work, namely upper and lower error bounds. A residual type error estimator is proposed, rigorously analysed, and shown to be efficient and reliable (for appropriate anisotropic discretizations). The analysis is partly similar to the isotropic counterpart of [BMV03].

## 2 Singularly perturbed model problem

Our exposition features a singularly perturbed reaction diffusion model problem. Its classical formulation reads:

$$-\varepsilon\Delta u + cu = f \quad \text{in } \Omega, \quad u = 0 \quad \text{on } \Gamma_D = \partial\Omega, \quad (1)$$

where  $\varepsilon > 0$  and  $c \geq 0$  are constant, and  $\Omega$  is a bounded domain of  $\mathbb{R}^d$ ,  $d = 2$  or  $3$ , with a polygonal or polyhedral boundary.

When the perturbation parameter is small,  $\varepsilon \ll c$ , then frequently the solution  $u$  of (1) is anisotropic, i.e. it exhibits boundary layers. Then anisotropic discretizations are particularly useful or even mandatory, cf. also the left part of Figure 1.

Note further that (1) contains the Poisson problem as a special case (via  $\varepsilon = 1, c = 0$ ). For 3D domains with concave edges, anisotropic edge singularities occur, and anisotropic discretizations similar to that of Figure 1 (right) are appropriate, cf. [Ape99].

Often the classical formulation (1) is too restrictive to describe the underlying physical phenomena properly. Then the so-called *weak formulation* is more appropriate. To this end denote by  $X = H_0^1(\Omega)$  the usual Sobolev space of functions whose first derivative is in

---

<sup>1</sup>This observation coincides with the fact that there are certain inherent links between the finite volume method and the finite element method.

$L^2(\Omega)$ , and whose trace on  $\partial\Omega$  vanishes. The weak formulation now reads:

$$\left. \begin{array}{l} \text{Find } u \in H_0^1(\Omega) : \quad b(u, v) = \langle f, v \rangle \quad \forall v \in H_0^1(\Omega) \\ \text{with } \quad b(u, v) := \int_{\Omega} \varepsilon \nabla u \nabla v + c u v \quad \langle f, v \rangle := \int_{\Omega} f v \quad . \end{array} \right\} \quad (2)$$

Let  $\|\cdot\|_{\omega}$  denote the  $L^2(\omega)$  norm. For the whole domain  $\Omega$  the subscript will be omitted. The natural energy norm related to (2) is

$$\|v\|_{\omega}^2 := \varepsilon \|\nabla v\|_{\omega}^2 + c \|v\|_{\omega}^2.$$

The bilinear form  $b(\cdot, \cdot)$  is elliptic and coercive with respect to the energy norm. The Lax-Milgram lemma then ensures existence and uniqueness of the weak solution  $u$  of (2).

Although the energy norm is weaker than e.g. the maximum norm, it can be appropriate to design adaptive algorithms (depending of course on the aim of the numerical solution process). The numerical examples of [Kun02] impressively show how an energy norm based error estimator drives the adaptive solution process towards (quasi) optimal meshes, cf. also [PV00].

Finally we remark that reaction diffusion problems (or more general convection diffusion problems) have been investigated to a large extend. In particular there exist *a posteriori* error estimators on anisotropic meshes for the *finite element method* [Kun01b, Kun01c, FPZ01, Kun03]. Some of their techniques can be employed here in modified forms.

## 3 Finite Volume discretization

### 3.1 Notation

For some domain  $\omega$  let  $|\omega| = \text{meas}(\omega)$  be its measure. Define  $\mathbb{P}^k(\omega)$  to be the space of polynomials of degree at most  $k$  on  $\omega$ . We use the shorthand notation  $x \lesssim y$  and  $x \sim y$  for  $x \leq c_1 y$  and  $c_1 x \leq y \leq c_2 x$ , respectively, with positive constants independent of  $x, y$  and  $\varepsilon, \mathcal{T}$ .

Let us start with the two dimensional case. Consider an admissible triangulation of  $\Omega$  into triangles  $T$  which form the *primal mesh*  $\mathcal{T} := \{T\}$ . The vertices (or nodes) of the triangles are denoted by  $x$  or  $x_i$ . Let  $E$  be a triangle edge, and define the set of interior edges by  $\mathcal{E} := \{E\}$ . The thick lines of Figure 2 illustrate this.

In order to construct the dual mesh  $\mathcal{T}^*$ , connect for all triangles  $T$  the barycentre of  $T$  with the midpoints of its edges. In conjunction with the boundary  $\partial\Omega$  this forms *control volumes*  $V_i$  around each vertex  $x_i$ , cf. Figure 2. The *dual mesh* is then  $\mathcal{T}^* := \{V\}$ . The set of the boundaries of the control volumes is denoted by  $\mathcal{E}^* := \{\partial V\}$ .

The (non-empty) intersection of a triangle  $T$  and a control volume  $V$  forms a *quadrilateral* denoted by  $Q$ . This quadrilateral  $Q$  can be splitted by drawing the line from the vertex  $x$  to the barycentre of  $T$ . This yields two small *sub-triangles* denoted by  $K$ . Finally

