

Scuola Internazionale Superiore di Studi Avanzati

---

Group “Mathematical Analysis, Modeling and Applications”

Final dissertation for Galileian School of Higher Education

**A-POSTERIORI ERROR ANALYSIS FOR THE DISCONTINUOUS  
GALERKIN METHOD FOR A GENERAL  
CONVECTION-REACTION PROBLEM**

Supervisor:  
Prof. Andrea Cangiani

Graduating student:  
Alessandro Vici

Serial number: 1204186

---

Academic year 2023/2024  
20<sup>th</sup> November 2024



# Contents

<b>Introduction</b>	<b>i</b>
<b>1 A posteriori error analysis for the Discontinuous Galerkin method in a simple diffusion problem</b>	<b>1</b>
1.1 The model problem . . . . .	1
1.2 The approximating structure . . . . .	1
1.2.1 The mesh . . . . .	1
1.2.2 Euristicis . . . . .	3
1.2.3 The energy space $E_h$ . . . . .	3
1.2.4 The space of approximating functions . . . . .	4
1.3 Formulation of the discrete problem . . . . .	5
1.3.1 Error estimates for problem 1.3.1 . . . . .	6
1.4 A posteriori error estimates . . . . .	6
<b>2 Discontinuous Galerkin method for a general convection-reaction problem, a posteriori error analysis</b>	<b>9</b>
2.1 Set up . . . . .	9
2.2 The mesh . . . . .	11
2.3 The broken $H^1$ space . . . . .	12
2.4 The space of approximating maps . . . . .	12
2.5 The discrete problem . . . . .	13
2.6 A new norm and seminorm . . . . .	14
2.7 A robust a-posteriori error estimator . . . . .	15
2.8 Proof of Theorem 2.7.2 . . . . .	17
2.8.1 Defining some auxiliary operators . . . . .	17
2.8.2 Properties of the auxiliary operators . . . . .	18
2.8.3 Interpolation operator . . . . .	21
2.8.4 Conclusion of the proof . . . . .	21



# Introduction

The study of most systems which are governed by Partial Differential Equations is too complicated. That means that it may be close to impossible to find the exact analytical solution to a given problem. The strategy which is employed then is to leave the exact solution unknown and to attack the problem numerically, that is to devise a method which takes as inputs the parameters of the Partial Differential Equation and which is able to return as output a computable solution which shall approximate the exact one.

There is a multitude of numerical methods which have been devised, such as the finite difference method, the finite element method, the finite volume method, and many others. Each method has its pros and cons, it might be suitable for certain types of problems but not working for other types of problems.

Nevertheless there are some theoretical aspects which every method needs to have:

- it must be well posed, i.e. it must have one and only one solution.
- its solution must be proven to converge to the exact solution of the original differential problem.
- it shall be accompanied by a theorem about its rate of convergence or by an estimate on its error (the criterion to measure the error is not unique, for every different method one could devise a different suitable criterion).

Denote by  $(V, \|\cdot\|)$  a normed space (in practise  $V$  is usually  $C^2$  or  $H^1$ ) in which the exact solution  $u$  to our problem is supposed to exist. Denote by  $V_h$  a finite dimensional subspace of  $V$  (e.g. the space of functions which are polynomial or piecewise polynomial with a fixed maximum degree) in which we wish to find a good approximation of  $u$ . Denote by  $u_h$  the function inside the subspace  $V_h$  which is returned by our numerical method. Then most estimates on the error are of the type

$$\|u - u_h\| \leq Ch\|u\|,$$

where  $h$  is a parameter which depends on the choice of the subspace  $V_h$  and  $C > 0$  is a positive constant which is independent of  $h$ .

One concrete example of numerical method is the finite element method, used to deal with a wide variety of differential problems. One such problem would be the

homogeneous Poisson problem on an open and bounded set  $\Omega \in \mathbb{R}^2$ : find  $u \in C^2(\Omega) \cap C^0(\overline{\Omega})$  such that

$$\begin{cases} \Delta u = f & \text{in } \Omega \\ u = 0 & \text{on } \partial\Omega \end{cases}$$

In the finite element method one partitions  $\Omega$  into small subsets (usually triangles), which we shall denote by  $K_1, \dots, K_m$ , fixes a degree of approximation  $p$ , and then looks for an approximating map  $u_h$  in the set of maps which are overall continuous and polynomial on each sets  $K_i$  with degree at most  $p$  (i.e. one look for an approximating map  $u_h \in C^0(\overline{\Omega})$  such that  $u_h|_{K_i}$  is a polynomial of degree at most  $p$  for all  $i$ ).

Denote by

$$h_i := \text{diam}(K_i) \quad \forall i = 1, \dots, m$$

the diameters of all subsets  $K_i$ , and define

$$h := \max\{h_1, \dots, h_m\}.$$

A typical error estimate for the finite element method is of the type

$$\|u - u_h\|_{H^1(\Omega)} \leq Ch^{p-1} \|u\|_{H^1(\Omega)}.$$

This estimate shows indeed the convergence of the method, but presents two inconveniences:

1. The solution  $u$  is unknown, so we are not able to compute an actual estimate on our error;
2. The estimate is usually found after deducing an estimate of the type

$$\|u - u_h\|_{H^1(\Omega)} \leq C \sum_{i=1}^m h_{K_i} \|u\|_{H^1(\Omega)},$$

and then using  $h = \max\{h_1, \dots, h_m\}$  in place of the summation. This means that the final estimate forgets about the local contributions for the error.

The information about local contributions to the error is therefore lost for these two reasons. If we had it though it could be exploited to refine the partition of  $\Omega$  in an optimal way and increase the rate of convergence. This is impossible when we have estimates which are expressed in terms of  $\|u\|$  (all these sort of estimates are called *a-priori error estimates*).

To move toward the study of local contributions on the error, other types of estimates have been found, and they are estimates of the type

$$\|u - u_h\|_{H^1(\Omega)} \leq C \sum_{i=1}^m h_{K_i} \|u_h\|_{H^1(K_i)}.$$

These estimates are called *a-posteriori error estimates* and are valuable precisely because they give information about local error contributions and are indeed computable.

First ideas and results of a-posteriori error analysis are present in [3] and [2]. There, the authors introduce the a-posteriori error analysis for the finite element method.

The interested reader may see [11] for a wide discussion on the variety of methods which arose from the aforementioned [3] and [2].

This thesis will be presenting the Discontinuous Galerkin method. It is similar to the finite element method, as it also partitions the domain into polygons and approximates the exact solution with maps which are piecewise polynomial, yet it differs on a substantial feature: in the Discontinuous Galerkin method the approximating maps are not required to be continuous, and are instead allowed to have discontinuities along the boundary of the polygons  $K_i$ . This has two immediate consequences:

1. There is the need to devise a way to quantify the jump which occurs at the discontinuity points;
2. The method is less restrictive on the choice of the partition. While in the finite element method (for a series of reasons) one can only partition  $\Omega$  into triangles and parallelograms, with the Discontinuous Galerkin method it is potentially possible to employ any type of polygon.

After describing the method the focus will be put on the description and analysis of an a-posteriori error estimator.

More specifically, in chapter 1 we will be introducing the discontinuous Galerkin method and its a-posteriori error analysis in a simple case: the homogeneous Poisson problem. We will adopt and follow the approach of Ohannes A. Karakashian and Frederic Pascal in [8]. Given its introductory purpose, in this chapter we will omit all proofs, which (unless differently specified) can be found in [8].

In chapter 2 we will be following along the lines of Dominik Schötzau and Liang Zhu in [10] and describe their Discontinuous Galerkin approach to a general convection-reaction problem, i.e. a problem of the type

$$a \cdot \nabla u + bu = f.$$

In their work they derive a robust a-posteriori error estimator. Here by “robust” we mean that the possibility of having a (small) diffusion term is taken into account. Formally, this means that the considered problem will be of the type

$$-\varepsilon \Delta u + a \cdot \nabla u + bu = f, \quad \text{with } \varepsilon \ll 1,$$

but it does not compromise the reliability of the estimator.

See [3] and [2] for an introduction on the A-posteriori error analysis for the finite element method and for adaptive finite element method.

See [11] for a wide discussion on the variety of methods which arose from the aforementioned [3] and [2].



# Chapter 1

## A posteriori error analysis for the Discontinuous Galerkin method in a simple diffusion problem

### 1.1 The model problem

We consider as model problem for a first analysis the homogeneous Poisson problem:

#### 1.1.1. Definition, Poisson problem

Given an open and bounded set  $\Omega \subset \mathbb{R}^n$  and  $f \in C^0(\Omega)$ , the homogeneous Poisson problem on  $\Omega$  associated to  $f$  is the problem

Find  $u \in C^2(\Omega) \cap C^0(\bar{\Omega})$  such that

$$\begin{aligned} -\Delta u &= f & \text{in } \Omega \\ u &= 0 & \text{on } \partial\Omega \end{aligned}$$

The Poisson problem admits a weak formulation, which is

#### 1.1.2. Definition, weak Poisson problem

Given an open set  $\Omega \subset \mathbb{R}^n$  and  $f \in L^2(\Omega)$ , the weak homogeneous Poisson problem on  $\Omega$  associated to  $f$  is the problem

Find  $u \in H_0^1(\Omega)$  such that

$$\int_{\Omega} \nabla u \cdot \nabla v \, d\mathcal{L}^n = \int_{\Omega} f v \quad \forall v \in H^1(\Omega).$$

### 1.2 The approximating structure

#### 1.2.1 The mesh

##### 1.2.1. Definition

We fix now some terminology and notation for the upcoming analysis:

**Mesh** : in order to approximate the exact solution to a problem, we will partition the domain  $\Omega$  into a finite number of subsets (which will usually be disjoint, except for their boundaries) and then approximate locally on each one. The *mesh* is the set  $\mathcal{T}_h$  of all subsets (which are usually chosen to be polyhedra) in which the problem's domain  $\Omega$  is fragmented. The value  $h$  appearing as pedice in the symbol  $\mathcal{T}_h$  is the quantity  $\max\{\text{diam}(K) : K \in \mathcal{T}_h\}$ .

**Set of meshes** : Depending on  $h$ , the approximation will result to be more or less precise. In order to study the convergence of the approximating solutions to the exact one, we will consider a fixed set  $\{\mathcal{T}_h\}_{h \in (0,1)}$  of partitions.

