

MOX-Report No. 13/2007

**Discontinuous Galerkin methods based on  
weighted interior penalties for second order  
PDEs with non-smooth coefficients**

PAOLO ZUNINO

MOX, Dipartimento di Matematica "F. Brioschi"  
Politecnico di Milano, Via Bonardi 29 - 20133 Milano (Italy)

[mox@mate.polimi.it](mailto:mox@mate.polimi.it)

<http://mox.polimi.it>



# Discontinuous Galerkin methods based on weighted interior penalties for second order PDEs with non-smooth coefficients

Paolo Zunino

June 12, 2007

‡ MOX– Modellistica e Calcolo Scientifico  
Dipartimento di Matematica “F. Brioschi”  
Politecnico di Milano  
via Bonardi 9, 20133 Milano, Italy  
`paolo.zunino@polimi.it`

## Abstract

We develop and analyze a Discontinuous Galerkin (DG) method based on weighted interior penalties (WIP) applied to second order PDEs and in particular to advection-diffusion-reaction equations featuring non-smooth and possibly vanishing diffusivity. First of all, looking at the derivation of a DG scheme with a bias to domain decomposition methods, we carefully discuss the set up of the discretization scheme in a general framework putting into evidence the helpful role of the weights and the connection with the well known Local Discontinuous Galerkin schemes (LDG). Then, we address the a-priori and the a-posteriori error analysis of the method, recovering optimal error estimates in suitable norms. By virtue of the introduction of the weighted penalties, these results turn out to be robust with respect to the diffusion parameter. Furthermore, we discuss the derivation of an a-posteriori local error indicator suitable for advection-diffusion-reaction problems with highly variable, locally small diffusivity. Finally, all the theoretical results are illustrated and discussed by means of numerical experiments.

## 1 Introduction

Although the DG methods are usually defined by means of the so called numerical fluxes between neighboring mesh cells, see [2], for most of the interior penalty methods for second order problems it is possible to correlate the expression of the numerical fluxes with a corresponding set of local interface conditions that are weakly enforced on each inter-element boundary. Such conditions are suitable to couple elliptic PDEs with smooth coefficients and it seems that a

little attention is paid to the case of problems with discontinuous data or to the limit case where the diffusivity vanishes in some parts of the computational domain.

To address these topics, we look at the derivation of a DG scheme with a bias to domain decomposition methods. Indeed, we discuss the set up of an interior penalty DG method arising from a set of generalized interface conditions, introduced in [10] to couple both elliptic and hyperbolic problems, which give rise to the so called heterogeneous domain decomposition methods, see [14]. Our purpose is to obtain a DG scheme that joins the efficacy of heterogeneous domain decomposition methods for the approximation of problems that vary in character from one part of the domain to another with the flexibility of the discontinuous approximation spaces, which allow to approximate a possibly discontinuous function without knowing a priori the location of the discontinuity. As a result of that, our method will turn out to be very effective for problems whose solution features very sharp internal layers.

In order to obtain such method, it is necessary to modify the numerical fluxes of a standard interior penalty scheme, replacing the arithmetic mean with suitably weighted averages where the weights depend on the coefficients of the problem. We will show that this technique provides several advantages with respect to standard IP schemes. First of all, it improves the capability of the method to regulate the degree of smoothness of the approximate solution by tuning local variation of the penalty weighting function with respect to the local variation of the coefficients of the problem. As pointed out in the seminal work by D.N. Arnold, [1], this is a fundamental feature for IP methods and we will see that it remarkably improves the stability of the scheme. Furthermore, we will show that the numerical fluxes corresponding to the methods based on the weighted interior penalties (WIP) are correlated with the ones defining the family of the local discontinuous Galerkin methods (LDG), see [8, 6, 7]. Interestingly, this provides a new recipe to choose the parameters needed to set up LDG methods applied to advection-diffusion-reaction problems. Secondly, we observe that the introduction of WIP can be seen as a way to incorporate into the definition of the scheme some partial knowledge of the solution. We will analyze how this feature turns out to be useful in the a-posteriori error analysis of the scheme and in particular in the definition of a local error indicator.

We point out that the idea of WIP is not completely new. In the framework of mortar finite-element methods, different authors have already highlighted the possibility of using an average with weights that differ from one half, see [17, 12]. These works present several mortaring techniques to match conforming finite elements on possibly non conforming computational meshes. However, these works do not consider any connection between the averaging weights and the coefficients of the problem. More recently, see [5], Burman and Zunino have introduced this dependence for an advection-diffusion-reaction problem with discontinuous diffusivity, and they have shown that the stability of the numerical solution obtained with  $H^1$ -conformal finite elements can be improved by means of a mortar scheme inspired to interior penalty methods where the penalty parameter is proportional to the harmonic mean of the diffusivity on

the interface of discontinuity. Later on, see [9], the WIP technique has been extended from mortars to the case of DG methods for the approximation of problems with anisotropic diffusion.

In this work, we upgrade to a more general framework the ideas proposed in [5] and [9]. The derivation of the method is carefully discussed together with its correlation to other families of DG schemes, in the framework of the unifying formulations proposed in [2] and [4]. Then, the benefits of the WIP scheme are highlighted in the a-priori and a-posteriori error analysis and confirmed by means of numerical experiments.

## 2 Problem set up and numerical approximation

Following [1], we observe that the IP methods are particularly interesting for the approximation of time dependent advection-diffusion-reaction equations with non-smooth and possibly vanishing (nonnegative) diffusivity. More precisely, we focus on problems governed by the following equation,

$$\partial_t u + \nabla \cdot (-\epsilon(t, u) \nabla u + \beta u) + \mu u = f$$

complemented by suitable boundary and initial conditions. However, we observe that the numerical treatment of the time dependence and of the nonlinearity resorts to consider a sequence of linear and stationary problems. Consequently, for the set up of a discretization scheme, we focus for simplicity on this type of equations.

First of all, let be  $x \in \mathbb{R}^d$  with  $d = 2, 3$  and let  $\Omega \subset \mathbb{R}^d$  be a bounded open set with a Lipschitz continuous boundary, for simplicity we assume that  $\Omega$  is a polygon/polyhedron. When  $\epsilon \in L^\infty(\Omega)$  is a nonnegative function that vanishes on a given part of the domain, an appropriate setting to approach advection-diffusion-reaction equations is the so called heterogeneous domain decomposition method developed in [10, 14]. Let  $\mu \in L^\infty(\Omega)$  be a positive function representing the reaction coefficient, let  $\beta \in [W^{1,\infty}(\Omega)]^d$  be the advection field such that  $\mu + \frac{1}{2} \nabla \cdot \beta \geq \mu_0 > 0$  and  $f \in L^2(\Omega)$ . In order to set up the heterogeneous method, we assume that  $\epsilon$  has some additional regularity, in particular there exist two open and measurable but not necessarily connected sets with Lipschitz continuous boundary,  $\Omega_{el}$  and  $\Omega_{hy}$ , such that

$$\Omega_{el} = \{x \in \Omega : \epsilon(x) > 0\}, \quad \Omega_{hy} = \Omega \setminus \bar{\Omega}_{el}.$$

Then, let  $\Gamma$  be the interface between the elliptic and hyperbolic regions, precisely  $\Gamma := \partial\Omega_{el} \cap \partial\Omega_{hy}$ . Furthermore, let  $n_{hy}$  be the outward unit normal with respect to  $\Omega_{hy}$  and be  $\partial\Omega_{in} := \{x \in \partial\Omega_{hy} : \beta \cdot n_{hy} < 0\}$ . The heterogeneous method

consists in finding a couple of functions  $u_{el}, u_{hy}$  such that,

$$\begin{cases} \nabla \cdot (-\epsilon \nabla u_{el} + \beta u_{el}) + \mu u_{el} = f & \text{in } \Omega_{el}, \\ \nabla \cdot (\beta u_{hy}) + \mu u_{hy} = f & \text{in } \Omega_{hy}, \\ -\epsilon \nabla u_{el} \cdot n + \beta \cdot n u_{el} = \beta \cdot n u_{hy} & \text{on } \Gamma, \\ u_{el} = u_{hy} & \text{on } \Gamma \cap \partial\Omega_{in}, \\ u_{el} = 0 & \text{on } \partial\Omega \cap \partial\Omega_{el}, \\ u_{hy} = 0 & \text{on } \partial\Omega \cap \partial\Omega_{in}, \end{cases} \quad (1)$$

where  $n$  is the unit normal vector associated to  $\Gamma$ . The weak formulation of this problem has been analyzed in [10] and it will be addressed later on. If the location of the interface  $\Gamma$  is known a-priori (typically all the cases where this line or surface is stationary and is described by a simple geometry) the numerical approximation of problem (1) is achieved in the framework of the domain decomposition methods. Furthermore, the numerical approximation of (1) has been recently reformulated in [5] to treat the case where  $\Gamma$  is any interface of discontinuity of  $\epsilon$ , and not only the one between the elliptic and the hyperbolic subregions. This leads to a mortar method that automatically adapts to the heterogeneous problem when  $\epsilon$  vanishes on some subregions of  $\Omega$ . In this work, we aim to apply the flexibility of DG methods in order to extend the ideas proposed in [5] to the case where the position of  $\Gamma$  is a-priori unknown.

To start with, we consider  $T_h$ , a shape regular triangulation of the domain  $\Omega$ , and we denote with  $K$  an element in  $T_h$  and with  $n_{\partial K}$  its outward unit normal. Let  $e$  be an edge (or face) of the element  $K \in T_h$ , which is an open simplex in  $\Omega$ . We say that  $e$  is an interior edge of the mesh if there are  $K^-(e)$  and  $K^+(e)$  in  $T_h$  such that  $e = \partial K^-(e) \cap \partial K^+(e)$ . We set  $K(e) = \{K^-(e), K^+(e)\}$  and let  $n_e$  be the unit normal vector to  $e$  pointing from  $K^-(e)$  towards  $K^+(e)$ . The analysis hereafter does not depend on the arbitrariness of this choice. Similarly, we say that  $e$  is a boundary face of the mesh if  $e = \partial K(e) \cap \partial\Omega$ . We denote with  $F_h, F_h^i$  and  $F_h^{\partial\Omega}$  the collections of all edges, of all the internal edges and of all the edges on  $\partial\Omega$  respectively. Finally, let  $h_e$  be the size of an edge and  $h_K$  be the one of an element.

Let  $H^s(T_h)$  be the broken Sobolev space of degree  $s > 0$  on  $T_h$ . For any  $v \in H^2(T_h)$  we introduce the following definitions,

$$\begin{aligned} v(x)|_e^\mp &:= \lim_{\delta \rightarrow 0^+} v(x \mp \delta n_e) \text{ for a.e. } x \in e, \\ v(x)|_{\partial K}^\mp &:= \lim_{\delta \rightarrow 0^+} v(x \mp \delta n_{\partial K}) \text{ for a.e. } x \in \partial K \setminus \partial\Omega. \end{aligned}$$

When not otherwise indicated, the values  $v|_{\partial K}^-, v|_e^-$  are implied. Moreover, the subscripts  $e$  and  $\partial K$  will be omitted if not necessary to the context. In this case, we will use the abridged notation  $v^\mp(x)$ . The jump over edges is defined as  $[[v(x)]]_e := v(x)|_e^- - v(x)|_e^+$ , while  $[[v(x)]]_{\partial K} := v(x)|_{\partial K}^- - v(x)|_{\partial K}^+$  is the jump with respect to element's outward normal vectors and we denote the arithmetic mean with  $\{v(x)\} := \frac{1}{2}v^-(x) + \frac{1}{2}v^+(x)$ .

To develop an approximation scheme for advection-diffusion-reaction equations with non-smooth and possibly vanishing diffusivity, we replace  $\epsilon$  with  $\epsilon_h$ , a discrete approximation of  $\epsilon$  onto a totally discontinuous approximation space defined on  $T_h$ . We notice that  $\epsilon_h$  must be either uniformly strictly positive either null on any  $K \in T_h$ , more precisely  $\epsilon_h(x) > 0$  or  $\epsilon_h(x) = 0$  for all  $x \in K$ , for all  $K \in T_h$ . This assumption is essential for the development of a numerical scheme, because we need that any elliptic/hyperbolic interface coincides with the edges of  $T_h$ . In other words, the case of an element  $K \in T_h$  where  $\epsilon_h > 0$  and  $\epsilon_h = 0$  simultaneously, is not admissible. We also observe that the introduction of a discrete vector field  $\beta_h$  may be useful from the practical point of view, but it is not strictly necessary.