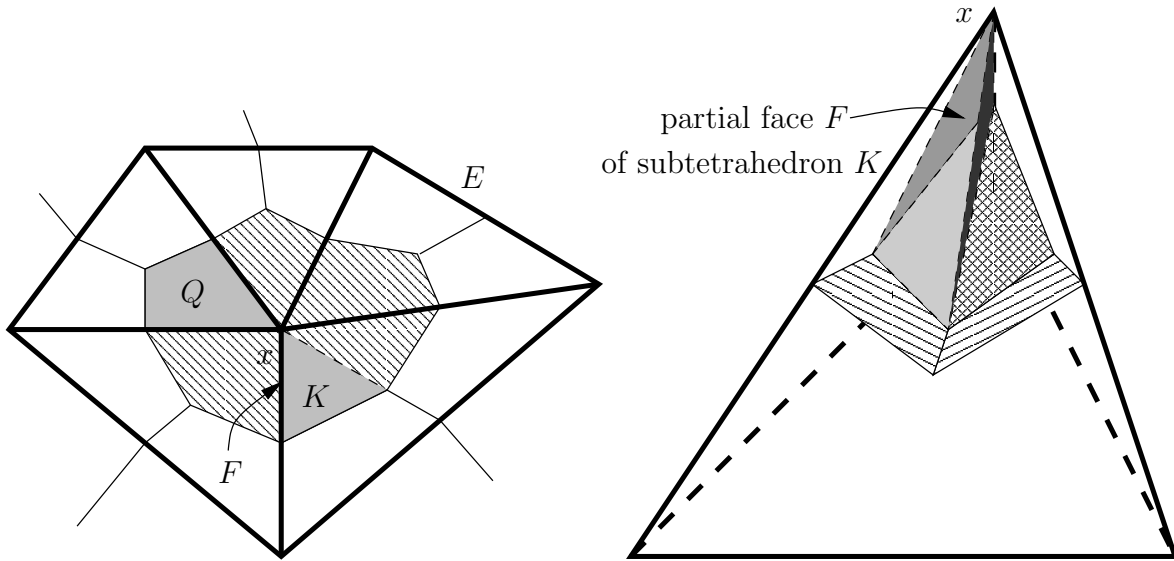


Figure 2: Primal and dual mesh: 2D case (left) and 3D case (right)

define a *partial edge*  $F$  to be the intersection of the boundary of  $K$  with the corresponding primal edge  $E \in \mathcal{E}$ . Figure 2 illustrates this.

In three spatial dimensions the *primal mesh*  $\mathcal{T}$  consists of tetrahedra  $T$ , and  $\mathcal{E}$  becomes the set of interior *faces*  $E$ .

The *control volumes*  $V$  are obtained by connecting the barycentre of a tetrahedron  $T$  with the barycentre of its faces and edges, cf. Figure 2. Set again  $\mathcal{E}^* := \{\partial V\}$ . The intersection of a tetrahedron  $T$  and a control volume  $V$  is a *hexahedron*  $Q$  which is depicted in the top part of the tetrahedron in Figure 2. The six *sub-tetrahedra*  $K \subset Q$  are each formed by the vertex  $x$ , the barycentre of  $T$  and appropriate face and edge barycentres, cf. the figure again. Define similarly the *partial face*  $F$  to be the intersection of  $\partial K$  and a primal face  $E \in \mathcal{E}$ .

From now on, most of the exposition describes the more challenging 3D case. The 2D case is treated only where it differs significantly.

### 3.2 Discretization

For some domain let  $\partial_n$  denote the directional derivative with respect to the outer unit normal vector. Next, consider an arbitrary control volume  $V \in \mathcal{T}^*$ . Green's formula gives the identity

$$\int_V cu + \int_{\partial V} -\varepsilon \partial_n u = \int_V f \quad \forall V \in \mathcal{T}^*. \quad (3)$$

The discrete approximation of the solution is sought in the space  $X_{0h}$  of continuous,  $\mathcal{T}$ -piecewise linear functions that satisfy homogeneous boundary conditions,

$$X_{0h} := \{v \in H_0^1(\Omega) : v|_T \in \mathbb{P}^1(T) \text{ for all } T \in \mathcal{T}\} \subset X = H_0^1(\Omega).$$

We employ a basis  $\{\varphi_i\}$  that consists of the usual hat functions, i.e.  $\varphi_i(x_j) = \delta_{i,j}$  for all interior nodes  $x_j$ .

The discrete solution  $u_h \in X_{0h}$  is determined by the discrete analogue to (3); namely the condition

$$\int_V cu_h + \int_{\partial V} -\varepsilon \partial_n u_h = \int_V f \quad (4)$$

has to be satisfied for all *interior* control volumes  $V$  (i.e.  $\partial V \cap \partial\Omega = \emptyset$ ).

For the practical implementation we express  $u_h$  uniquely as a sum over all interior nodes,

$$u_h = \sum_i u_i \varphi_i \in X_{0h}.$$

This leads to a system of equation for the unknown coefficients  $u_i \in \mathbb{R}$ . The next section investigates the matrix and discusses existence and uniqueness of the discrete solution  $u_h$ .

### 3.3 Existence and uniqueness of a discrete solution

The above discretization results in a matrix  $A$  with entries

$$a_{i,j} = -\varepsilon \int_{\partial V_i} \partial_n \varphi_j + c \int_{V_i} \varphi_j \quad \forall i \text{ with } V_i \in \mathcal{T}^*, \forall j$$

which is splitted into  $a_{i,j} = a_{i,j}^\varepsilon + a_{i,j}^c$ , with  $a_{i,j}^\varepsilon := -\varepsilon \int_{\partial V_i} \partial_n \varphi_j$  and  $a_{i,j}^c := c \int_{V_i} \varphi_j$ . In general  $A$  is not an M-matrix since its off diagonal entries can be positive. Hence one resorts to show certain equivalences to (positive definite) matrices that arise from a finite element discretization, cf. also [BR87, Bey97]. Therefore recall the standard finite element bilinear form  $b(u, v) := \int_\Omega \varepsilon \nabla u \nabla v + cuv$  introduced in (2). Together with the standard basis  $\{\varphi_j\}$  of linear finite elements this gives rise to the finite element stiffness matrix  $B$  with entries

$$b_{i,j} = \int_\Omega \varepsilon \nabla \varphi_i \nabla \varphi_j + c \varphi_i \varphi_j.$$

Again a splitting  $b_{i,j} = b_{i,j}^\varepsilon + b_{i,j}^c$  is employed with  $b_{i,j}^\varepsilon := \varepsilon \int_\Omega \nabla \varphi_i \nabla \varphi_j$  and  $b_{i,j}^c := c \int_\Omega \varphi_i \varphi_j$ .

An easy calculation confirms that the diffusion related FVM matrix part  $A^\varepsilon := (a_{i,j}^\varepsilon)_{i,j}$  coincides with the corresponding FEM matrix part  $B^\varepsilon := (b_{i,j}^\varepsilon)_{i,j}$ , cf. also [BR87], [Bey97, Corollary 4.2.20]. Since the FEM matrix  $B^\varepsilon$  is known to be positively definite, the same holds true for the FVM matrix  $A^\varepsilon \equiv B^\varepsilon$ .

