

Higher-order adaptive virtual element methods with contraction properties

*Original*

Higher-order adaptive virtual element methods with contraction properties / Canuto, Claudio; Fassino, Davide. - In: MATHEMATICS IN ENGINEERING. - ISSN 2640-3501. - ELETTRONICO. - 5:6(2023), pp. 1-33.  
[10.3934/mine.2023101]

*Availability:*

This version is available at: 11583/2984223 since: 2023-11-30T13:08:08Z

*Publisher:*

AIMS Press

*Published*

DOI:10.3934/mine.2023101

*Terms of use:*

This article is made available under terms and conditions as specified in the corresponding bibliographic description in the repository

*Publisher copyright*

(Article begins on next page)



---

*Research article*

## Higher-order adaptive virtual element methods with contraction properties<sup>†</sup>

Claudio Canuto\* and Davide Fassino

Department of Mathematical Sciences “Giuseppe Luigi Lagrange”, Politecnico di Torino, Corso Duca degli Abruzzi 24, Torino, 10129, Italy

<sup>†</sup> **This contribution is part of the Special Issue:** Advancements in Polytopal Element Methods

Guest Editors: Michele Botti; Franco Dassi; Lorenzo Mascotto; Ilario Mazzieri

Link: <https://www.aimspress.com/mine/article/6538/special-articles>

\* **Correspondence:** Email: [claudio.canuto@polito.it](mailto:claudio.canuto@polito.it).

**Abstract:** The realization of a standard Adaptive Finite Element Method (AFEM) preserves the mesh conformity by performing a completion step in the refinement loop: In addition to elements marked for refinement due to their contribution to the global error estimator, other elements are refined. In the new perspective opened by the introduction of Virtual Element Methods (VEM), elements with hanging nodes can be viewed as polygons with aligned edges, carrying virtual functions together with standard polynomial functions. The potential advantage is that all activated degrees of freedom are motivated by error reduction, not just by geometric reasons. This point of view is at the basis of the paper [L. Beirão da Veiga et al., “Adaptive VEM: stabilization-free a posteriori error analysis and contraction property”, SIAM Journal on Numerical Analysis, vol. 61, 2023], devoted to the convergence analysis of an adaptive VEM generated by the successive newest-vertex bisections of triangular elements without applying completion, in the lowest-order case (polynomial degree  $k = 1$ ). The purpose of this paper is to extend these results to the case of VEMs of order  $k \geq 2$  built on triangular meshes. The problem at hand is a variable-coefficient, second-order self-adjoint elliptic equation with Dirichlet boundary conditions; the data of the problem are assumed to be piecewise polynomials of degree  $k - 1$ . By extending the concept of global index of a hanging node, under an admissibility assumption of the mesh, we derive a stabilization-free a posteriori error estimator. This is the sum of residual-type terms and certain virtual inconsistency terms (which vanish for  $k = 1$ ). We define an adaptive VEM of order  $k$  based on this estimator, and we prove its convergence by establishing a contraction result for a linear combination of (squared) energy norm of the error, (squared) residual estimator, and (squared) virtual inconsistency estimator.

**Keywords:** diffusion-reaction problems; Virtual Element Methods; global index of a hanging node; a posteriori error analysis; stabilization-free estimator; adaptivity; contraction property; convergence

---

## 1. Introduction

Adaptive Finite Element Methods (AFEM) for self-adjoint coercive problems written in the form

$$u \in \mathbb{V} : \mathcal{B}(u, v) = F(v), \quad \forall v \in \mathbb{V},$$

iterate the sequence

$$\text{SOLVE} \rightarrow \text{ESTIMATE} \rightarrow \text{MARK} \rightarrow \text{REFINE}$$

to produce better and better approximations of  $u$ . Their practical efficiency is corroborated by sound theoretical results of convergence, complexity, and optimality, which in various cases (such as, e.g., conforming  $h$ -versions) completely explain the behaviour of the adaptive algorithms [11, 13–15, 18].

The standard AFEM realization preserves the conformity of the initial mesh, at the expense of performing a completion step in **REFINE**: In addition to elements marked for refinement due to their contribution to the global error estimator, other elements are refined. Without this step, one would obtain nonconforming meshes, containing elements with hanging nodes.

In the new perspective opened by the introduction of Virtual Element Methods (VEM) [3, 4], elements with hanging nodes can be viewed as polygons with aligned edges, carrying virtual (i.e., non-accessible) functions together with standard polynomial functions. The potential advantage is that all activated degrees of freedom are motivated by error reduction, not just by geometric reasons. On the other hand, in this transformation of an adaptive FEM into an adaptive VEM, one loses the availability of a general convergence theory, which so far is lacking (although results on a posteriori error estimates [8, 12] have been obtained, together with efficient practical recipes for refining polytopal meshes [2, 9, 10]).

Such a shift in perspective inspired the recent papers [5, 6], devoted to the analysis of an adaptive VEM generated by the successive newest-vertex bisections of triangular elements without applying completion, in the lowest-order case (polynomial degree  $k = 1$ ). Despite the simple geometric setup, the investigation faced some VEM-specific obstacles in the analysis, giving answers that could prove useful in the study of more general adaptive VEM discretizations. For instance, a VEM solution  $u_{\mathcal{T}} \in \mathbb{V}_{\mathcal{T}} \subset \mathbb{V}$ , defined by the Galerkin projection

$$u_{\mathcal{T}} \in \mathbb{V}_{\mathcal{T}} : \mathcal{B}_{\mathcal{T}}(u_{\mathcal{T}}, v_{\mathcal{T}}) = F_{\mathcal{T}}(v_{\mathcal{T}}), \quad \forall v_{\mathcal{T}} \in \mathbb{V}_{\mathcal{T}},$$

satisfies an a posteriori error bound of the type

$$\|u - u_{\mathcal{T}}\|_{\mathbb{V}}^2 \lesssim \eta_{\mathcal{T}}^2(u_{\mathcal{T}}) + S_{\mathcal{T}}(u_{\mathcal{T}}, u_{\mathcal{T}}),$$

where  $\eta_{\mathcal{T}}(u_{\mathcal{T}})$  is a residual-type error estimator,  $S_{\mathcal{T}}(u_{\mathcal{T}}, u_{\mathcal{T}})$  is the stabilization term that makes the discrete bilinear form  $\mathcal{B}_{\mathcal{T}}(u_{\mathcal{T}}, v_{\mathcal{T}})$  coercive in  $\mathbb{V}$ , and for simplicity we assume piecewise constant data on the mesh  $\mathcal{T}$ . Unfortunately, the term  $S_{\mathcal{T}}(u_{\mathcal{T}}, u_{\mathcal{T}})$  need not reduce under a mesh refinement, as  $\eta_{\mathcal{T}}^2(u_{\mathcal{T}})$  does: This makes the convergence analysis problematic. However, one of the key results obtained in [5] states that  $S_{\mathcal{T}}(u_{\mathcal{T}}, u_{\mathcal{T}})$  is dominated by  $\eta_{\mathcal{T}}^2(u_{\mathcal{T}})$ , i.e.,

$$S_{\mathcal{T}}(u_{\mathcal{T}}, u_{\mathcal{T}}) \lesssim \eta_{\mathcal{T}}^2(u_{\mathcal{T}}),$$

provided an assumption of *admissibility* of the non-conforming meshes generated by successive refinements is fulfilled; such a restriction, which appears to have little practical impact, amounts to

requiring the uniform boundedness of the *global index* of all hanging node, a useful concept introduced in [5] to hierarchically organize the set of hanging nodes. Once the a posteriori error bound is reduced to

$$\|u - u_{\mathcal{T}}\|_{\mathbb{V}}^2 \lesssim \eta_{\mathcal{T}}^2(u_{\mathcal{T}}),$$

the convergence analysis becomes feasible, and a contraction property is proven to hold for a linear combination of the (squared) energy norm of the error and the (squared) residual estimator.

The purpose of this paper is to extend the results in [5] to the case of VEMs of order  $k \geq 2$  built on triangular meshes. The problem at hand is again a variable-coefficient, second-order self-adjoint elliptic equation with Dirichlet boundary conditions. The geometric concept of hanging node (a vertex for some elements, contained inside an edge of some other elements) is replaced by a functional one, referring to the degrees of freedom associated with the node; once the meaning of hanging node is clarified, the definition of *global index* of a node, and its role in the analysis, is similar to the one given in [5].

A significant difference with respect to the content of that paper concerns the control of the stabilization term, which does not involve only the residual estimator, but a new term, called the *virtual inconsistency estimator* and denoted by  $\Psi_{\mathcal{T}}(u_{\mathcal{T}})$ . It measures the projection error, upon local spaces of polynomials, of certain expressions depending on the operator coefficients and the discrete solution; it vanishes when  $k = 1$  or when the coefficients are constant. The new stabilization bound, which we derive under an admissibility assumption of the mesh, takes the form

$$S_{\mathcal{T}}(u_{\mathcal{T}}, u_{\mathcal{T}}) \lesssim \eta_{\mathcal{T}}^2(u_{\mathcal{T}}) + \Psi_{\mathcal{T}}^2(u_{\mathcal{T}}),$$

which leads to the a posteriori, stabilization-free error control

$$\|u - u_{\mathcal{T}}\|_{\mathbb{V}}^2 \lesssim \eta_{\mathcal{T}}^2(u_{\mathcal{T}}) + \Psi_{\mathcal{T}}^2(u_{\mathcal{T}}).$$

Correspondingly, we obtain the convergence of the adaptive VEM of order  $k$  by proving a contraction result for a linear combination of (squared) energy norm of the error, (squared) residual estimator, and (squared) virtual inconsistency estimator.

Similarly to [5], we assume here that the data  $\mathcal{D}$  of our boundary-value problem are piecewise polynomials of degrees related to  $k - 1$ , on the initial mesh  $\mathcal{T}_0$  and consequently on each mesh  $\mathcal{T}$  derived by newest-vertex bisection. This is not a restriction, since we propose to insert the adaptive VEM procedure just described, which we now consider as a module GALERKIN, into an outer loop AVEM of the form

```
[ $\mathcal{T}, u_{\mathcal{T}}$ ] = AVEM( $\mathcal{T}_0, \epsilon_0, \omega, \text{tol}$ )
j = 0
while  $\epsilon_j > \frac{1}{2}\text{tol}$  do
  [ $\hat{\mathcal{T}}_j, \hat{\mathcal{D}}_j$ ] = DATA( $\mathcal{T}_j, \mathcal{D}, \omega\epsilon_j$ )
  [ $\mathcal{T}_{j+1}, \mathcal{D}_{j+1}$ ] = GALERKIN( $\hat{\mathcal{T}}_j, \hat{\mathcal{D}}_j, \epsilon_j$ )
   $\epsilon_{j+1} \leftarrow \frac{1}{2}\epsilon_j$ 
  j  $\leftarrow$  j + 1
end while
return
```

where the module **DATA** produces, via greedy-type iterations, a piecewise polynomial approximation of the input data with prescribed accuracy, defined on a suitable refinement of the input partition. Manifestly, the target accuracy is matched after a finite number of calls to **DATA** and **GALERKIN**. Properties of complexity and quasi-optimality of this two-loop algorithm are investigated in [6] in the linear case  $k = 1$ . We plan to do the same for the case  $k \geq 2$  in a forthcoming paper.

The outline of this paper is as follows. In Sections 2 and 3, we introduce the model boundary-value problem, and its discretization by an enhanced version of the VEM ([1]). In Section 4 we define the global index of a node, and we formulate the admissibility assumption on the mesh. Two essential properties for bounding the stabilization term are established in Section 5. The a posteriori error estimators are defined in Section 6, whereas stabilization-free a posteriori error estimates are proven in Section 7. In Section 8, we investigate how the a posteriori error estimators are reduced under mesh refinement. These properties are needed to justify the refinement strategy in our adaptive module **GALERKIN**, which is described in Section 9. In Section 10, we discuss the proof of convergence of the loop **GALERKIN**. The paper ends with some numerical experiments, reported in Section 11.