**Element of the mesh** : we will call *element* of the mesh each single polihedron (or subset)  $K$  such that  $K \in \mathcal{T}_h$ .

**Subelement of the mesh** : we will call *subelement* of a mesh any face, edge or vertex of any polihedron of the mesh (when the mesh is made of more general subsets the definition of subelement becomes slightly more delicate, but shall be clear from the context or shape of the elements themselves).

**Adjacent elements** : we say that  $K_1, K_2 \in \mathcal{T}_h$  are *adjacent* if<sup>(1)</sup>  $\mathcal{H}^{n-1}(K_1 \cap K_2) > 0$

**Conforming mesh** : it is a mesh in which the intersection of any two elements is either a subelement or the empty set.

### 1.2.2. Definition, internal and boundary edges

We shall define the set  $\mathcal{E}^I$  of internal edges and the set  $\mathcal{E}^B$  of boundary edges as

$$\begin{aligned}\mathcal{E}^I &:= \{\ell = \partial K \cap \partial K' : K, K' \in \mathcal{T}_h, K \neq K', \mathcal{H}^{n-1}(\ell) > 0\}, \\ \mathcal{E}^B &:= \{\ell = \partial K \cap \partial \Omega : K \in \mathcal{T}_h, \mathcal{H}^{n-1}(\ell) > 0\}.\end{aligned}$$

where  $\mathcal{H}^{n-1}$  is the  $(n-1)$ -dimensional Hausdorff measure.

**Note:** one or more elements of  $\mathcal{E}^B$  might be curved, depending on the shape of  $\Omega$ .

### 1.2.3. Assumprions on the mesh

Throughout this thesis we shall always assume that the following properties hold for the set  $\{\mathcal{T}_h\}_{h \in (0,1)}$ :

- (i) for all  $h \in (0, 1)$ ,  $\mathcal{T}_h$  is made of poligons (if  $n = 2$ ) or polihedra (if  $n = 3$ ) and it is conforming.
- (ii) there exists  $\theta_0 > 0$  such that  $h_k/\rho_K \geq \theta_0 \forall K \in \mathcal{T}_h, \forall h \in (0, 1)$ . Here  $h_k$  and  $\rho_K$  denote respectively the radius of the circumscribed circle and the radius of the inscribed one for the element  $K$ . This assumption is usually called *Shape regularity* or *minimal angle condition*.

---

<sup>1</sup>Here  $\mathcal{H}^{n-1}$  is the  $(n-1)$ -dimensional Hausforff measure.

(iii) the mesh is *quasi uniform*, i.e.  $\exists C_*, C^* > 0$  such that if  $K_1, K_2 \in \mathcal{T}_h$  are adjacent, then

$$C_* \text{diam}(K_1) \leq \text{diam}(K_2) \leq C^* \text{diam}(K_1).$$

### 1.2.2 Heuristics

We will now spend a few words on the idea behind the upcoming definitions and the upcoming method. The discontinuous Galerkin method here presented will, in a sense, replicate/mimic the action of a map  $u$  when integrated against the laplacian of a test function  $v \in C^2(\Omega) \cap C^0(\bar{\Omega})$ <sup>(2)</sup>. To start grasping the idea behind the method, we make the following observation: define the space  $C^2(\mathcal{T}_h)$  as

$$C^2(\mathcal{T}_h) := \{u \in L^2(\Omega) : u|_K \in C^2(K) \cap C^0(\bar{K}) \forall K \in \mathcal{T}_h\}.$$

Then if  $v \in C^2(\Omega) \cap C^0(\bar{\Omega})$  and  $u \in C^2(\mathcal{T}_h)$ , then

$$\begin{aligned} \int_{\Omega} u \Delta v \, dx &= \sum_{K \in \mathcal{T}_h} \int_K u \Delta v \, dx \\ &= \sum_{K \in \mathcal{T}_h} \left( - \int_K \nabla u \cdot \nabla v \, dx + \int_{\partial K} u \nabla v \cdot \mathbf{n}_K \, ds \right). \end{aligned}$$

The actual space which we will be dealing with is the space of maps which are locally  $H^2$ -regular on the mesh, i.e. maps  $u \in L^2(\Omega)$  such that  $u|_K \in H^2(K) \forall K \in \mathcal{T}_h$ . Therefore we will define a bilinear form which much recalls the summation above.

### 1.2.3 The energy space $E_h$

#### 1.2.4. Definition, the energy space $E_h$

Given a mesh  $\mathcal{T}_h$  on the domain  $\Omega$  we shall define the *energy space*  $E_h$  on  $\mathcal{T}_h$  as the set

$$E_h := \{u \in L^2(\Omega) : u|_K \in H^2(K) \forall K \in \mathcal{T}_h\}$$

endowed with an appropriate *energy norm*, which is defined below.

Notice that the energy space  $E_h$  allows for jumps along the edges of the elements of the mesh. To formulate an analogous of the weak Poisson problem defined in 1.1.2 we will need to add term which quantify such jumps.

#### 1.2.5. Definition, jump and normal derivative along an edge

Fix an order  $K_1, K_2, \dots, K_m$  in which we shall enumerate the elements of  $\mathcal{T}_h$ .

Let  $1 \leq n_1 < n_2 \leq m$  be positive integers.

Let  $K_{n_1}, K_{n_2} \in \mathcal{T}_h$  be adjacent elements of the mesh.

Let  $\ell := \overline{K_{n_1}} \cap \overline{K_{n_2}}$  be the edge shared by  $K_{n_1}$  and  $K_{n_2}$ .

For later use, denote  $K_{n_1}, K_{n_2}$  as  $K^- := K_{n_1}, K^+ := K_{n_2}$ .

<sup>2</sup>This concept is treated rigorously in the field of distributions and distributional derivatives. For a formal introduction to the notion of distribution and distributional derivative we refer to [9].

Let  $u \in E_h$ .

Let  $u^-$  be the extension by continuity of  $u|_{K^-}$  on  $\overline{K^-}$ .

Let  $u^+$  be the extension by continuity of  $u|_{K^+}$  on  $\overline{K^+}$ .

Let  $\mathbf{n}^-(\cdot)$  and  $\mathbf{n}^+(\cdot)$  be the exterior normals of  $K^-$  and  $K^+$  respectively (indeed only defined on their boundaries).

Then define the following functions and quantities:

- $\langle v, w \rangle_\ell := \int_\ell uv \, d\mathcal{H}^{n-1} \quad \forall v, w \in E_h.$
- $[u]_\ell := u^+|_\ell - u^-|_\ell.$
- $\{\partial_n u\}|_\ell := \frac{\partial u^+}{\partial \mathbf{n}^+}.$

**1.2.6. Remark:** the choice of the ordering  $K_1, \dots, K_m$  plays a role in determining the sign of  $[u]$  and the value of  $\{\partial_n u\}|_\ell$ .

**1.2.7. Remark:** Arnold, in [1], formulates the Discontinuous Galerkin method in an analogous way, defining  $\{\partial_n u\}|_\ell$  as

$$\{\partial_n u\}|_\ell := \frac{1}{2} \left( \frac{\partial u^+}{\partial \mathbf{n}^+} + \frac{\partial u^-}{\partial \mathbf{n}^+} \right).$$

All results presented here remain valid under Arnold's alternative formulation.

### 1.2.8. Definition, normal derivative along a boundary edge

Let  $\ell \in \mathcal{E}^B$  be a boundary edge of some element  $K^+ \in \mathcal{T}_h$ .

Then, similarly as in 1.2.5, define

- $\partial_n u|_\ell := \frac{\partial u}{\partial \mathbf{n}^+}.$

### 1.2.9. Definition, Energy norm

We shall endow the energy space  $E_h$  with the energy norm  $\|\cdot\|_{1,h}$  defined as

$$\begin{aligned} \|u\|_{1,h} := & \left( \sum_{K \in \mathcal{T}_h} \|\nabla u\|_{H^2(K)}^2 + \sum_{\ell \in \mathcal{E}^I} \left[ h_\ell \|\{\partial_n u\}\|_{L^2(\ell)}^2 + \frac{1}{h_\ell} \|[u]\|_{L^2(\ell)}^2 \right] + \right. \\ & \left. + \sum_{\ell \in \mathcal{E}^B} \left[ h_\ell \|\partial_n u\|_{L^2(\ell)}^2 + \frac{1}{h_\ell} \|u\|_{L^2(\ell)}^2 \right] \right)^{1/2}. \end{aligned}$$

### 1.2.4 The space of approximating functions

The set of approximating functions that we will consider is the finite dimensional subspace  $V_h$  of  $L^2(\Omega)$  defined below:

#### 1.2.10. Definition

The set of approximating maps will be

$$V_h^{(k)} := \{u_h \in E_h : u_h|_{\text{int}(K)} \in P_k(K) \forall K \in \mathcal{T}_h\},$$

where  $P_k(K) := \{\text{polynomial functions on } K \text{ of degree } \leq k\}$ .

This will be endowed naturally with the energy norm defined in 1.2.9.

When not necessary we will omit the degree  $k$  and we will write simply  $V_h$  in place of  $V_h^{(k)}$ .

### 1.3 Formulation of the discrete problem

We have set all the necessary notation and we can finally define a bilinear form on the finite dimensional space  $V_h$  and use it to formulate the weak problem yielding the approximate solutions of the Discontinuous Galerkin method.

#### 1.3.1. Problem (weak problem with jumps)

Let  $V_h$  be as in 1.2.10.

Let  $\gamma \in \mathbb{R}_{>0}$ .

Let  $(u, v)_K := \int_K uv \, dx \quad \forall K \in \mathcal{T}_h$  denote the inner product in  $L^2(K)$ .