The reaction related parts of both discretization matrices are different but closely linked. In the *two dimensional* case one obtains

$$\begin{aligned} \text{FVM part:} \quad a_{i,j}^c &= c \sum_{T \subset \text{supp} \varphi_i \cap \text{supp} \varphi_j} |T| \cdot \begin{cases} 22/108 & \text{for } i = j, \\ 7/108 & \text{for } i \neq j, \end{cases} \\ \text{FEM part:} \quad b_{i,j}^c &= c \sum_{T \subset \text{supp} \varphi_i \cap \text{supp} \varphi_j} |T| \cdot \begin{cases} 2/12 & \text{for } i = j, \\ 1/12 & \text{for } i \neq j. \end{cases} \end{aligned}$$

With the temporary notation  $A^c := (a_{i,j}^c)_{i,j}$  and  $B^c := (B_{i,j}^c)_{i,j}$  one obtains

$$A^c = \frac{7}{9}B^c + \frac{4}{9}\text{diag}(B^c).$$

Assume first  $c > 0$ . Again from standard FEM theory one knows  $B^c$  to be positively definite; the same holds for  $\text{diag}(B^c)$  because all matrix entries are positive. Hence  $A^c$  is positively definite, and with the previous result for  $A^\varepsilon$  the positive definiteness extends to the FVM matrix  $A = A^\varepsilon + A^c$ . For  $c = 0$  one concludes  $A = A^\varepsilon = B^\varepsilon$  and the positive definiteness of  $A$  directly follows from the one of  $B^\varepsilon$ . In both cases, this ensures unique solvability of the discrete system and existence and uniqueness of the discrete solution  $u_h$ .

In *three spatial dimension* the precise formulas for  $A^c$  and  $B^c$  change to

$$\begin{aligned} \text{FVM part:} \quad a_{i,j}^c &= c \sum_{T \subset \text{supp}\varphi_i \cap \text{supp}\varphi_j} |T| \cdot \begin{cases} 75/576 & \text{for } i = j, \\ 23/576 & \text{for } i \neq j, \end{cases} \\ \text{FEM part:} \quad b_{i,j}^c &= c \sum_{T \subset \text{supp}\varphi_i \cap \text{supp}\varphi_j} |T| \cdot \begin{cases} 2/20 & \text{for } i = j, \\ 1/20 & \text{for } i \neq j, \end{cases} \end{aligned}$$

which implies

$$A^c = \frac{115}{144}B^c + \frac{145}{288}\text{diag}(B^c).$$

The remainder of the exposition is exactly the same.

## 4 Anisotropic meshes

The introduction of Section 1 illustrates that anisotropic discretizations can be very advantageous or, in certain situations, be even mandatory. More information and arguments can be found [Ape99].

The previous exposition is independent of the shape of the elements. In this section we now introduce and describe anisotropic elements in detail, present their notation, basic properties, and some weak mesh assumptions.

Start with an arbitrary (anisotropic) tetrahedron  $T \in \mathcal{T}$ . Enumerate its vertices such that  $P_0P_1$  is the longest edge,  $\text{meas}_2(\triangle P_0P_1P_2) \geq \text{meas}_2(\triangle P_0P_1P_3)$ , and  $\text{meas}_1(P_1P_2) \geq \text{meas}_1(P_0P_2)$ . Further, introduce three orthogonal vectors  $\mathbf{p}_{i,T}$  of length

$$h_{i,T} := |\mathbf{p}_{i,T}|,$$

cf. Figure 3.

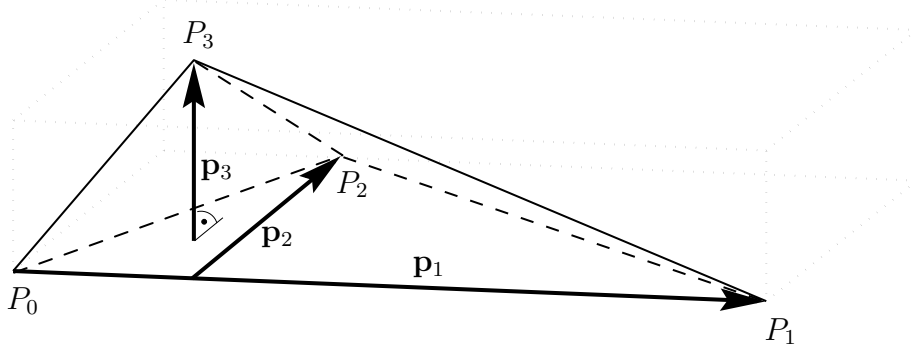
The minimal element size is particularly important; thus define

$$h_{\min,T} := h_{3,T}.$$

The three main anisotropic directions  $\mathbf{p}_{i,T}$  play an important role in several proofs. They are comprised in the orthogonal matrix

$$C_T := (\mathbf{p}_1, \mathbf{p}_2, \mathbf{p}_3) \in \mathbb{R}^{3 \times 3}.$$

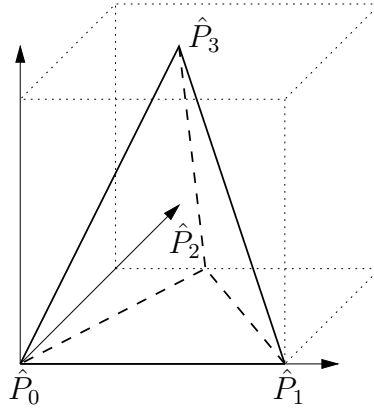


Figure 3: Notation of tetrahedron  $T$ 

Then  $C_T$  can be considered as a transformation matrix which defines implicitly the *reference element*  $\hat{T}$  via

$$\hat{T} := C_T^{-1}(T - \vec{P}_0),$$

cf. Figure 4. In order to facilitate the understanding of this mapping, the circumscribing box of  $T$  has been drawn in Figure 3. This box is mapped onto the unit cube given in Figure 4. Note in particular that the reference element  $\hat{T}$  is of size  $\mathcal{O}(1)$ .

Figure 4: Reference tetrahedron  $\hat{T}$ 

In 2D the notation is similar. The enumeration of the triangle  $T$  is as in the bottom triangle  $P_0P_1P_2$  of Figure 3. Furthermore  $h_{min,T} := h_{2,T}$ , and  $C_T$  becomes a  $2 \times 2$  matrix.