## 2. VEM spaces of order $k \geq 2$

We consider the following Dirichlet boundary value problem in a polygonal domain  $\Omega$ ,

$$\begin{cases} -\nabla \cdot (A \nabla u) + cu = f & \text{in } \Omega, \\ u = 0 & \text{on } \partial\Omega, \end{cases} \quad (2.1)$$

where  $A \in (L^\infty(\Omega))^{2 \times 2}$  is symmetric and uniformly positive definite in  $\Omega$ ,  $c \in L^\infty(\Omega)$  and non-negative in  $\Omega$ ,  $f \in L^2(\Omega)$ . Data will be denoted by  $\mathcal{D} = (A, c, f)$ . The variational formulation of this problem is written as

$$\begin{cases} \text{find } u \in \mathbb{V} := H_0^1(\Omega) & \text{such that} \\ \mathcal{B}(u, v) = (f, v), & \forall v \in \mathbb{V}, \end{cases} \quad (2.2)$$

where  $(\cdot, \cdot)$  is the scalar product in  $L^2(\Omega)$  and  $\mathcal{B}(u, v) := a(u, v) + m(u, v)$  is the bilinear form associated with Problem (2.1), i.e.,

$$a(u, v) := (A \nabla u, \nabla v), \quad m(u, v) := (c u, v).$$

We denote the energy norm as  $\|\cdot\| = \sqrt{\mathcal{B}(\cdot, \cdot)}$ , which satisfies

$$c_{\mathcal{B}} |v|_{1, \Omega}^2 \leq \|v\|^2 \leq c^{\mathcal{B}} |v|_{1, \Omega}^2, \quad \forall v \in \mathbb{V}, \quad (2.3)$$

for suitable  $0 < c_{\mathcal{B}} \leq c^{\mathcal{B}}$ .

*Remark 2.1.* For the sake of simplicity, in (2.1) we consider the Poisson problem with vanishing Dirichlet conditions on the whole boundary domain. The extension to generic Dirichlet and/or Neumann/Robin boundary conditions does not pose conceptual difficulties. In the numerical examples, we actually provide experiments with more general Dirichlet boundary conditions.

In order to find a discrete approximation of the solution of Problem (2.2), we firstly introduce a fixed initial partition  $\mathcal{T}_0$  on the domain  $\bar{\Omega}$  made of triangular elements  $E$ . We will denote by  $\mathcal{T}$  any refinement of  $\mathcal{T}_0$  obtained by a finite number of newest-vertex element bisections. We underline that we are not requiring  $\mathcal{T}$  to be a conforming mesh, since hanging nodes may arise in the refinement. The classification of nodes, which will play a crucial role in the proofs presented in this paper, is postponed in Section 4.

According to the Virtual Element theory [3], an element  $E$  of the triangulation can be viewed as a polygon with more than three edges, if some hanging nodes are sitting on its boundary. We can then denote by  $\mathcal{E}_E$  the set of edges  $e$  of element  $E$  and  $\mathcal{E} := \bigcup_{E \in \mathcal{T}} \mathcal{E}_E$ . We finally define the diameter of an element  $E$  as  $h_E = |E|^{1/2}$  and  $h = \max_{E \in \mathcal{T}} \{h_E\}$ .

We introduce the functional spaces needed to apply the VEM. We start by defining the space of functions on the boundary of  $E$ ,  $\mathbb{V}_{\partial E, k}$ , which is constituted by the functions that are continuous on the boundary of  $E$  and that, when restricted to any edge of  $\partial E$ , are polynomials of degree  $k > 0$ , i.e.,

$$\mathbb{V}_{\partial E, k} := \{v \in C^0(\partial E) : v|_e \in \mathbb{P}_k(e), \forall e \subset \partial E\}.$$

Then, we define the “enhanced” VEM space in  $E$ , as done in [1], such that

$$\mathbb{V}_{E, k} := \left\{ v \in H^1(E) : v|_{\partial E} \in \mathbb{V}_{\partial E, k}, \Delta v \in \mathbb{P}_k(E), (v - \Pi_E^\nabla v, q)_E = 0 \forall q \in \mathbb{P}_k(E) \setminus \mathbb{P}_{k-2}(E) \right\}, \quad (2.4)$$

where  $\mathbb{P}_k(E) \setminus \mathbb{P}_{k-2}(E)$  is the space spanned by the monomials of degree equal to  $k$  and  $k - 1$ , and  $\Pi_E^\nabla : H^1(E) \rightarrow \mathbb{P}_k(E)$  is the projector defined by

$$(\nabla(v - \Pi_E^\nabla v), \nabla q)_E = 0, \quad \forall q \in \mathbb{P}_k(E), \quad \int_{\partial E} (v - \Pi_E^\nabla v) = 0.$$

We remark that  $\mathbb{V}_{E, k}$  contains the polynomial space of degree  $k$  on  $E$  and its dimension is

$$\dim(\mathbb{V}_{E, k}) = n_e^E k + \frac{k(k-1)}{2}, \quad (2.5)$$

where  $n_e^E$  is the number of edges of  $E$ . We notice that in the case  $k > 1$  a function  $v$  in  $\mathbb{V}_{E, k}$  is uniquely defined by

- the set of the values at the vertices of  $E$ ;
- the set of the values at the  $k - 1$  equally-spaced internal points on each edge of  $\partial E$ ;
- the set of the moments  $\frac{1}{|E|} \int_E v(\mathbf{x}) m(\mathbf{x}) d\mathbf{x} \forall m \in \mathcal{M}_{k-2}(E)$ ,

where the set  $\mathcal{M}_p(E)$ ,  $p \geq 0$ , is defined as

$$\mathcal{M}_p(E) = \left\{ \left( \frac{\mathbf{x} - \mathbf{x}_E}{h_E} \right)^s, |s| \leq p \right\}. \quad (2.6)$$

We will denote by  $\boldsymbol{\mu}_p(E, v) = \left( \frac{1}{|E|} \int_E v(\mathbf{x}) m(\mathbf{x}) d\mathbf{x} : m \in \mathcal{M}_p(E) \setminus \mathcal{M}_{p-1}(E) \right)$  the vector of the moments of  $v$  of order  $p$ . By  $|\boldsymbol{\mu}_p(E, v)|$  we will denote the  $l^2$ -norm of this vector.

We can now introduce the global discrete space as

$$\mathbb{V}_{\mathcal{T}} := \{v \in \mathbb{V} : v|_E \in \mathbb{V}_{E, k} \forall E \in \mathcal{T}\}.$$

On  $\mathcal{T}$  we need also to give the definition of the space of piecewise polynomial functions on  $\mathcal{T}$

$$\mathbb{W}_{\mathcal{T}}^k := \{w \in L^2(\Omega) : w|_E \in \mathbb{P}_k(E) \forall E \in \mathcal{T}\}, \quad (2.7)$$

and its subspace

$$\mathbb{V}_{\mathcal{T}}^0 := \mathbb{V}_{\mathcal{T}} \cap \mathbb{W}_{\mathcal{T}}^k, \quad (2.8)$$

which plays a crucial role in the forthcoming analysis.

We now introduce a series of projectors that will be used in the rest of the paper. For any  $E \in \mathcal{T}$ , we denote by  $\Pi_{p,E}^0 : L^2(E) \rightarrow \mathbb{P}_p(E)$  the  $L^2(E)$ -orthogonal projector onto the space of polynomial of degree  $p$  on  $E$ . Thanks to the choice of the enhanced space  $\mathbb{V}_{E,k}$  (2.4), we remark that  $\Pi_{k,E}^0 v$  and  $\Pi_{k-1,E}^0 \nabla v$  can be computed for any function  $v \in \mathbb{V}_{E,k}$ , see [1] for the details. To simplify the notation, in the following we will drop the symbol  $E$  from  $\Pi_{k,E}^0$  when no confusion arises. The global  $L^2$ -orthogonal projector is denoted by  $\Pi_{p,\mathcal{T}}^0 : L^2(\Omega) \rightarrow \mathbb{W}_{\mathcal{T}}^p$ .

We can also define the Lagrange interpolation operator  $\mathcal{I}_E : \mathbb{V}_{E,k} \rightarrow \mathbb{P}_k(E)$  on  $E$ , which builds a polynomial of degree  $k$  using the  $3k$  degrees of freedom on the boundary of  $E$  and the moments of order  $\leq k-3$ , since

$$\dim(\mathbb{P}_k(E)) = 3k + \frac{(k-1)(k-2)}{2}.$$

Moreover, we will denote by  $\mathcal{I}_{\mathcal{T}} : \mathbb{V}_{\mathcal{T}} \rightarrow \mathbb{W}_{\mathcal{T}}^k$  the Lagrange interpolation operator that restricts to  $\mathcal{I}_E$  on each  $E \in \mathcal{T}$ .

### 3. Discretization with data of degree $k-1$

In the rest of this paper, we assume that data  $\mathcal{D} = (A, c, f)$  are piecewise polynomials of degree  $k-1$  on the initial partition  $\mathcal{T}_0$ , hence on each partition  $\mathcal{T}$  obtained by newest-vertex refinement. Their values on each element of the triangulation will be denoted by

$$(A_E, c_E, f_E) \in (\mathbb{P}_{k-1}(E))^{2 \times 2} \times \mathbb{P}_{k-1}(E) \times \mathbb{P}_{k-1}(E).$$

We here define the bilinear forms that we need for the Galerkin discretization problem, starting from  $a_E, m_E : \mathbb{V}_{E,k} \times \mathbb{V}_{E,k} \rightarrow \mathbb{R}$ , such that

$$\begin{aligned} a_{\mathcal{T}}(v, w) &:= \sum_{E \in \mathcal{T}} \int_E (A_E \Pi_{k-1}^0 \nabla v) \cdot (\Pi_{k-1}^0 \nabla w) =: \sum_{E \in \mathcal{T}} a_E(v, w), \\ m_{\mathcal{T}}(v, w) &:= \sum_{E \in \mathcal{T}} \int_E c_E \Pi_k^0 v \Pi_k^0 w =: \sum_{E \in \mathcal{T}} m_E(v, w). \end{aligned}$$

We also introduce the symmetric bilinear form  $s_E : \mathbb{V}_E \times \mathbb{V}_E \rightarrow \mathbb{R}$  as

$$s_E(v, w) := \sum_{i=1}^{\bar{N}_E} v(\mathbf{x}_i) w(\mathbf{x}_i),$$

where  $\{\mathbf{x}_i\}_{i=1}^{\overline{N}_E}$  indicates the set of the degrees of freedom on the boundary of  $E$ . Indeed, we remark that in this case the stabilization term can be built without using the internal degrees of freedom, as shown in [7]. We assume for  $s_E$  the existence of two positive constant  $c_s$  and  $C_s$  independent on  $E$ , such that

$$c_s |v|_{1,E}^2 \leq s_E(v, v) \leq C_s |v|_{1,E}^2, \quad \forall v \in \mathbb{V}_E \setminus \mathbb{P}_k(E). \quad (3.1)$$

We define the local stabilizing form as

$$S_E(v, w) = s_E(v - \mathcal{I}_E v, w - \mathcal{I}_E w), \quad \forall v, w \in \mathbb{V}_E,$$

and the global stabilization form

$$S_{\mathcal{T}}(v, w) := \sum_{E \in \mathcal{T}} S_E(v, w), \quad \forall v, w \in \mathbb{V}_{\mathcal{T}}.$$

From (3.1), we get

$$S_{\mathcal{T}}(v, v) \simeq |v - \mathcal{I}_{\mathcal{T}} v|_{1,\mathcal{T}}^2, \quad \forall v \in \mathbb{V}_{\mathcal{T}},$$

where  $|\cdot|_{1,\mathcal{T}}$  denotes the broken  $H^1$ -seminorm over  $\mathcal{T}$ . Thus, we can now define the bilinear form  $\mathcal{B}_{\mathcal{T}}(\cdot, \cdot), \mathcal{B}_{\mathcal{T}} : \mathbb{V}_{\mathcal{T}} \times \mathbb{V}_{\mathcal{T}} \rightarrow \mathbb{R}$ , as

$$\mathcal{B}_{\mathcal{T}}(v, w) = a_{\mathcal{T}}(v, w) + m_{\mathcal{T}}(v, w) + \gamma S_{\mathcal{T}}(v, w), \quad (3.2)$$

with  $\gamma$  independent of  $\mathcal{T}$  satisfying  $\gamma \geq \gamma_0$  for some fixed  $\gamma_0 > 0$ . For the loading term we introduce  $\mathcal{F}_{\mathcal{T}} : \mathbb{V}_{\mathcal{T}} \rightarrow \mathbb{R}$  as

$$\mathcal{F}_{\mathcal{T}}(v) := \sum_{E \in \mathcal{T}} \int_E f_E \Pi_k^0 v = \sum_{E \in \mathcal{T}} \int_E f_E v, \quad \forall v \in \mathbb{V}_{\mathcal{T}}, \quad (3.3)$$

since  $f_E$  has been already approximated with a polynomial of degree  $k - 1$ . Note that the equality in (3.3) remains true if  $f_E$  is an approximation of  $f$  of degree  $k$  on  $E$ .