In order to develop a suitable method to approximate a problem that may change in character from element to element, we start from a local problem formulation, already proposed in [4] for the Laplace equation. The idea consists in splitting problem (1) in subproblems localized to each element  $K \in T_h$ , complementing them with suitable matching conditions on inter-element interfaces and with boundary conditions on  $\partial\Omega$ . The key point is to set up suitable interface conditions that automatically adapt to the variations of  $\epsilon_h$ . As shown in problem (1), we need two types of interface conditions, one enforcing the continuity of the fluxes and one enforcing the continuity of the solution, when necessary. Let us focus on the latter. If  $\epsilon_h$  were positive and quasi-uniform on  $\Omega$ , the standard condition to enforce the continuity of the solution would be,

$$\left[\frac{1}{2}|\beta \cdot n_{\partial K}| + \{\epsilon_h\}\right][u]_{\partial K} = 0 \quad \text{on } \partial K \setminus \partial\Omega.$$

In order to correct this condition in the case where  $\epsilon_h$  changes from element to element, we introduce an *heterogeneity factor*, which quantifies the variation of  $\epsilon_h$  on each inter-element interface,  $\lambda_h(x)|_{\partial K} : \partial K \setminus \partial\Omega \rightarrow [-1, 1]$  such that

$$\lambda_h(x)|_{\partial K} := \begin{cases} \frac{1}{2} \frac{[\epsilon_h(x)]_{\partial K}}{\{\epsilon_h(x)\}}, & \text{if } \{\epsilon_h(x)\} > 0, \\ 0, & \text{if } \{\epsilon_h(x)\} = 0. \end{cases}$$

Then, starting from the case of uniform diffusivity considered above, we propose the following generalized interface conditions for the continuity of the solution,

$$\left[\frac{1}{2}|\beta \cdot n_{\partial K}|(1 - \text{sign}(\beta \cdot n_{\partial K})\varphi_{\partial K}(\lambda_h)) + \{\epsilon_h\}(1 - \chi_{\partial K}(\lambda_h))\right][u]_{\partial K} = 0 \quad \text{on } \partial K \setminus \partial\Omega,$$

where  $\varphi_{\partial K}(\lambda_h)$  and  $\chi_{\partial K}(\lambda_h)$  are chosen to make each local problem to be well posed. To this purpose, we assume that they satisfy  $|\chi_{\partial K}(\lambda_h)| \leq 1$ ,  $|\varphi_{\partial K}(\lambda_h)| \leq 1$  and,

$$\begin{aligned} \chi_{\partial K}(\lambda_h) &= 0 & \text{if } \lambda_h|_{\partial K} &= 0, \\ \chi_{\partial K}(\lambda_h) &= 1 & \text{if } \lambda_h|_{\partial K} &= \pm 1 \\ \varphi_{\partial K}(\lambda_h) &= 0 & \text{if } \lambda_h|_{\partial K} &= 0, \\ \varphi_{\partial K}(\lambda_h) &= \mp 1 & \text{if } \lambda_h|_{\partial K} &= \pm 1. \end{aligned} \tag{2}$$

According to these properties, we further assume that  $\chi_{\partial K}(\lambda_h)$  is a symmetric function while  $\varphi_{\partial K}(\lambda_h)$  is skewsymmetric. Then, setting  $\sigma(v) = -\epsilon_h \nabla v + \beta v$ ,

we reformulate problem (1) as follows: we look for a function  $u$  such that for all  $K \in T_h$ ,

$$\begin{cases} \nabla \cdot (\sigma(u)) + \mu u = f & \text{in } K, \\ \llbracket \sigma(u) \rrbracket_{\partial K} \cdot n_{\partial K} = 0 & \text{on } \partial K \setminus \partial \Omega, \\ \left[ \frac{1}{2} (|\beta \cdot n| - \beta \cdot n) + \epsilon_h \right] u = 0 & \text{on } \partial K \cap \partial \Omega, \\ \left[ \frac{1}{2} |\beta \cdot n_{\partial K}| (1 - \text{sign}(\beta \cdot n_{\partial K})) \varphi_{\partial K}(\lambda_h) \right. \\ \quad \left. + \{\epsilon_h\} (1 - \chi_{\partial K}(\lambda_h)) \right] \llbracket u \rrbracket_{\partial K} = 0 & \text{on } \partial K \setminus \partial \Omega. \end{cases} \quad (3)$$

We notice that problem (3) is equivalent to problem (1) when  $\Gamma \subset F_h^i$ . This makes problem (3) the right starting point to approximate second order PDEs with nonnegative characteristic form and non-smooth coefficients, without exploiting any a-priori knowledge of the interface between the elliptic and hyperbolic subregions. For instance, this is particularly interesting in those cases where the interface evolves in time, because it overrides the application of computationally expensive interface tracking and re-meshing procedures.

We are now ready to introduce the weak formulation of problem (3). Let  $\Omega_{el}$  be the subset of  $\Omega$  where  $\epsilon_h > 0$  and let be  $\Omega_{hy} := \Omega \setminus \Omega_{el}$ . We denote with  $\partial \Omega_{el}$  and  $\partial \Omega_{hy}$  their boundary with  $n_{el}$  and  $n_{hy}$  their outward normal vectors and with  $\Gamma_h := \partial \Omega_{el} \cap \partial \Omega_{hy}$  their interface. Here, the subscript  $h$  reminds that  $\Gamma_h$  lies on the edges of  $T_h$  because  $\epsilon_h$  is either positive or null on its elements. We also define  $\Gamma_h^{in} := \{x \in \Gamma_h \mid \beta \cdot n_{hy} < 0\}$ ,  $\Gamma_h^{out} := \Gamma_h \setminus \Gamma_h^{in}$  and  $\partial \Omega_{hy}^{in} := \{x \in \partial \Omega_{hy} \cap \partial \Omega \mid \beta \cdot n_{hy} < 0\}$ ,  $\partial \Omega_{hy}^{out} := \{x \in \partial \Omega_{hy} \cap \partial \Omega \mid \beta \cdot n_{hy} > 0\}$ . In this setting, we introduce the following functional spaces,

$$\begin{aligned} V^{el} &:= H^1(\Omega_{el}), & V^{hy} &:= \{v \in L^2(\Omega_{hy}), \beta \cdot \nabla v \in L^2(\Omega_{hy})\}, \\ V_0^{el} &:= H_{\partial \Omega \cap \partial \Omega_{el}}^1(\Omega_{el}), & V_{\beta,0}^{hy} &:= \{v \in V^{hy}, \beta \cdot n v|_{\partial \Omega_{hy}^{in}} = 0\}. \end{aligned}$$

We also define  $V := V^{el} \times V^{hy}$  and  $V_{\beta,0} := V_0^{el} \times V_{\beta,0}^{hy}$ . Following the analysis pursued in [10], for any  $u := (u_{el}, u_{hy})$ ,  $v := (v_{el}, v_{hy}) \in V_{\beta,0}$  we introduce the bilinear form relative to problem (3),

$$\begin{aligned} a(u, v) &:= \int_{\Omega_{el}} \left( (\epsilon_h \nabla u_{el} - \beta u_{el}) \cdot \nabla v_{el} + \mu u_{el} v_{el} \right) + \int_{\Omega_{hy}} \left( -\beta u_{hy} \cdot \nabla v_{hy} + \mu u_{hy} v_{hy} \right) \\ &\quad + \int_{\Gamma_h^{in}} \beta \cdot n_e u_{el} \llbracket v \rrbracket_e + \int_{\Gamma_h^{out}} \beta \cdot n_e u_{hy} \llbracket v \rrbracket_e + \int_{\partial \Omega_{hy}^{out}} \beta \cdot n u_{hy} v_{hy}, \end{aligned}$$

and the corresponding right hand side  $\mathcal{F}(v) := \int_{\Omega_{el}} f v_{el} + \int_{\Omega_{hy}} f v_{hy}$ . Then, the weak formulation of problem (3) reads as follows: find  $u \in V_{\beta,0}$  such that,

$$a(u, v) = \mathcal{F}(v), \quad \forall v \in V_{\beta,0}. \quad (4)$$

For the analysis of problem (4) we refer to [10, 14]. Conversely, we aim to briefly discuss here some regularity properties that arise from the multi-domain formulation (3). Let us consider the case of  $\epsilon_h > 0$  for any  $K \in T_h$  and assume

that the coefficients  $\beta$  and  $\mu$  are regular enough to ensure that the operator  $\nabla \cdot \sigma(u) + \mu u$  is an isomorphism between  $H^{2+s}(K) \cap H_0^1(K)$  and  $H^s(K)$ , see [11]. Let us denote with  $\mathcal{H}_K(\lambda)$  the lifting on  $K$  of the Dirichlet data  $\lambda|_{\partial K}$ , obtained by means of the operator  $\nabla \cdot \sigma(u) + \mu u$ . By virtue of the regularity of  $\nabla \cdot \sigma(u) + \mu u$  and owing to the trace theorem, we assert that the operator  $\sigma(\mathcal{H}_K(\lambda)) \cdot n_e$  maps  $H^{\frac{3}{2}+s}(e)$  into  $H^{\frac{1}{2}+s}(e)$  for any  $e \in \partial K$ . Recalling now the fundamentals of domain-decomposition methods, see for instance [15], we observe that problem (3) is equivalent to determine a function  $\lambda \in H^{\frac{1}{2}}(\cup_{F_h^i} e)$ , such that  $\sum_{e \in F_h^i} [\sigma(\mathcal{H}_K(\lambda))] \cdot n_e = \mathcal{G}$ , where  $\mathcal{G} \in H^{\frac{1}{2}+s}(F_h^i)$  is a given right hand side depending on  $f$ . Thanks to the regularity of  $\mathcal{G}$  and of the Dirichlet to Neumann map  $\lambda \rightarrow \sigma(\mathcal{H}_K(\lambda)) \cdot n_e$ , we assert that  $\lambda \in H^{\frac{1}{2}}(\cup_{F_h^i} e) \cap H^{\frac{3}{2}+s}(F_h^i)$ . Finally, we reconstruct the solution  $u$  of problem (3) by means of the lifting operator  $\mathcal{H}_K(\lambda)$  on each element  $K \in T_h$ , and by virtue of its regularity we conclude that  $u \in H^{2+s}(T_h) \cap H_0^1(\Omega)$ . These observations justify the property  $u \in H^{p+1}(T_h)$  with  $p > 0$  that will be assumed later on for the analysis of our numerical scheme.

## 2.1 Derivation of the numerical method

To set up our discretization scheme we could exploit the general framework proposed in [4] where a discrete solution is obtained by weakly enforcing that the residuals corresponding to equations (3)<sub>1</sub>, (3)<sub>2</sub>, (3)<sub>3</sub> and (3)<sub>4</sub> are equal to zero. However, to put into evidence the role of weighted averages to ensure the consistency of the method, we prefer to derive the scheme following a more classical approach, while the interpretation on the framework of [4] will be discussed in the following section.

To start with, for any  $v \in H^2(T_h)$  we introduce the weighted averages,

$$\begin{aligned} \{v(x)\}_w &:= w_e^-(x)v^-(x) + w_e^+(x)v^+(x), \quad \forall x \in e, \forall e \in F_h^i, \\ \{v(x)\}^w &:= w_e^+(x)v^-(x) + w_e^-(x)v^+(x), \quad \forall x \in e, \forall e \in F_h^i, \end{aligned}$$

where the weights necessarily satisfy  $w_e^-(x) + w_e^+(x) = 1$ . We say that these averages are conjugate, because they fulfill the following identity,

$$[[uv]] = \{u\}_w [[v]] + \{v\}^w [[u]], \quad \forall u, v \in H^2(T_h). \quad (5)$$

The role of  $\{\cdot\}_w$  and  $\{\cdot\}^w$  can also be interchanged, but for symmetry this choice does not affect the final setting of the method. Let  $u \in V_{\beta,0}$  be the solution of problem (4), let  $T_h^{el}$  be the collection of elements  $K \in T_h$  such that  $K \subset \Omega_{el}$  and let  $T_h^{hy} := T_h \setminus T_h^{el}$ , we assume that  $u \in H^2(T_h^{el}) \cap V_{\beta,0}$ . Let us denote for simplicity  $\sigma := \sigma(u)$ . Then, starting from (3), for any  $v \in H^1(T_h)$  we obtain,

$$\begin{aligned} \int_{\Omega} f v &= \int_{\Omega} (\nabla \cdot \sigma v + \mu v) = \sum_{K \in T_h} \int_K (\nabla \cdot \sigma v + \mu v) \\ &= \sum_{K \in T_h} \left[ \int_K (-\sigma \cdot \nabla v + \mu v) + \int_{\partial K} \sigma \cdot n_{\partial K} v \right], \quad \forall v \in H^1(T_h). \quad (6) \end{aligned}$$

Then, considering the identity,

$$\sum_{K \in T_h} \int_{\partial K} \sigma \cdot n_{\partial K} v = \sum_{e \in F_h^i} \int_e \llbracket \sigma v \rrbracket_e \cdot n_e + \sum_{e \in F_h^{\partial \Omega}} \int_e (\sigma v) \cdot n,$$

and replacing it into (6), owing to (5) we obtain,

$$\begin{aligned} \sum_{e \in F_h^i} \int_e \left( \{\sigma\}_w \cdot n_e \llbracket v \rrbracket_e + \llbracket \sigma \rrbracket_e \cdot n_e \{v\}^w \right) + \sum_{e \in F_h^{\partial \Omega}} \int_e \sigma \cdot n v \\ + \sum_{K \in T_h} \int_K \left( -\sigma \cdot \nabla v + \mu u v \right) = \int_{\Omega} f v, \quad \forall v \in H^1(T_h). \end{aligned} \quad (7)$$

Now, we need to enforce conditions (3)<sub>2</sub>, (3)<sub>3</sub>, (3)<sub>4</sub> on each inter-element interface and on the boundary of the domain. To this aim, we introduce  $\lambda_h|_e := \llbracket \epsilon_h \rrbracket_e / 2 \{\epsilon_h\}_e$ . Then, we remind that the function  $\chi_{\partial K}(\cdot)$  is assumed to be symmetric, consequently for any  $e \in F_h^i$  with  $e = \partial K^+ \cap \partial K^-$ , we define  $\chi_e(\lambda_h|_e) := \chi_{\partial K}(\lambda_h|_{\partial K^+}) = \chi_{\partial K}(\lambda_h|_{\partial K^-})$ , although the sign of  $\lambda_h|_e$  is arbitrarily determined. Analogously, since  $\varphi_{\partial K}$  have to be skewsymmetric we set,

$$\beta \cdot n_e \varphi_e(\lambda_h|_e) := \beta \cdot n_{\partial K^+} \varphi_{\partial K}(\lambda_h|_{\partial K^+}) = \beta \cdot n_{\partial K^-} \varphi_{\partial K}(\lambda_h|_{\partial K^-}).$$

Then, on any edge  $e \in F_h^i$  we introduce,

$$\gamma_{h,e}(\epsilon_h, \beta) := \frac{1}{2} (|\beta \cdot n_e| - \beta \cdot n_e \varphi_e(\lambda_h|_e)) + \{\epsilon_h\} (1 - \chi_e(\lambda_h|_e) \xi h_e^{-1}). \quad (8)$$

Accordingly, for the boundary we set,  $\gamma_{h,\partial \Omega}(\epsilon_h, \beta) := \frac{1}{2} (|\beta \cdot n| - \beta \cdot n) + \epsilon_h \xi h_e^{-1}$ . We notice that we have adjusted the scaling between the advective and diffusive terms with the introduction of a factor  $\xi h_e^{-1}$ . This will lead to convenient error estimates in the energy norm. Furthermore, from now on the heterogeneity factor on each edge,  $\lambda_h|_e$  will be simply denoted with  $\lambda_h$ . Expression (8) allows us to rewrite (3)<sub>2</sub>, (3)<sub>3</sub> and (3)<sub>4</sub> in order to be applied on  $F_h^i$ ,

$$\begin{aligned} \llbracket \sigma \rrbracket_e \cdot n_e &= 0, & \text{on } e \in F_h^i, \\ \gamma_{h,e}(\epsilon_h, \beta) \llbracket u \rrbracket_e &= 0, & \text{on } e \in F_h^i, \\ \gamma_{h,\partial \Omega}(\epsilon_h, \beta) u &= 0, & \text{on } e \in F_h^{\partial \Omega}. \end{aligned} \quad (9)$$

The first of (9) is enforced weakly by replacing it into (7), while the second and the third of (9) are enforced at the discrete level by means of a penalty method.