For a sub-element  $K$  (cf. Section 3.1) introduce analogously the anisotropic directions  $\mathbf{p}_{i,K}$ , the matrix  $C_K$ , and the minimal dimension  $h_{min,K}$ .

The next lemma states that the tetrahedron  $T$  and the sub-tetrahedron  $K \subset T$  have similar anisotropic directions  $\mathbf{p}_{i,T}$  and  $\mathbf{p}_{i,K}$ . Mathematically this is expressed as a certain norm equivalence.

**Lemma 4.1 (Equivalence of  $C_T$  and  $C_K$ )** *For any subelement  $K \subset T$ , one has*

$$|C_T^\top y|_{\mathbb{R}^{d \times d}} \sim |C_K^\top y|_{\mathbb{R}^{d \times d}} \quad \forall y \in \mathbb{R}^d.$$

**Proof:** Start with the mapping  $C_K$  which maps the reference element  $\hat{K}$  onto the element  $K$  (up to translation). Similarly  $C_T^{-1}$  maps  $T$  onto  $\hat{T}$ . Since  $K \subset T$  we further conclude  $C_T^{-1} : K \rightarrow C_T^{-1}K \subset C_T^{-1}T \equiv \hat{T}$ . The superposition of both mappings then maps  $C_T^{-1}C_K : \hat{K} \rightarrow C_T^{-1}K \subset \hat{T}$ .

The spectral norm of the corresponding transformation matrix can be bounded by  $|C_T^{-1}C_K|_{\mathbb{R}^{d \times d}} \lesssim \text{diam}(C_T^{-1}K)/\rho(\hat{K})$ , where  $\rho(\hat{K})$  denotes the diameter of the largest inscribed sphere of  $\hat{K}$ .

Both  $\hat{T}$  and  $\hat{K}$  are isotropic reference elements of size  $\mathcal{O}(1)$ . Thus  $\rho(\hat{K}) \sim 1$  and  $\text{diam}(C_T^{-1}K) \leq \text{diam}(C_T^{-1}T) \equiv \text{diam}(\hat{T}) \sim 1$ . With the help of these relations one obtains  $|C_K^\top C_T^{-\top}|_{\mathbb{R}^{d \times d}} = |C_T^{-1}C_K|_{\mathbb{R}^{d \times d}} \lesssim 1$  and  $|C_K^\top y|_{\mathbb{R}^{d \times d}} \leq |C_K^\top C_T^{-\top}|_{\mathbb{R}^{d \times d}} |C_T^\top y|_{\mathbb{R}^{d \times d}} \lesssim |C_T^\top y|_{\mathbb{R}^{d \times d}}$ .

The converse inequality is shown similarly. The technicality there is to prove the relation  $\text{diam}(C_K^{-1}T) \lesssim 1$ . In the 2D case, a magnification of the sub-triangle  $K$  by a factor of 4 contains  $T$  (this corresponds to  $4 \times$ identity operator+translation). Hence  $\text{diam}(C_K^{-1}T) \leq 4 \text{diam}(C_K^{-1}K) \equiv 4 \text{diam}(\hat{K}) \sim 1$ , and the proof finishes as above. In the 3D case the magnification factor is 8; all other arguments are the same.  $\blacksquare$

The primal mesh  $\mathcal{T}$  has to meet the usual conformity conditions, cf. [Cia78, Chapter 2]. Moreover, two further mild conditions have to be satisfied.

**Anisotropic mesh requirements:**

1. The number of tetrahedra containing a node  $x$  is bounded uniformly.
2. The dimensions of adjacent tetrahedra must not change rapidly, i.e.

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

At some places of our exposition it is advantageous to replace the minimal anisotropic dimension  $h_{min,T}$  (which is related to an element  $T$ ) by certain average values. For a control volume  $V \in \mathcal{T}^*$  thus define

$$h_{min,V} := \frac{\sum_{T:T \cap V \neq \emptyset} h_{min,T}}{\sum_{T:T \cap V \neq \emptyset} 1}$$

For the common face  $E$  of two elements  $T_1$  and  $T_2$  set

$$h_{min,E} := \frac{h_{min,T_1} + h_{min,T_2}}{2}.$$

For a partial face  $F$  of two sub-tetrahedra  $K_1, K_2$  define similarly  $h_{min,F} := (h_{min,K_1} + h_{min,K_2})/2$ .

The definitions are modified in the obvious way for boundary faces. In the 2D case the notation is analogous.

Note that the original term  $h_{min,T}$  is of comparable size as the average values since the dimensions of adjacent tetrahedra do not change rapidly, see above. More precisely,

$$\begin{aligned} h_{min,T} &\sim h_{min,V} \sim h_{min,K} \sim h_{min,E} \sim h_{min,F} \\ &\text{for } V \cap T \neq \emptyset, K \subset T, E \subset \partial T, F \subset \partial T. \end{aligned} \quad (5)$$

## 5 Anisotropic interpolation error estimates

In order to obtain an accurate discrete solution  $u_h$ , it is obvious to align the anisotropic tetrahedra  $T$  of the primal mesh according to the anisotropy of the solution. It turns out that this intuitive alignment is also necessary to prove sharp upper error bounds. In particular the proof employs specific interpolation error estimates. However, these interpolation estimates do *not* hold for general meshes; instead the mesh has to have the aforementioned anisotropic alignment with the function to be interpolated.

In order to quantify this alignment, we introduce a so-called *alignment measure*  $m_1(v, \mathcal{T})$  which originates from [Kun00].

**Definition 5.1 (Alignment measure)** *Let  $v \in H^1(\Omega)$ , and  $\mathcal{F} = \{\mathcal{T}\}$  be a family of triangulations of  $\Omega$ . Define the alignment measure  $m_1 : H^1(\Omega) \times \mathcal{F} \mapsto \mathbb{R}$  by*

$$m_1(v, \mathcal{T}) := \left( \sum_{T \in \mathcal{T}} h_{min,T}^{-2} \|C_T^\top \nabla v\|_T^2 \right)^{1/2} / \|\nabla v\|. \quad (6)$$

By definition one has  $m_1(v, \mathcal{T}) \geq 1$ . For arbitrary *isotropic* meshes one obtains that  $m_1(v, \mathcal{T}) \sim 1$ . The same is achieved for *anisotropic* meshes  $\mathcal{T}$  that are *aligned* with the anisotropic function  $v$ . However, if the mesh  $\mathcal{T}$  is *not aligned* then  $m_1$  can be arbitrarily large, cf. [Kun01a, Section 4.2].