We have now defined all the forms that appear in the discrete formulation of the Problem (2.2). It reads as

$$\begin{cases} \text{find } u_{\mathcal{T}} \in \mathbb{V}_{\mathcal{T}} \text{ such that} \\ \mathcal{B}_{\mathcal{T}}(u_{\mathcal{T}}, v) = \mathcal{F}_{\mathcal{T}}(v), \quad \forall v \in \mathbb{V}_{\mathcal{T}}. \end{cases} \quad (3.4)$$

The bilinear form  $\mathcal{B}_{\mathcal{T}}$  is continuous and coercive, hence, there exists a unique and stable solution of the Problem (3.4). Furthermore, the following result extends Lemma 2.6 in [5].

**Lemma 3.1** (Galerkin quasi-orthogonality). *For any  $v \in \mathbb{V}_{\mathcal{T}}$  and  $w \in \mathbb{V}_{\mathcal{T}}^0$ , it holds*

$$\begin{aligned} a_{\mathcal{T}}(v, w) &= a(v, w) - \sum_{E \in \mathcal{T}} \int_E (A_E (I - \Pi_{k-1}^0) \nabla v) \cdot \nabla w, \\ m_{\mathcal{T}}(v, w) &= m(v, w) - \sum_{E \in \mathcal{T}} \int_E c_E ((I - \Pi_k^0) v) w, \\ S_{\mathcal{T}}(v, w) &= 0. \end{aligned}$$

Consequently,

$$|\mathcal{B}(u - u_{\mathcal{T}}, w)| \lesssim S_{\mathcal{T}}(u_{\mathcal{T}}, u_{\mathcal{T}})^{1/2} |w|_{1,\Omega},$$

where  $u$  is the solution of (2.2) and  $u_{\mathcal{T}}$  the solution of (3.4).



#### 4. The index of a node

A crucial concept, firstly introduced in [5] for the case  $k = 1$ , is the global index of a node: It will be used in the proofs of Section 5. In order to extend its definition to the case  $k > 1$ , we preliminarily introduce some useful definitions.

Let

$$\hat{E} := \{(x, y) \in \mathbb{R}^2 : x \geq 0, y \geq 0, x + y \leq 1\}$$

be the reference element and denote by  $\hat{R}_{\hat{E},k}$  the  $k$ -lattice built on  $\hat{E}$ , i.e.,

$$\hat{R}_{\hat{E},k} := \left\{ \left( \frac{i}{k}, \frac{j}{k} \right) \in \mathbb{R}^2 : i \geq 0, j \geq 0, i + j \leq k \right\}.$$

Considering the affine function  $F_E : \hat{E} \rightarrow E$  mapping the reference element onto an element  $E \in \mathcal{T}$ , we define the physical lattice on  $E$  by

$$R_{E,k} := F_E(\hat{R}_{\hat{E},k}),$$

and the set of *proper nodes* of  $E$  as the points of the physical lattice sitting on the boundary of  $E$ , i.e.,

$$\mathcal{P}_E := R_{E,k} \cap \partial E.$$

Observe that we implicitly assume that  $k \geq 2$  is sufficiently small so that interpolation on equally spaced nodes is numerically stable.

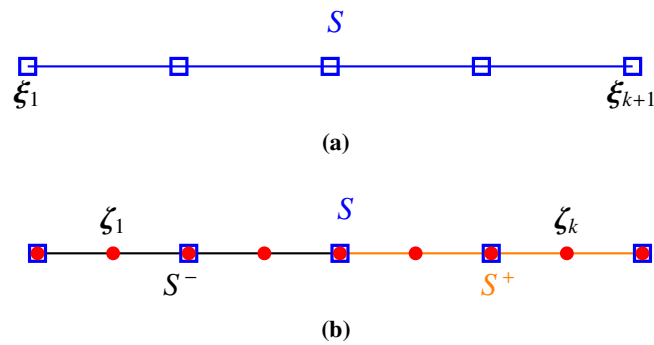
Next, we denote by  $\mathcal{H}_E$  the set of *hanging nodes* of  $E$ , i.e., the set of points  $\mathbf{x} \in \partial E$  that are not proper nodes of  $E$ , but that are proper nodes of some other element  $E'$ , i.e.,

$$\mathcal{H}_E := \{\mathbf{x} \in \partial E : \exists E' \in \mathcal{T} \text{ such that } \mathbf{x} \in \mathcal{P}_{E'}\} \setminus \mathcal{P}_E.$$

Finally, let  $\mathcal{N}_E := \mathcal{P}_E \cup \mathcal{H}_E$  be the set of all nodes sitting on  $E$ .

At the global level,  $\mathcal{N} := \bigcup_{E \in \mathcal{T}} \mathcal{N}_E$  will be the set of all nodes of the triangulation  $\mathcal{T}$ , which we split into the set  $\mathcal{P} := \{\mathbf{x} \in \mathcal{N} : \mathbf{x} \in \mathcal{P}_E \forall E \text{ containing } \mathbf{x}\}$  of the *proper nodes of  $\mathcal{T}$* , and the set  $\mathcal{H} := \mathcal{N} \setminus \mathcal{P}$  of the *hanging nodes of  $\mathcal{T}$* .

Next, let us clarify what happens when a hanging node is created. Let  $S$  be an element edge that is being refined, i.e., split into two contiguous edges  $S^-$  and  $S^+$ . Before the refinement,  $S$  contains  $k + 1$  equally-spaced nodes  $\xi_n$ ,  $n = 1, \dots, k + 1$ : The endpoints and the  $k - 1$  internal ones. After the refinement,  $S$  contains  $2k + 1$  nodes, precisely  $k + 1$  equally-spaced nodes on each sub-edge  $S^\pm$ , with the midpoint in common; see Figure 1. The spacing of the ‘old’ nodes on  $S$  was  $\frac{|S|}{k}$  (where  $|S|$  denotes the length of  $S$ ), whereas the spacing of the ‘new’ nodes is  $\frac{|S|}{2k}$ . Consequently,  $k + 1$  of these nodes coincide with those initially on  $S$ , and the new nodes introduced in the refinement are only  $k$ . We will denote these latter by  $\zeta_i$ ,  $i = 1, \dots, k$ .



**Figure 1.** Blue squares represent the  $k + 1$  equally-spaced nodes  $\xi_n$  on the edge  $S$  before refinement. Red circles represent the  $2k + 1$  nodes that arise after refinement. We have denoted by  $\zeta_i$  the new nodes that do not coincide with any  $\xi_n$ .

This suggests the following definition.

**Definition 4.1** (closest neighbors of a node). *With the previous notation, if  $x := \zeta_i$  is created as the midpoint of the segment  $[x', x''] := [\xi_{n_i}, \xi_{n_i+1}]$  for some  $n_i$ , we define the set  $\mathcal{B}(x) := \{x', x''\}$ .*

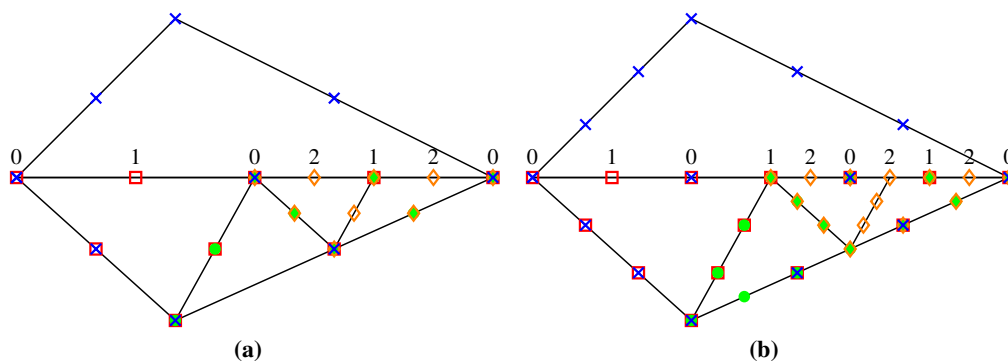
We are ready to give the announced definition of global index of a node of the triangulation  $\mathcal{T}$ .

**Definition 4.2** (global index of a node). *Given a node  $x \in \mathcal{N}$ , we define its global index  $\lambda$  recursively as follows:*

- If  $x$  is a proper node, then  $\lambda(x) := 0$ ;
- If  $x$  is a hanging node, with  $x', x'' \in \mathcal{B}(x)$ , then set

$$\lambda(x) := \max\{\lambda(x'), \lambda(x'')\} + 1.$$

Figure 2 shows the evolution of the global index after three refinements in the cases  $k = 2$  (a) and  $k = 3$  (b). We remark that, for instance, the midpoint of the horizontal edge is a proper node in case (a), and a hanging node in case (b).



**Figure 2.** Triangulation after the three refinements in the case  $k = 2$  (a) and in the case  $k = 3$  (b). Blue crosses represent the original degrees of freedom. Red squares, green circles and orange triangles are used for the degrees of freedom of the first, second and third refinement, respectively. All nodes are proper, except those on the horizontal line, whose global index is reported.

The largest global index in  $\mathcal{T}$  will be denoted by  $\Lambda_{\mathcal{T}} := \max_{\mathbf{x} \in \mathcal{N}} \{\lambda(\mathbf{x})\}$ . In this paper, as in [5], we will consider sequences of successively refined triangulations  $\{\mathcal{T}\}$  whose global index does not blow up.

**Assumption 4.3.** *There exists a constant  $\Lambda > 0$  such that, for any triangulation  $\mathcal{T}$  generated by successive refinements of  $\mathcal{T}_0$ , it holds*

$$\Lambda_{\mathcal{T}} \leq \Lambda.$$

Any such triangulation will be called  $\Lambda$ -admissible.

## 5. Two key properties

In this section we discuss the validity of some results for the degree  $k > 1$  that will be used in the rest of the paper. We will highlight in particular the differences from the case  $k = 1$ .

**Proposition 5.1** (scaled Poincaré inequality in  $\mathbb{V}_{\mathcal{T}}$ ). *There exists a constant  $C_P > 0$ , independent of  $\mathcal{T}$ , such that*

$$\sum_{E \in \mathcal{T}} h_E^{-2} \|v\|_{0,E}^2 \leq C_P |v|_{1,\Omega}^2, \quad \forall v \in \mathbb{V}_{\mathcal{T}} \text{ such that } v(\mathbf{x}) = 0, \forall \mathbf{x} \in \mathcal{P}. \quad (5.1)$$

*Proof.* Let  $E \in \mathcal{T}$  be an element of the triangulation. If  $E$  is an element of the original partition  $\mathcal{T}_0$ , all its vertices are proper nodes. Otherwise,  $E$  has been generated after some refinements by splitting an element  $\tilde{E}$  into two elements,  $E$  and  $E'$ . Let  $L$  be the common edge shared by  $E$  and  $E'$ . If  $L$  is not further refined, then all the nodes on  $L$  are proper because they are shared by  $E$  and  $E'$ . If  $L$  is refined and  $k$  is even, then the midpoint of  $L$  is a proper node.