For any integer  $p \geq 0$ , we define the classical totally discontinuous approximation spaces,

$$V_h^p := \{v_h \in L^2(\Omega); \forall K \in T_h, v_h|_K \in \mathbb{P}^p\}.$$

Starting from (7) and (9), we aim to find  $u_h \in V_h^p$  such that,

$$\begin{aligned}
& \sum_{K \in T_h} \int_K \left( -\sigma_h \cdot \nabla v_h + \mu u_h v_h \right) + \sum_{e \in F_h^{\partial\Omega}} \int_e \left( \sigma_h \cdot n v_h - \epsilon_h \nabla v_h \cdot n u_h \right) \\
& + \sum_{e \in F_h^i} \int_e \left( \{\sigma_h\}_w \cdot n_e \llbracket v_h \rrbracket_e - \{\epsilon_h \nabla v_h\}_w \cdot n_e \llbracket u_h \rrbracket_e \right) \\
& + \sum_{e \in F_h^i} \int_e \left[ \frac{1}{2} (|\beta \cdot n_e| - \beta \cdot n_e \varphi_e(\lambda_h)) + \{\epsilon_h\} (1 - \chi_e(\lambda_h)) \xi h_e^{-1} \right] \llbracket u_h \rrbracket_e \llbracket v_h \rrbracket_e \\
& + \sum_{e \in F_h^{\partial\Omega}} \int_e \left[ \frac{1}{2} (|\beta \cdot n| - \beta \cdot n) + \epsilon_h \xi h_e^{-1} \right] u_h v_h = \int_{\Omega} f v_h, \forall v_h \in V_h^p. \quad (10)
\end{aligned}$$

where  $\sigma_h := -\epsilon_h \nabla u_h + \beta u_h$  for simplicity. We notice that we have added the new terms  $\{\epsilon_h \nabla v_h\}_w \cdot n_e \llbracket u_h \rrbracket_e$  on  $F_h^i$  and  $\epsilon_h \nabla v_h \cdot n u_h$  on  $F_h^{\partial\Omega}$  to preserve symmetry. The left hand side of equation (10) can be split in two parts. The former corresponds to the symmetric terms and it reads as follows,

$$\begin{aligned}
a_h^s(u_h, v_h) & := \sum_{K \in T_h} \int_K \left[ \epsilon_h \nabla u_h \cdot \nabla v_h + \left( \mu + \frac{1}{2} \nabla \cdot \beta \right) u_h v_h \right. \\
& + \sum_{e \in F_h^i} \int_e \left[ -\{\epsilon_h \nabla u_h\}_w \cdot n_e \llbracket v_h \rrbracket_e - \{\epsilon_h \nabla v_h\}_w \cdot n_e \llbracket u_h \rrbracket_e \right. \\
& \quad \left. + \left( \frac{1}{2} |\beta \cdot n_{\partial K}| + \{\epsilon_h\} (1 - \chi_e(\lambda_h)) \right) \xi h_e^{-1} \llbracket u_h \rrbracket_e \llbracket v_h \rrbracket_e \right] \\
& + \sum_{e \in F_h^{\partial\Omega}} \int_e \left[ -\epsilon_h \nabla u_h \cdot n v_h - \epsilon_h \nabla v_h \cdot n u_h + \left( \frac{1}{2} |\beta \cdot n| + \epsilon_h \xi h_e^{-1} \right) u_h v_h \right].
\end{aligned}$$

The remaining part of the bilinear form is,

$$\begin{aligned}
a_h^r(u_h, v_h) & := - \sum_{K \in T_h} \int_K \left[ (\beta u_h) \cdot \nabla v_h + \frac{1}{2} (\nabla \cdot \beta) u_h v_h \right] \\
& + \sum_{e \in F_h^i} \int_e \left[ \{\beta u_h\}_w \cdot n_e \llbracket v_h \rrbracket_e - \frac{1}{2} \beta \cdot n_e \varphi_e(\lambda_h) \llbracket u_h \rrbracket_e \llbracket v_h \rrbracket_e \right] \\
& + \sum_{e \in F_h^{\partial\Omega}} \int_e \frac{1}{2} \beta \cdot n u_h v_h.
\end{aligned}$$

Then, setting  $a_h(u_h, v_h) := a_h^s(u_h, v_h) + a_h^r(u_h, v_h)$  and  $F(v_h) := \int_{\Omega} f v_h$ , our prototype of method reads as follows: find  $u_h \in V_h^p$  such that,

$$a_h(u_h, v_h) = F(v_h), \quad \forall v_h \in V_h^p. \quad (11)$$

We point out that (11) corresponds to a family of methods that may differ for the definition of the weights and for the expression of  $\chi_e$  and  $\varphi_e$  into (8). We will consider their choice in the following section.

## 2.2 Definition of the weights and of the scaling functions

In this section, we discuss how to identify specific choices of the weights  $w_e^-, w_e^+$  and of the scaling functions  $\chi_e$  and  $\varphi_e$  that basically influence the behavior of the method that we have proposed.

First of all, we aim to point out suitable recipes to define the weights on each edge. Our strategy is to introduce a suitable weighing function  $\phi$  that will be applied to construct the tilted weights depending on the heterogeneity factor,  $\lambda_h$ . Precisely, we set  $w_e^\pm := \phi(\pm\lambda_h)$ . Observing that  $\lambda_h(x) \in [-1, 1]$  we require that  $\phi$  satisfies the following general properties,

$$\phi(\cdot) \in C^0([-1, 1]), \quad \phi([-1, 1]) = [0, 1], \quad \phi(-t) + \phi(t) = 1, \quad \forall t \in [-1, 1],$$

with the following particular assumption, that is necessary to enable the method to satisfy (2),

$$\phi(-1) = 0, \quad \text{and} \quad \phi(1) = 1, \quad \text{or vice versa} \quad \phi(-1) = 1, \quad \text{and} \quad \phi(1) = 0. \quad (12)$$

By virtue of the partition of unity theorem, there are infinitely many functions satisfying these requirements. For instance, we propose the following family of weighing functions,

$$\phi(t) := \frac{1}{2}(1 \pm \text{sign}(t)|t|^\alpha), \quad (13)$$

where  $\alpha \in \mathbb{R}^+$  plays the role of *tilting factor*. We notice that at this level the choice of the sign into  $\phi(t)$  is arbitrary, it will be fixed later exploiting (2). Moreover, we observe that the smaller is  $\alpha$ , the more the weights  $w_e^\pm$  differ from the standard value  $w_e^\pm = \frac{1}{2}$ . For instance, when  $\phi(t) = \frac{1}{2}(1 + \text{sign}(t)|t|^\alpha)$ , we have the following limit cases,

$$\lim_{\alpha \rightarrow 0} \phi(t) = \begin{cases} 0, & \text{if } t \in [-1, 0), \\ \frac{1}{2}, & \text{if } t = 0, \\ 1, & \text{if } t \in (0, 1], \end{cases} \quad \lim_{\alpha \rightarrow \infty} \phi(t) = \begin{cases} 0, & \text{if } t = -1, \\ \frac{1}{2}, & \text{if } t \in (-1, 1), \\ 1, & \text{if } t = 1. \end{cases}$$

We will see later that another case of particular interest is  $\alpha = 1$ .

Secondly, we focus our attention on  $a^r(\cdot, \cdot)$ , the advective part of the bilinear form. We aim to find a suitable expression for  $\varphi_e$  that makes  $a^r(\cdot, \cdot)$  to be skewsymmetric. By means of integration by parts, equation (5) and exploiting the regularity of  $\beta \in [W_\infty^1(\Omega)]^d$ , we have that

$$\begin{aligned} a_h^r(u_h, v_h) &= \sum_{K \in \mathcal{T}_h} \int_K [(\beta v_h) \cdot \nabla u_h + \frac{1}{2}(\nabla \cdot \beta) v_h u_h] - \sum_{e \in F_h^{\partial\Omega}} \int_e \frac{1}{2} \beta \cdot n v_h u_h \\ &\quad - \sum_{e \in F_h^i} \int_e \left[ \beta \cdot n_e \{v_h\}^w \llbracket u_h \rrbracket_e + \frac{1}{2} \beta \cdot n_e \varphi_e(\lambda_h) \llbracket v_h \rrbracket_e \llbracket u_h \rrbracket_e \right]. \quad (14) \end{aligned}$$

For any  $e \in F_h^i$  we remind that  $\{v_h\}^w = \{v_h\}_w - (w_e^- - w_e^+) \llbracket v_h \rrbracket_e$ . Then, by means of the specific definition,

$$\varphi_e(\lambda_h) := (w_e^- - w_e^+) \quad (15)$$

we obtain the following identity,

$$\{v_h\}^w + \frac{1}{2}\varphi_e(\lambda_h)\llbracket v_h \rrbracket_e = \{v_h\}_w - \frac{1}{2}\varphi_e(\lambda_h)\llbracket v_h \rrbracket_e. \quad (16)$$

Replacing (16) into (14) we conclude that  $a_h^r(u_h, v_h) = -a_h^r(v_h, u_h)$ . Furthermore, owing to (15) we obtain the identity,

$$\beta \cdot n_e \{v_h\}_w - \frac{1}{2}\beta \cdot n_e (w_e^- - w_e^+) \llbracket v_h \rrbracket_e = \beta \cdot n_e \{v_h\},$$

which shows that the advective flux on each inter-element interface coincides with the *standard upwind flux*.

Finally, we study how to choose  $\chi_e$ . We propose a particular choice of  $\chi_e$  that allows us to rewrite our scheme in the framework introduced in [4] for the set up and the analysis of DG methods. In particular, we require that  $\chi_e$  satisfies the following property,

$$\{\epsilon_h\}(1 - \chi_e(\lambda_h)) = \{\epsilon_h\}_w, \quad \text{i.e.} \quad \chi_e(\lambda_h) = (w_e^+ - w_e^-)\lambda_h. \quad (17)$$

Let us now consider (15) and (17) and replace  $\phi(\cdot)$  with (13). We obtain the following expressions,

$$\varphi_e(\lambda_h) = \mp \text{sign}(\lambda_h) |\lambda_h|^\alpha, \quad \chi_e(\lambda_h) = \pm |\lambda_h|^{\alpha+1},$$

and we observe that only the choice  $\phi(\lambda_h) = \frac{1}{2}(1 + \text{sign}(\lambda_h)|\lambda_h|^\alpha)$  allows to satisfy (2) for both  $\varphi_e$  and  $\chi_e$ . By consequence the ambiguity in definition (13) is resolved by choosing the positive sign.

Then, to rewrite (10) in the framework proposed in [4], we set  $\omega_e^\pm = w_e^\pm \epsilon_h^\pm$  and we exploit the identity

$$\{\epsilon_h v\}_w = \{\epsilon_h\}_w \{v\} + \frac{1}{2} \llbracket \omega \rrbracket_e \llbracket v \rrbracket_e, \quad \forall v \in H^2(T_h),$$

leading to the following equivalence,

$$\{\epsilon_h \nabla v_h\}_w \cdot n_e \llbracket u_h \rrbracket_e = \{\epsilon_h\}_w \llbracket u_h \rrbracket_e \left( \{\nabla v_h\} \cdot n_e + \frac{\llbracket \omega \rrbracket_e}{4\{\omega\}} \llbracket \nabla v_h \rrbracket_e \cdot n_e \right). \quad (18)$$

Applying integration by parts into (10) over the term  $\int_K -\sigma_h \cdot \nabla v_h$  and exploit-

ing (17) and (18), we obtain the following equation,

$$\begin{aligned}
& \sum_{K \in T_h} \int_K \left( \nabla \cdot \sigma_h + \mu u_h - f \right) v_h + \sum_{e \in F_h^i} \int_e \llbracket \sigma_h \rrbracket_e \cdot n_e \{v_h\}^w \\
& + \sum_{e \in F_h^i} \{\epsilon_h\}_w \llbracket u_h \rrbracket_e \left( h_e^{-1} \llbracket v_h \rrbracket_e - \{\nabla v_h\} \cdot n_e - \frac{\llbracket \omega \rrbracket_e}{4\{\omega\}} \llbracket \nabla v_h \rrbracket_e \cdot n_e \right) \\
& + \sum_{e \in F_h^i} \int_e \frac{1}{2} (|\beta \cdot n_e| - \beta \cdot n_e \varphi_e(\lambda_h)) \llbracket u_h \rrbracket_e \llbracket v_h \rrbracket_e \\
& + \sum_{e \in F_h^{\partial\Omega}} \int_e (\epsilon_h u_h) (h_e^{-1} v_h - \nabla v_h \cdot n) \\
& + \sum_{e \in F_h^{\partial\Omega}} \int_e \frac{1}{2} (|\beta \cdot n| - \beta \cdot n) u_h v_h = 0, \quad \forall v_h \in V_h^p. \tag{19}
\end{aligned}$$

According to [4], we observe that equation (19) can be regarded as the weak counterpart of (3) obtained by enforcing that a linear combination of the residuals associated to equations (3)<sub>1</sub>-(3)<sub>4</sub> is equal to zero. In other words, equation (19) can be rewritten as,

$$\begin{aligned}
& \sum_{K \in T_h} \int_K R_0(u_h) W_0(v_h) + \sum_{e \in F_h^i} \left[ R_1(u_h) W_1(v_h) + R_2(u_h) W_2(v_h) + R_3(u_h) W_3(v_h) \right] \\
& \sum_{e \in F_h^{\partial\Omega}} \left[ R_4(u_h) W_4(v_h) + R_5(u_h) W_5(v_h) \right] = 0, \quad \forall v_h \in V_h^p, \tag{20}
\end{aligned}$$

being  $R_i(u_h)$  the residuals associated to (3)<sub>1</sub>-(3)<sub>4</sub> in the discrete case, while  $W_i(v_h)$  are suitable test functions. Precisely, they read as follows,

$$\begin{aligned}
R_0(u_h) &= \nabla \cdot \sigma_h + \mu u_h - f, & W_0(v_h) &= v_h, \\
R_1(u_h) &= \llbracket \sigma_h \rrbracket_e \cdot n_e, & W_1(v_h) &= \{v_h\}^w, \\
R_2(u_h) &= \{\epsilon_h\}_w \llbracket u_h \rrbracket_e, & W_2(v_h) &= h_e^{-1} \llbracket v_h \rrbracket_e \\
& & & - (\{\nabla v_h\} + \frac{\llbracket \omega \rrbracket_e}{4\{\omega\}} \llbracket \nabla v_h \rrbracket_e) \cdot n_e, \\
R_3(u_h) &= \frac{1}{2} (|\beta \cdot n_e| - \beta \cdot n_e \varphi_e(\lambda_h)) \llbracket u_h \rrbracket_e, & W_3(v_h) &= \llbracket v_h \rrbracket_e, \\
R_4(u_h) &= \epsilon_h u_h, & W_4(v_h) &= h_e^{-1} v_h - \nabla v_h \cdot n, \\
R_5(u_h) &= \frac{1}{2} (|\beta \cdot n| - \beta \cdot n) u_h, & W_5(v_h) &= v_h, \tag{21}
\end{aligned}$$