In practice, one almost always obtains  $m_1(v, \mathcal{T}) \sim 1$  for ‘sensible’ anisotropic meshes, i.e. the alignment measure is no obstacle for reliable error estimation. Since the focus of our work is on *a posteriori* error estimation, we refer to [Kun00] for more details.

Next we require a global,  $\mathcal{T}^*$ -piecewise constant interpolation operator:  $I_m : L^2(\Omega) \rightarrow L^2(\Omega)$  which is defined by

$$I_m v := \begin{cases} |V|^{-1} \int_V v & \text{for interior volumes } V \\ 0 & \text{for boundary volumes } V \text{ (i.e. } \partial V \cap \Gamma_D \neq \emptyset) \end{cases}$$

Note that  $I_m$  satisfies homogeneous Dirichlet boundary conditions, i.e.  $I_m v = 0$  on  $\Gamma_D$ .

**Lemma 5.2 (Local interpolation error bounds for  $I_m$ )** *Let  $v \in H_0^1(\Omega)$ . Then*

$$\|v - I_m v\|_V \leq \|v\|_V \quad (7)$$

$$\|v - I_m v\|_V^2 \lesssim \sum_{K \subset V} \|C_K^\top \nabla v\|_K^2 \quad (8)$$

$$\sum_{K \subset V} \|C_K^\top \nabla(v - I_m v)\|_K^2 = \sum_{K \subset V} \|C_K^\top \nabla v\|_K^2. \quad (9)$$

**Proof:** The first and last relation are trivial.

The remaining inequality (8) has been proven for an interior control volume  $V$  in [Kun00, Lemma 4]. In order to treat also boundary control volumes  $V$  we utilize the intermediate result  $\| |V|^{-1} \int_V v \|_V^2 \lesssim \sum_{K \subset V} \|C_K^\top \nabla v \|_K^2$  from the proof of [Kun00, Theorem 1]. The triangle inequality  $\|v - I_m v \|_V = \|v \|_V \leq \|v - |V|^{-1} \int_V v \|_V + \| |V|^{-1} \int_V v \|_V$  in conjunction with the previous bound  $\|v - |V|^{-1} \int_V v \|_V \lesssim \sum_{K \subset V} \|C_K^\top \nabla v \|_K^2$  finishes the proof for a boundary control volume  $V$ . Alternatively one can apply the Poincaré inequality and scaling arguments.  $\blacksquare$

**Lemma 5.3** *Let  $v \in H_0^1(\Omega)$ . Then*

$$\|v - I_m v \| \leq \|v \| \quad (10)$$

$$\sum_{V \in \mathcal{T}^*} h_{min,V}^{-2} \|v - I_m v \|_V^2 \lesssim m_1(v, \mathcal{T})^2 \|\nabla v \|^2 \quad (11)$$

$$\sum_{V \in \mathcal{T}^*} h_{min,V}^{-2} \sum_{K \subset V} \|C_K^\top \nabla(v - I_m v) \|_K^2 \lesssim m_1(v, \mathcal{T})^2 \|\nabla v \|^2. \quad (12)$$

**Proof:** The first interpolation error bound is obvious again.

For the second estimate, start with some control volume  $V$  and recall

$$\|v - I_m v \|_V^2 \lesssim \sum_{K \subset V} \|C_K^\top \nabla v \|_K^2.$$

With the help of Lemma 4.1 one obtains  $\|C_K^\top \nabla v \|_K \sim \|C_T^\top \nabla v \|_K \leq \|C_T^\top \nabla v \|_T$  and concludes (11) by

$$\begin{aligned} \sum_{V \in \mathcal{T}^*} h_{min,V}^{-2} \|v - I_m v \|_V^2 &\lesssim \sum_{V \in \mathcal{T}^*} \sum_{T: T \cap V \neq \emptyset} h_{min,V}^{-2} \|C_T^\top \nabla v \|_T^2 \\ &\lesssim \sum_{T \in \mathcal{T}} h_{min,T}^{-2} \|C_T^\top \nabla v \|_T^2 = m_1(v, \mathcal{T})^2 \|\nabla v \|^2. \end{aligned}$$

Here we have also employed  $h_{min,V} \sim h_{min,T}$  for  $T \cap V \neq \emptyset$ .

For the last interpolation bound proceed similarly and derive

$$\begin{aligned} \sum_{V \in \mathcal{T}^*} h_{min,V}^{-2} \sum_{K \subset V} \|C_K^\top \nabla(v - I_m v) \|_K^2 &\lesssim \sum_{V \in \mathcal{T}^*} h_{min,V}^{-2} \sum_{T: T \cap V \neq \emptyset} \|C_T^\top \nabla v \|_T^2 \\ &\lesssim m_1(v, \mathcal{T})^2 \|\nabla v \|^2 \end{aligned}$$

which finishes the proof.  $\blacksquare$

Later on, a certain factor which is closely linked to the actual PDE (1) plays a crucial role to derive error estimates. This term  $\alpha_V$  is defined by

$$\alpha_V := \min\{c^{-1/2}, \varepsilon^{-1/2} h_{min,V}\}. \quad (13)$$

Introduce  $\alpha_E$ ,  $\alpha_K$  and  $\alpha_F$  in an analogous fashion. With their help we derive some specific interpolation estimates which are closely related to the error estimator to be analysed.