Define  $a_h^\gamma : V_h \times V_h \rightarrow \mathbb{R}$  to be the bilinear form

$$\begin{aligned} a_h^\gamma(u_h, v_h) := & \sum_{K \in \mathcal{T}_h} (\nabla u_h, \nabla v_h)_K + \\ & - \sum_{\ell \in \mathcal{E}^I} \left[ \langle \{\partial_n u_h\}_\ell, [v_h] \rangle_\ell + \langle \{\partial_n v_h, [u_h]\}_\ell - \frac{\gamma}{h_\ell} \langle [u_h], [v_h] \rangle_\ell \right] \\ & - \sum_{\ell \in \mathcal{E}^B} \left[ \langle \partial_n u_h, v_h \rangle_\ell + \langle \partial_n v_h, u_h \rangle_\ell - \frac{\gamma}{h_\ell} \langle u_h, v_h \rangle_\ell \right]. \end{aligned}$$

Find  $u_h^\gamma \in V_h$  such that

$$a_h^\gamma(u_h^\gamma, v_h) = (f, v_h) \quad \forall v_h \in V_h.$$

#### 1.3.2. Lemma, continuity and coercivity of $a_h^\gamma$

The following two statements hold:

(a)  $a_h^\gamma(u, v) \leq (1 + \gamma) \|u\|_{1,h} \|v\|_{1,h} \quad \forall u, v \in E_h.$

(b) There exist constants  $\gamma_0 > 0$  and  $c_a > 0$  such that

$$a_h^\gamma(u, u) \geq c_a \|u\|_{1,h}^2 \quad \forall \gamma \geq \gamma_0, \forall u \in V_h^{(r)}.$$

Here the constant  $\gamma_0$  depends on  $r$  and on the ratios  $h_K/\rho_K$ .

**1.3.3. Remark:** lemma 1.3.2 implies, by means of Lax-Milgram's lemma, that for  $\gamma$  sufficiently large, problem 1.3.1 is well posed.

### 1.3.1 Error estimates for problem 1.3.1

#### 1.3.4. Definition, approximation hypothesis (AH)

Let  $r \in \mathbb{N} \setminus \{0\}$  be a positive integer.

Let  $m \in \{0, 1, 2, \dots, r\}$  be an integer between zero and  $r$ .

Let  $|u|_{j,D}$  denote the  $j$ -th sobolev seminorm of  $u$  for any  $u \in H^j(D)$ .

Denote by (AH) the following statement:

“There exists a constant  $c > 0$ , independent of  $h$ , such that  $\forall u \in H^m(\Omega)$ ,  $\forall K \in \mathcal{T}_h$ ,  $\exists \chi \in P_{k-1}(K)$  satisfying

$$|u - \chi|_{j,K} \leq ch_K^{m-j} |u|_{m,K} \quad \forall j \in \{0, 1, \dots, m\}.”$$

**1.3.5. Remark:** (AH) is an hypothesis which concerns the geometry of the mesh. Its core detail is that  $c$  is independent of  $h$  <sup>(3)</sup>.

#### 1.3.6. Theorem, error estimates for problem 1.3.1

Assume that (AH) holds.

Let  $u$  and  $u_h^\gamma$  be the solutions, respectively, to problems 1.1.2 and 1.3.1.

Assume that  $u \in H^r(\Omega) \cap H_0^1(\Omega)$ , with  $r \geq 2$ .

Then there exists a positive constant  $c$ , independent of  $h$  and  $u$ , such that

$$\|u - u_h^\gamma\|_{1,h} \leq c \left( \sum_{K \in \mathcal{T}_h} h_K^{2(r-1)} |u|_{r,K}^2 \right)^{1/2}, \quad (1.1)$$

$$\|u - u_h^\gamma\| \leq ch^r |u|_{r,\Omega}. \quad (1.2)$$

## 1.4 A posteriori error estimates

### 1.4.1. Theorem

Let  $e := u - u_h^\gamma$ .

Then

$$\begin{aligned} \sum_{K \in \mathcal{T}_h} \|\nabla e\|_{L^2(K)}^2 &\leq c \left( \sum_{K \in \mathcal{T}_h} h_K^2 \|f + \Delta u_h^\gamma\|_{L^2(K)}^2 + \sum_{\ell \in \mathcal{E}^B} h_\ell \|\partial_n u_h^\gamma\|_{L^2(\ell)}^2 + \right. \\ &\quad \left. + \gamma \sum_{\ell \in \mathcal{E}^I} \frac{1}{h_\ell} \|[u_h^\gamma]\|_{L^2(\ell)}^2 + \gamma \sum_{\ell \in \mathcal{E}^B} \frac{1}{h_\ell} \|u_h^\gamma\|_{L^2(\ell)}^2 \right). \end{aligned}$$

### 1.4.2. Proposition, approximation with maps in $H_0^1(\Omega)$

If  $\{\mathcal{T}_h\}_h$  is a set of conforming meshes made of triangles (tetrahedra if  $d = 3$ ) which satisfies the hypotheses in 1.2.3, then there exists  $C \in \mathbb{R}_{>0}$ , independent of  $h$ , such that

$$\inf_{\chi \in V_h^{(p)} \cap H_0^1(\Omega)} \sum_{K \in \mathcal{T}_h} \|\nabla(v - \chi)\|_{L^2(K)}^2 \leq C \left( \sum_{\ell \in \mathcal{E}^I} \frac{1}{h_\ell} \|[v]\|_{L^2(\ell)}^2 + \sum_{\ell \in \mathcal{E}^B} \frac{1}{h_\ell} \|v\|_{L^2(\ell)}^2 \right). \quad (1.3)$$

<sup>3</sup>Remember that we are considering a set  $\{\mathcal{T}_h\}_{h \in (0,1)}$  of meshes. (AH) requires that all meshes  $\mathcal{T}_h$  satisfy those inequalities with the same constant  $c$ .

**1.4.3. Remark:** proposition 1.4.2 allows to define a family  $\{\mathcal{A}_h\}_h$  of interpolation operators  $\mathcal{A}_h : V_h \rightarrow V_h \cap H_0^1(\Omega)$  such that  $\mathcal{A}_h(v_h)$  satisfies (1.3) without the “inf” for all  $h$  and  $v_h \in V_h$ .

**1.4.4. Theorem**

Suppose that  $f$  is piecewise polynomial on  $\mathcal{T}_h$ . Then

- (i)  $h_K^2 \|f + \Delta u_h^\gamma\|_{L^2(K)}^2 \leq c \|\nabla e\|_{L^2(K)}^2 \quad \forall K \in \mathcal{T}_h.$
- (ii)  $h_\ell \|\partial_n u_h^\gamma\|_{L^2(\ell)}^2 \leq c \left( \|\nabla e\|_{L^2(K^+)}^2 + \|\nabla e\|_{L^2(K^-)}^2 \right) \quad \forall \ell = K^+ \cap K^- \in \mathcal{E}^I.$





## Chapter 2

# Discontinuous Galerkin method for a general convection-reaction problem, a posteriori error analysis

### 2.1 Set up

Let  $\Omega \subset \mathbb{R}^2$  be a bounded lipschitz polygon.

Let  $\Gamma$  denote its topological boundary ( $\Gamma := \partial\Omega$ ).

Let  $f \in L^2(\Omega)$ .

Let  $0 < \epsilon \ll 1$ .

Let  $a(\cdot) = \begin{pmatrix} a_1(\cdot) \\ a_2(\cdot) \end{pmatrix} \in [W^{1,\infty}(\Omega)]^2$ .

Let  $b(\cdot) \in L^\infty(\Omega)$ .

In this chapter we will study a numerical approach to the following problem:

#### 2.1.1. Problem, convection-diffusion problem

Find  $u \in C^2(\Omega) \cap C^0(\bar{\Omega})$  such that

$$\begin{cases} -\epsilon\Delta u + a(x) \cdot \nabla u + b(x)u = f(x). & \text{in } \Omega \\ u = 0 & \text{on } \Gamma \end{cases}$$

This problem admits a weak formulation, derived from applying integration by parts to the LHS of the equation  $\int_{\Omega} \left( -\epsilon\Delta u + a(x) \cdot \nabla u + b(x)u \right) v \, dx = \int_{\Omega} f v \, dx$ , which is

**2.1.2. Problem, weak convection-diffusion problem**

For any  $u, v \in H_0^1(\Omega)$ , define

$$A(u, v) := \int_{\Omega} \left( \epsilon \nabla u \cdot \nabla v + (a(x) \cdot \nabla u)v + b(x)uv \right) dx$$

Find  $u \in H_0^1(\Omega)$  such that

$$A(u, v) = \int_{\Omega} f v dx \quad \forall v \in H_0^1(\Omega).$$

**2.1.3. Remark**

Consider  $A(\cdot, \cdot)$  defined above in 2.1.2. An equivalent expression for  $A(\cdot, \cdot)$  can be found integrating by parts also the convective term, and one gets

$$A(u, v) = \int_{\Omega} \left( \epsilon \nabla u \cdot \nabla v + (a(x) \cdot \nabla v)u + (b(x) - \nabla \cdot a(x))uv \right) dx. \quad (2.1)$$

**2.1.4. Lemma**

Assume that  $a$  and  $\Omega$  are of order 1.

Assume that  $-\frac{1}{2}\nabla \cdot a + b \geq \beta \quad \exists \beta \geq 0$ .

Assume that  $\|-\nabla \cdot a + b\|_{L^\infty(\Omega)} \leq c_* \beta \quad \exists c_* \geq 0$ .

Then a solution for problem 2.1.2 exists and it is unique.

*Proof of lemma 2.1.4.*

The integration by parts formula for Sobolev maps gives, for any  $u, v \in H^1(\Omega)$ , that

$$\begin{aligned} \int_{\Omega} (av) \cdot \nabla u dx &= - \int_{\Omega} u \nabla \cdot (av) dx + \int_{\partial\Omega} (av) \cdot \mathbf{n}_{\Omega} u ds \\ &= - \int_{\Omega} (\nabla \cdot a)uv dx - \int_{\Omega} (a \cdot \nabla v)u + \int_{\partial\Omega} (a \cdot \mathbf{n}_{\Omega})uv ds. \end{aligned}$$

These identities, together with the fact that the boundary term becomes zero when  $u, v \in H_0^1(\Omega)$ , allow to write  $A(\cdot, \cdot)$  in the form

$$A(u, v) = \int_{\Omega} \left( \epsilon \nabla u \cdot \nabla v + \frac{1}{2}(av) \cdot \nabla u - \frac{1}{2}(au) \cdot \nabla v + (b - \frac{1}{2}\nabla \cdot a)uv \right) dx \quad \forall u, v \in H_0^1(\Omega) \quad (2.2)$$

Using (2.2) and the hypothesis  $(b - \frac{1}{2}\nabla \cdot a) \geq \beta$  it is immediate to deduce that  $\exists c_1 \in \mathbb{R}_{>0}$  such that

$$A(u, u) \geq c_1 \|u\|_{H^1(\Omega)}^2 \quad \forall u \in H_0^1(\Omega).$$

In the case in which  $\beta = 0$  one needs to employ the Poincaré inequality.