In conclusion, we notice that the specific choice of  $\chi_e$  proposed in (17) is particularly interesting because it allows us to rewrite the method (11) in the general framework (20). An interesting feature is to obtain a formulation where the test functions  $W_i(v_h)$  are independent from the coefficients of the problem. This property is not satisfied in our case, since both  $W_1(v_h)$  and  $W_2(v_h)$  depend on the weights  $w_e^\pm$  and thus on  $\epsilon_h$ . However, we notice that the exception

of  $W_1(v_h)$  is of minor importance with respect to our purpose, because it can be easily verified that the term involving  $W_1(v_h)$  cancels out going back from (20) to (10). Moreover, in the particular case  $[\![\omega]\!]_e = 0$ , we obtain that the test function  $W_2(v_h)$  becomes  $W_2(v_h) = h_e^{-1}[\![v_h]\!]_e - \{\nabla v_h\} \cdot n_e$  and thus it is independent on  $\epsilon_h$ . In this case the weights have to satisfy  $w_e^+ \epsilon_h^+ = w_e^- \epsilon_h^-$  and by consequence they are defined as follows,

$$\begin{aligned} w_e^- &= \frac{\epsilon_h^+}{\epsilon_h^- + \epsilon_h^+}, \quad w_e^+ = \frac{\epsilon_h^-}{\epsilon_h^- + \epsilon_h^+}, \quad \text{if } \epsilon_h^- + \epsilon_h^+ > 0, \\ w_e^- &= w_e^+ = \frac{1}{2}, \quad \text{if } \epsilon_h^- = \epsilon_h^+ = 0, \end{aligned} \quad (22)$$

which correspond to choose  $w_e^\pm = \frac{1}{2}(1 \pm \lambda_h)$  or equivalently  $\phi(t) = \frac{1}{2}(1 + t)$  that means  $\alpha = 1$  into (13). Equation (22) implies that the scaling factor of the penalty term, namely  $\{\epsilon_h\}_w$ , is equivalent to the harmonic average of the values  $\epsilon_h^\pm$ . Consequently, the mortar between elements is proportional to the stiffness of two sequential springs of modulus  $\epsilon_h^\pm$  respectively. This interpretation suggests that (22), and thus  $\alpha = 1$ , seems to be a very natural choice for problems with discontinuous coefficients.

### 2.3 Reinterpretation by means of numerical fluxes

Several DG methods are set up by means of numerical fluxes, denoted here with  $\tilde{u}_h$  and  $\tilde{\sigma}_h \cdot n_e$ , which represent suitable approximations of  $u$  and  $\sigma(u) \cdot n_e$  on the edges of  $T_h$ . In order to highlight the relationship between the method proposed here and other DG approximation schemes, we aim to identify the numerical fluxes corresponding to our scheme. Following the paradigm presented in [2], we start from the discretization of the governing equations  $\nabla \cdot \sigma + \mu u = f$  and  $\sigma = -\epsilon_h \nabla u + \beta u$ , and exploiting the fluxes  $\tilde{u}_h$  and  $\tilde{\sigma}_h \cdot n_e$  we obtain,

$$\begin{aligned} \sum_{K \in T_h} \int_K f v_h &= \sum_{K \in T_h} \int_K \left( -\sigma_h \cdot \nabla v_h + \mu u_h v_h \right) \\ &+ \sum_{e \in F_h^i} \left( \{\tilde{\sigma}_h\}_w \cdot n_e [\![v_h]\!]_e + [\![\tilde{\sigma}_h]\!]_e \cdot n_e \{v_h\}^w \right) + \sum_{e \in F_h^{\partial\Omega}} \tilde{\sigma}_h \cdot n v_h, \end{aligned} \quad (23)$$

$$\begin{aligned} &- \sum_{K \in T_h} \int_K \sigma_h \cdot \nabla v_h = - \sum_{K \in T_h} \int_K \left( \epsilon_h u_h \Delta v_h + (\beta u_h) \cdot \nabla v_h \right) \\ &+ \sum_{e \in F_h^i} \left( \{\epsilon_h \nabla v_h\}_w \cdot n_e [\![\tilde{u}_h]\!]_e + [\![\epsilon_h \nabla v_h]\!]_e \cdot n_e \{\tilde{u}_h\}^w \right) \\ &+ \sum_{e \in F_h^{\partial\Omega}} \epsilon_h \tilde{u}_h \nabla v_h \cdot n. \end{aligned} \quad (24)$$

Moreover, by means of integration by parts we obtain,

$$\begin{aligned} & - \sum_{K \in \mathcal{T}_h} \int_K \epsilon_h u_h \Delta v_h = \sum_{K \in \mathcal{T}_h} \int_K \epsilon_h \nabla u_h \cdot \nabla v_h \\ & - \sum_{e \in F_h^i} \left( \{\epsilon_h \nabla v_h\}_w \cdot n_e \llbracket u_h \rrbracket_e + \llbracket \epsilon_h \nabla v_h \rrbracket_e \cdot n_e \{u_h\}^w \right) - \sum_{e \in F_h^{\partial\Omega}} \epsilon_h u_h \nabla v_h \cdot n, \end{aligned}$$

which can be replaced into (24) in order to give,

$$\begin{aligned} & - \sum_{K \in \mathcal{T}_h} \int_K \sigma_h \cdot \nabla v_h = \sum_{K \in \mathcal{T}_h} \int_K \left( \epsilon_h \nabla u_h - \beta u_h \right) \cdot \nabla v_h + \sum_{e \in F_h^{\partial\Omega}} \epsilon_h (\tilde{u}_h - u_h) \nabla v_h \cdot n \\ & + \sum_{e \in F_h^i} \left( \{\epsilon_h \nabla v_h\}_w \cdot n_e \llbracket \tilde{u}_h - u_h \rrbracket_e + \llbracket \epsilon_h \nabla v_h \rrbracket_e \cdot n_e \{\tilde{u}_h - u_h\}^w \right). \quad (25) \end{aligned}$$

Combining equations (23) and (25) and reminding that  $\sigma_h = -\epsilon_h \nabla u_h + \beta u_h$ , we conclude that the bilinear form corresponding to the fluxes  $\tilde{u}_h$  and  $\tilde{\sigma}_h \cdot n_e$  is,

$$\begin{aligned} \tilde{a}_h(u_h, v_h) & := \sum_{K \in \mathcal{T}_h} \int_K \left( -\sigma_h \cdot \nabla v_h + \mu u_h v_h \right) \\ & + \sum_{e \in F_h^i} \left( \{\tilde{\sigma}_h\}_w \cdot n_e \llbracket v_h \rrbracket_e + \llbracket \tilde{\sigma}_h \rrbracket_e \cdot n_e \{v_h\}^w \right) \\ & + \sum_{e \in F_h^i} \left( \{\epsilon_h \nabla v_h\}_w \cdot n_e \llbracket \tilde{u}_h - u_h \rrbracket_e + \llbracket \epsilon_h \nabla v_h \rrbracket_e \cdot n_e \{\tilde{u}_h - u_h\}^w \right) \\ & + \sum_{e \in F_h^{\partial\Omega}} \left( \tilde{\sigma}_h \cdot n v_h + \epsilon_h (\tilde{u}_h - u_h) \nabla v_h \cdot n \right). \end{aligned}$$

Now we compare  $\tilde{a}_h(u_h, v_h)$  with  $a(u_h, v_h)$  and we observe that these two forms are equivalent provided that,

$$\tilde{\sigma}_h \cdot n_e = \{\sigma_h\}_w \cdot n_e + \gamma_{h,e}(\epsilon_h, \beta) \llbracket u_h \rrbracket_e, \quad \tilde{u}_h = \{u_h\}^w. \quad (26)$$

Indeed, owing to (26) we obtain,

$$\begin{aligned} & \{\tilde{\sigma}_h\}_w \cdot n_e = \tilde{\sigma}_h \cdot n_e \quad \text{and} \quad \llbracket \tilde{\sigma}_h \rrbracket_e \cdot n_e = 0, \\ & \{\tilde{u}_h - u_h\}^w = \{\tilde{u}_h\}^w - \{u_h\}^w = 0 \quad \text{and} \quad \llbracket \tilde{u}_h - u_h \rrbracket_e = \llbracket \tilde{u}_h \rrbracket_e - \llbracket u_h \rrbracket_e = -\llbracket u_h \rrbracket_e, \end{aligned}$$

which directly imply  $\tilde{a}_h(u_h, v_h) = a(u_h, v_h)$ . In conclusion, the numerical fluxes associated to our weighted interior penalty scheme are given by (26). According to the terminology introduced by [2] such fluxes are *conservative*. Finally, a simple manipulation of (26) puts into evidence a close connection between our scheme and the LDG methods. Indeed, reminding (15), for any  $v_h \in V_h^p$  we obtain,

$$\{v_h\}_w = \{v_h\} + \frac{1}{2} \varphi_e \llbracket v_h \rrbracket_e, \quad \{v_h\}^w = \{v_h\} - \frac{1}{2} \varphi_e \llbracket v_h \rrbracket_e,$$

and consequently (26) can be rewritten as,

$$\begin{aligned}\tilde{\sigma}_h \cdot n_e &= \{\sigma_h\} \cdot n_e + \frac{1}{2}\varphi_e \llbracket \sigma_h \rrbracket_e \cdot n_e + \gamma_{h,e} \llbracket u_h \rrbracket_e, \\ \tilde{u}_h &= \{u_h\} - \frac{1}{2}\varphi_e \llbracket u_h \rrbracket_e.\end{aligned}$$

On the edges where  $\lambda_h \neq 0$  and thus  $\varphi_e(\lambda_h) \neq 0$ , these fluxes turn out to be equivalent to the fluxes of the LDG method, proposed in [6], if we set  $C_{11} = \gamma_{h,e}$ ,  $C_{12} = \frac{1}{2}\varphi_e$  and  $C_{22} = 0$ , according to the terminology defined there. A specific choice for the parameter  $C_{12}$  together with the expression  $\tilde{u}_h = w_e^+ u_h^+ + w_e^- u_h^-$  with  $w_e^+, w_e^- \neq \frac{1}{2}$  has already been applied in [7], in order to prove a superconvergence property for the approximation of the Poisson problem. However, the ideas underlying [7] are very different from ours. Indeed, definition (26) can be seen as a new recipe to properly scale the characteristic parameters  $C_{11}$  and  $C_{12}$  of the LDG method with respect to the coefficients of an advection-diffusion-reaction problem. In other words, the WIP scheme is capable to automatically adapt its numerical fluxes to the degree of smoothness of the diffusion coefficient. In particular, WIP coincides with the symmetric interior penalty scheme on the edges where  $\epsilon_h$  is continuous, while it switches to the LDG formulation where  $\epsilon_h$  is discontinuous and thus  $\lambda_h \neq 0$ .

### 3 A-priori error analysis

The goal of this section is to establish an error estimate being robust with respect to locally vanishing diffusion. The analysis is performed in the spirit of Strang's Second Lemma by addressing the consistency, coercivity and continuity properties of the bilinear form. To obtain such result, we follow the analysis proposed in [9], with the necessary modifications to adapt it to the method at hand.

For any  $v \in V$  we consider the energy norm associated to problem (11),

$$\begin{aligned}\|v\|_a^2 &:= \|\epsilon_h^{\frac{1}{2}} \nabla v\|_{0,T_h}^2 + \|\mu_0^{\frac{1}{2}} v\|_{0,T_h}^2 + \|(\frac{1}{2}|\beta \cdot n| + \epsilon_h h_e^{-1})^{\frac{1}{2}} v\|_{0,F_h^{\partial\Omega}}^2 \\ &\quad + \|(\frac{1}{2}|\beta \cdot n_e| + \{\epsilon_h\}_w h_e^{-1})^{\frac{1}{2}} \llbracket v \rrbracket_e\|_{0,F_h^i}^2\end{aligned}$$

and for any  $v \in V \cap H^2(T_h^{el})$  we introduce the following augmented norm,

$$\begin{aligned}\|v\|_{\perp}^2 &:= \|v\|_a^2 + \| |\beta \cdot n_e|^{\frac{1}{2}} v \|_{0,F_h^i}^2 + \|(\epsilon_h h_e)^{\frac{1}{2}} \nabla v \cdot n\|_{0,F_h^{\partial\Omega}}^2 \\ &\quad + \|(\epsilon_h h_e)^{\frac{1}{2}} \nabla v \cdot n_e\|_{0,F_h^i}^2,\end{aligned}$$

where  $\|\cdot\|_{0,K}$ ,  $\|\cdot\|_{0,e}$  are the  $L^2$ -norms on  $K$ ,  $e$  respectively, and,

$$\|v\|_{0,T_h}^2 := \sum_{K \in T_h} \|v\|_{0,K}^2, \quad \|v\|_{0,F_h^i}^2 := \sum_{e \in F_h^i} \|v\|_{0,e}^2, \quad \|v\|_{0,F_h^{\partial\Omega}}^2 := \sum_{e \in F_h^{\partial\Omega}} \|v\|_{0,e}^2.$$

In the analysis, we will make use of the following inverse inequalities, see [1], that hold true for all  $K \in T_h$  and for all  $v_h \in V_h^p$ , provided that  $T_h$  is shape regular,

$$h_e^{\frac{1}{2}} \|v_h\|_{0,e} \lesssim \|v_h\|_{0,K}, \quad (27)$$

$$h_K \|\nabla v_h\|_{0,K} \lesssim \|v_h\|_{0,K}. \quad (28)$$

Here and in the sequel, the symbol  $\lesssim$  denotes an inequality involving a positive constant  $C$  independent of the size of the mesh family and of the diffusion parameter. Moreover, we will make use of the following approximation properties of  $V_h^p$ . Let  $\pi_h : H^{p+1}(T_h) \rightarrow V_h^p$  be the  $L^2$ -projection operator onto  $V_h^p$ . Then, it satisfies,

$$\|v - \pi_h v\|_{0,K} + h_K \|\nabla(v - \pi_h v)\|_{0,K} \lesssim h_K^{p+1} |v|_{p+1,K}, \quad (29)$$

$$\|v - \pi_h v\|_{0,e} + h_e \|\nabla(v - \pi_h v) \cdot n_e\|_{0,e} \lesssim h_K^{p+\frac{1}{2}} |v|_{p+1,K}, \quad (30)$$

where  $|v|_{s,K}$  is the seminorm in  $H^s(K)$ .