**Lemma 5.4 (Interpolation error bounds for  $\mathbf{I}_m$ )** *Let  $v \in H_0^1(\Omega)$*

$$\sum_{V \in \mathcal{T}^*} \alpha_V^{-2} \|v - \mathbf{I}_m v\|_V^2 \lesssim m_1(v, \mathcal{T})^2 \|v\|^2 \quad (14)$$

$$\varepsilon^{1/2} \sum_{E \in \mathcal{E}} \alpha_E^{-1} \|v - \mathbf{I}_m v\|_E^2 \lesssim m_1(v, \mathcal{T})^2 \|v\|^2. \quad (15)$$

**Proof:** The ideas are similar to [Kun01b, Lemma 3.11]. For self-containment, we repeat major steps here. Observe first  $\alpha_V^{-2} = \max\{c, \varepsilon/h_{\min, V}^2\}$ . The previous lemma implies

$$\begin{aligned} \sum_{V \in \mathcal{T}^*} \alpha_V^{-2} \|v - \mathbf{I}_m v\|_V^2 &= \sum_{\substack{V \in \mathcal{T}^* \\ c \geq \varepsilon h_{\min, V}^{-2}}} c \|v - \mathbf{I}_m v\|_V^2 + \sum_{\substack{V \in \mathcal{T}^* \\ c < \varepsilon h_{\min, V}^{-2}}} \varepsilon h_{\min, V}^{-2} \|v - \mathbf{I}_m v\|_V^2 \\ &\leq c \|v - \mathbf{I}_m v\|^2 + \varepsilon \sum_{V \in \mathcal{T}^*} h_{\min, V}^{-2} \|v - \mathbf{I}_m v\|_V^2 \\ &\lesssim \|v\|^2 + \varepsilon m_1(v, \mathcal{T})^2 \|\nabla v\|^2 \leq m_1(v, \mathcal{T})^2 \|v\|^2 \end{aligned}$$

In order to show interpolation bound (15), start with some partial face  $F \subset \partial K$ . The anisotropic trace inequality then reads

$$\frac{|K|}{|F|} \|v - \mathbf{I}_m v\|_F^2 \lesssim \|v - \mathbf{I}_m v\|_K \left( \|v - \mathbf{I}_m v\|_K + \|C_K^\top \nabla(v - \mathbf{I}_m v)\|_K \right).$$

for  $v \in H^1(K)$ , cf. [Kun01b, Lemma 3.5]. Next, switch from the matrix  $C_K$  to  $C_T$  via Lemma 4.1. Recall further  $\alpha_F \sim \alpha_V$  for any partial face  $F \cap V \neq \emptyset$ , see (5). The previous Lemma 5.3 and the geometrical property  $h_{\min, T} \lesssim |K|/|F|$  then lead to the desired bound (15).  $\blacksquare$

## 6 Residual error estimation

The discretization error can be bounded by the residual in the  $H^{-1}(\Omega)$  norm. Since this norm is difficult to evaluate, one commonly tries to approximate it from the data at hand, cf. the standard textbook [Ver96].

To this end introduce the *exact element residual*  $R_{\mathcal{T}}$  in a  $\mathcal{T}$ -piecewise fashion by

$$R_{\mathcal{T}}|_T = f - (-\varepsilon \Delta u_h + c u_h) \quad \forall T \in \mathcal{T}.$$

This exact residual involves the fixed but otherwise arbitrary function  $f$ . In order to derive a lower error bound, this function is commonly replaced by some approximation  $f_{\mathcal{T}}$  from a finite dimensional space. The crucial condition is  $f_{\mathcal{T}}|_K \in \mathbb{P}^1(K)$  for each sub-tetrahedron  $K$ .<sup>2</sup>

<sup>2</sup>A practical choice could be an approximation  $f_{\mathcal{T}}$  that is linear over each element  $T$ , or that is linear over each control volume  $V$ . The actual choice will certainly depend on the numerical implementation and its data structure.

With the help of  $f_{\mathcal{T}}$  the *approximate element residual*  $r_{\mathcal{T}} \in L^2(\Omega)$  is given for each hexahedron  $Q$  by

$$r_{\mathcal{T}}|_Q := f_{\mathcal{T}} - (-\varepsilon \Delta u_h + cu_h) \in \mathbb{P}^1(Q) \quad \text{on each } Q.$$

Next, introduce the face residual  $r_E$  which is basically the  $\varepsilon$  scaled gradient jump. To this end consider a face  $E$  and fix one of the two unit normal vectors  $n_E$ . Define the jump  $[[v]]_E$  of a function  $v \in L^2(\Omega)$  across a face  $E$  (with respect to the orientation of  $n_E$ ) by

$$[[v]]_E(y) := \lim_{t \rightarrow +0} v(y - tn_E) - v(y + tn_E) \quad y \in E.$$

The *face residual*  $r_E$  is given by

$$r_E := \begin{cases} [[\partial_{n_E} u_h]]_E & \text{for } E \in \mathcal{E} \\ 0 & \text{for } E \subset \Gamma_D. \end{cases}$$

Note that  $r_E$  is independent of the orientation of  $n_E$ .

Now we are ready to define the error estimator.

**Definition 6.1 (Error estimator)** *For a control volume  $V$  define the error estimator by*

$$\eta_V^2 := \alpha_V^2 \|r_{\mathcal{T}}\|_V^2 + \varepsilon^{-1/2} \alpha_V \sum_{F: F \cap V \neq \emptyset} \|r_E\|_F^2$$

(recall the notation for the partial face  $F$  from Section 3.1). The approximation term is given by

$$\zeta_V := \alpha_V \|f - f_V\|_V$$

Finally introduce the global terms

$$\eta^2 := \sum_{V \in \mathcal{T}^*} \eta_V^2 \quad \zeta^2 := \sum_{V \in \mathcal{T}^*} \zeta_V^2.$$

**Theorem 6.2 (Residual error estimation)** *Let  $u \in H_0^1(\Omega)$  be the exact solution and  $u_h \in X_{0h}$  be the finite element solution. Then the error is bounded locally from below by*

$$\eta_V \lesssim \| \|u - u_h\| \|_V + \zeta_V \quad (16)$$

for all  $V \in \mathcal{T}^*$ . The error is bounded globally from above by

$$\| \|u - u_h\| \| \lesssim m_1(u - u_h, \mathcal{T}) [\eta^2 + \zeta^2]^{1/2} . \quad (17)$$