Now we shall use the Holder inequality for maps in  $L^2(\Omega)$ , the formulation (2.1) of  $A(\cdot, \cdot)$ , the hypothesis on  $\|b - \nabla \cdot a\|_{L^\infty(\Omega)}$  and the lipschitzianity of  $a$  (from which follows that  $\|a\|_{L^\infty(\Omega)} < \infty$ ), in order to deduce that  $\exists c_2 \in \mathbb{R}_{>0}$  such that

$$|A(u, v)| \leq c_2 \|u\|_{H^1(\Omega)} \|v\|_{H^1(\Omega)} \quad \forall u, v \in H_0^1(\Omega).$$

These two properties, which are respectively the coercivity and the continuity of  $A(\cdot, \cdot)$  with respect to the Sobolev  $H^1$  norm, allow to apply Lax-Milgram's lemma, which gives both existence and uniqueness of a solution.  $\square$

**2.1.5. Remark:** The use of Lax-Milgram's lemma in the proof above also gives a bound on  $\|u\|_{H^1(\Omega)}$ , which is explicit when  $\beta > 0$ , and reads

$$\|u\|_{H^1(\Omega)} \leq \frac{1}{\min\{\varepsilon, \beta\}} \|f\|_{L^2(\Omega)}.$$

**2.1.6. Core assumptions:** from this point onward we shall assume that all hypotheses of Lemma 2.1.4 hold.

## 2.2 The mesh

We will use again the notation from section 1.2.1.

We will assume that our domain is polygonal, so that the mesh can be entirely polygonal as well (including the elements which are adjacent to  $\partial\Omega$ ).

We will also assume that:

- $\mathcal{T}_h$  is made of only triangles and parallelograms;
- $\{\mathcal{T}_h\}_{h>0}$  satisfies the conditions stated in 1.2.3.

We shall use, as in chapter 1, the following notation:

- $\mathcal{E}$  denotes the set of all edges;
- $\mathcal{E}^I$  denotes the set of internal edges;
- $\mathcal{E}^B$  denotes the set of boundary edges;
- $h_K$  denotes the diameter of the element  $K$ ;
- $\rho_K$  denotes the radius of (one of) the biggest circle(s) that can be inscribed in the element  $K$ ;
- $\mathbf{n}_K(\cdot)$  denotes the outward normal unitary vector of the element  $K$  (indeed  $\mathbf{n}_K$  is only defined on the set  $\partial K$ ).

Finally fix an  $h > 0$  and consider:

- two adjacent elements of  $\mathcal{T}_h$ ,  $K^-$  and  $K^+$ , meeting along the edge  $\ell := K^- \cap K^+$ ;
- a map  $u : \Omega \rightarrow \mathbb{R}$  such that  $u|_{\text{int}(K)} \in C^\infty(\text{int}(K))$ .

Denote by  $u^-$  the extension by continuity of  $u$  to  $K^-$ .

Denote by  $u^+$  the extension by continuity of  $u$  to  $K^+$ .

We shall quantify the behaviour of  $u$  along the edge  $\ell$  using the following quantities, which are slightly different from the ones used in chapter 1:

- $\{u\}(x) := \frac{1}{2}(u^-(x) + u^+(x))$ , the average of  $u$  along  $\ell$ ;

- $\llbracket u \rrbracket(x) := u^-(x)\mathbf{n}_{K^-}(x) + u^+(x)\mathbf{n}_{K^+}(x)$ , the jump of  $u$  along  $\ell$ ;

If  $q(\cdot) = \begin{pmatrix} q_1(\cdot) \\ q_2(\cdot) \end{pmatrix}$  is a vector field on  $\Omega$  such that  $q|_{\text{int}(K)} \in C^1(\text{int}(K), \mathbb{R}^2) \forall K \in \mathcal{T}_h$ ,

then we shall adopt an analogous notation for  $q^-$  and  $q^+$  and define the quantities

- $\{\!\{q\}\!\}(x) := \frac{1}{2}(q^-(x) + q^+(x))$ , the average of  $u$  along  $\ell$ ;
- $\llbracket q \rrbracket(x) := q^-(x) \cdot \mathbf{n}_{K^-}(x) + q^+(x) \cdot \mathbf{n}_{K^+}(x)$ , the jump of  $u$  along  $\ell$ ;

We shall give similar definitions on the boundary edges  $\ell \in \mathcal{E}^B$ : if  $u \in C^1(\mathcal{T}_h)$ ,  $q(\cdot) = \begin{pmatrix} q_1(\cdot) \\ q_2(\cdot) \end{pmatrix} \in C^1(\mathcal{T}_h, \mathbb{R}^2)$ , and  $\ell = \partial K \cap \Gamma \in \mathcal{E}^B$ , then we shall define

- $\{\!\{u\}\!\}(x) := u(x)$ ;
- $\{\!\{q\}\!\}(x) := q(x)$ ;
- $\llbracket u \rrbracket(x) := u(x)\mathbf{n}_K(x)$ ;
- $\llbracket q \rrbracket(x) := q(x) \cdot \mathbf{n}_K(x)$ .

Finally we shall define the inward and outward flow portions of  $\Gamma$  and of  $\partial K$  as:

- $\Gamma_{\text{in}} := \{x \in \Gamma : a(x) \cdot \mathbf{n}_\Omega(x) < 0\}$ ;
- $\Gamma_{\text{out}} := \{x \in \Gamma : a(x) \cdot \mathbf{n}_\Omega(x) \geq 0\}$ ;
- $\partial K_{\text{in}} := \{x \in \partial K : a(x) \cdot \mathbf{n}_K(x) < 0\}$ ;
- $\partial K_{\text{out}} := \{x \in \partial K : a(x) \cdot \mathbf{n}_K(x) \geq 0\}$ ;

## 2.3 The broken $H^1$ space

Given a fixed  $h$  and the corresponding mesh  $\mathcal{T}_h$ , we can naturally associate to  $\mathcal{T}_h$  the broken  $H^1$  space  $H^1(\mathcal{T}_h)$  defined below.

### 2.3.1. Definition, broken $H^1$ space

We define the *broken  $H^1$  space* associated to  $\mathcal{T}_h$  as

$$H^1(\mathcal{T}_h) := \{u \in L^2(\Omega) : u|_K \in H^1(K) \forall K \in \mathcal{T}_h\}.$$

## 2.4 The space of approximating maps

Let  $p \in \mathbb{N} \setminus \{0\}$  be the desired order of approximation.

Since we are admitting both triangles and parallelograms in our mesh, we shall differentiate the class of approximating polynomials accordingly. Define then, for any  $K \in \mathcal{T}_h$ ,

$$\mathcal{S}_p(K) := \begin{cases} P_p(K) & \text{if } K \text{ is a triangle} \\ Q_p(K) & \text{if } K \text{ is a parallelogram} \end{cases}$$

where  $P_p(K)$  is the set of polynomials  $p(x, y)$  on  $K$  having total degree  $\leq p$ , i.e.

$$P_p(K) := \left\{ \sum_{\substack{i, j \geq 0 \\ i+j \leq p}} a_{ij} x^i y^j : a_{ij} \in \mathbb{R} \right\},$$

while  $Q_p(K)$  is the set of polynomials  $p(x, y)$  on  $K$  having the degree of each variable  $\leq p$ , i.e

$$Q_p(K) := \left\{ \sum_{\substack{0 \leq i \leq p \\ 0 \leq j \leq p}} a_{ij} x^i y^j : a_{ij} \in \mathbb{R} \right\}.$$

#### 2.4.1. Definition, the approximating space $V_h^p$

The space from which the Discontinuous Galerkin method will extract an approximating map will be the space  $V_h^p$ , defined as

$$V_h^p := \left\{ u_h \in L^2(\Omega) : u|_{\text{int}(K) \in S_p(K)} \right\}.$$

In the following, unless needed, we will omit the degree  $p$  and simply write  $V_h$  in place of  $V_h^p$ .

## 2.5 The discrete problem

Putting together all the objects which have been defined so far we are able to formulate the following problem:

#### 2.5.1. Problem, the approximating problem

Let  $\gamma \in \mathbb{R}_{>0}$  be an *interior penalty parameter*.

Let  $f$  be the same  $f$  of problem 2.1.2.

Define the bilinear form  $A_h^\gamma(\cdot, \cdot)$  on  $V_h \times V_h$  as

$$\begin{aligned} A_h^\gamma(u, v) := & \sum_{K \in \mathcal{T}_h} \int_K (\varepsilon \nabla u \cdot \nabla v + (a \cdot \nabla v)u + buv) dx + \\ & - \sum_{\ell \in \mathcal{E}} \int_\ell \varepsilon \{\{\nabla u\}\} \cdot \llbracket v \rrbracket ds - \sum_{\ell \in \mathcal{E}} \int_\ell \varepsilon \{\{\nabla v\}\} \cdot \llbracket u \rrbracket ds + \\ & + \sum_{\ell \in \mathcal{E}} \frac{\varepsilon \gamma}{h_\ell} \int_\ell \varepsilon \llbracket u \rrbracket \cdot \llbracket v \rrbracket ds + \\ & + \sum_{K \in \mathcal{T}_h} \int_{\partial K_{\text{in}} \cap \Gamma_{\text{in}}} (a \cdot \mathbf{n}_K) uv ds + \sum_{K \in \mathcal{T}_h} \int_{\partial K_{\text{in}} \setminus \Gamma} (a \cdot \mathbf{n}_K) (u^e - u) v ds. \end{aligned}$$

Find  $u_h^\gamma \in V_h$  such that

$$A_h^\gamma(u_h, v_h) = \int_\Omega f v_h dx \quad \forall v_h \in V_h.$$

## 2.6 A new norm and seminorm

### 2.6.1. Definition, a norm on $H^1(\mathcal{T}_h)$

Let  $u \in H^1(\mathcal{T}_h)$  <sup>(1)</sup>.