**Lemma 1.** *For any  $v \in H^{p+1}(T_h)$  with  $p > 0$  we have,*

$$\| \|v - \pi_h v\| \|_* \lesssim \sum_{K \in T_h} (h_K^p \|\epsilon_h\|_{L^\infty(K)} + h_K^{p+\frac{1}{2}} \|\beta\|_{W_\infty^1(K)} + h_K^{p+1} \mu_0) |v|_{p+1,K},$$

where  $\| \cdot \|_*$  represents both  $\| \cdot \|_a$  and  $\| \cdot \|_\perp$ .

*Proof.* The result follows immediately from the application of (29) and (30) into  $\| \|v - \pi_h v\| \|_a$  and  $\| \|v - \pi_h v\| \|_\perp$ .  $\square$

**Lemma 2.** *Let  $u \in V_{\beta,0} \cap H^2(T_h^{el})$  be the solution of problem (4). Then,  $a_h(u, v) = F(v)$ ,  $\forall v \in (V \cap H^2(T_h^{el})) \oplus V_h^p$  and  $a_h(u - u_h, v_h) = 0 \forall v_h \in V_h^p$ .*

*Proof.* Let us consider  $a_h(u, v)$ . First of all, we observe that  $u \in H^2(T_h^{el})$  implies that  $\llbracket u \rrbracket_e = 0$  in the sense of traces for any edge such that  $K(e) \in T_h^{el}$ . Analogously,  $\{\epsilon_h \nabla v\}_w \cdot n_e = 0$  for any edge such that  $K(e) \in T_h^{hy}$ . Finally, by virtue of the definition of the weights, see (12), we obtain that  $\{\epsilon_h \nabla v\}_w \cdot n_e = 0$  for all  $e \in \Gamma_h$ . Then, the term  $\{\epsilon_h \nabla v\}_w \cdot n_e \llbracket u \rrbracket_e$  is equal to zero on any  $e \in F_h^i$ , because either  $\llbracket u \rrbracket_e = 0$  and  $\{\epsilon_h \nabla v\}_w \neq 0$  or  $\llbracket u \rrbracket_e \neq 0$  and  $\{\epsilon_h \nabla v\}_w = 0$ . Proceeding similarly, we obtain that  $\epsilon_h \nabla v \cdot nu = 0$ . For the same reasons, we assert that  $\gamma_{h,e}(\epsilon_h, \beta) \llbracket u \rrbracket_e \llbracket v \rrbracket_e = 0$  for any internal edge of  $T_h^{el}$  while for the edges inside  $T_h^{hy}$  we may have  $\llbracket u \rrbracket_e \neq 0$  but  $\beta \cdot n_e \llbracket u \rrbracket_e = 0$ , since  $u \in V$  and  $\beta \in W_\infty^1(\Omega)$ . Moreover, since  $\llbracket u \rrbracket_e = 0$  on  $\Gamma_h^{in}$ , and  $\gamma_{h,e}(\epsilon_h, \beta) = 0$  on  $\Gamma_h^{out}$ , we conclude that  $\gamma_{h,e}(\epsilon_h, \beta) \llbracket u \rrbracket_e = 0$ , for any  $e \in F_h^i$ . Finally, we observe that  $u \in V \cap H^2(T_h^{el})$  ensures that  $\llbracket \sigma \rrbracket_e \cdot n_e = 0$  in the sense of traces for any  $e \in F_h^i$ . By consequence,  $a(u, v)$  is equivalent to the second row of (6). Then, integrating by parts and applying the equivalence  $\nabla \cdot \sigma + \mu u = f$  in  $L^2(\Omega)$ , we obtain  $a(u, v) = F(v)$  for all  $v \in (V \cap H^2(T_h^{el})) \oplus V_h^p$ . Finally, setting  $v = v_h$ , the Galerkin orthogonality  $a_h(u - u_h, v_h) = 0 \forall v_h \in V_h^p$  follows immediately.  $\square$

**Lemma 3.** *Choosing  $\varphi_e$  as in equation (15), namely  $\varphi_e(\lambda_h) := (w_e^- - w_e^+)$ , the bilinear form  $a_h^r(\cdot, \cdot)$  skew-symmetric, more precisely  $a_h^r(u, v) = -a^r(v, u)$  for all  $u, v \in (V \cap H^2(T_h^{el})) \oplus V_h^p$ .*

*Proof.* The proof follows immediately from the combination of equations (14) and (16).  $\square$

**Lemma 4.** *Choosing  $\xi$  large enough in the definition of  $\gamma_{h,e}(\epsilon_h, \beta)$  and  $\gamma_{h,\partial\Omega}(\epsilon_h, \beta)$ , see equation (8), the bilinear form  $a_h^s(\cdot, \cdot)$  is coercive in the energy norm, precisely  $\|v_h\|_a \lesssim a_h^s(v_h, v_h)$  for all  $v_h \in V_h^p$ .*

*Proof.* First of all, it is straightforward to verify that,

$$\begin{aligned}
& \sum_{K \in T_h} \int_K \left[ \epsilon_h (\nabla v_h)^2 + \left( \mu + \frac{1}{2} \nabla \cdot \beta \right) v_h^2 \right] \\
& + \sum_{e \in F_h^i} \int_e \left( \frac{1}{2} |\beta \cdot n_e| + \{\epsilon_h\}_w \xi h_e^{-1} \right) \llbracket v_h \rrbracket_e^2 + \sum_{e \in F_h^{\partial\Omega}} \int_e \left( \frac{1}{2} |\beta \cdot n| + \epsilon_h \xi h_e^{-1} \right) v_h^2 \\
& \geq \|\epsilon_h^{\frac{1}{2}} \nabla v_h\|_{0, T_h}^2 + \|\mu_0^{\frac{1}{2}} v_h\|_{0, T_h}^2 + \left\| \left( \frac{1}{2} |\beta \cdot n_e| + \{\epsilon_h\}_w \xi h_e^{-1} \right)^{\frac{1}{2}} \llbracket v_h \rrbracket_e \right\|_{0, F_h^i}^2 \\
& + \left\| \left( \frac{1}{2} |\beta \cdot n| + \epsilon_h \xi h_e^{-1} \right)^{\frac{1}{2}} v_h \right\|_{0, F_h^{\partial\Omega}}^2. \tag{31}
\end{aligned}$$

For the remaining terms of  $a_h^s(\cdot, \cdot)$ , reminding that  $\omega_e^\pm = w_e^\pm \epsilon_h^\pm$  and  $\omega_e^+ + \omega_e^- = \{\epsilon_h\}_w$ , we obtain the following bounds,

$$\begin{aligned}
& 2 \sum_{e \in F_h^i} \int_e \{\epsilon_h \nabla v_h\}_w \cdot n_e \llbracket v_h \rrbracket_e + 2 \sum_{e \in F_h^{\partial\Omega}} \int_e \epsilon_h \nabla v_h \cdot n v_h \\
& = 2 \sum_{e \in F_h^i} \int_e (\omega_e^- \nabla v_h^- + \omega_e^+ \nabla v_h^+) \cdot n_e \llbracket v_h \rrbracket_e + 2 \sum_{e \in F_h^{\partial\Omega}} \int_e \epsilon_h \nabla v_h \cdot n v_h \\
& \leq \sum_{e \in F_h^i} \left[ \alpha h_e \left( \|(\epsilon_h^-)^{\frac{1}{2}} \nabla v_h^- \cdot n_e\|_{0,e}^2 + \|(\epsilon_h^+)^{\frac{1}{2}} \nabla v_h^+ \cdot n_e\|_{0,e}^2 \right) + \frac{1}{\alpha h_e} \|\{\epsilon_h\}_w^{\frac{1}{2}} \llbracket v_h \rrbracket_e\|_{0,e}^2 \right] \\
& + \sum_{e \in F_h^{\partial\Omega}} \left[ \alpha h_e \|\epsilon_h^{\frac{1}{2}} \nabla v_h \cdot n\|_{0,e}^2 + \frac{1}{\alpha h_e} \|\epsilon_h^{\frac{1}{2}} v_h\|_{0,e}^2 \right].
\end{aligned}$$

By means of inequality (27), we get,

$$\begin{aligned}
& 2 \sum_{e \in F_h^i} \int_e \{\epsilon_h \nabla v_h\}_w \cdot n_e \llbracket v_h \rrbracket_e + 2 \sum_{e \in F_h^{\partial\Omega}} \int_e \epsilon_h \nabla v_h \cdot n v_h \\
& \lesssim \alpha \|\epsilon_h^{\frac{1}{2}} \nabla v_h\|_{0, T_h}^2 + \frac{1}{\alpha} \|\{\epsilon_h\}_w^{\frac{1}{2}} h_e^{-\frac{1}{2}} \llbracket v_h \rrbracket_e\|_{0, F_h^i}^2 + \frac{1}{\alpha} \|\epsilon_h^{\frac{1}{2}} h_e^{-\frac{1}{2}} v_h\|_{0, F_h^{\partial\Omega}}^2. \tag{32}
\end{aligned}$$

The coercivity of  $a_h^s(\cdot, \cdot)$  in the norm  $\|\cdot\|_a$  directly follows from the combination of (31) and (32) provided  $\alpha$  and  $\xi$  are such that  $\alpha \lesssim 1$  and  $\frac{1}{\alpha} \lesssim \xi$ .  $\square$

**Corollary 1.** *Under the assumptions of Lemma 3 and Lemma 4 the bilinear form  $a_h(\cdot, \cdot)$  is coercive in the energy norm  $|||\cdot|||_a$  for all  $v_h \in V_h^p$ .*

**Corollary 2.** *Corollary 1 ensures that problem (11) is well posed.*

Now, to suitably define the continuity of the bilinear form  $a_h(\cdot, \cdot)$ , we introduce the space,

$$V_\perp := \{v \in V \cap H^2(T_h^{el}); \int_K v v_h = 0, \forall v_h \in V_h^p\}.$$

**Lemma 5.** *The bilinear form  $a_h(\cdot, \cdot)$  is continuous, precisely it satisfies*

$$a_h(u, v_h) \lesssim |||u|||_\perp |||v_h|||_a \quad \forall u \in V_\perp, \quad \forall v_h \in V_h^p.$$

*Proof.* We proceed by bounding each term of  $a_h(u, v_h)$  with respect to  $|||u|||_\perp$  and  $|||v_h|||_a$ . First, we consider  $a_h^s(u, v_h)$  and we observe that,

$$\begin{aligned} & \sum_{e \in F_h^i} \int_e \{\epsilon_h \nabla u\}_w \cdot n_e \llbracket v_h \rrbracket_e + \sum_{e \in F_h^{\partial\Omega}} \int_e \epsilon_h \nabla u \cdot n v_h \\ & \lesssim \|(\epsilon_h h_e)^{\frac{1}{2}} \nabla u \cdot n_e\|_{0, F_h^i} \|(\frac{1}{2} \{\epsilon_h\}_w h_e^{-1})^{\frac{1}{2}} \llbracket v_h \rrbracket_e\|_{0, F_h^i} \\ & \quad + \|(\epsilon_h h_e)^{\frac{1}{2}} \nabla u \cdot n\|_{0, F_h^{\partial\Omega}} \|(\epsilon_h h_e^{-1})^{\frac{1}{2}} v_h\|_{0, F_h^{\partial\Omega}}, \end{aligned}$$

and from (32) we get,

$$\begin{aligned} & \sum_{e \in F_h^i} \int_e \{\epsilon_h \nabla v_h\}_w \cdot n_e \llbracket u \rrbracket_e + \sum_{e \in F_h^{\partial\Omega}} \int_e \epsilon_h \nabla v_h \cdot n u \\ & \lesssim \|\epsilon_h^{\frac{1}{2}} \nabla v_h\|_{0, T_h} \left( \|(\{\epsilon_h\}_w h_e^{-1})^{\frac{1}{2}} \llbracket u \rrbracket_e\|_{0, F_h^i} + \|\epsilon_h^{\frac{1}{2}} h_e^{-\frac{1}{2}} u\|_{0, F_h^{\partial\Omega}} \right). \end{aligned}$$

Exploiting the previous inequalities, the bilinear form  $a_h^s(u, v_h)$  can be easily estimated as follows,

$$a_h^s(u, v_h) - \int_\Omega (\frac{1}{2} \nabla \cdot \beta) u v_h \lesssim |||u|||_\perp |||v_h|||_a.$$

where the term  $\frac{1}{2} \nabla \cdot \beta u v_h$  cancels out with the opposite one in  $a_h^r(u, v_h)$ .