So, let us consider the case  $k$  odd and let us assume that  $L$  is refined  $M \geq 1$  times. We focus in particular on the internal node  $\bar{\mathbf{x}}$  of  $L$  is at distance  $\frac{|L|}{k}$  from one of the endpoints, Figure 3 shows the case  $k = 3$ . This point belongs to one of the  $M + 1$  intervals in which  $L$  is refined, having width  $|L|/2^s$ , for some  $1 \leq s \leq M$ . We remark that  $s$  depends on how  $L$  has been refined (in the case of uniform refinements of  $L$ , one has  $2^s = M + 1$ ). We localize the chosen node  $\bar{\mathbf{x}}$  in  $L$  by defining an  $m \geq 0$  such that

$$\frac{|L| m}{2^s} \leq \frac{|L|}{k} \leq \frac{|L|(m+1)}{2^s},$$

or, equivalently,

$$k m \leq 2^s \leq k (m + 1). \quad (5.2)$$

The interval going from  $\frac{|L|m}{2^s}$  to  $\frac{|L|(m+1)}{2^s}$  is an edge for a smaller element  $E'$ , thus it contains  $k - 1$  internal nodes. Since they are equi-spaced, their positions are at

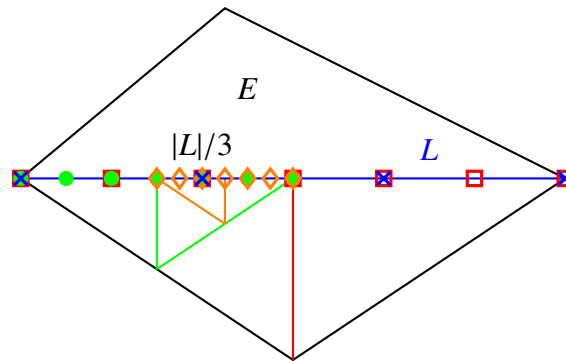
$$\frac{|L|}{2^s} \left( m + \frac{n}{k} \right) \quad \text{with } n = 0, \dots, k.$$

By taking  $n = 2^s - m k$ , which is compatible with conditions (5.2), we conclude that one of the internal nodes of  $E'$  coincides with  $\bar{\mathbf{x}}$ .

This guarantees that  $E$  has at least one proper node  $\mathbf{x}$  on its boundary. By hypothesis  $v(\mathbf{x}) = 0$ , and so we can apply the classical Poincaré inequality,

$$h_E^{-2} \|v\|_{0,E}^2 \lesssim |v|_{1,E}^2,$$

that concludes the proof.  $\square$



**Figure 3.** The case  $k = 3$  with 3 refinements of the edge  $L$  (in blue) is shown. Red, green and orange lines are the lines needed to refine  $L$  the first, the second and the third time respectively. Blue crosses are the degrees of freedom on  $L$  of the function living on  $E$ . Red squares, green circles, orange diamonds are the degrees of freedom on  $L$  generated after the first, the second and the third refinement of  $L$ .

*Remark 5.2.* The previous proof exploits the fact that when  $k > 1$ , each element of the triangulation contains at least a proper node. This differs from the case  $k = 1$  in which the edges do not contain internal nodes, and then elements with all hanging nodes as vertices are admissible. As a further difference from the case  $k = 1$ , we highlight that in Proposition 5.1 the constant  $C_P$  does not depend on the constant  $\Lambda$ , whose existence has been introduced in Assumption 4.3.

The next result we are going to establish is a hierarchical representation of the interpolation error  $v - \mathcal{I}_{EV}$  on the boundary  $\partial E$  of an element  $E \in \mathcal{T}$ . Assume that  $v \in \mathbb{V}_{E,k}$ , and let  $L$  be a side of the triangle  $E$ ; for simplicity, in the sequel the restriction of  $v$  to  $L$ , which is a piecewise polynomial of degree  $k$ , will be still denoted by  $v$ . The subsequent bisections of  $L$  which generate the nodes in  $\mathcal{N}_E \cap L$  allow us to write the difference  $(v - \mathcal{I}_{EV})|_L$  telescopically as

$$(v - \mathcal{I}_{EV})|_L = \sum_{j=1}^{J_L} (\mathcal{I}_j - \mathcal{I}_{j-1})v; \quad (5.3)$$

here,  $\mathcal{I}_0 = \mathcal{I}_{E|L}$ ,  $\mathcal{I}_{J_L}$  is the identity operator, whereas  $\mathcal{I}_j v$  for  $1 \leq j \leq J_L - 1$  is the piecewise polynomial of degree  $k$  which interpolates  $v$  on the partition of  $L$  of level  $j$ , namely the partition formed by sub-edges of length  $\leq \frac{|L|}{2^j}$ .

In order to understand the structure of the detail  $(\mathcal{I}_j - \mathcal{I}_{j-1})v$ , assume that  $S$  is a sub-edge of  $L$  of length  $= \frac{|L|}{2^{j-1}}$ , which is split into two sub-edges  $S^\pm$  of length  $= \frac{|L|}{2^j}$  (see Figure 1). On  $S$  we have two interpolation operators, namely

$$\mathcal{I} := \mathcal{I}_{j-1|S} : C^0(S) \rightarrow \mathbb{P}_k(S)$$

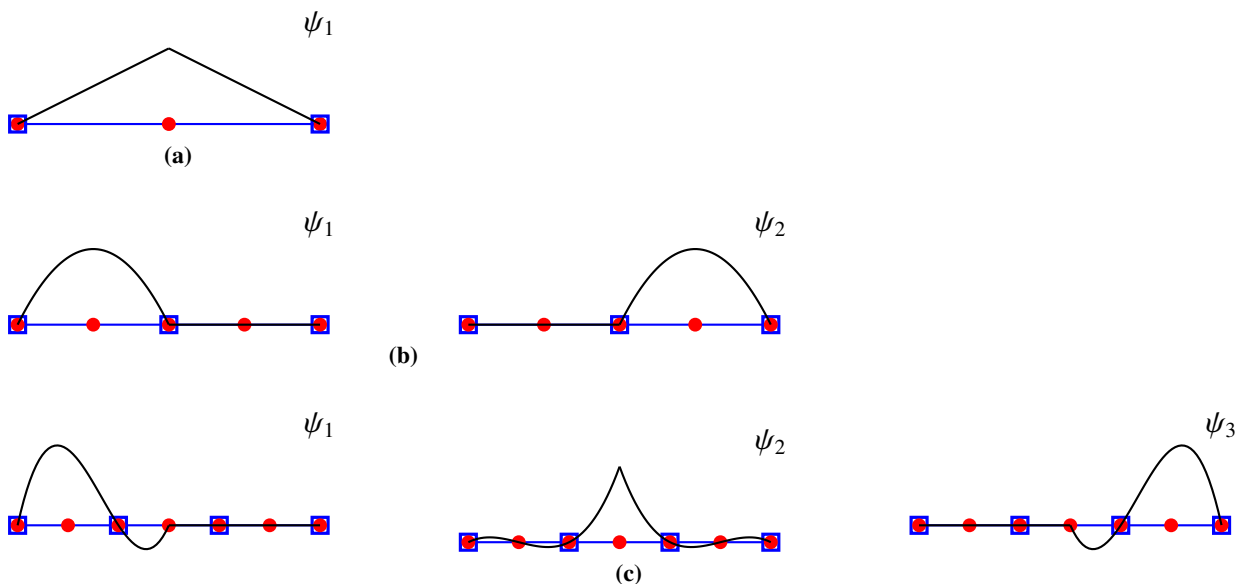
and

$$\mathcal{I}^\vee := \mathcal{I}_{j|_S} : C^0(S) \rightarrow \mathbb{P}_k(S^-, S^+) = \{v \in C^0(S) : v|_{S^-} \in \mathbb{P}_k(S^-) \text{ and } v|_{S^+} \in \mathbb{P}_k(S^+)\},$$

which coincides with the interpolation operator  $\mathcal{I}^- : C^0(S^-) \rightarrow \mathbb{P}_k(S^-)$  when restricted to  $S^-$  and with the analogous operator  $\mathcal{I}^+$  when restricted to  $S^+$ . With the notation introduced just before Definition 4.1, we can quantify the discrepancy between the two interpolation operators by defining the  $k$  basis functions

$$\psi_i \in \mathbb{P}_k(S^-, S^+) \text{ such that } \psi_i(x) = \begin{cases} 1 & \text{if } x = \zeta_i, \\ 0 & \text{if } x = \zeta_j, j \neq i, \\ 0 & \text{if } x = \xi_n, n = 1, \dots, k+1, \end{cases} \quad 1 \leq i \leq k.$$

See Figure 4 for a graphical representation of these functions in the cases  $k = 1$  (a),  $k = 2$  (b),  $k = 3$  (c).



**Figure 4.** Blue squares are the  $k + 1$  equi-spaced original nodes on the blue edge. Red points represent the nodes added after the refinement of the interval. Black lines show the shapes of the basis  $\psi_i$ ,  $i = 1, \dots, k$ , in the case  $k = 1$  (a),  $k = 2$  (b),  $k = 3$  (c).

Hence, the difference between the two interpolation operators on  $S$  can be written as

$$\mathcal{I}^\vee v - \mathcal{I}v = \sum_{i=1}^k d(v, \zeta_i) \psi_i,$$

where  $d$  is defined as

$$d(v, \zeta_i) := (\mathcal{I}^\vee v - \mathcal{I}v)(\zeta_i) = (v - \mathcal{I}v)(\zeta_i). \quad (5.4)$$

The values of  $\mathcal{I}v$  at the  $k$  nodes  $\zeta_i$  are a linear combination of the values of  $\mathcal{I}v$  at the  $k + 1$  nodes  $\zeta_n$ , where  $\mathcal{I}v$  coincides with  $v$ . Thus, there exist coefficients  $\alpha_{i,n}$  such that

$$(\mathcal{I}v)(\zeta_i) = \sum_{n=1}^{k+1} \alpha_{i,n} v(\zeta_n), \quad i = 1, \dots, k. \quad (5.5)$$

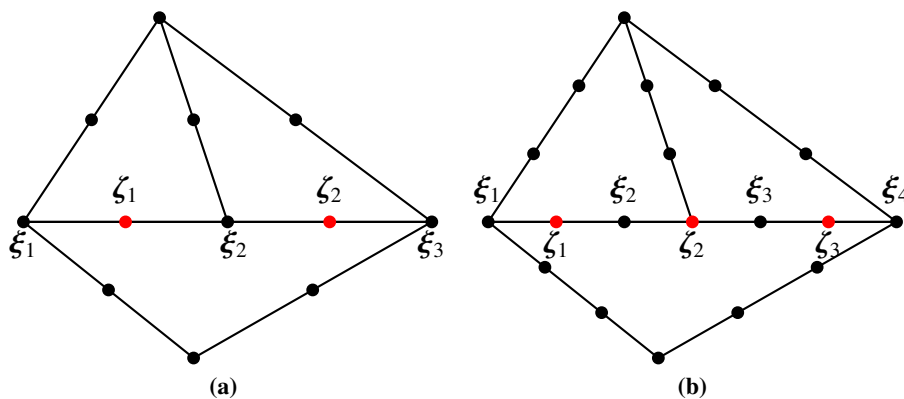
The explicit values of these coefficients in the case  $k = 2$  for the two new nodes  $\zeta_1$  and  $\zeta_2$  are given in these expressions:

$$\begin{aligned} (\mathcal{I}v)(\zeta_1) &= \frac{3}{8}v(\xi_1) + \frac{3}{4}v(\xi_2) - \frac{1}{8}v(\xi_3), \\ (\mathcal{I}v)(\zeta_2) &= -\frac{1}{8}v(\xi_1) + \frac{3}{4}v(\xi_2) + \frac{3}{8}v(\xi_3), \end{aligned}$$

where  $\xi_i \leq \zeta_i \leq \xi_{i+1}$ ,  $i = 1, 2$ . Similarly, in the case  $k = 3$ , we get

$$\begin{aligned} (\mathcal{I}v)(\zeta_1) &= \frac{5}{16}v(\xi_1) + \frac{15}{16}v(\xi_2) - \frac{5}{16}v(\xi_3) + \frac{1}{16}v(\xi_4), \\ (\mathcal{I}v)(\zeta_2) &= -\frac{1}{16}v(\xi_1) + \frac{9}{16}v(\xi_2) + \frac{9}{16}v(\xi_3) - \frac{1}{16}v(\xi_4), \\ (\mathcal{I}v)(\zeta_3) &= \frac{1}{16}v(\xi_1) - \frac{5}{16}v(\xi_2) + \frac{15}{16}v(\xi_3) + \frac{5}{16}v(\xi_4), \end{aligned}$$

where again  $\xi_i \leq \zeta_i \leq \xi_{i+1}$ ,  $i = 1, 2, 3$ . Figure 5 shows both cases. We notice that the coefficients  $\alpha_{i,n}$  depend only on the relative positions of the nodes on  $S$ , not on the level  $j$  of refinement.