Then define the norm

$$\|u\|^2 := \sum_{K \in \mathcal{T}_h} (\varepsilon \|\nabla u\|_{L^2(K)}^2 + \beta \|u\|_{L^2(K)}^2) + \sum_{\ell \in \mathcal{E}} \frac{\gamma \varepsilon}{h_\ell} \|[[u]]\|_{L^2(K)}^2.$$

### 2.6.2. Definition, a seminorm on $L^2(\Omega)^2$

Let  $q = \begin{pmatrix} q_1 \\ q_2 \end{pmatrix} \in L^2(\Omega)^2$ .

Define then the seminorm

$$|q|_\star := \sup_{v \in H_0^1(\Omega) \setminus \{0\}} \frac{\int_\Omega q \cdot \nabla v \, dx}{\|v\|}.$$

### 2.6.3. Lemma, characterization of $|\cdot|_\star$

Let  $q \in L^2(\Omega)^2$ .

Let  $q = \nabla \varphi + q_0$  be a Helmholtz decomposition<sup>(2)</sup> of  $q$ .

Then

$$|q|_\star = 0 \Leftrightarrow q = q_0.$$

### 2.6.4. Definition, a norm on $H_0^1(\Omega)$

Let  $\varphi \in H_0^1(\Omega)$ .

Then define the norm

$$\|\varphi\|_\star := \sup_{v \in H_0^1(\Omega) \setminus \{0\}} \frac{\int_\Omega \nabla \varphi \cdot \nabla v \, dx}{\|v\|}.$$

**2.6.5. Remark:** If  $q = \nabla \varphi + q_0$  is the Helmholtz decomposition of  $q \in L^2(\Omega)^2$ , then

$$|q|_\star = \|\varphi\|_\star.$$

### 2.6.6. Definition, seminorm associated with the form $A(\cdot, \cdot)$

Consider  $A_h^\gamma(\cdot, \cdot)$  defined as in 2.5.1.

Then, for any  $u \in V_h$ , define

$$|u|_A := \left( |au|_\star^2 + \sum_{\ell \in \mathcal{E}} \left( \beta h_\ell + \frac{h_\ell}{\varepsilon} \right) \|[[u]]\|_{L^2(\ell)}^2 \right)^{1/2}.$$

<sup>1</sup>Recall the definition 2.3.1 of broken  $H^1$  space

<sup>2</sup>**Helmholtz's Theorem (weak form):** for any  $q \in L^2(\Omega)^2$ , there exists exactly one couple  $(\varphi, q_0)$  such that

- $\varphi \in H_0^1(\Omega)$  satisfies  $\int_\Omega \nabla \varphi \cdot \nabla v \, dx = \int_\Omega q \cdot \nabla v \, dx \quad \forall \varphi \in H_0^1(\Omega)$ ;
- $q_0 = q - \nabla \varphi$  is divergence free, in the sense that  $\int_\Omega q_0 \cdot \nabla v \, dx = 0 \quad \forall v \in H_0^1(\Omega)$ .

This produces a unique decomposition of  $q$  given by  $q = \nabla \varphi + q_0$  and this is called *Helmholtz decomposition*.

Moreover, this decomposition is orthogonal in  $L^2(\Omega)^2$ .

## 2.7 A robust a-posteriori error estimator

Before defining the estimator, we need to take into account the presence of round-up errors. We will formalize this aspect of the analysis by defining a machine-computable problem. This problem will be formulated just like problem 2.5.1, but it will not take the original  $a, b, f$  as input parameters, but approximated versions  $a_h, b_h, f_h$ . Here we will not need to specify any approximation criterion. We will only assume that such a criterion exists and that its precision can be quantified in terms of the differences  $a - a_h, b - b_h$  and  $f - f_h$ , as we will write explicitly below when defining the parameter  $\Theta$ .

### 2.7.1. Problem, machine-computable problem

Let  $a, b, f$  be as in section 2.1.

Let  $a_h, b_h, f_h \in V_h$  be approximations of  $a, b, f$  respectively.

Let  $V_h$  be as in definition 2.4.1.

Let  $\gamma \in \mathbb{R}_{>0}$  be the interior penalty parameter.

Define  $\overline{A}_h^\gamma(u_h, v_h)$  as in problem 2.5.1, but using the approximations  $a_h, b_h, f_h$  in place of  $a, b, f$ .

Find  $u_h^\gamma \in V_h$  such that

$$\overline{A}_h^\gamma(u_h^\gamma, v_h) = \int_{\Omega} f_h v_h dx \quad \forall v_h \in V_h.$$

In order to define our estimator in a constructive way we shall define some distinct quantities which will be then put together to define the final estimator.

First, we will define two families of local parameters  $\{\rho_K\}_{K \in \mathcal{T}_h}$  and  $\{\rho_\ell\}_{\ell \in \mathcal{E}}$ , one associated to elements of the mesh and one associated to the set of edges:

$$\rho_K := \begin{cases} \min \left\{ \frac{h_K}{\sqrt{\varepsilon}}, \frac{1}{\sqrt{\beta}} \right\} & \text{if } \beta > 0 \\ \frac{h_K}{\sqrt{\varepsilon}} & \text{if } \beta = 0 \end{cases}$$

$$\rho_\ell := \begin{cases} \min \left\{ \frac{h_\ell}{\sqrt{\varepsilon}}, \frac{1}{\sqrt{\beta}} \right\} & \text{if } \beta > 0 \\ \frac{h_\ell}{\sqrt{\varepsilon}} & \text{if } \beta = 0 \end{cases}$$

Second, we will define three estimators, each one estimating a different source of error in the method:

**Local interior residual** :  $\eta_{R_K}^2 := \rho_K^2 \|f_h + \varepsilon \Delta u_h - a_h \cdot \nabla u_h - b_h u_h\|_{L^2(K)}^2$ .

**Local edge residual** :  $\eta_{E_K}^2 := \sum_{\substack{\ell \in \mathcal{E}^I \\ \ell \subset \partial K}} \frac{\rho_\ell}{2\sqrt{\varepsilon}} \|[\varepsilon \nabla u_h]\|_{L^2(\ell)}^2$ .

$$\begin{aligned} \text{Local jump term : } \eta_{J_K}^2 := & \frac{1}{2} \sum_{\substack{\ell \in \mathcal{E}^I \\ \ell \subset \partial K}} \left( \frac{\gamma \varepsilon}{h_\ell} + h_\ell \beta + \frac{h_\ell}{\varepsilon} \right) \| \llbracket u_h \rrbracket \|_{L^2(\ell)} + \\ & + \sum_{\substack{\ell \in \mathcal{E}^B \\ \ell \subset \partial K}} \left( \frac{\gamma \varepsilon}{h_\ell} + h_\ell \beta + \frac{h_\ell}{\varepsilon} \right) \| \llbracket u_h \rrbracket \|_{L^2(\ell)} \end{aligned}$$

Third, we will put together the quantities above to define two local estimators:

$$\text{Local error estimators : } \eta_K^2 := \eta_{R_K}^2 + \eta_{E_K}^2 + \eta_{J_K}^2.$$

**Local data approximation errors :**

$$\Theta_K^2 := \rho_K^2 \left( \|f - f_h\|_{L^2(K)}^2 + \|(a - a_h) \cdot \nabla u_h\|_{L^2(K)}^2 + \|(b - b_h)u_h\|_{L^2(K)}^2 \right).$$

Fourth, finally, we can define two global estimators:

$$\text{Global error estimator : } \eta^2 := \sum_{K \in \mathcal{T}_h} \eta_K^2.$$

$$\text{Global data approximation error : } \Theta^2 := \sum_{K \in \mathcal{T}_h} \Theta_K^2.$$

Below we will use the symbol “ $\lesssim$ ” to denote an inequality which holds up to multiplication by a positive constant which is independent of the mesh parameter  $h$ , independent of the diffusion coefficient  $\varepsilon$  and independent of the parameter  $\gamma$ . In other words

$$f(\cdot) \lesssim g(\cdot) \Leftrightarrow \begin{array}{l} \exists C \in \mathbb{R}_{>0} \text{ such that } \forall h, \forall \varepsilon, \forall \gamma \text{ it} \\ \text{holds that } f(\cdot) \leq Cg(\cdot). \end{array}$$

Using the notation and quantities described so far, we can state the two main results of this thesis.

### 2.7.2. Theorem, reliability of the estimator $\eta$

Let  $u$  be the exact solution of problem 2.1.2.

Let  $u_h$  be the solution to problem 2.7.1.

Then

$$\| \|u - u_h\| \| + |u - u_h|_A \lesssim \eta + \Theta.$$

### 2.7.3. Theorem, efficiency of the estimator $\eta$

Let  $u$  be the exact solution of problem 2.1.2.

Let  $u_h$  be the solution to problem 2.7.1.

Then

$$\eta \lesssim \| \|u - u_h\| \| + |u - u_h|_A + \Theta.$$

### 2.7.4. Remark, robustness of $\eta$

The two constants implied by the use of  $\lesssim$  in Theorems 2.7.2 and 2.7.3 are independent of the diffusion parameter  $\varepsilon$ . We shall refer at this property as robustness of the estimator, since its application and effectiveness does not depend on the actual value of  $\varepsilon$ .



In this thesis we will prove only Theorem 2.7.2. The interested reader is referred to [10] for the proof of Theorem 2.7.3. The following section is dedicated entirely to the proof of Theorem 2.7.2.

## 2.8 Proof of Theorem 2.7.2

### 2.8.1 Defining some auxiliary operators

We will decompose this proof in a sequence of lemmas.

To start we shall define some auxiliary forms defined on  $V_h \times V_h$ :

**Diffusion and reaction form :**

$$D_h(u, v) := \sum_{K \in \mathcal{T}_h} \int_K (\varepsilon \nabla u \cdot \nabla v + (b - \nabla \cdot a)uv) dx.$$

**Convection and flux form :**

$$\begin{aligned} O_h(u, v) := & - \sum_{K \in \mathcal{T}_h} \int_K (a \cdot \nabla v)u dx + \sum_{K \in \mathcal{T}_h} \int_{\partial K_{\text{out}} \cap \Gamma} (a \cdot \mathbf{n}_K)uv ds + \\ & + \sum_{K \in \mathcal{T}_h} \int_{\partial K_{\text{out}} \setminus \Gamma} (a \cdot \mathbf{n}_K)u(v - v^e) ds \end{aligned}$$

where  $v$  represents the polynomial value of  $v$  in  $K$  and  $v^e$  represents the polynomial value of  $v$  in the adjacent element.