Second, we consider  $a_h^r(u, v_h)$ . To bound the term  $\sum_{K \in T_h} \int_K (\beta u) \cdot \nabla v_h$ , let  $\beta_h := \Pi_h^0 \beta$  be the piecewise constant vector-valued field equal to the mean value of  $\beta$  on each  $K \in T_h$ . Then,

$$\int_K (\beta u) \cdot \nabla v_h = \int_K u \beta_h \cdot \nabla v_h + \int_K u (\beta - \beta_h) \cdot \nabla v_h = \int_K u (\beta - \beta_h) \cdot \nabla v_h, \quad (33)$$

since  $\beta_h \cdot \nabla v_h \in V_h^p$  and  $u \in V_\perp$ . Moreover,  $\beta \in [W^{1, \infty}(\Omega)]^d$  thus we have  $\|\beta - \beta_h\|_{[W_\infty^0(K)]^d} \lesssim h_K$  for all  $K \in T_h$ , so that the inverse inequality (28) yields

$$\sum_{K \in T_h} \int_K |(\beta u) \cdot \nabla v_h| \lesssim \|u\|_{0, T_h} \|v_h\|_{0, T_h} \leq |||u|||_a |||v|||_a.$$

For the remaining terms of  $a_h^r(u, v_h)$  we get,

$$\begin{aligned} & \sum_{e \in F_h^i} \int_e \left[ \{\beta u\}_w \cdot n_e \llbracket v_h \rrbracket_e - \frac{1}{2} \beta \cdot n_e \varphi_e(\epsilon_h) \llbracket u \rrbracket_e \llbracket v_h \rrbracket_e \right] + \frac{1}{2} \sum_{e \in F_h^{\partial\Omega}} \int_e \beta \cdot n e u v_h \\ & \lesssim \|\beta \cdot n_e\|^{\frac{1}{2}} \|u\|_{0, F_h^i} \|\beta \cdot n_e\|^{\frac{1}{2}} \|\llbracket v_h \rrbracket_e\|_{0, F_h^i} + \|\beta \cdot n_e\|^{\frac{1}{2}} \|\llbracket u \rrbracket_e\|_{0, F_h^i} \|\beta \cdot n_e\|^{\frac{1}{2}} \|\llbracket v_h \rrbracket_e\|_{0, F_h^i} \\ & + \|\beta \cdot n\|^{\frac{1}{2}} \|u\|_{0, F_h^{\partial\Omega}} \|\beta \cdot n\|^{\frac{1}{2}} \|v_h\|_{0, F_h^{\partial\Omega}}. \end{aligned}$$

From the two previous inequalities we obtain that,

$$a_h^r(u, v_h) + \int_{\Omega} \left( \frac{1}{2} \nabla \cdot \beta \right) u v_h \lesssim \|u\|_{\perp} \|v_h\|_a. \quad (34)$$

Finally, the result directly follows from the combination of (33) and (34).  $\square$

We are now ready to prove an a-priori error estimate in the energy norm for method (11).

**Lemma 6.** *Let  $u \in V_{\beta,0} \cap H^2(T_h^{el})$  be the solution of (3), let  $\pi_h u$  be the  $L^2$ -projection of  $u$  onto  $V_h^p$  and let  $u_h$  be the solution of (11). Under assumptions of lemmas 2, 4, 3 and 5 we obtain,*

$$\|u - u_h\|_a \lesssim \|u - \pi_h u\|_{\perp}. \quad (35)$$

Moreover, under the assumptions of lemma 1 we get,

$$\begin{aligned} & \|u - u_h\|_a \\ & \lesssim \sum_{K \in T_h} \left[ (h_K^p \|\epsilon_h\|_{L^\infty(K)} + h_K^{p+\frac{1}{2}} \|\beta\|_{W_\infty^1(K)} + h_K^{p+1} \mu_0) |v|_{p+1, K} \right]. \end{aligned} \quad (36)$$

*Proof.* Lemmas 2, 4, 3 and 5 imply that

$$\begin{aligned} \|u_h - \pi_h u\|_a & \lesssim \frac{a_h(u_h - \pi_h u, u_h - \pi_h u)}{\|u_h - \pi_h u\|_a} \lesssim \frac{a_h(u - \pi_h u, u_h - \pi_h u)}{\|u_h - \pi_h u\|_a} \\ & \lesssim \|u - \pi_h u\|_{\perp}, \end{aligned} \quad (37)$$

owing to the fact that  $u - \pi_h u \in V_{\perp}$ . We complete the proof of (35) by applying the triangle inequality and using the fact that  $\|\cdot\|_a \leq \|\cdot\|_{\perp}$ , precisely,

$$\|u - u_h\|_a \lesssim \|u_h - \pi_h u\|_a + \|u - \pi_h u\|_a \lesssim \|u - \pi_h u\|_{\perp}.$$

Moreover, (36) directly follows from the combination of (35) and lemma 1.  $\square$

## 4 Duality based a-posteriori error analysis

In this section we aim to put into evidence some peculiar advantages of the weighted interior penalty method in the derivation of a local error estimator.

First of all, we observe that there are different strategies to obtain a local error estimator, which can be mainly classified into residual based or duality based error estimates. In our context, the derivation of a residual based error estimate is particularly challenging, because we aim to obtain a result that is robust with respect to the diffusion parameter, see [18] for a recent discussion of these topics. Conversely, a duality based error estimate can be derived following the framework proposed in [3] and references therein, and can be also easily adapted to the case of nonconforming elements, see for instance [13]. In particular, we notice that the duality based approach straightforwardly preserves the robustness of the weighted interior penalty method with respect to locally vanishing diffusivities. By consequence, for this preliminary study, we focus on the duality based a-posteriori error analysis.

#### 4.1 Duality based error representation

To start with, we introduce a linear functional  $J(\cdot) : V \rightarrow \mathbb{R}$  that is the output functional for which we aim to control the error. The definition of  $J(\cdot)$  will be made precise later on. Let  $\hat{\beta}$  be the *dual* advection field, namely  $\hat{\beta} = -\beta$ . We denote with the same superscript all the quantities, depending on the new advection field, that are involved in the definition of problem (4). Furthermore, we introduce the dual space  $V_{\hat{\beta},0} := V_0^{el} \times V_{\hat{\beta},0}^{hy}$  and for any  $u, v \in V$  we define the dual bilinear form,

$$\begin{aligned} \hat{a}(u, v) := & \int_{\Omega_{el}} \left( (\epsilon_h \nabla u_{el} - \hat{\beta} u_{el}) \cdot \nabla v_{el} + (\mu - \nabla \cdot \hat{\beta}) u_{el} v_{el} \right) \\ & + \int_{\Omega_{hy}} \left( -\hat{\beta} u_{hy} \cdot \nabla v_{hy} + (\mu - \nabla \cdot \hat{\beta}) u_{hy} v_{hy} \right) \\ & + \int_{\hat{\Gamma}_h^{in}} \hat{\beta} \cdot nu_{el} \llbracket v \rrbracket + \int_{\hat{\Gamma}_h^{out}} \hat{\beta} \cdot nu_{hy} \llbracket v \rrbracket + \int_{\partial\Omega_{hy}^{out}} \hat{\beta} \cdot nu_{hy} v_{hy}, \end{aligned}$$

where  $\hat{\Gamma}_h^{in}$ ,  $\hat{\Gamma}_h^{out}$  and  $\partial\Omega_{hy}^{out}$  refer to  $\hat{\beta} = -\beta$ . It is easily verified that  $\hat{a}(u, v) = a(v, u) \forall u, v \in V_{\hat{\beta},0}$  or vice versa  $a(u, v) = \hat{a}(v, u) \forall u, v \in V_{\hat{\beta},0}$ . In this framework, we introduce the dual problem with respect to (4): find  $z \in V_{\hat{\beta},0}$  such that,

$$\hat{a}(z, \varphi) = J(\varphi), \quad \forall \varphi \in V_{\hat{\beta},0}. \quad (38)$$

The analysis of the dual problem is analogous to the primal one, namely problem (4). To proceed, we set up a discretization method for the dual problem. Mimicking the derivation of the WIP method, we obtain the bilinear form  $\hat{a}_h(u_h, v_h)$ , which differs from  $a_h(u_h, v_h)$ , because the advection field  $\beta$  is replaced by  $\hat{\beta} = -\beta$ . We observe that  $\hat{a}_h(u_h, v_h)$  satisfies the following properties,

$$\begin{aligned} \hat{a}_h(u_h, v_h) &:= \hat{a}_h^s(u_h, v_h) + \hat{a}_h^r(u_h, v_h), \\ \hat{a}_h^s(u_h, v_h) &:= a_h^s(u_h, v_h), \\ \hat{a}_h^r(u_h, v_h) &:= a_h^r(\hat{\beta}; u_h, v_h), \end{aligned}$$

where in the last row we have highlighted the dependence of  $a_h^r(u_h, v_h)$  from  $\beta$ . Owing to the symmetry of  $a_h^s(u_h, v_h)$  and to the skew-symmetry of  $a_h^r(u_h, v_h)$  we immediately obtain that,

$$\begin{aligned}\hat{a}_h(u_h, v_h) &= a_h^s(u_h, v_h) - a_h^r(\beta; u_h, v_h) \\ &= a_h^s(v_h, u_h) + a_h^r(\beta; v_h, u_h) = a_h(v_h, u_h).\end{aligned}\quad (39)$$

Then, we consider the discrete dual problem that consists in finding  $z_h \in V_h^q$  such that,

$$\hat{a}_h(z_h, \varphi_h) = J(\varphi_h), \quad \forall \varphi_h \in V_h^q, \quad (40)$$

where the discrete dual space  $V_h^q$  is generally richer than  $V_h^p$ , i.e.  $q > p$ .

We notice that mimicking Lemma 3 and 4 as well as Corollary 1 and 2 we immediately conclude that problem (40) is coercive and thus well-posed. Finally, in analogy with Lemma 2, we obtain the following result that states the consistency of problem (40) with respect to (38).

**Lemma 7.** *Let  $z \in V_{\beta,0} \cap H^2(T_h^{el})$  be the solution of problem (38). Then,  $\hat{a}_h(z, \varphi) = J(\varphi)$ ,  $\forall \varphi \in (V \cap H^2(T_h^{el})) \oplus V_h^p$  and  $a_h(z - z_h, \varphi_h) = 0 \forall \varphi_h \in V_h^p$ .*

The proof is omitted since it is analogous to one of lemma 2.

Now, let  $e := u - u_h$  be the error relative to our numerical method, where  $u \in V_{\beta,0} \cap H^2(T_h^{el})$  is the solution of (4) and  $u_h \in V_h^p$  satisfies (11). We easily conclude that  $e \in (V \cap H^2(T_h^{el})) \oplus V_h^p$ . Lemma 7 allows to rewrite the error on the output functional  $J(e) = J(u) - J(u_h)$  in terms of the residuals of the numerical method, more precisely we obtain the following error representation formula.

**Lemma 8.** *Let  $z \in V_{\beta,0} \cap H^2(T_h^{el})$  be the solution of problem (38) and  $u_h \in V_h^p$  the one of (11). Then, for any  $\zeta := (z - v_h) \in V \cap H^2(T_h^{el}) \oplus V_h^p$ ,*

$$\begin{aligned}J(e) &= - \sum_{K \in T_h} \int_K R_0(u_h) W_0(\zeta) \\ &\quad - \sum_{e \in F_h^i} \int_e \left( R_1(u_h) W_1(\zeta) + R_2(u_h) W_2(\zeta) + R_3(u_h) W_3(\zeta) \right) \\ &\quad - \sum_{e \in F_h^{\partial\Omega}} \int_e \left( R_4(u_h) W_4(\zeta) + R_5(u_h) W_5(\zeta) \right),\end{aligned}\quad (41)$$

where  $R_i(\cdot)$  and  $W_i(\cdot)$  are defined in (21).

*Proof.* Since  $e \in (V \cap H^2(T_h^{el})) \oplus V_h^p$ , we apply Lemma 7 with  $\varphi = e$ . Owing to (39) and the Galerkin orthogonality of the primal problem, we get,

$$J(e) = \hat{a}_h(z, e) = a_h(e, z) = a_h(e, \zeta).$$

Now, starting from the expression of  $a_h(e, \zeta)$ , mimicking the steps that lead from (10) to (19), and taking into account that  $u$  weakly satisfies problem (3), we exactly obtain (41).  $\square$

## 4.2 Definition of the local error indicator

We are interested on the error representation identity (41) because it is the starting point to develop an adaptive finite element method. To this purpose, the first step is to define a local error indicator, that is a discrete function that quantifies at which extent the local approximation on each element  $K \in T_h$  contributes to the global error  $J(e)$ . Looking at equation (41), it is clear that the quantities  $R_i(u_h)W_i(\zeta)$ ,  $i = 0, \dots, 5$  are the natural bricks to build up the local error indicator. In this perspective, we observe that it is necessary to set up a suitable strategy for the repartition on each element  $K \in T_h$  of the error indicators  $R_i(u_h)W_i(\zeta)$  with  $i = 1, 2, 3$ , lying on the mesh skeleton, namely  $F_h^i$ . Moreover, this is particularly significant in the case of DG methods, since they are suited to the use of non-conformingly refined meshes, where a non smooth distribution of the error indicator directly reflects into the pattern of the refined mesh.

We also observe that, in the specific case of problems with highly variable diffusivity, the residual  $R_1(u_h) = \llbracket \sigma_h \rrbracket_e \cdot n_e$  is likely to be one of the leading contributions to determine the error. In this framework, we will see that different strategies to split  $R_1(u_h)W_1(\zeta)$  over the neighboring elements lead to remarkably different local error indicators.

In general, the most common and natural strategy to break up the weighted residuals on each edge, precisely  $\int_e R_i(u_h)W_i(\zeta)$ ,  $i = 1, 2, 3$ ,  $e \in F_h^i$ , is to equally divide them into the elements  $K^\pm \in K(e)$ . This seems to be the only possibility to treat  $R_2(u_h)W_2(\zeta)$  and  $R_3(u_h)W_3(\zeta)$ , since no information on the dual solution  $z$  is a-priori available. This is not the case for  $\int_e R_1(u_h)W_1(\zeta)$ . Indeed,  $W_1(\zeta)$  is the only weighting function that depends on the heterogeneity factor, more precisely  $W_1(\zeta) = \{\zeta\}^w = w_e^- \zeta^+ + w_e^+ \zeta^-$ . By consequence, exploiting the information of  $w_e^\pm$ , we can conceive different options to separate  $W_1(\zeta)$  into  $K^\pm$ . Let us denote with  $\int_e^- R_1(u_h)W_1(\zeta) := \int_e R_1(u_h)W_1^*(\zeta)$  the contribution of the error indicator that falls on  $K^-$ . We consider the following alternative splitting strategies,

$$W_1^*(\zeta) := \begin{cases} \frac{1}{2}W_1(\zeta) & (a) \\ w_e^- \zeta^+ & (b) \\ w_e^+ \zeta^- & (c) \end{cases} \quad (42)$$

Then, we notice that equation (41) can be rewritten as follows, where now  $\pm$  refer to the normal vector  $n_{\partial K}$ ,

$$\begin{aligned} J(e) = & - \sum_{K \in T_h} \left[ \int_K R_0(u_h)W_0(\zeta) \right. \\ & + \int_{\partial K \setminus \partial \Omega} \left( R_1(u_h)W_1^*(\zeta) + \frac{1}{2}R_2(u_h)W_2(\zeta) + \frac{1}{2}R_3(u_h)W_3(\zeta) \right) \\ & \left. + \int_{\partial K \cap \partial \Omega} \left( R_4(u_h)W_4(\zeta) + R_5(u_h)W_5(\zeta) \right) \right]. \end{aligned} \quad (43)$$

There are several ways to derive from (43) a local error indicator  $\eta_K^*(u_h, \zeta)$  such

that

$$J(e) \leq \sum_{K \in T_h} \eta_K^*(u_h, \zeta).$$

In our case, we privilege the efficacy of the indicator to set up a mesh refinement strategy rather than its sharpness to represent the global error. Then, disregarding the criteria aiming to minimize the *effectivity index*, see [3], we propose the following local estimator,

$$\begin{aligned} \eta_K^*(u_h, \zeta) &= \int_{\partial K \setminus \partial \Omega} \left( |R_1(u_h)W_1^*(\zeta)| + \frac{1}{2}|R_2(u_h)W_2(\zeta)| + \frac{1}{2}|R_3(u_h)W_3(\zeta)| \right) \\ &+ \int_{\partial K \cap \partial \Omega} \left( |R_4(u_h)W_4(\zeta)| + |R_5(u_h)W_5(\zeta)| \right) + \int_K |R_0(u_h)W_0(\zeta)|. \end{aligned} \quad (44)$$

The difference between the three options proposed in (42) and the most appropriate choice will be pointed out later on by means of numerical experiments.