**Figure 5.** Black points are the proper nodes. Red points represent the hanging nodes generated after a refinement. In (a) the case  $k = 2$  is shown,  $\zeta_1$  is the hanging node obtained after the refinement of  $\xi_1$  and  $\xi_3$  and it is the midpoint of  $\xi_1$  and  $\xi_2$ . We notice that if we have called the other red point  $\zeta_2$ ,  $\xi_1$  and  $\xi_3$  would have been switched. Analogously, (b) represents the case  $k = 3$ .

Summarizing, at the level  $j$  of refinement of the edge  $L$ , we get

$$(\mathcal{I}_j - \mathcal{I}_{j-1})v = \sum_{x \in \mathcal{H}_{L,j}} d(v, \mathbf{x})\psi_x,$$

where  $\mathcal{H}_{L,j}$  is the set of hanging nodes on  $L$  created at the level  $j$  of refinement, whereas

$$d(v, \mathbf{x}) = (\mathcal{I}_j v - \mathcal{I}_{j-1} v)(\mathbf{x}) = (v - \mathcal{I}_{j-1} v)(\mathbf{x}).$$

Summing-up over the levels and recalling (5.3), we obtain

$$(v - \mathcal{I}_E v)|_L = \sum_{x \in \mathcal{H}_L} d(v, \mathbf{x})\psi_x.$$

where  $\mathcal{H}_L = \mathcal{H}_E \cap L$ , whence

$$(v - \mathcal{I}_E v)|_{\partial E} = \sum_{\mathbf{x} \in \mathcal{H}_E} d(v, \mathbf{x}) \psi_{\mathbf{x}}.$$

We now introduce the subspace of  $\mathbb{V}_{E,k}$

$$X_E := \left\{ w \in \mathbb{V}_{E,k} : w(\mathbf{x}) = 0 \ \forall \mathbf{x} \in \mathcal{P}_E, \text{ and } \boldsymbol{\mu}_p(w, E) = \mathbf{0}, \ 0 \leq p \leq k-3 \right\},$$

which contains  $v - \mathcal{I}_E v$  by definition of  $\mathcal{I}_E$ . On  $X_E$ , we have two norms, namely the seminorm  $|w|_{1,E}$  (which is a norm on  $X_E$  due to the vanishing of  $w$  at the three vertices of  $E$ ) and the norm

$$\|w\|_{X_E} := \left( \sum_{\mathbf{x} \in \mathcal{H}_E} d^2(w, \mathbf{x}) + |\boldsymbol{\mu}_{k-2}(E, w)|^2 \right)^{1/2}.$$

Note that, due to Assumption 4.3, the dimension of  $X_E$  is uniformly bounded by a constant depending on  $\Lambda$ ; furthermore, the number of possible patterns of hanging nodes on  $\partial E$ , which determines the details  $d(w, \mathbf{x})$ , is also bounded in terms of  $\Lambda$ . As a consequence, the two norms are equivalent, with equivalence constants depending on  $\Lambda$ . Therefore,

$$\sum_{\mathbf{x} \in \mathcal{H}_E} d^2(w, \mathbf{x}) \leq \|w\|_{X_E}^2 \simeq |w|_{1,E}^2, \quad \forall w \in X_E.$$

Since  $v - \mathcal{I}_E v \in X_E$  and  $d(v - \mathcal{I}_E v, \mathbf{x}) = d(v, \mathbf{x})$  for any  $\mathbf{x} \in \mathcal{H}_E$ , we obtain

$$\sum_{\mathbf{x} \in \mathcal{H}_E} d^2(v, \mathbf{x}) \lesssim |v - \mathcal{I}_E v|_{1,E}^2.$$

Summing-up over all the elements of the triangulation, we arrive at the following result.

**Lemma 5.3** (global interpolation error vs hierarchical errors). *There exists a constant  $C_D > 0$  depending on  $\Lambda$  but independent of the triangulation  $\mathcal{T}$  such that*

$$\sum_{\mathbf{x} \in \mathcal{H}} d^2(v, \mathbf{x}) \leq C_D |v - \mathcal{I}_{\mathcal{T}} v|_{1,\mathcal{T}}^2, \quad \forall v \in \mathbb{V}_{\mathcal{T}}. \quad (5.6)$$

Next, we introduce the interpolation operator

$$\mathcal{I}_{\mathcal{T}}^0 : \mathbb{V}_{\mathcal{T}} \rightarrow \mathbb{V}_{\mathcal{T}}^0, \quad (5.7)$$

where  $\mathbb{V}_{\mathcal{T}}^0$  is defined in (2.8), by the following conditions:

- $(\mathcal{I}_{\mathcal{T}}^0 v)(\mathbf{x}) = v(\mathbf{x})$  for all  $\mathbf{x} \in \mathcal{P}$ ,
- $\boldsymbol{\mu}_p(E, \mathcal{I}_{\mathcal{T}}^0 v) = \boldsymbol{\mu}_p(E, v)$  for all  $0 \leq p \leq k-3$  and for all  $E \in \mathcal{T}$ .

These conditions uniquely identify  $\mathcal{I}_{\mathcal{T}}^0 v$ . Indeed, if  $\mathbf{x} \in \mathcal{H}$  is generated by a refinement of level  $j$  of an edge  $L$  (say,  $\mathbf{x} = \zeta_i$  with the notation introduced before Definition 4.1), then  $(\mathcal{I}_{\mathcal{T}}^0 v)(\mathbf{x})$  can be expressed in terms of the values of  $\mathcal{I}_{\mathcal{T}}^0 v$  at the  $k+1$  nodes (say,  $\boldsymbol{\xi}_n$ ) created at the previous levels of refinement of  $L$ , using the same coefficients as in formula (5.5), i.e.,

$$(\mathcal{I}_{\mathcal{T}}^0 v)(\zeta_i) = \sum_{n=1}^{k+1} \alpha_{i,n} (\mathcal{I}_{\mathcal{T}}^0 v)(\boldsymbol{\xi}_n), \quad i = 1, \dots, k; \quad (5.8)$$

and so on recursively.

The following result provides a representation of the error  $\mathcal{I}_{\mathcal{T}} v - \mathcal{I}_{\mathcal{T}}^0 v$ .

**Lemma 5.4.** *It holds*

$$|\mathcal{I}_{\mathcal{T}}v - \mathcal{I}_{\mathcal{T}}^0v|_{1,\mathcal{T}}^2 \simeq \sum_{\mathbf{x} \in \mathcal{H}} \delta^2(v, \mathbf{x}), \quad \forall v \in \mathbb{V}_{\mathcal{T}},$$

where  $\delta(v, \mathbf{x}) := v(\mathbf{x}) - (\mathcal{I}_{\mathcal{T}}^0v)(\mathbf{x})$ .

*Proof.* Consider an element  $E \in \mathcal{T}$ . Recall that by construction it holds  $\boldsymbol{\mu}_p(E, \mathcal{I}_E v) = \boldsymbol{\mu}_p(E, v) = \boldsymbol{\mu}_p(E, \mathcal{I}_{\mathcal{T}}^0v)$ , whence  $\boldsymbol{\mu}_p(\mathcal{I}_E v - \mathcal{I}_{\mathcal{T}}^0v, E) = \mathbf{0}$  for all  $0 \leq p \leq k-3$ . Consequently,

$$|\mathcal{I}_E v - \mathcal{I}_{\mathcal{T}}^0v|_{1,E}^2 \simeq \sum_{\mathbf{x} \in \mathcal{P}_E} |(\mathcal{I}_E v - \mathcal{I}_{\mathcal{T}}^0v)(\mathbf{x})|^2.$$

If  $\mathbf{x} \in \mathcal{P}_E$ ,  $(\mathcal{I}_E v)(\mathbf{x}) = v(\mathbf{x})$ , hence

$$|\mathcal{I}_E v - \mathcal{I}_{\mathcal{T}}^0v|_{1,E}^2 \simeq \sum_{\mathbf{x} \in \mathcal{P}_E} |(v - \mathcal{I}_{\mathcal{T}}^0v)(\mathbf{x})|^2.$$

Summing on all the elements of the partition, we get

$$\sum_{E \in \mathcal{T}} |\mathcal{I}_E v - \mathcal{I}_{\mathcal{T}}^0v|_{1,E}^2 \simeq \sum_{\mathbf{x} \in \mathcal{N}} |(v - \mathcal{I}_{\mathcal{T}}^0v)(\mathbf{x})|^2 \simeq \sum_{\mathbf{x} \in \mathcal{H}} |(v - \mathcal{I}_{\mathcal{T}}^0v)(\mathbf{x})|^2,$$

since if  $\mathbf{x} \in \mathcal{P}$ ,  $(\mathcal{I}_{\mathcal{T}}^0v)(\mathbf{x}) = v(\mathbf{x})$ . This concludes the proof.  $\square$

Concatenating Lemmas 5.3 and 5.4, we can prove the second key property of this section.

**Proposition 5.5** (comparison between interpolation operators). *Let  $\mathcal{T}$  be  $\Lambda$ -admissible. Then, there exists a constant  $C_I > 0$ , depending on  $\Lambda$ , but independent of  $\mathcal{T}$ , such that*

$$|v - \mathcal{I}_{\mathcal{T}}^0v|_{1,\Omega} \leq C_I |v - \mathcal{I}_{\mathcal{T}}v|_{1,\mathcal{T}}, \quad \forall v \in \mathbb{V}_{\mathcal{T}}.$$

*Proof.* Given a function  $v \in \mathbb{V}_{\mathcal{T}}$ , by the triangle inequality

$$|v - \mathcal{I}_{\mathcal{T}}^0v|_{1,\Omega} = |v - \mathcal{I}_{\mathcal{T}}^0v|_{1,\mathcal{T}} \leq |v - \mathcal{I}_{\mathcal{T}}v|_{1,\mathcal{T}} + |\mathcal{I}_{\mathcal{T}}v - \mathcal{I}_{\mathcal{T}}^0v|_{1,\mathcal{T}},$$

so it is enough to bound the last norm on the right-hand side. To this end, considering the vectors

$$\boldsymbol{\delta} = (\delta(\mathbf{x}))_{\mathbf{x} \in \mathcal{H}} := (\delta(v, \mathbf{x}))_{\mathbf{x} \in \mathcal{H}}, \quad \mathbf{d} = (d(\mathbf{x}))_{\mathbf{x} \in \mathcal{H}} := (d(v, \mathbf{x}))_{\mathbf{x} \in \mathcal{H}},$$

and recalling the two Lemmas, the proof can be concluded if we show that

$$\|\boldsymbol{\delta}\|_{\ell^2(\mathcal{H})} \lesssim \|\mathbf{d}\|_{\ell^2(\mathcal{H})}.$$