**Proof: Upper error bound:** Let  $v \in H_0^1(\Omega)$  and recall the geometrical setting as described in Section 3.1. Integration by parts then yields

$$\begin{aligned}
a(u - u_h, v) &= \int_{\Omega} \varepsilon \nabla(u - u_h) \nabla v + c(u - u_h)v \\
&= \int_{\Omega} \varepsilon \nabla u \nabla v - \sum_{V \in \mathcal{T}^*} \sum_{Q \subset V} \int_Q \varepsilon \nabla u_h \nabla(v - \mathbf{I}_m v) + \int_{\Omega} c(u - u_h)v \\
&= \int_{\Omega} (-\varepsilon \Delta u + cu - cu_h)v + \int_{\partial\Omega} \varepsilon \partial_n uv \\
&\quad - \sum_{V \in \mathcal{T}^*} \sum_{Q \subset V} \left( \int_Q -\varepsilon \Delta u_h(v - \mathbf{I}_m v) + \int_{\partial Q} \varepsilon \partial_n u_h(v - \mathbf{I}_m v) \right) \\
&= \int_{\Omega} R_{\mathcal{T}}(v - \mathbf{I}_m v) + \int_{\Omega} (f - cu_h) \mathbf{I}_m v \\
&\quad - \sum_{V \in \mathcal{T}^*} \sum_{Q \subset V} \left( \int_{\partial Q \cap \mathcal{E}} \varepsilon \partial_n u_h(v - \mathbf{I}_m v) + \int_{\partial Q \cap \mathcal{E}^*} \varepsilon \partial_n u_h(v - \mathbf{I}_m v) \right).
\end{aligned}$$

Figure 5 illustrates the last two boundary integrals (2D case; same domain as in Figure 2).

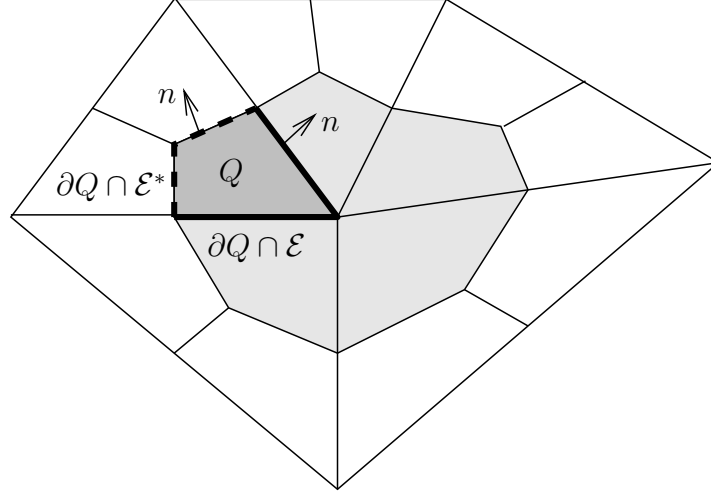


Figure 5: Boundary parts  $\partial Q \cap \mathcal{E}$  (thick solid line) and  $\partial Q \cap \mathcal{E}^*$  (thick dashed line)

For interior volumes  $V$  the finite volume discretization (4) implies

$$\int_V (f - cu_h) \mathbf{I}_m v + \int_{\partial V} \varepsilon \partial_n u_h \mathbf{I}_m v = 0 \quad \forall v \in L^2(\Omega)$$

since  $\mathbf{I}_m v$  is constant on  $V$ . For boundary volumes the same identity holds since  $\mathbf{I}_m v = 0$  there by definition.

Furthermore, since  $v \in H^1(\Omega)$ , the trace of  $v$  on a face  $E^* \in \mathcal{E}^*$  is in  $L^2(E^*)$ . Furthermore,  $\partial_{n_{E^*}} u_h$  is continuous across  $E^* \in \mathcal{E}^*$ , and thus one can simplify<sup>3</sup>

$$\sum_{V \in \mathcal{T}^*} \sum_{Q \subset V} \int_{\partial Q \cap \mathcal{E}^*} \varepsilon \partial_n u_h v = \sum_{E^* \in \mathcal{E}^*} \int_{E^*} \varepsilon \partial_{n_{E^*}} u_h \llbracket v \rrbracket_{E^*} = 0.$$

Hence we end up with

$$\begin{aligned} a(u - u_h, v) &= \int_{\Omega} R_{\mathcal{T}}(v - \mathbf{I}_m v) - \sum_{E \in \mathcal{E}} \int_E \varepsilon \llbracket \partial_{n_E} u_h \rrbracket_E (v - \mathbf{I}_m v) \\ &\leq \left( \sum_{V \in \mathcal{T}^*} \alpha_V^2 \|R_{\mathcal{T}}\|_V^2 \right)^{1/2} \left( \sum_{V \in \mathcal{T}^*} \alpha_V^{-2} \|v - \mathbf{I}_m v\|_V^2 \right)^{1/2} \\ &\quad + \left( \sum_{E \in \mathcal{E}} \varepsilon^{-1/2} \alpha_E \|r_E\|_E^2 \right)^{1/2} \left( \sum_{E \in \mathcal{E}} \varepsilon^{1/2} \alpha_E^{-1} \|v - \mathbf{I}_m v\|_E^2 \right)^{1/2}. \end{aligned}$$

Now the interpolation error estimates of Lemma 5.4 are applied to bound  $v - \mathbf{I}_m v$ . One concludes

$$a(u - u_h, v) \lesssim \left( \sum_{V \in \mathcal{T}^*} \alpha_V^2 \|R_{\mathcal{T}}\|_V^2 + \sum_{E \in \mathcal{E}} \varepsilon^{-1/2} \alpha_E \|r_E\|_E^2 \right)^{1/2} m_1(v, \mathcal{T}) \|v\|.$$

Next, insert  $v = u - u_h$  and observe  $a(u - u_h, u - u_h) = \|v\|^2$ . A simple triangle inequality provides the change from  $R_{\mathcal{T}}$  to  $r_{\mathcal{T}}$  and finishes the proof of the upper error bound (17).

**Lower error bound:** This part of the proof is analogous to [Kun01b, Theorem 4.3] for a finite element method. As it is commonly the case, the lower error bounds is not specific to a certain discretization but relies heavily on the definition of the residuals. Therefore we do not repeat the proof itself but present the two main auxiliary inequalities instead.