Finally, for the sake of completeness, starting from the error representation formula (41), we prove an a-posteriori error estimate in the  $L^2$  norm. We consider  $\epsilon_h > 0$  for any  $K \in T_h$  and by consequence, provided that the data of the primal problem are regular enough, we assume that  $u \in H_0^1(\Omega) \cap H^{p+1}(T_h)$ , which directly implies  $e \in H^{p+1}(T_h)$ . To proceed, we define  $J(\varphi)$  as follows,

$$J(\varphi) := \int_{\Omega} j\varphi \quad \text{with} \quad j := \frac{e}{\|e\|_{0,\Omega}} \in H^{p+1}(T_h) \subset L^2(\Omega),$$

that immediately gives  $J(e) = \|e\|_{0,\Omega}$ . Moreover, since the functional  $J(\varphi)$  admits an  $L^2$  representation that is locally regular, we conclude that the solution of the dual problem satisfies  $z \in H_0^1(\Omega) \cap H^{p+1}(T_h)$  too. In this framework, we obtain the following result.

**Lemma 9.** *Under the regularity assumption  $z \in H^{p+1}(T_h)$  the following a-posteriori error estimate holds,*

$$\|e\|_{0,\Omega} \lesssim \sum_{K \in T_h} h_K^{p+1} \bar{\eta}_K^*(u_h) |z|_{p+1,K},$$

where  $\bar{\eta}_K^*(u_h)$  is defined as follows,

$$\begin{aligned} \bar{\eta}_K^*(u_h) &:= \|R_0(u_h)\|_{0,K} + \sum_{e \in \partial K \cap \partial \Omega} \left( \frac{1}{2} h_K^{-\frac{3}{2}} \|R_4(u_h)\|_{0,e} + \frac{1}{2} h_K^{-\frac{1}{2}} \|R_5(u_h)\|_{0,e} \right) \\ &+ \sum_{e \in \partial K \setminus \partial \Omega} \left( h_K^{-\frac{1}{2}} w^* \|R_1(u_h)\|_{0,e} + \frac{1}{2} h_K^{-\frac{3}{2}} \|R_2(u_h)\|_{0,e} + \frac{1}{2} h_K^{-\frac{1}{2}} \|R_3(u_h)\|_{0,e} \right), \end{aligned}$$

where  $w^* = w^a = \frac{1}{2}$ ,  $w^* = w^b = w_e^-$  or  $w^* = w^c = w_e^+$  according to (42).

*Proof.* The desired result directly follows from (41) by choosing  $\zeta = z - \pi_h z$  and replacing the following inequalities into the definition of  $W_i(\zeta)$ ,  $i = 0, \dots, 5$ ,

$$\begin{aligned} \|\zeta\|_{0,K} &\lesssim h_K^{p+1} |z|_{p+1,K}, \\ h_e^{-1} \|\zeta\|_{0,e} + \|\nabla \zeta \cdot n_e\|_{0,e} &\lesssim h_K^{p-\frac{1}{2}} |z|_{p+1,K}. \end{aligned}$$

□

## 5 Numerical results

In the previous sections we have built up a family of weighted interior penalty (WIP) methods that depend on the value of a scalar parameter, the tilting factor  $\alpha$ . The aim of this section is to highlight the benefits they provide to the approximation of advection - diffusion - reaction equations with non-smooth and possibly locally vanishing diffusivity.

We will first compare the WIP scheme with several variants of the interior penalty DG methods, trying to point out how the tilting factor influences the accuracy and the robustness of the method. Secondly, we will focus on the a-posteriori error estimate and we will compare the three local error indicators arising from (42), with the perspective to set up an adaptive mesh refinement strategy.

### 5.1 A test case with non-smooth coefficients

In order to pursue a quantitative comparison between our scheme and the standard interior penalty method, we aim to set up a test problem, featuring discontinuous coefficients, that allows us to analytically compute the exact solution. To this aim, we consider the following test case, already proposed in [5, 9]. For the sake of clarity, we remind it here.

Let be  $\Omega \subset \mathbb{R}^2$  and let  $x, y$  be the space coordinates. We split the domain  $\Omega$  into two subregions,  $\Omega_1 = (x_0, x_{\frac{1}{2}}) \times (y_0, y_1)$ ,  $\Omega_2 = (x_{\frac{1}{2}}, x_1) \times (y_0, y_1)$  and we choose for simplicity  $x_0 = 0$ ,  $x_{\frac{1}{2}} = \frac{1}{2}$ ,  $x_1 = 1$  while  $y_0 = 0$ ,  $y_1 = 1$ . The viscosity  $\epsilon(x, y)$  is a discontinuous function across the interface  $x = x_{\frac{1}{2}}$ , for any  $y \in (y_0, y_1)$ . Precisely, we will consider a constant  $\epsilon(x, y)$  in each subregion with several values of  $\epsilon_1$  in  $\Omega_1$  and a fixed  $\epsilon_2 = 1.0$  in  $\Omega_2$ . Moreover, we set  $\beta = [\beta_x = 1, \beta_y = 0]$ ,  $\mu = 0$ ,  $f = 0$  and the boundary conditions  $u_1(x_0, y) = u_0 = 1$ ,  $u_2(x_1, y) = u_1 = 0$ ,  $\partial_y u(x, y_0) = \partial_y u(x, y_1) = 0$ . Then, the exact solution of the problem on each subregion  $\Omega_1, \Omega_2$  can be expressed as an exponential function with respect to  $x$  independently from  $y$ . The global solution  $u(x, y)$  is provided by choosing the value at the interface  $x = x_{\frac{1}{2}}$  in order to ensure the following matching conditions,

$$\begin{aligned} \lim_{x \rightarrow x_{\frac{1}{2}}^-} u(x, y) &= \lim_{x \rightarrow x_{\frac{1}{2}}^+} u(x, y), \\ \lim_{x \rightarrow x_{\frac{1}{2}}^-} -\epsilon(x, y) \partial_x u(x, y) + \beta_x u(x, y) &= \lim_{x \rightarrow x_{\frac{1}{2}}^+} -\epsilon(x, y) \partial_x u(x, y) + \beta_x u(x, y). \end{aligned}$$

More precisely, introducing the Péclet numbers  $pe_1 := |\beta_x|(x_{\frac{1}{2}} - x_0)/\epsilon_1$ ,  $pe_2 := |\beta_x|(x_1 - x_{\frac{1}{2}})/\epsilon_2$ , by consequence of the matching conditions we obtain that  $u_{\frac{1}{2}} := u(x_{\frac{1}{2}}, y)$  reads as follows,

$$u_{\frac{1}{2}} = \left[ \frac{u_0 \exp(pe_1)}{1 - \exp(pe_1)} + \frac{u_1}{1 - \exp(pe_2)} \right] \left[ \frac{\exp(pe_1)}{1 - \exp(pe_1)} + \frac{1}{1 - \exp(pe_2)} \right]^{-1}.$$

As a result of that, the exact solution in each subdomain can be expressed as,

$$u_1(x, y) = \frac{u_{\frac{1}{2}} - \exp(pe_1)u_0 + [u_0 - u_{\frac{1}{2}}] \exp(\beta(x - x_0)/\epsilon_1)}{1 - \exp(pe_1)}$$

$$u_2(x, y) = \frac{u_1 - \exp(pe_2)u_{\frac{1}{2}} + [u_{\frac{1}{2}} - u_1] \exp(\beta(x - x_{\frac{1}{2}})/\epsilon_2)}{1 - \exp(pe_2)}.$$

It is easy to see that when  $0 \simeq \epsilon_1 \ll \epsilon_2 = 1$  the global solution,  $u$ , features a very sharp internal layer upwind to the discontinuity of  $\epsilon$ , located at  $x = x_{\frac{1}{2}}$ .

## 5.2 Comparison of different interior penalty methods

In the numerical simulations that will follow, our reference standard interior penalty method (IP) is obtained by setting  $\lambda_h = 0$  and choosing  $w_e^\pm$ ,  $\varphi_e(\lambda_h)$  and  $\chi_e(\lambda_h)$  according to (13), (15) and (17) respectively. We consider both the symmetric interior penalty method (SIP) and the skew-symmetric version (SSIP), also known as NIPG, [16]. The latter variant has the advantage that it only requires the condition  $\xi > 0$  to prove lemma 4. Consequently, we will set  $\xi = 2 \cdot 10^{-2}$  for SSIP while  $\xi = 2$  for SIP and WIP, and we will study how this parameter influences the accuracy when  $\epsilon$  is vanishing. For all test cases we consider a uniform triangulation  $T_h$  with  $h = 0.05$  and we apply piecewise linear elements. In this setting, we perform a quantitative comparison based on several indicators. In particular, we consider the energy norm, the  $L^2$  norm and the norm of the advective derivative that is defined as  $\|v\|_{h,\beta}^2 := \sum_{K \in T_h} h_K \|\beta \cdot \nabla v\|_{0,K}^2$  and it has been analyzed in [9] for a case similar to the present one. Finally we introduce the following indicator,

$$\Delta := \max(|\max_{\Omega}(u_h) - \max_{\Omega}(u)|, |\min_{\Omega}(u_h) - \min_{\Omega}(u)|),$$

which quantifies to which extent the numerical solution exceeds the extrema of the exact one.

The results, reported in figure 1, put into evidence that the WIP scheme performs better than the standard IP methods, particularly in those cases where the solution is non smooth and at the same time the computational mesh is not completely adequate to capture the singularities. From the analysis of figure 1, it is possible to identify three regimens where the numerical methods behave differently. The first one consists on the diffusive region, namely  $2^{-4} < \epsilon_1 \leq 1$ , where all the methods provide similar results. For  $\epsilon_1 \leq 2^{-4}$  a transition takes place and all the error indicators increase, because the computational mesh is not adequate any more to capture the very sharp internal layer that originates upwind to the discontinuity of  $\epsilon$ . In this case, we notice that the error is quite sensitive to the choice of the tilting factor  $\alpha$ . More precisely, it seems that the tilting factor influences the tradeoff between the accuracy of the method and its robustness. The smaller is  $\alpha$  the more the method is robust with respect to a discontinuity of the diffusivity, as it is suggested by the behavior of the  $L^2$  norm and of the indicator  $\Delta$ . However, for the smallest values of  $\alpha$  we notice that

Figure 1: The norms  $\|\cdot\|_{0,\Omega}$  (top-left),  $\|\cdot\|$  (bottom-left),  $\|\cdot\|_{h,\beta}$  (bottom-right) and the indicator  $\Delta$  (top-right) are plotted for the values  $\epsilon_1 = 2^{-i}$ ,  $i = 0, \dots, 16$ . Several schemes are compared with respect to these indicators.

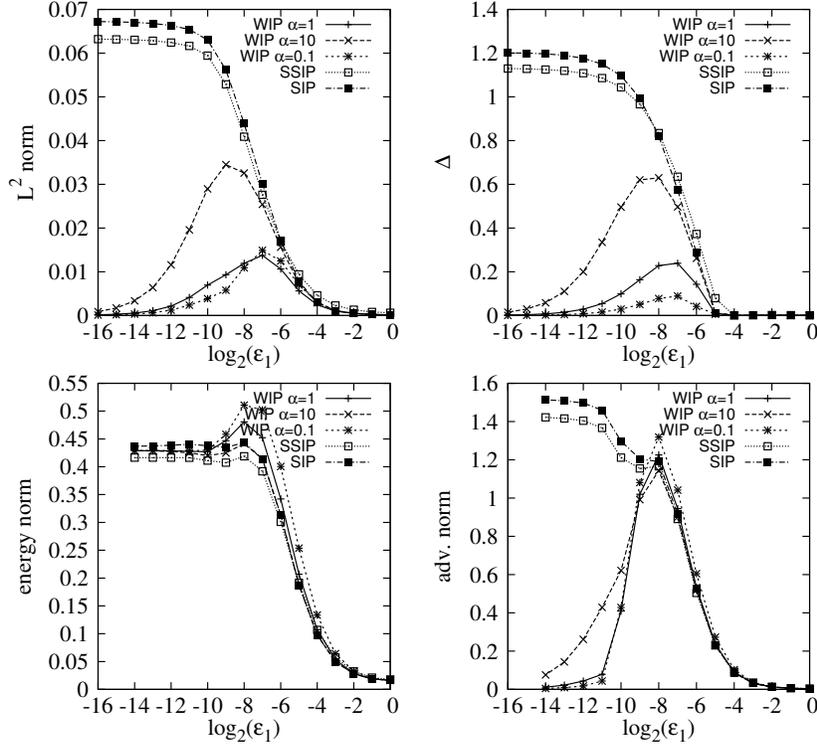
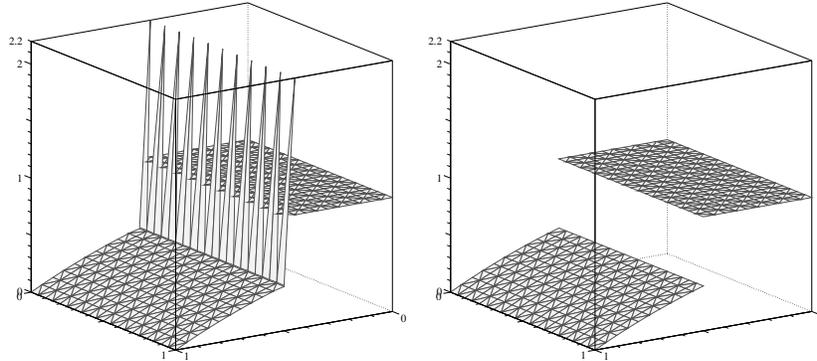


Figure 2: The solution of the SIP and WIP methods for  $\epsilon_1 = 10^{-6}$  and  $\epsilon_2 = 1$ .



the accuracy of the scheme slightly decreases in the energy and in the advective norms. Finally, the smallest values of  $\epsilon_1$ , namely  $\epsilon_1 < 2^{-12}$ , correspond to the hyperbolic regimen. In this case, the global solution  $u$  can be seen as a discontinuous function in correspondence of the jump of  $\epsilon$ . However, we observe that the standard interior penalty schemes compute a solution that is almost continuous, as reported in figure 2. This behavior promotes the instability of the approximate solution in the neighborhood of the boundary layer, because the computational mesh is not adequate in order to smoothly approximate the very high gradients across the interface. The quantity  $\Delta$  shows that the spurious oscillations generated in this case reach the 100% of the maximum of the exact solution. This is true both for the symmetric and the skew-symmetric variants, thus we conclude that the magnitude of the penalty parameter  $\xi$  is less significant than the application of the weighted averages. Conversely, the WIP methods are very effective for any value of  $\alpha$ , because the scheme is consistent with the elliptic-hyperbolic limit case.