Given  $\mathbf{x} \in \mathcal{H}$ , assume that it is generated by a refinement of level  $j$  of an edge  $L$  (say,  $\mathbf{x} = \boldsymbol{\zeta}_i$  with the notation introduced before Definition 4.1). Writing  $v^* := \mathcal{I}_{\mathcal{T}}^0v$  for short, and exploiting formulas (5.4) and (5.5), we get

$$\begin{aligned} \delta(\boldsymbol{\zeta}_i) &= v(\boldsymbol{\zeta}_i) - v^*(\boldsymbol{\zeta}_i) = v(\boldsymbol{\zeta}_i) - \sum_{n=1}^{k+1} \alpha_{i,n} v^*(\boldsymbol{\xi}_n) \\ &= v(\boldsymbol{\zeta}_i) - \sum_{n=1}^{k+1} \alpha_{i,n} v(\boldsymbol{\xi}_n) - \sum_{n=1}^{k+1} \alpha_{i,n} (v^*(\boldsymbol{\xi}_n) - v(\boldsymbol{\xi}_n)) \\ &= d(\boldsymbol{\zeta}_i) + \sum_{n=1}^{k+1} \alpha_{i,n} \delta(\boldsymbol{\xi}_n). \end{aligned} \tag{5.9}$$



Thus, we can build a matrix  $\mathbf{W} : \ell^2(\mathcal{H}) \rightarrow \ell^2(\mathcal{H})$  such that  $\boldsymbol{\delta} = \mathbf{W}\mathbf{d}$ , and we just need to prove that

$$\|\mathbf{W}\|_2 \lesssim 1.$$

We now organize the hanging nodes with respect to the global index  $\lambda \in [1, \Lambda_{\mathcal{T}}]$ . Calling  $\mathcal{H}_\lambda = \{\mathbf{x} \in \mathcal{H} : \lambda(\mathbf{x}) = \lambda\}$ , and  $\mathcal{H} = \bigcup_{1 \leq \lambda \leq \Lambda_{\mathcal{T}}} \mathcal{H}_\lambda$ , the matrix  $\mathbf{W}$  can be factorized in lower triangular matrices  $\mathbf{W}_\lambda$ , that change the nodes of level  $\lambda$ , leaving the others unchanged. In particular,

$$\mathbf{W} = \mathbf{W}_{\Lambda_{\mathcal{T}}} \mathbf{W}_{\Lambda_{\mathcal{T}}-1} \dots \mathbf{W}_2 \mathbf{W}_1,$$

where  $\mathbf{W}_1$  is just the identity matrix  $\mathbf{I}$ , whereas each other matrix  $\mathbf{W}_\lambda$  differs from the identity only in the rows of block  $\lambda$ . In each of these rows, all entries are zero, but the entries  $\alpha_{i,n}$  in the off-diagonal part and 1 on the diagonal. In order to estimate  $\mathbf{W}_\lambda$ , we use the Hölder inequality  $\|\mathbf{W}_\lambda\|_2^2 \leq \|\mathbf{W}_\lambda\|_1 \|\mathbf{W}_\lambda\|_\infty$ . From the construction of  $\mathbf{W}_\lambda$  have that

$$\|\mathbf{W}_\lambda\|_\infty \leq \max_n \left\{ \sum_{i=1}^{k+1} |\alpha_{i,n}| \right\} + 1 =: \beta_1, \quad \|\mathbf{W}_\lambda\|_1 \leq 5k \max_{i,n} |\alpha_{i,n}| + 1 =: \beta_2,$$

where in the last inequality it has been used the fact that a hanging node of global index  $< \lambda$  may appear at most 5 times on the right-hand side of (5.9), since at most five edges meet at a node [5, Proposition 3.2]. These bring us to the following bound

$$\|\mathbf{W}\|_2 \leq \prod_{2 \leq \lambda \leq \Lambda_{\mathcal{T}}} \|\mathbf{W}_\lambda\|_2 \leq (\beta_1 \cdot \beta_2)^{\frac{\Lambda-1}{2}}$$

and the proof is concluded.  $\square$

## 6. A posteriori error estimator

With the aim of discussing the a posteriori error analysis, and following [12], we define the a posteriori error estimators, starting from the internal residual over an element  $E$ , i.e.,

$$r_{\mathcal{T}}(E; v, \mathcal{D}) := f_E + \nabla \cdot (A_E \Pi_{k-1}^0 \nabla v) - c_E \Pi_k^0 v, \quad (6.1)$$

for any  $v \in \mathbb{V}_{E,k}$ . We highlight that in the case  $k = 1$ , with piecewise constant data, the diffusion term in the residual vanishes. Furthermore, we define the jump residual over  $e$ , where  $e$  is an edge shared by two elements  $E_1$  and  $E_2$  of the partition  $\mathcal{T}$ , as

$$j_{\mathcal{T}}(e; v, \mathcal{T}) := [[A \Pi_{k-1}^0 \nabla v]]_e = (A_{E_1} \Pi_{k-1}^0 \nabla v|_{E_1}) \cdot \mathbf{n}_1 + (A_{E_2} \Pi_{k-1}^0 \nabla v|_{E_2}) \cdot \mathbf{n}_2,$$

where  $\mathbf{n}_i$  denotes the unit normal vector to  $e$  pointing outward with respect to  $E_i$ ; we set  $j_{\mathcal{T}}(e; v) = 0$  of  $e \in \partial\Omega$ . Then, let the local residual estimator associated with  $E$  be

$$\eta_{\mathcal{T}}^2(E; v, \mathcal{D}) := h_E^2 \|r_{\mathcal{T}}(E; v, \mathcal{D})\|_{0,E}^2 + \frac{1}{2} \sum_{e \in \mathcal{E}_E} h_E \|j_{\mathcal{T}}(e; v, \mathcal{D})\|_{0,e}^2, \quad (6.2)$$

and the global residual estimator as the sum of the local residuals

$$\eta_{\mathcal{T}}^2(v, \mathcal{D}) := \sum_{E \in \mathcal{T}} \eta_{\mathcal{T}}^2(E; v, \mathcal{D}).$$

In contrast to what has been done for the case  $k = 1$ , we also need to introduce the *virtual inconsistency terms*, defined by

$$\begin{aligned} \Psi_{\mathcal{T},A}^2(E; v, \mathcal{D}) &:= \|(I - \Pi_{k-1}^0)(A_E \Pi_{k-1}^0 \nabla v)\|_{0,E}^2, \\ \Psi_{\mathcal{T},c}^2(E; v, \mathcal{D}) &:= h_E^2 \|(I - \Pi_k^0)(c_E \Pi_k^0 v)\|_{0,E}^2, \end{aligned} \quad (6.3)$$

as well as their sum

$$\Psi_{\mathcal{T}}^2(v, \mathcal{D}) := \sum_{E \in \mathcal{T}} \Psi_{\mathcal{T}}^2(E; v, \mathcal{D}) := \sum_{E \in \mathcal{T}} \Psi_{\mathcal{T},A}^2(E; v, \mathcal{D}) + \Psi_{\mathcal{T},c}^2(E; v, \mathcal{D}). \quad (6.4)$$

## 7. A posteriori error estimates

In this section we present one of the main results of this paper, a stabilization-free a posteriori error bound. In this view, we firstly start by introducing the classical Clément operator upon the space  $\mathbb{V}_{\mathcal{T}}^0$ ,  $\tilde{\mathcal{I}}_{\mathcal{T}}^0 : \mathbb{V} \rightarrow \mathbb{V}_{\mathcal{T}}^0$ ; it is defined at the proper nodes on the skeleton of  $\mathcal{T}$  as the average of the target function on the support of the associated basis functions, whereas the internal moments (if any) coincide with those of the target function.

The scaled Poincaré inequality (Proposition 5.1) and Proposition 5.5 guarantee the validity of the error estimate for  $\tilde{\mathcal{I}}_{\mathcal{T}}^0$ . Given these propositions, its proof does not involve the polynomial degree  $k$ , hence, it does not change with respect to the one presented in [5].

**Lemma 7.1** (Clément interpolation estimate).  $\forall v \in \mathbb{V}$ , it holds

$$\sum_{E \in \mathcal{T}} h_E^{-2} \|v - \tilde{\mathcal{I}}_{\mathcal{T}}^0 v\|_{0,E}^2 \lesssim |v|_{1,\Omega}^2,$$

where the hidden constant depends on  $\Lambda$  but not on  $\mathcal{T}$ .

We can now prove the following results, which is similar to Theorem 13 in [12], but with a slightly modified proof.

**Proposition 7.2** (upper bound). *There exists a constant  $C_{apost} > 0$ , independent of  $u$ ,  $\mathcal{T}$ ,  $u_{\mathcal{T}}$  and  $\gamma$ , such that*

$$|u - u_{\mathcal{T}}|_{1,\Omega}^2 \leq C_{apost} \left( \eta_{\mathcal{T}}^2(u_{\mathcal{T}}, \mathcal{D}) + S_{\mathcal{T}}(u_{\mathcal{T}}, u_{\mathcal{T}}) \right). \quad (7.1)$$

*Proof.* For any  $v \in \mathbb{V}$ , using the definition of Problem (2.2), we have that

$$\begin{aligned} \mathcal{B}(u - u_{\mathcal{T}}, v) &= \mathcal{B}(u, v) - \mathcal{B}(u_{\mathcal{T}}, v) - (f, v_{\mathcal{T}}) + \mathcal{B}(u, v_{\mathcal{T}}) \\ &= (f, v - v_{\mathcal{T}}) - \mathcal{B}(u_{\mathcal{T}}, v) + \mathcal{B}(u, v_{\mathcal{T}}) - \mathcal{B}(u_{\mathcal{T}}, v_{\mathcal{T}}) + \mathcal{B}(u_{\mathcal{T}}, v_{\mathcal{T}}) \\ &= ((f, v - v_{\mathcal{T}}) - \mathcal{B}(u_{\mathcal{T}}, v - v_{\mathcal{T}})) + \mathcal{B}(u - u_{\mathcal{T}}, v_{\mathcal{T}}) =: I + II, \end{aligned}$$

where  $v_{\mathcal{T}} := \tilde{I}_{\mathcal{T}}^0 v \in \mathbb{V}_{\mathcal{T}}^0$ . The first term can be written as

$$\begin{aligned} I &= \sum_{E \in \mathcal{T}} \left\{ \int_E f_E(v - v_{\mathcal{T}}) - \int_E A_E \nabla u_{\mathcal{T}} \cdot \nabla(v - v_{\mathcal{T}}) - \int_E c_E u_{\mathcal{T}}(v - v_{\mathcal{T}}) \right\} \\ &= \sum_{E \in \mathcal{T}} \left\{ \int_E f_E(v - v_{\mathcal{T}}) - \int_E (A_E \Pi_{k-1}^0 \nabla u_{\mathcal{T}}) \cdot \nabla(v - v_{\mathcal{T}}) - \int_E (c_E \Pi_k^0 u_{\mathcal{T}})(v - v_{\mathcal{T}}) \right\} \\ &\quad + \sum_{E \in \mathcal{T}} \left\{ \int_E (A_E(\Pi_{k-1}^0 - I) \nabla u_{\mathcal{T}}) \cdot \nabla(v - v_{\mathcal{T}}) + \int_E (c_E(\Pi_k^0 - I) u_{\mathcal{T}})(v - v_{\mathcal{T}}) \right\} =: I_1 + I_2. \end{aligned}$$

The addend  $I_1$  can be expressed as

$$\begin{aligned} I_1 &= \sum_{E \in \mathcal{T}} \left\{ \int_E (f_E + \nabla \cdot (A_E \Pi_{k-1}^0 \nabla u_{\mathcal{T}}) - c_E \Pi_k^0 u_{\mathcal{T}})(v - v_{\mathcal{T}}) \right\} \\ &\quad + \sum_{E \in \mathcal{T}} \int_{\partial E} \mathbf{n} \cdot (A_E \Pi_{k-1}^0 \nabla u_{\mathcal{T}})(v - v_{\mathcal{T}}), \end{aligned}$$

which can be bounded by using Lemma 7.1,

$$|I_1| \lesssim \eta_{\mathcal{T}}(u_{\mathcal{T}}, \mathcal{D}) |v|_{1,\Omega}.$$