Notice that in we extend  $v$  outside of the domain  $\Omega$  by defining it to be identically 0, then along boundary edges we have that  $v - v^e = v - 0 = v$ , coherently with the definition of  $O_h$ .

**Edge diffusion form :**

$$K_h(u, v) := - \sum_{\ell \in \mathcal{E}} \int_{\ell} \{\{\varepsilon \nabla u\}\} \cdot \llbracket v \rrbracket ds - \sum_{\ell \in \mathcal{E}} \int_{\ell} \{\{\varepsilon \nabla v\}\} \cdot \llbracket u \rrbracket ds.$$

Notice that  $\{\{\varepsilon \nabla u\}\}$  is by definition the average of  $\varepsilon \nabla u$  along  $\ell$ , while  $\llbracket v \rrbracket$  is the jump along  $\ell$ , which mimics  $-\nabla v$  in a distributional sense. This means that overall the term  $\{\{\varepsilon \nabla u\}\} \cdot \llbracket v \rrbracket$  mimics the term  $-\varepsilon \nabla u \cdot \nabla v$  along edge  $\ell$ , and this explains the name ‘‘Edge diffusion norm’’. The second summation is analogous to the first one and it is added to make the form symmetric.

**Jump diffusion form :**

$$J_h(u, v) := \sum_{\ell \in \mathcal{E}} \frac{\varepsilon \gamma}{h_\ell} \int_{\ell} \llbracket u \rrbracket \cdot \llbracket v \rrbracket ds.$$

**Null trace form :**  $\tilde{A}_h(u, v) := D_h(u, v) + O_h(u, v) + J_h(u, v)$ .

We shall see below that if  $u, v \in H_0^1(\Omega)$ , then  $A_h(u, v) = \tilde{A}_h(u, v)$ . This motivates the name of  $\tilde{A}_h$ .

### 2.8.1. Remark, alternative form of $\tilde{A}_h$

When  $u|_K \in C^2(\overline{K}) \forall K \in \mathcal{T}_h$ , which happens in our case since  $u|_K$  is a polynomial, then integration by parts can be applied on single elements, and it allows to rewrite  $\tilde{A}_h(u, v)$  as

$$\begin{aligned} \tilde{A}_h(u, v) = & \sum_{K \in \mathcal{T}_h} \int_K (-\varepsilon v \Delta u + (a \cdot \nabla u)v + buv) dx + \\ & + \sum_{K \in \mathcal{T}_h} \int_{\partial K} \varepsilon v \nabla u \cdot \mathbf{n}_K ds + \\ & - \sum_{K \in \mathcal{T}_h} \int_{\partial K_{\text{in}} \cap \Gamma} (a \cdot \mathbf{n}_K) uv ds + \\ & - \sum_{K \in \mathcal{T}_h} \int_{\partial K_{\text{in}} \setminus \Gamma} (a \cdot \mathbf{n}_K) (u - u^e) v ds. \end{aligned}$$

which shows the resemblance between  $\tilde{A}_h$  and our original weak problem.

**2.8.2. Remark:** recalling remarks 2.1.3 and ?? one can easily see that  $A_h^\gamma = D_h + O_h + K_h + J_h$ .

**2.8.3. Remark:** The forms  $D_h, O_h, K_h, J_h$  are also well defined for any  $u, v \in V_h + H_0^1(\Omega)$ , and not only for  $u, v \in V_h$ . More precisely:

- Take any  $u \in H_0^1(\Omega)$ . For any  $K \in \mathcal{T}_h$  the trace operator allows to associate to  $u|_K$  its trace map  $T_K u \in L^2(\partial K)$ .
- Using the trace maps  $T_K u$  one can consider some edge  $\ell = \overline{K^-} \cap \overline{K^+}$  and define  $\llbracket u \rrbracket|_\ell$  as  $\llbracket u \rrbracket|_\ell := (T_{K^+} u) \mathbf{n}_{K^+} + (T_{K^-} u) \mathbf{n}_{K^-}$ .
- If  $u \in H^1(\Omega)$  and  $\ell = \overline{K^-} \cap \overline{K^+}$ , then  $(T_{K^+} u)|_\ell = (T_{K^-} u)|_\ell$ , so that  $\llbracket u \rrbracket|_\ell \equiv 0$ .
- With a similar argument one can see that if  $u \in H_0^1(\Omega)$  and  $\ell \in \mathcal{E}^B$ , then  $\llbracket u \rrbracket|_\ell \equiv 0$ .
- If  $u \in H_0^1(\Omega)$ , then for  $L^2$ -a.e. point  $x \in \Omega$  it holds that

$$\nabla u(x) = \lim_{r \rightarrow 0} \frac{1}{|B_r(x)|} \int_{B_r(x)} \nabla u dx.$$

### 2.8.4. Observation

If  $u, v \in H_0^1(\Omega)$ , then  $K_h(u, v) = 0$ .

## 2.8.2 Properties of the auxiliary operators

### 2.8.5. Lemma, coercivity of $\tilde{A}_h$

If  $u \in H_0^1(\Omega)$ , then  $\tilde{A}_h(u, u) \geq \lll u \rrr^2$ .

*Proof of lemma 2.8.5.*

Through integration by parts and a proper rearrangement of the boundary integrals one can rewrite  $\tilde{A}_h(u, v)$  as

$$\begin{aligned} \tilde{A}_h(u, v) = & \sum_{K \in \mathcal{T}_h} \int_K \left( \varepsilon \nabla u \cdot \nabla v + \left( b - \frac{1}{2} \nabla \cdot a \right) uv \right) dx + \\ & - \frac{1}{2} \sum_{K \in \mathcal{T}_h} \left( \int_K (au) \cdot \nabla v dx - \int_K (av) \cdot \nabla u dx \right) + \\ & - \frac{1}{2} \sum_{K \in \mathcal{T}_h} \int_{\partial K_{\text{in}}} (a \cdot \mathbf{n}_K) uv ds + \frac{1}{2} \sum_{K \in \mathcal{T}_h} \int_{\partial K_{\text{out}}} (a \cdot \mathbf{n}_K) uv ds + \\ & + \sum_{K \in \mathcal{T}_h} \int_{\partial K_{\text{in}} \setminus \Gamma} (a \cdot \mathbf{n}_K) u^e v ds + \\ & + \sum_{\ell \in \mathcal{E}} \frac{\varepsilon \gamma}{h_\ell} \int_\ell \llbracket u \rrbracket \llbracket v \rrbracket ds. \end{aligned}$$

Now the thesis follows because:

- when  $u, v \in H_0^1(\Omega)$  the integrals on the boundary edges are all zero;
- when  $u = v$  the second line cancels itself;
- when  $u = v$  and  $u \in H_0^1(\Omega)$  line 3 cancels precisely line 4.

Now the thesis of the lemma follows from the assumptions made at the very beginning on  $a$  and  $b$  (see section 2.1).  $\square$

### 2.8.6. Lemma, continuity of $D_h, J_h, O_h$

The following inequalities hold:

- $|D_h(u, v)| \leq \|u\| \|v\| \quad \forall u, v \in V_h + H_0^1(\Omega).$
- $|J_h(u, v)| \leq \|u\| \|v\| \quad \forall u, v \in V_h + H_0^1(\Omega).$
- $|O_h(u, v)| \leq |au|_* \|v\| \quad \forall u \in V_h + H_0^1(\Omega), \forall v \in V_h.$

*Proof of lemma 2.8.6.*

The first inequality follows from Cauchy-Schwartz and 2.1.6.

The second inequality follows from Cauchy-Schwartz.

The third inequality follows from the definition of  $|\cdot|_*$ .  $\square$

### 2.8.7. Lemma, inequality for $K_h$

For any  $u \in V_h$  and any  $v \in H_0^1(\Omega) \cap V_h$  it holds that

$$K_h(u, v) \leq \frac{1}{\gamma} \left( \sum_{\ell \in \mathcal{E}} \frac{\varepsilon \gamma}{h_\ell} \|\llbracket u \rrbracket\|_{L^2(\ell)}^2 \right)^{1/2} \|v\|.$$

*Proof of lemma 2.8.7.*

When  $v \in H_0^1(\Omega)$ , the jumps  $\llbracket v \rrbracket$  are identically zero, and therefore if  $v \in H_0^1(\Omega) \cap V_h$  we have

$$K_h(u, v) = - \sum_{\ell \in \mathcal{E}} \int_{\ell} \{\{\varepsilon \nabla v\}\} \cdot \llbracket u \rrbracket ds.$$

The summation on edges can be rewritten as a summation on elements as

$$\sum_{\ell \in \mathcal{E}} \int_{\ell} \{\{\varepsilon \nabla v\}\} \cdot \llbracket u \rrbracket ds = \frac{1}{2} \sum_{K \in \mathcal{T}_h} \int_{\partial K \setminus \Gamma} \varepsilon \nabla v \cdot \llbracket u \rrbracket ds + \sum_{K \in \mathcal{T}_h} \int_{\partial K \cap \Gamma} \varepsilon \nabla v \cdot \llbracket u \rrbracket ds.$$

It also holds that

$$\begin{aligned} \left| \frac{1}{2} \sum_{K \in \mathcal{T}_h} \int_{\partial K \setminus \Gamma} \varepsilon \nabla v \cdot \llbracket u \rrbracket ds + \sum_{K \in \mathcal{T}_h} \int_{\partial K \cap \Gamma} \varepsilon \nabla v \cdot \llbracket u \rrbracket ds \right| &\leq \\ &\leq \frac{3}{2} \sum_{K \in \mathcal{T}_h} \int_{\partial K} |\varepsilon \nabla v| |\llbracket u \rrbracket| ds. \end{aligned}$$

Now using Cauchy-Schwartz, the inverse inequality  $\|v\|_{L^2(\partial K)} \lesssim h_K^{-1/2} \|v\|_{L^2(K)} \quad \forall v \in S_p(K)$  (which holds thanks to shape regularity), and then Cauchy-Schwartz again, one can further deduce that

$$\frac{3}{2} \sum_{K \in \mathcal{T}_h} \int_{\partial K} |\varepsilon \nabla v| |\llbracket u \rrbracket| ds \lesssim \frac{1}{\gamma} \left( \sum_{K \in \mathcal{T}_h} \varepsilon h_K \|v\|_{L^2(K)}^2 \right)^{1/2} \left( \frac{\varepsilon \gamma}{h_K} \|\llbracket u \rrbracket\|_{L^2(K)}^2 \right)^{1/2}.$$