### 5.3 Approximation and comparison of the local error indicators

In order to apply the local error indicators  $\eta_K^*(u_h, \zeta)$  to the test case of section 5.1, first of all we have to provide a precise definition for the output functional  $J(u)$ . For simplicity, we consider  $J(u) := \int_{\Omega} u$ .

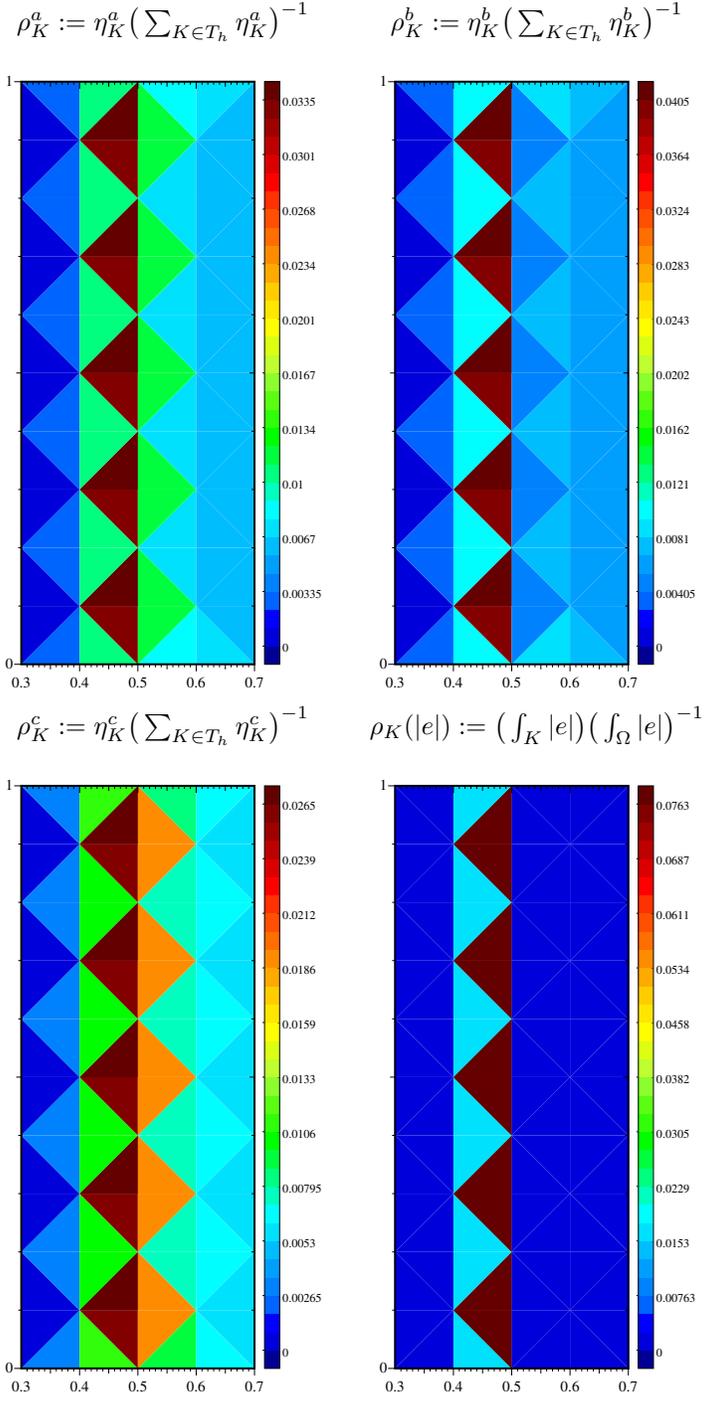
To proceed, we approximate the exact solution,  $z$ , of the dual problem (38) by means of its discretization through the WIP method. Because of the arbitrariness of  $v_h$  in the definition of  $\zeta$ , it is clear that we have to compute  $z_h$  choosing  $q > p$  into (41). The simplest case is to choose  $p = 1$  and  $q = 2$ . We denote with  $z_h^{(2)}$  the solution of (40) into  $V_h^2$  and with  $u_h^{(1)}$  the solution of (11) into  $V_h^1$ . Then, our approximate local error indicator is given by  $\eta_K^*(u_h^{(1)}, z_h^{(2)})$  on each element  $K \in T_h$ .

For the comparison of the indicators  $\eta_K^{a,b,c}$  we rescale each of them with respect to the estimated global error. This results in a piecewise constant function that quantifies to which percentage the error on each element contributes to the global one. More precisely, we introduce the *relative* local error estimators,

Table 1: Quantitative comparison between the relative error indicators.

indicator	$\rho_K( e ) - \rho_K^a$	$\rho_K( e ) - \rho_K^b$	$\rho_K( e ) - \rho_K^c$
$\max_{K \in T_h}$	0.0428329	0.0358393	0.0498264
$\min_{K \in T_h}$	-0.0107334	-0.00820335	-0.017727
indicator	$\rho_K( e )/\rho_K^a$	$\rho_K( e )/\rho_K^b$	$\rho_K( e )/\rho_K^c$
$\max_{K \in T_h}$	2.27904	1.88532	2.88062
$\min_{K \in T_h}$	0.0476816	0.0476819	0.0476812

Figure 3: Comparison between the relative local error estimators  $\rho_K^{a,b,c}$ .



defined by

$$\rho_K^* := \eta_K^* \left( \sum_{K \in T_h} \eta_K^* \right)^{-1},$$

and we compare them with the quantity

$$\rho_K(|e|) := (\int_K |e|) (\int_\Omega |e|)^{-1} = J(|e|)|_K (J(|e|))^{-1}.$$

The relative local error indicators  $\rho_K^*$  and  $\rho_K(|e|)$  are reported in figure 3, where they are plotted in a region neighboring the interface,  $x = \frac{1}{2}$ , where  $\epsilon$  jumps from  $\epsilon_1 = 5 \cdot 10^{-3}$  to  $\epsilon_2 = 1$ . For their comparison, we take the exact indicator  $\rho_K(|e|)$  as a reference. From the plot of  $\rho_K(|e|)$  in figure 3, we immediately notice that the elements that mostly contribute to the error are the ones on the left of the interface where  $\epsilon$  is discontinuous, because the exact solution of the problem at hand features a very sharp internal layer in this region. We observe that the indicators  $\rho_K^{a,b,c}$  are significantly different on those elements where the heterogeneity factor differs from zero and thus the averaging weights  $w_e^\pm$  differ from  $\frac{1}{2}$ . Moreover, we notice that  $\rho_K^b$  is the one that mostly resembles to  $\rho_K(|e|)$ . More precisely,  $\rho_K^b$  favorably clusters the error on the elements that lay on the left with respect to the interface, while  $\rho_K^a$  and in particular  $\rho_K^c$  promote the dispersion of the local error on both sides.

A more quantitative comparison is pursued in table 1, where we consider the indicators  $\rho_K(|e|) - \rho_K^{a,b,c}$  and  $\rho_K(|e|) / \rho_K^{a,b,c}$ , which can be seen as two alternative ways to define a *local and relative* effectivity index. For both cases, we conclude that  $\rho_K^b$  is the best relative indicator with respect to  $\rho_K(|e|)$  and it is definitely more effective than  $\rho_K^{a,c}$ . This is achieved by means of the averaging weights  $w_e^\pm$  that have been suitably exploited to better gather the local residuals of the numerical solution into the local error indicator.

To sum up, figure 3 suggests that the strategies (a), (b) and (c) proposed in (42) to build up a local error indicator may lead to considerably different adaptively fitted computational meshes, especially when non conforming refinements are allowed, thanks to the flexibility of the DG method. Indeed, by means of the simple fixed error reduction strategy applied iteratively, or through one single step of a mesh optimization strategy, see [3], the piecewise constant (and thus discontinuous) relative error indicators  $\rho_K^{a,b,c}$  can be translated into a function that prescribes how to refine the mesh in order to satisfy a suitable tolerance on the global error. Furthermore, this is not only applicable in the context of mesh refinement, but it also fits to the case of error reduction by means of hierarchical basis functions.

## 6 Conclusions

We have proposed a family of DG methods that extends the standard IP schemes by means of weighted averages. This generalization does not increase the computational cost of the scheme but remarkably improves its robustness for the approximation of advection-diffusion-reaction problems that vary in character

from one part of the domain to another, because the diffusivity coefficient may be discontinuous and locally vanishing. These benefits emerge from the a-priori error analysis of the method and are also confirmed by numerical experiments. Finally, we have considered the a-posteriori error analysis of the scheme, putting into evidence that the introduction of weighted interior penalties also helps to improve the effectivity of a local error estimator.

## References

- [1] Douglas N. Arnold. An interior penalty finite element method with discontinuous elements. *SIAM J. Numer. Anal.*, 19(4):742–760, 1982.
- [2] Douglas N. Arnold, Franco Brezzi, Bernardo Cockburn, and L. Donatella Marini. Unified analysis of discontinuous Galerkin methods for elliptic problems. *SIAM J. Numer. Anal.*, 39(5):1749–1779 (electronic), 2001/02.
- [3] Wolfgang Bangerth and Rolf Rannacher. *Adaptive finite element methods for differential equations*. Lectures in Mathematics ETH Zürich. Birkhäuser Verlag, Basel, 2003.
- [4] F. Brezzi, B. Cockburn, L. D. Marini, and E. Süli. Stabilization mechanisms in discontinuous Galerkin finite element methods. *Comput. Methods Appl. Mech. Engrg.*, 195(25-28):3293–3310, 2006.
- [5] Erik Burman and Paolo Zunino. A domain decomposition method based on weighted interior penalties for advection-diffusion-reaction problems. *SIAM J. Numer. Anal.*, 44(4):1612–1638 (electronic), 2006.
- [6] Paul Castillo, Bernardo Cockburn, Ilaria Perugia, and Dominik Schötzau. An a priori error analysis of the local discontinuous Galerkin method for elliptic problems. *SIAM J. Numer. Anal.*, 38(5):1676–1706 (electronic), 2000.
- [7] Bernardo Cockburn, Guido Kanschat, Ilaria Perugia, and Dominik Schötzau. Superconvergence of the local discontinuous Galerkin method for elliptic problems on Cartesian grids. *SIAM J. Numer. Anal.*, 39(1):264–285 (electronic), 2001.
- [8] Bernardo Cockburn and Chi-Wang Shu. The local discontinuous Galerkin method for time-dependent convection-diffusion systems. *SIAM J. Numer. Anal.*, 35(6):2440–2463 (electronic), 1998.
- [9] Alexandre Ern, Annette Stephansen, and Paolo Zunino. A Discontinuous Galerkin method with weighted averages for advection-diffusion equations with locally vanishing and anisotropic diffusivity. Technical Report 103, MOX, Department of Mathematics, Politecnico di Milano, 2007. Submitted.

- [10] Fabio Gastaldi and Alfio Quarteroni. On the coupling of hyperbolic and parabolic systems: analytical and numerical approach. *Appl. Numer. Math.*, 6(1-2):3–31, 1989/90. Spectral multi-domain methods (Paris, 1988).
- [11] P. Grisvard. *Elliptic problems in nonsmooth domains*, volume 24 of *Monographs and Studies in Mathematics*. Pitman (Advanced Publishing Program), Boston, MA, 1985.
- [12] B. Heinrich and K. Pönitz. Nitsche type mortaring for singularly perturbed reaction-diffusion problems. *Computing*, 75(4):257–279, 2005.
- [13] G. Kanschat and R. Rannacher. Local error analysis of the interior penalty discontinuous Galerkin method for second order elliptic problems. *J. Numer. Math.*, 10(4):249–274, 2002.
- [14] Alfio Quarteroni, Franco Pasquarelli, and Alberto Valli. Heterogeneous domain decomposition: principles, algorithms, applications. In *Fifth International Symposium on Domain Decomposition Methods for Partial Differential Equations (Norfolk, VA, 1991)*, pages 129–150. SIAM, Philadelphia, PA, 1992.
- [15] Alfio Quarteroni and Alberto Valli. *Domain decomposition methods for partial differential equations*. Numerical Mathematics and Scientific Computation. The Clarendon Press Oxford University Press, New York, 1999. Oxford Science Publications.
- [16] B. Rivière, M. Wheeler, and V. Girault. Improved energy estimates for interior penalty, constrained and discontinuous Galerkin methods for elliptic problems. Part I. *Comput. Geosci.*, 3:337–360, 1999.
- [17] Rolf Stenberg. Mortaring by a method of J. A. Nitsche. In *Computational mechanics (Buenos Aires, 1998)*. Centro Internac. Métodos Numér. Ing., Barcelona, 1998.
- [18] R. Verfürth. Robust a posteriori error estimates for stationary convection-diffusion equations. *SIAM J. Numer. Anal.*, 43(4):1766–1782 (electronic), 2005.