On the other hand, noting that

$$\begin{aligned} \|(I - \Pi_{k-1}^0) \nabla u_{\mathcal{T}}\|_{0,E} &= \|(I - \Pi_{k-1}^0) \nabla (I - \Pi_k^0) u_{\mathcal{T}}\|_{0,E} \\ &\leq \|\nabla (I - \Pi_k^0) u_{\mathcal{T}}\|_{0,E} \end{aligned} \quad (7.2)$$

and applying again Lemma 7.1 and the stability of the Clément operator in the  $H^1$  norm, the addend  $I_2$  can be bounded as follows:

$$\begin{aligned} |I_2| &\leq \left( \sum_{E \in \mathcal{T}} \|A_E (I - \Pi_{k-1}^0) \nabla u_{\mathcal{T}}\|_{0,E}^2 \|\nabla(v - v_{\mathcal{T}})\|_{0,E}^2 + \sum_{E \in \mathcal{T}} h_E^2 \|c_E (I - \Pi_k^0) u_{\mathcal{T}}\|_{0,E}^2 h_E^{-2} \|v - v_{\mathcal{T}}\|_{0,E}^2 \right)^{1/2} \\ &\lesssim \left( \sum_{E \in \mathcal{T}} \|A_E (I - \Pi_{k-1}^0) \nabla u_{\mathcal{T}}\|_{0,E}^2 + h_E^2 \|c_E (I - \Pi_k^0) u_{\mathcal{T}}\|_{0,E}^2 \right)^{1/2} |v|_{1,\Omega} \\ &\lesssim \left( \sum_{E \in \mathcal{T}} \|\nabla(u_{\mathcal{T}} - \Pi_k^0 u_{\mathcal{T}})\|_{0,E}^2 + h_E^2 \|(u_{\mathcal{T}} - \Pi_k^0 u_{\mathcal{T}})\|_{0,E}^2 \right)^{1/2} |v|_{1,\Omega} \\ &\lesssim (S_{\mathcal{T}}(u_{\mathcal{T}}, u_{\mathcal{T}}))^{1/2} |v|_{1,\Omega}. \end{aligned}$$

Looking now at the term  $II$ , we have by Lemma 3.1

$$|\mathcal{B}(u - u_{\mathcal{T}}, v)| \lesssim S_{\mathcal{T}}(u_{\mathcal{T}}, u_{\mathcal{T}})^{1/2} |v|_{1,\Omega}.$$

Finally, by taking  $v := u - u_{\mathcal{T}} \in \mathbb{V}$ , we get

$$\mathcal{B}(u - u_{\mathcal{T}}, u - u_{\mathcal{T}}) \lesssim (\eta_{\mathcal{T}}(u_{\mathcal{T}}, \mathcal{D}) + S_{\mathcal{T}}(u_{\mathcal{T}}, u_{\mathcal{T}})^{1/2}) |u - u_{\mathcal{T}}|_{1,\Omega},$$

which, using the coercivity of  $\mathcal{B}$ , concludes the proof.  $\square$

We now report a bound for the local residual estimator, proved in [12, Theorem 16].

**Proposition 7.3** (local lower bound). *There holds*

$$\eta_{\mathcal{T}}^2(E; u_{\mathcal{T}}, \mathcal{D}) \lesssim \sum_{E' \in w_E} (|u - u_{\mathcal{T}}|_{1,E'}^2 + S_{E'}(u_{\mathcal{T}}, u_{\mathcal{T}}))$$

where  $w_E := \{E' : |\partial E \cap \partial E'| \neq \emptyset\}$ . The hidden constant is independent of  $\gamma$ ,  $h$ ,  $u$  and  $u_{\mathcal{T}}$ .

Summing on all the elements of the partition, we get the following corollary.

**Corollary 7.4** (global lower bound). *There exists a constant  $c_{\text{apost}} > 0$ , independent of  $u$ ,  $\mathcal{T}$ ,  $u_{\mathcal{T}}$  and  $\gamma$ , such that*

$$c_{\text{apost}} \eta_{\mathcal{T}}^2(u_{\mathcal{T}}, \mathcal{D}) \leq |u - u_{\mathcal{T}}|_{1,\Omega}^2 + S_{\mathcal{T}}(u_{\mathcal{T}}, u_{\mathcal{T}}).$$

In the following proposition we present a bound of the stabilization term. We remark that in the case  $k = 1$  the inconsistency term does not appear.

**Proposition 7.5** (bound of the stabilization term). *There exists a constant  $C_B > 0$  independent of  $\mathcal{T}$ ,  $u_{\mathcal{T}}$  and  $\gamma$ , such that*

$$\gamma^2 S_{\mathcal{T}}(u_{\mathcal{T}}, u_{\mathcal{T}}) \leq C_B (\eta_{\mathcal{T}}^2(u_{\mathcal{T}}, \mathcal{D}) + \Psi_{\mathcal{T}}^2(u_{\mathcal{T}}, \mathcal{D})). \quad (7.3)$$

*Proof.* From the Definition (3.2) of the form  $\mathcal{B}_{\mathcal{T}}$  and from (3.4),  $\forall w \in \mathbb{V}_{\mathcal{T}}^0$  it holds

$$\begin{aligned} \gamma S_{\mathcal{T}}(u_{\mathcal{T}}, u_{\mathcal{T}}) &= \gamma S_{\mathcal{T}}(u_{\mathcal{T}}, u_{\mathcal{T}} - w) \\ &= \mathcal{B}_{\mathcal{T}}(u_{\mathcal{T}}, u_{\mathcal{T}} - w) - a_{\mathcal{T}}(u_{\mathcal{T}}, u_{\mathcal{T}} - w) - m_{\mathcal{T}}(u_{\mathcal{T}}, u_{\mathcal{T}} - w) \\ &= \mathcal{F}(u_{\mathcal{T}} - w) - a_{\mathcal{T}}(u_{\mathcal{T}}, u_{\mathcal{T}} - w) - m_{\mathcal{T}}(u_{\mathcal{T}}, u_{\mathcal{T}} - w). \end{aligned}$$

Defining  $e_{\mathcal{T}} := u_{\mathcal{T}} - w$ , we get

$$\gamma S_{\mathcal{T}}(u_{\mathcal{T}}, e_{\mathcal{T}}) = \sum_{E \in \mathcal{T}} \left\{ \int_E f e_{\mathcal{T}} - \int_E (A_E \Pi_{k-1}^0(\nabla u_{\mathcal{T}})) \cdot \Pi_{k-1}^0(\nabla e_{\mathcal{T}}) - \int_E c_E \Pi_k^0 u_{\mathcal{T}} \Pi_k^0 e_{\mathcal{T}} \right\}. \quad (7.4)$$

We notice that

$$\begin{aligned} \int_E (A_E \Pi_{k-1}^0(\nabla u_{\mathcal{T}})) \cdot \Pi_{k-1}^0(\nabla e_{\mathcal{T}}) &= \int_E (\Pi_{k-1}^0(A_E \Pi_{k-1}^0 \nabla u_{\mathcal{T}})) \cdot \nabla e_{\mathcal{T}} \\ &= \int_E (\Pi_{k-1}^0 - I)(A_E \Pi_{k-1}^0 \nabla u_{\mathcal{T}}) \cdot \nabla e_{\mathcal{T}} + \int_E (A_E \Pi_{k-1}^0 \nabla u_{\mathcal{T}}) \cdot \nabla e_{\mathcal{T}} \end{aligned} \quad (7.5)$$

and

$$\begin{aligned} \int_E c_E \Pi_k^0 u_{\mathcal{T}} \Pi_k^0 e_{\mathcal{T}} &= \int_E \Pi_k^0(c_E \Pi_k^0 u_{\mathcal{T}}) e_{\mathcal{T}} \\ &= \int_E (\Pi_k^0 - I)(c_E \Pi_k^0 u_{\mathcal{T}}) e_{\mathcal{T}} + \int_E c_E (\Pi_k^0 u_{\mathcal{T}}) e_{\mathcal{T}}. \end{aligned} \quad (7.6)$$

By substituting (7.5) and (7.6) into (7.4), it results

$$\begin{aligned}
\gamma S_{\mathcal{T}}(u_{\mathcal{T}}, u_{\mathcal{T}}) &= \sum_{E \in \mathcal{T}} \int_E (f + \nabla \cdot (A_E \Pi_{k-1}^0 \nabla u_{\mathcal{T}}) - c_E \Pi_k^0 u_{\mathcal{T}}) e_{\mathcal{T}} - \sum_{E \in \mathcal{T}} \int_{\partial E} \mathbf{n} \cdot \nabla (A_E \Pi_{k-1}^0 \nabla u_{\mathcal{T}}) e_{\mathcal{T}} \\
&\quad + \sum_{E \in \mathcal{T}} \int_E (I - \Pi_{k-1}^0)(A_E \Pi_{k-1}^0 \nabla u_{\mathcal{T}}) \cdot \nabla e_{\mathcal{T}} + \sum_{E \in \mathcal{T}} \int_E (I - \Pi_k^0)(c_E \Pi_k^0 u_{\mathcal{T}}) e_{\mathcal{T}} \\
&\leq \sum_{E \in \mathcal{T}} h_E \|r_{\mathcal{T}}(E; u_{\mathcal{T}}, \mathcal{D})\|_{0,E} h_E^{-1} \|e_{\mathcal{T}}\|_{0,E} + \frac{1}{2} \sum_{e \in \mathcal{E}} h_e^{1/2} \|j_{\mathcal{T}}(e; u_{\mathcal{T}}, \mathcal{D})\|_{0,e} h_e^{-1/2} \|e_{\mathcal{T}}\|_{0,e} \\
&\quad + \sum_{E \in \mathcal{T}} \|(I - \Pi_{k-1}^0)(A_E \Pi_{k-1}^0 \nabla u_{\mathcal{T}})\|_{0,E} \|\nabla e_{\mathcal{T}}\|_{0,E} + \sum_{E \in \mathcal{T}} h_E \|(I - \Pi_k^0) c_E \Pi_k^0 u_{\mathcal{T}}\|_{0,E} h_E^{-1} \|e_{\mathcal{T}}\|_{0,E} \\
&\leq \sum_{E \in \mathcal{T}} h_E \|r_{\mathcal{T}}(E; u_{\mathcal{T}}, \mathcal{D})\|_{0,E} h_E^{-1} \|e_{\mathcal{T}}\|_{0,E} + \frac{1}{2} \sum_{e \in \mathcal{E}} h_e^{1/2} \|j_{\mathcal{T}}(e; u_{\mathcal{T}}, \mathcal{D})\|_{0,e} h_e^{-1/2} \|e_{\mathcal{T}}\|_{0,e} \\
&\quad + C_{\text{inv}} \sum_{E \in \mathcal{T}} \Psi_A(E; u_{\mathcal{T}}, \mathcal{D}) h_E^{-1} \|e_{\mathcal{T}}\|_{0,E} + \sum_{E \in \mathcal{T}} \Psi_c(E; u_{\mathcal{T}}, \mathcal{D}) h_E^{-1} \|e_{\mathcal{T}}\|_{0,E}.
\end{aligned}$$

With the same strategy used in [5], for any  $\delta > 0$ , we get

$$\gamma S_{\mathcal{T}}(u_{\mathcal{T}}, u_{\mathcal{T}}) \leq \frac{1}{2\delta} \left( \eta_{\mathcal{T}}^2(u_{\mathcal{T}}, \mathcal{D}) + \Psi_{\mathcal{T}}^2(u_{\mathcal{T}}, \mathcal{D}) \right) + \frac{\delta}{2} \Phi_{\mathcal{T}}(e_{\mathcal{T}}),$$

where

$$\Phi_{\mathcal{T}}(e_{\mathcal{T}}) = \sum_{E \in \mathcal{T}} \left\{ \max\{C_{\text{inv}}^2, 1\} h_E^{-2} \|e_{\mathcal{T}}\|_{0,E}^2 + \frac{1}{2} \sum_{e \in \mathcal{E}_E} h_e^{-1} \|e_{\mathcal{T}}\|_{0,e}^2 \right\}.$$

Posing now  $w = \mathcal{I}_{\mathcal{T}}^0 u_{\mathcal{T}}$  and applying Proposition 5.1, we get

$$\Phi_{\mathcal{T}}(u_{\mathcal{T}} - \mathcal{I}_{\mathcal{T}}^0 u_{\mathcal{T}}) \lesssim |u_{\mathcal{T}} - \mathcal{I}_{\mathcal{T}}^0 u_{\mathcal{T}}|_{1,\Omega}^2,$$

whereas Proposition 5.5 yields

$$|u_{\mathcal{T}} - \mathcal{I}_{\mathcal{T}}^0 u_{\mathcal{T}}|_{1,\Omega}^2 \lesssim |u_{\mathcal{T}} - \mathcal{I}_{\mathcal{T}} u_{\mathcal{T}}|_{1,\Omega}^2 \simeq S_{\mathcal{T}}(u_{\mathcal{T}}, u_{\mathcal{T}}),$$

so we obtain

$$\gamma^2 S_{\mathcal{T}}(u_{\mathcal{T}}, u_{\mathcal{T}}) \leq C_B \left( \eta_{\mathcal{T}}^2(u_{\mathcal{T}}, \mathcal{D}) + \Psi_{\mathcal{T}}^2(u_{\mathcal{T}}, \mathcal{D}) \right),$$

for a suitable constant  $C_B > 0$ . □

Combining Propositions 7.2 and 7.5, we arrive at the following key result.