Chaining up all the equalities and inequalities written so far in this proof we get

$$\begin{aligned} |K_h(u, v)| &\lesssim \frac{1}{\gamma} \left( \sum_{K \in \mathcal{T}_h} \varepsilon h_K \|v\|_{L^2(K)}^2 \right)^{1/2} \left( \frac{\varepsilon \gamma}{h_K} \|\llbracket u \rrbracket\|_{L^2(K)}^2 \right)^{1/2} \\ &\lesssim \frac{1}{\gamma} \left( \sum_{K \in \mathcal{T}_h} \varepsilon \|v\|_{L^2(K)}^2 \right)^{1/2} \left( \frac{\varepsilon \gamma}{h_K} \|\llbracket u \rrbracket\|_{L^2(K)}^2 \right)^{1/2}, \end{aligned}$$

from which the thesis of the lemma follows immediately.  $\square$

### 2.8.8. Lemma, inf-sup condition on $\tilde{A}_h$

There exists  $C \in \mathbb{R}_{>0}$  such that

$$\inf_{u \in H_0^1(\Omega) \setminus \{0\}} \sup_{v \in H_0^1(\Omega) \setminus \{0\}} \frac{\tilde{A}_h(u, v)}{(\|u\| + |au|_{\star}) \|v\|} \geq C.$$

*Proof of lemma 2.8.8.*

Let  $u \in H_0^1(\Omega)$  and  $\theta \in (0, 1)$ .

Recall that, by definition of  $|\cdot|_{\star}$ ,  $|au|_{\star} = \sup_{v \in H_0^1(\Omega) \setminus \{0\}} \frac{\int_{\Omega} (a \cdot \nabla v) u dx}{\|v\|}$ .

Then there exists  $w_{\theta} \in H_0^1 \setminus \{0\}$  such that

$$\|w_{\theta}\| = 1 \quad \text{and} \quad O_h(u, w_{\theta}) = - \int_{\Omega} (a \cdot \nabla w_{\theta}) u dx \geq \theta |au|_{\star}.$$

Also recall, from lemma 2.8.6, that  $\exists C_1 > 0$  such that  $D_h(u, w_\theta) + J_h(u, w_\theta) \leq C_1 \| \|u\| \|w_\theta\|$ .

Then

$$\begin{aligned} \tilde{A}_h(u, w_\theta) &= (D_h + J_h + O_h)(u, w_\theta) \\ &\geq |au|_\star - C_1 \| \|u\| \|w_\theta\| \\ &= |au|_\star - C_1 \| \|u\|. \end{aligned}$$

Define  $v_\theta := u + \frac{\| \|u\|}{1 + C_1} w_\theta$ .

Indeed  $\| \|v_\theta\| \leq \left(1 + \frac{1}{1 + C_1}\right) \| \|u\|$ .

Recall that, by lemma 2.8.5,  $\tilde{A}_h(u, u) \geq \| \|u\|^2$ .

Then, to conclude the proof of the lemma, we see that

$$\begin{aligned} \sup_{v \in H_0^1(\Omega) \setminus \{0\}} \frac{\tilde{A}_h(u, v)}{\| \|v\|} &\geq \frac{\tilde{A}_h(u, v_\theta)}{\| \|v_\theta\|} \\ &\geq \frac{\| \|u\| + (1 + C_1)^{-1} \| \|u\| (\theta |au|_\star - C_1 \| \|u\|)}{(1 + 1/(1 + C_1)) \| \|v_\theta\|} \\ &= \frac{1}{2 + C_1} (\theta |u|_\star + \| \|u\|). \end{aligned}$$

This inequality was proven for a generic  $v \in H_0^1(\Omega)$  and a generic  $\theta$ , and therefore the proof is finished.  $\square$

### 2.8.3 Interpolation operator

Let  $V_h^c := V_h \cap H_0^1(\Omega)$ .

Let  $\mathcal{A}_h : V_h \rightarrow V_h^c$  be as in remark 1.4.3.

Let  $\mathcal{I}_h$  be as in [12], Lemma 3.3, i.e.

$$\mathcal{I}_h : H_0^1(\Omega) \rightarrow \{\varphi \in C^0(\bar{\Omega}) : \varphi|_K \in S_1(K) \forall K \in \mathcal{T}_h \text{ and } \varphi|_\Gamma \equiv 0\}$$

such that

- $\| \mathcal{I}_h v \| \lesssim \| \|v\| \quad \forall v \in H_0^1(\Omega)$ .
- $\left( \sum_{K \in \mathcal{T}_h} \frac{1}{\rho_K^2} \| \|v - \mathcal{I}_h v\|_{L^2(K)}^2 \right)^{1/2} \lesssim \| \|v\| \quad \forall v \in H_0^1(\Omega)$ .
- $\left( \sum_{\ell \in \mathcal{E}} \frac{\sqrt{\varepsilon}}{\rho_\ell} \| \|v - \mathcal{I}_h v\|_{L^2(\ell)}^2 \right)^{1/2} \lesssim \| \|v\| \quad \forall v \in H_0^1(\Omega)$ .

### 2.8.4 Conclusion of the proof

For any  $u_h \in V_h$  let  $u_h^c := \mathcal{A}_h u_h \in V_h^c$  be its conforming part and let  $u_h^r := u_h - u_h^c$  be the remaining part, so that we have in fact decomposed  $u_h$  as  $u_h = u_h^c + u_h^r$ .

Now, by an elementary application of the triangular inequality we can write the estimate

$$\| \|u - u_h\| \| + |u - u_h|_A \leq \| \|u - u_h^c\| \| + |u - u_h^c|_A + \| \|u_h^r\| \| + |u_h^r|_A. \quad (2.3)$$

We will prove the theorem by estimating separately the lements appearing at the RHS.

### 2.8.9. Lemma, estimating the $u_h^r$ part

It holds that  $\| \|u_h^r\| \| + |u_h^r|_A \lesssim \eta$ .

*Proof of lemma 2.8.9.*

By construction,  $\llbracket u_h^r \rrbracket = \llbracket u_h \rrbracket$ .

Therefore

$$\begin{aligned} \| \|u_h^r\| \|^2 + |u_h^r|_A^2 &= \sum_{K \in \mathcal{T}_h} (\varepsilon \|\nabla u\|_{L^2(K)}^2 + \beta \|u\|_{L^2(K)}^2) + |au_h^r|_{\star}^2 + \\ &\quad + \sum_{\ell \in \mathcal{E}} \left( \frac{\varepsilon \gamma}{h_\ell} + \beta h_\ell + \frac{h_\ell}{\varepsilon} \right) \|\llbracket u_h \rrbracket\|_{L^2(\ell)}^2 \\ &\lesssim \sum_{K \in \mathcal{T}_h} (\varepsilon \|\nabla u\|_{L^2(K)}^2 + \beta \|u\|_{L^2(K)}^2) + |au_h^r|_{\star}^2 + \\ &\quad + \sum_{K \in \mathcal{T}_h} \eta_{J_K} \end{aligned}$$

By proposition 1.4.2 we have

$$\sum_{K \in \mathcal{T}_h} \varepsilon \|\nabla u_h^r\|_{L^2(K)}^2 \lesssim \frac{1}{\gamma} \sum_{K \in \mathcal{T}_h} \frac{\varepsilon \gamma}{h_\ell} \|\llbracket u \rrbracket\|_{L^2(K)}^2 \lesssim \frac{1}{\gamma} \sum_{K \in \mathcal{T}_h} \eta_{J_K}^2.$$

In [7] (see proposition 5.2, at section 5.4), the authors prove that the operator  $\mathcal{A}$  also satisfies

$$\sum_{K \in \mathcal{T}_h} \| \|u_h^r\| \|_{L^2(K)}^2 \lesssim \sum_{\ell \in \mathcal{E}} h_\ell \|\llbracket u_h \rrbracket\|_{L^2(\ell)}^2,$$

which implies

$$\beta \sum_{K \in \mathcal{T}_h} \| \|u_h^r\| \|_{L^2(K)}^2 \lesssim \frac{1}{\gamma} \sum_{\ell \in \mathcal{E}} \beta h_\ell \| \|u_h\| \|_{L^2(\ell)}^2 \lesssim \sum_{K \in \mathcal{T}_h} \beta h_\ell \eta_{J_K}^2.$$

Recall that  $|au|_{\star} = \sup_{v \in H_0^1(\Omega) \setminus \{0\}} \frac{\int_{\Omega} (a \cdot \nabla v) u \, dx}{\| \|v\| \|} \lesssim \| \|u\| \|_{L^2(\Omega)}$  (the inequality holds because  $a$  was assumed to be of order 1).

Finally, using again the properties of the operator  $\mathcal{A}$ , we have

$$|au_h^r|_{\star}^2 \lesssim \frac{1}{\varepsilon} \| \|u_h^r\| \|_{L^2(\Omega)}^2 \lesssim \sum_{\ell \in \mathcal{E}} \frac{h_\ell}{\varepsilon} \|\llbracket u_h \rrbracket\|_{L^2(\ell)}^2 \lesssim \frac{h_\ell}{\varepsilon} \eta_{J_K}^2,$$

and thus the proof is concluded.  $\square$

### 2.8.10. Lemma, estimating the error of $\tilde{A}_h$ and $\mathcal{I}_h$

For any  $v \in H_0^1(\Omega)$  we have that

$$\int_{\Omega} f(v - \mathcal{I}_h v) \, dx - \tilde{A}_h(u_h, v - \mathcal{I}_h v) \lesssim (\eta + \Theta) \| \|v\| \|.$$

*Proof of lemma 2.8.10.*

Let  $T := \int_{\Omega} f(v - \mathcal{I}_h v) dx - \tilde{A}_h(u_h, v - \mathcal{I}_h v)$ .