For an triangle/tetrahedron  $K \subset V$  one concludes

$$\alpha_K \|r_{\mathcal{T}}\|_K \lesssim \|u - u_h\|_K + \alpha_K \|f - f_V\|_K,$$

cf. [Kun01b, Theorem 4.3]. Consider next a partial face  $F$  of two sub-tetrahedra  $K_1$  and  $K_2$ . The corresponding bound then reads

$$\varepsilon^{-1/2} \alpha_F \|r_E\|_F^2 \lesssim \|u - u_h\|_{K_1 \cup K_2}^2 + \alpha_F^2 \|f - f_V\|_{K_1 \cup K_2}^2.$$

Equivalence (5) states that the smallest dimensions do not change rapidly. This means  $h_{\min, V} \sim h_{\min, F}$  and subsequently  $\alpha_V \sim \alpha_F$  for  $F \subset V$ .

Combining both inequality readily provides the lower error bound (16). ■

---

<sup>3</sup>Here  $n_{E^*}$  is one of the two unit normal vectors for a face  $E^* \in \mathcal{E}^*$ . Then  $\partial_{n_{E^*}} u_h \llbracket v \rrbracket_{E^*}$  is independent of the orientation of  $n_{E^*}$ .



## 7 Summary

We have proposed and rigorously analysed an *a posteriori* error estimator for the *finite volume method* on *anisotropic meshes*. While such estimators are known for the finite element method, no corresponding result has been available for the finite volume method.

The error estimation has been shown to be robust with respect to the small perturbation parameter  $\varepsilon$ . Robustness with respect to the mesh anisotropy has been also achieved for well aligned anisotropic meshes (which is frequently met in practice). Thus our novel error estimator is reliable and efficient.

Numerical experiments will be carried in a subsequent work.

## References

- [Ang95] L. Angermann. Balanced a-posteriori error estimates for finite volume type discretizations of convection-dominated elliptic problems. *Computing*, 55(4):305–323, 1995.
- [AO00] M. Ainsworth and J.T. Oden. *A posteriori error estimation in finite element analysis*. Wiley, 2000.
- [Ape99] Th. Apel. *Anisotropic finite elements: Local estimates and applications*. Advances in Numerical Mathematics. Teubner, Stuttgart, 1999.
- [Bey97] J. Bey. *Finite-Volumen- und Mehrgitterverfahren für elliptische Randwertprobleme*. PhD thesis, Universität Tübingen, 1997. also published in *Advances in Numerical Mathematics* series, Teubner, Stuttgart, 1998.
- [BM00] A. Bergam and Z. Mghazli. Estimateurs a posteriori d’un schéma de volumes finis pour un problème non linéaire. *C. R. Acad. Sci.*, 331(6):475–478, 2000.
- [BMV03] A. Bergam, Z. Mghazli, and R. Verfürth. Estimations a posteriori d’un schéma de volumes finis pour un problème non linéaire. *to appear in Numer. Math.*, 2003.
- [BR87] R.E. Bank and D. Rose. Some error estimates for the box method. *SIAM J. Numer. Anal.*, 24:777–787, 1987.
- [Cia78] P. G. Ciarlet. *The finite element method for elliptic problems*. North-Holland, Amsterdam, 1978.
- [FPZ01] L. Formaggia, S. Perotto, and P. Zunino. An anisotropic a-posteriori error estimate for a convection-diffusion problem. *Comput. Vis. Sci.*, 4(2):99–104, 2001.

- [KO00] D. Kröner and M. Ohlberger. A posteriori error estimates for upwind finite volume schemes for nonlinear conservation laws in multi dimensions. *Math. Comput.*, 69(229):25–39, 2000.
- [Kun00] G. Kunert. An a posteriori residual error estimator for the finite element method on anisotropic tetrahedral meshes. *Numer. Math.*, 86(3):471–490, 2000. DOI 10.1007/s002110000170.
- [Kun01a] G. Kunert. A local problem error estimator for anisotropic tetrahedral finite element meshes. *SIAM J. Numer. Anal.*, 39(2):668–689, 2001.
- [Kun01b] G. Kunert. Robust a posteriori error estimation for a singularly perturbed reaction–diffusion equation on anisotropic tetrahedral meshes. *Adv. Comp. Math.*, 15(1–4):237–259, 2001.
- [Kun01c] G. Kunert. Robust local problem error estimation for a singularly perturbed problem on anisotropic finite element meshes. *Math. Model. Numer. Anal.*, 35(6):1079–1109, 2001.
- [Kun02] G. Kunert. A note on the energy norm for a singularly perturbed model problem. *Computing*, 69(3):265–272, 2002.
- [Kun03] G. Kunert. A posteriori error estimation for convection dominated problems on anisotropic meshes. *Math. Methods Appl. Sci.*, 26(7):589–617, 2003.
- [LT02] R. Lazarov and S. Tomov. A posteriori error estimates for finite volume element approximations of convection-diffusion-reaction equations. *Comput. Geosci.*, 6(3-4):483–503, 2002.
- [MMM93] J.A. Mackenzie, D.F. Mayers, and A.J. Mayfield. Error estimates and mesh adaption for a cell vertex finite volume scheme. *Notes Numer. Fluid Mech.*, 44:290–310, 1993.
- [Mor96] K. W. Morton. *Numerical solution of convection-diffusion problems*. Chapman & Hall, London, 1996.
- [MOS96] J.J.H. Miller, E. O’Riordan, and G.I. Shishkin. *Fitted numerical methods for singularly perturbed problems. Error Estimates in the maximum norm for linear problems in one and two dimensions*. World Scientific Publications, Singapore, 1996.
- [Ohl01a] M. Ohlberger. A posteriori error estimate for finite volume approximations to singularly perturbed nonlinear convection-diffusion equations. *Numer. Math.*, 87(4):737–761, 2001.

- [Ohl01b] M. Ohlberger. A posteriori error estimates for vertex centered finite volume approximations of convection-diffusion-reaction equations. *Math. Model. Numer. Anal.*, 35(2):355–387, 2001.
- [PV00] A. Papastavrou and R. Verfürth. A posteriori error estimators for stationary convection–diffusion problems: A computational comparison. *Comput. Methods Appl. Mech. Eng.*, 189(2):449–462, 2000.
- [RST96] H.-G. Roos, M. Stynes, and L. Tobiska. *Numerical methods for singularly perturbed differential equations. Convection-diffusion and flow problems*. Springer, Berlin, 1996.
- [Ver96] R. Verfürth. *A review of a posteriori error estimation and adaptive mesh-refinement techniques*. Wiley-Teubner, Chichester; Stuttgart, 1996.