Integrating by parts as in remark 2.8.1 we get

$$\begin{aligned} T &= \sum_{K \in \mathcal{T}_h} \int_K (f + \varepsilon \Delta u_h - a \cdot \nabla u_h - b u_h)(v - \mathcal{I}_h v) dx + \\ &\quad - \sum_{K \in \mathcal{T}_h} \int_{\partial K} ((\varepsilon \nabla u_h \cdot \mathbf{n}_K)(v - \mathcal{I}_h v)) ds + \\ &\quad + \sum_{K \in \mathcal{T}_h} \int_{\partial K_{\text{in}} \setminus \Gamma} a \cdot \mathbf{n}_K (u_h - u_h^e)(v - \mathcal{I}_h v) ds \\ &=: T_1 + T_2 + T_3. \end{aligned}$$

**Claim 1.**  $T_1 \lesssim \left( \sum_{K \in \mathcal{T}_h} \eta_{R_K}^2 + \Theta_K^2 \right)^{1/2} \|v\|.$

*Proof of claim 1.*

First, we shall add and subtract the data approximation terms to  $T_1$ :

$$\begin{aligned} T_1 &= \sum_{K \in \mathcal{T}_h} \int_K (f_h + \varepsilon \Delta u_h - a_h \cdot \nabla u_h - b_h u_h)(v - \mathcal{I}_h v) dx + \\ &\quad + \sum_{K \in \mathcal{T}_h} \int_K ((f - f_h) + \varepsilon \Delta u_h - (a - a_h) \cdot \nabla u_h - (b - b_h) u_h)(v - \mathcal{I}_h v) dx \end{aligned}$$

Now using Cauchy-Schwartz and the properties of  $\mathcal{I}_h$ , we get

$$\begin{aligned} T_1 &\lesssim \left( \sum_{K \in \mathcal{T}_h} \eta_{R_K}^2 \right)^{1/2} \left( \sum_{K \in \mathcal{T}_h} \frac{1}{\rho_K^2} \|v - \mathcal{I}_h v\|_{L^2(K)}^2 \right)^{1/2} + \\ &\quad + \left( \sum_{K \in \mathcal{T}_h} \Theta_K^2 \right)^{1/2} \left( \sum_{K \in \mathcal{T}_h} \frac{1}{\rho_K^2} \|v - \mathcal{I}_h v\|_{L^2(K)}^2 \right)^{1/2} \\ &\lesssim \left( \sum_{K \in \mathcal{T}_h} \eta_{R_K}^2 + \Theta_K^2 \right)^{1/2} \|v\|. \end{aligned}$$

□

**Claim 2.**  $T_2 \lesssim \left( \sum_{K \in \mathcal{T}_h} \eta_{E_K}^2 \right)^{1/2} \|v\|.$

*Proof of claim 2.*

We shall rewrite  $T_2$  in terms of the jumps of  $\nabla u$  as

$$T_2 = - \sum_{\ell \in \mathcal{E}^I} \varepsilon \llbracket \nabla u \rrbracket (v - \mathcal{I}_h v) ds.$$

Now Cauchy-Schwartz and the properties of  $\mathcal{I}_h$  imply that

$$\begin{aligned} T_2 &\lesssim \left( \sum_{\ell \in \mathcal{E}^I} \frac{\rho_\ell}{\sqrt{\varepsilon}} \|\llbracket \varepsilon \nabla u \rrbracket\|_{L^2(\ell)}^2 \right)^{1/2} \left( \sum_{\ell \in \mathcal{E}} \frac{\sqrt{\varepsilon}}{\rho_\ell} \|v - \mathcal{I}_h v\|_{L^2(\ell)}^2 \right)^{1/2} \\ &\lesssim \left( \sum_{\ell \in \mathcal{E}^I} \eta_{E_K}^2 \right)^{1/2} \|v\|. \end{aligned}$$

□

**Claim 3.**  $T_3 \lesssim \left( \sum_{K \in \mathcal{T}_h} \eta_{J_K}^2 \right)^{1/2} \|v\|.$

*Proof of claim 2.*

Remembering that by the shape regularity assumption we have  $\rho_\ell \leq \frac{1}{\sqrt{\varepsilon}} h_K$ , we can use Cauchy-Schwartz and the properties of  $\mathcal{I}_h$  to deduce that

$$\begin{aligned} T_3 &\lesssim \left( \sum_{\ell \in \mathcal{E}} \frac{\rho_\ell}{\sqrt{\varepsilon}} \|u_h\|_{L^2(\ell)}^2 \right)^{1/2} \left( \sum_{\ell \in \mathcal{E}} \frac{\sqrt{\varepsilon}}{\rho_\ell} \|v - \mathcal{I}_h v\|_{L^2(\ell)}^2 \right)^{1/2} \\ &\lesssim \left( \sum_{K \in \mathcal{T}_h} \eta_{J_K}^2 \right)^{1/2} \|v\|. \end{aligned}$$

□

Claims 1, 2 and 3 together prove lemma 2.8.10. □

### 2.8.11. Lemma, estimating the $u - u_h^c$ part

There holds

$$\|u - u_h^c\| + |u - u_h^c|_A \lesssim \eta + \Theta.$$

*Proof of lemma 2.8.11.*

Applying lemma 2.8.8 applied to  $u - u_h^c$  gives

$$\|u - u_h^c\| + |a(u - u_h^c)|_* \lesssim \sup_{v \in H_0^1(\Omega) \setminus \{0\}} \frac{\tilde{A}_h(u - u_h^c, v)}{\|v\|}.$$

Notice that if  $u, v \in H_0^1(\Omega)$ , then

$$\tilde{A}_h(u, v) = A_h(u, v) = A(u, v). \quad (2.4)$$

Also notice, simply rewriting remark 2.8.2, that

$$A_h(u, v) = \tilde{A}_h(u, v) + K_h(u, v) \quad \forall u, v \in V_h. \quad (2.5)$$



Now, by definition  $u \in H_0^1(\Omega)$  and  $u_h^c \in V_h \cap H_0^1(\Omega)$ , so for all  $v \in H_0^1(\Omega)$  it holds that

$$\begin{aligned} \tilde{A}_h(u - u_h^c, v) &= A(u, v) - \tilde{A}_h(u_h^c, v) \\ &= \int_{\Omega} f v \, dx - \tilde{A}_h(u_h, v) + \tilde{A}_h(u_h^r, v) \\ &= \int_{\Omega} f v \, dx + \left( - \int_{\Omega} f \mathcal{I}_h v + A_h(u_h, \mathcal{I}_h v) \right) - \tilde{A}_h(u_h, v) + \tilde{A}_h(u_h^r, v) \\ &= \int_{\Omega} f(v - \mathcal{I}_h v) \, dx - \tilde{A}_h(u_h, v - \mathcal{I}_h v) + K_h(u_h, \mathcal{I}_h v) + \tilde{A}_h(u_h^r, v). \end{aligned}$$

So, in the end, we may define

$$\begin{aligned} U_1 &:= \int_{\Omega} f(v - \mathcal{I}_h v) \, dx - \tilde{A}_h(u_h, v - \mathcal{I}_h v) \\ U_2 &:= \tilde{A}_h(u_h^r, v) \\ U_3 &:= K_h(u_h, \mathcal{I}_h v) \end{aligned}$$

and write neatly  $\tilde{A}_h(u - u_h^c, v) = U_1 + U_2 + U_3$ .

Now lemma 2.8.10 gives

$$T_1 \lesssim (\eta + \Theta) \|v\|,$$

lemmas 2.8.6 and 2.8.9 give

$$T_2 \lesssim (\|u_h^r\| + |au_h^r|_*) \|v\| \lesssim \eta \|v\|,$$

and lemma 2.8.7, together with the properties of  $\mathcal{I}_h$ , gives

$$T_3 \lesssim \frac{1}{\gamma} \left( \sum_{K \in \mathcal{T}_h} \eta_{J_K}^2 \right)^{1/2} \|\mathcal{I}_h v\| \lesssim \frac{1}{\gamma} \left( \sum_{K \in \mathcal{T}_h} \eta_{J_K}^2 \right)^{1/2} \|v\|.$$

The proof of lemma 2.8.11 is concluded.  $\square$

The statement of Theorem 2.7.2 now follows immediately from (2.3) and from lemmas 2.8.9 and 2.8.11.



# Bibliography

- [1] Douglas N Arnold. An interior penalty finite element method with discontinuous elements. *SIAM journal on numerical analysis*, 19(4):742–760, 1982.
- [2] Ivo Babuška and Werner C Rheinboldt. A-posteriori error estimates for the finite element method. *International journal for numerical methods in engineering*, 12(10):1597–1615, 1978.
- [3] I Babuvška and Werner C Rheinboldt. Error estimates for adaptive finite element computations. *SIAM Journal on Numerical Analysis*, 15(4):736–754, 1978.
- [4] Roland Becker, Peter Hansbo, and Mats G Larson. Energy norm a posteriori error estimation for discontinuous galerkin methods. *Computer Methods in Applied Mechanics and Engineering*, 192(5-6):723–733, 2003.
- [5] Willy Dörfler. A convergent adaptive algorithm for poisson’s equation. *SIAM Journal on Numerical Analysis*, 33(3):1106–1124, 1996.
- [6] Xiaobing Feng and Ohannes A Karakashian. Two-level additive schwarz methods for a discontinuous galerkin approximation of second order elliptic problems. *SIAM Journal on Numerical Analysis*, 39(4):1343–1365, 2001.
- [7] Paul Houston, Dominik Schötzau, and Thomas P Wihler. Energy norm a posteriori error estimation of hp-adaptive discontinuous galerkin methods for elliptic problems. *Mathematical Models and Methods in Applied Sciences*, 17(01):33–62, 2007.
- [8] Ohannes A Karakashian and Frederic Pascal. A posteriori error estimates for a discontinuous galerkin approximation of second-order elliptic problems. *SIAM Journal on Numerical Analysis*, 41(6):2374–2399, 2003.
- [9] Steven G Krantz and Harold R Parks. *Geometric integration theory*. Springer Science & Business Media, 2008.
- [10] Dominik Schötzau and Liang Zhu. A robust a-posteriori error estimator for discontinuous galerkin methods for convection–diffusion equations. *Applied numerical mathematics*, 59(9):2236–2255, 2009.

- [11] Rüdiger Verfürth. A review of a posteriori error estimation and adaptive mesh-refinement techniques. *(No Title)*, 1996.
- [12] Rüdiger Verfürth. Robust a posteriori error estimates for stationary convection-diffusion equations. *SIAM journal on numerical analysis*, 43(4):1766–1782, 2005.