**Corollary 7.6** (stabilization-free a posteriori error upper bound). *It holds*

$$|u - u_{\mathcal{T}}|_{1,\Omega} \leq C_{U_1} \eta_{\mathcal{T}}^2(u_{\mathcal{T}}, \mathcal{D}) + C_{U_2} \Psi_{\mathcal{T}}^2(u_{\mathcal{T}}, \mathcal{D}),$$

where  $C_{U_1} = C_{\text{apost}} \left( \frac{C_B}{\gamma^2} + 1 \right)$  and  $C_{U_2} = C_{\text{apost}} \frac{C_B}{\gamma^2}$ .

*Remark 7.7.* Note that the chosen stabilization affects the value of the constant  $c_{\text{apost}}$ , which in principle may depend on the polynomial degree and the geometry of the mesh. However, this dependence is under control; indeed, (i) we are not proposing a  $p$ -method, so the polynomial degree  $k$  is fixed, (ii) the refinement procedure is obtained by newest-vertex bisection, which guarantees shape regularity on the refined elements, (iii) Assumption 4.3 enforces an upper bound on the number of hanging nodes on each edge.

## 8. The effect of a mesh refinement

In view of the convergence analysis of the adaptive algorithm GALERKIN, in this section we analyse the effect of refining the partition  $\mathcal{T}$  by applying one or more newest-vertex bisections to some of its elements. Specifically, in Sect. 8.1 we prove that the residual estimator (6.2) is reduced by a fixed fraction (up to an addend proportional to the stabilization term) when the element  $E$  is split into two elements by one bisection. We prove a similar result for the inconsistency term estimator (6.4), provided a suitable number of bisections is applied to  $E$ . Next, in Sect. 8.2 we establish a quasi-orthogonality property in the energy norm between the solutions on two partitions, one being a refinement of the other.

### 8.1. Reduction of estimators under refinement

Let us consider an element  $E$  in  $\mathcal{T}$  which is bisected into elements  $E_1$  and  $E_2$ ; the refined partition containing these two elements will be denoted by  $\mathcal{T}_*$ . Given  $v \in \mathbb{V}_{\mathcal{T}}$ , we notice that  $v$  is known on  $\partial E$ , and in particular at the new vertex of  $E_1$  and  $E_2$  produced by the bisection. Denoting by  $e = E_1 \cap E_2$  the new edge, we associate a function  $v_* \in \mathbb{V}_{\mathcal{T}_*}$  to  $v$  such that  $v_*|_{\partial E} = v|_{\partial E}$ ,  $v_*|_e \in \mathbb{P}_1(e)$ , and  $\boldsymbol{\mu}_p(E_i, v_*) = \boldsymbol{\mu}_p(E, v)$  for all  $0 \leq p \leq k-2$  and for  $i = 1, 2$ . In the following we will write  $v$  instead of  $v_*$  when no confusion arises.

#### 8.1.1. The residual estimator

Let  $\eta_{\mathcal{T}}(E; v, \mathcal{D})$  be defined in (6.2) and  $\eta_{\mathcal{T}_*}(E; v, \mathcal{D})$  be the sum of the local residual estimators on the two newly formed elements, defined as follows:

$$\eta_{\mathcal{T}_*}^2(E; v, \mathcal{D}) := \sum_{i=1}^2 \eta_{\mathcal{T}_*}^2(E_i; v, \mathcal{D}) = \sum_{i=1}^2 \left\{ h_{E_i}^2 \|r_{\mathcal{T}}(E_i; v, \mathcal{D})\|_{0, E_i}^2 + \frac{1}{2} \sum_{e \in \mathcal{E}_{E_i}} h_{E_i} \|j_{\mathcal{T}}(e; v, \mathcal{D})\|_{0, e}^2 \right\},$$

where we recall that  $h_{E_i} = \frac{1}{\sqrt{2}} h_E$ ,  $i = 1, 2$ . We notice that, since  $\mathcal{D}$  does not change under refinement, the functions  $f_{E_i} = f_E|_{E_i}$ ,  $c_{E_i} = c_E|_{E_i}$  and  $A_{E_i} = A_E|_{E_i}$  will be denoted again by  $f_E$ ,  $c_E$  and  $A_E$ , respectively.

**Lemma 8.1** (local residual estimator reduction). *There exist constants  $\mu_r \in (0, 1)$  and  $c_{er,1} > 0$  such that for any  $v \in \mathbb{V}_{\mathcal{T}}$*

$$\eta_{\mathcal{T}_*}(E; v, \mathcal{D}) \leq \mu_r \eta_{\mathcal{T}}(E; v, \mathcal{D}) + c_{er,1} S_{\mathcal{T}(E)}^{1/2}(v, v),$$

where  $S_{\mathcal{T}(E)}(v, v) := \sum_{E' \in \mathcal{T}(E)} S_{E'}(v, v)$  with  $\mathcal{T}(E) := \{E' \in \mathcal{T} : \mathcal{E}_E \cap \mathcal{E}_{E'} \neq \emptyset\}$ .

*Proof.* Recalling the Definition (6.1), we have the following residuals

$$\begin{aligned} r_E &:= f_E + \nabla \cdot (A_E \Pi_{k-1, E}^0 \nabla v) - c_E \Pi_{k, E}^0 v, \\ r_{E_i} &:= f_E + \nabla \cdot (A_E \Pi_{k-1, E_i}^0 \nabla v) - c_E \Pi_{k, E_i}^0 v. \end{aligned}$$

Writing

$$r_{E_i} = r_E - \nabla \cdot (A_E \Pi_{k-1, E}^0 \nabla v - A_E \Pi_{k-1, E_i}^0 \nabla v) + c_E \Pi_{k, E}^0 v - c_E \Pi_{k, E_i}^0 v,$$

we get, for any  $\epsilon > 0$ ,

$$\begin{aligned} \sum_{i=1}^2 h_{E_i}^2 \|r_{E_i}\|_{0,E_i}^2 &\leq \sum_{i=1}^2 h_{E_i}^2 (1 + \epsilon) \|r_E\|_{0,E_i}^2 \\ &\quad + 2 \sum_{i=1}^2 h_{E_i}^2 \left(1 + \frac{1}{\epsilon}\right) \|\nabla \cdot (A_E (\Pi_{k-1,E}^0 \nabla v - \Pi_{k-1,E_i}^0 \nabla v))\|_{0,E_i}^2 \\ &\quad + 2 \sum_{i=1}^2 h_{E_i}^2 \left(1 + \frac{1}{\epsilon}\right) \|c_E (\Pi_{k,E}^0 v - \Pi_{k,E_i}^0 v)\|_{0,E_i}^2. \end{aligned}$$

The second term can be bounded by using the inverse inequality and the minimality of  $\Pi_{k-1,E_i}^0$  as follows:

$$\begin{aligned} \sum_{i=1}^2 h_{E_i}^2 \|\nabla \cdot (A_E (\Pi_{k-1,E}^0 \nabla v - \Pi_{k-1,E_i}^0 \nabla v))\|_{0,E_i}^2 &\lesssim \sum_{i=1}^2 \|\Pi_{k-1,E}^0 \nabla v - \Pi_{k-1,E_i}^0 \nabla v\|_{0,E_i}^2 \\ &\leq 2 \|\nabla v - \Pi_{k-1,E}^0 \nabla v\|_{0,E}^2 + 2 \sum_{i=1}^2 \|\nabla v - \Pi_{k-1,E_i}^0 \nabla v\|_{0,E_i}^2 \\ &\leq 4 \|\nabla v - \Pi_{k-1,E}^0 \nabla v\|_{1,E}^2 \lesssim |v - \mathcal{I}_E v|_{1,E}^2 \lesssim S_E(v, v), \end{aligned}$$

while, for the last term, using the Poincaré inequality we have

$$\begin{aligned} \sum_{i=1}^2 h_{E_i}^2 \|c_E (\Pi_{k,E}^0 v - \Pi_{k,E_i}^0 v)\|_{0,E_i}^2 &\lesssim h_E^2 \sum_{i=1}^2 \|\Pi_{k,E}^0 v - \Pi_{k,E_i}^0 v\|_{0,E_i}^2 \\ &\leq h_E^2 \|v - \Pi_{k,E}^0 v\|_{0,E}^2 \lesssim h_E^2 |v - \Pi_{k,E}^0 v|_{1,E}^2 \lesssim h_E^2 S_E(v, v). \end{aligned}$$

Finally, taking an appropriate value of  $\epsilon$  and setting  $\mu := \frac{1+\epsilon}{2} \in (0, 1)$  (for instance, if  $\epsilon = \frac{1}{2}$ ,  $\mu = \frac{3}{4}$ ) we get

$$\sum_{i=1}^2 h_{E_i}^2 \|r_{E_i}\|_{0,E_i}^2 \leq \mu h_E^2 \|r_E\|_{0,E}^2 + C(1 + h_E^2) S_E(v, v),$$

where  $C > 0$  is a constant.

For the jump condition, we will essentially use the proof given in [5, Lemma 5.2]. In particular, we write  $j_{\mathcal{T}_*}(e; v) = j_{\mathcal{T}}(e; v) + (j_{\mathcal{T}_*}(e; v) - j_{\mathcal{T}}(e; v))$  and for any  $\epsilon > 0$

$$\sum_{j=1}^2 \sum_{e \in \mathcal{E}_{E_i}} h_{E_i} \|j_{\mathcal{T}_*}(e; v)\|_{0,e}^2 \leq (1 + \epsilon) T_1 + \left(1 + \frac{1}{\epsilon}\right) T_2,$$

with  $T_1 := \sum_{i=1}^2 \sum_{e \in \mathcal{E}_{E_i}} h_{E_i} \|j_{\mathcal{T}}(e; v)\|_{0,e}^2$  and  $T_2 := \sum_{i=1}^2 \sum_{e \in \mathcal{E}_{E_i}} h_{E_i} \|j_{\mathcal{T}_*}(e; v) - j_{\mathcal{T}}(e; v)\|_{0,e}^2$ . On the new edge we notice that  $j_{\mathcal{T}}(e; v) = 0$ , then,

$$T_1 \leq \frac{1}{\sqrt{2}} \sum_{e \in \mathcal{E}_E} h_E \|j_{\mathcal{T}}(e; v)\|_{0,e}^2.$$

We now define  $\mathcal{T}_*(E_i) := \{E' \in \mathcal{T}_* : \mathcal{E}_{E_i} \cap \mathcal{E}_{E'} \neq \emptyset\}$ ; for any edge  $e \in \mathcal{E}_{E_i}$ , we denote by  $E_{i,e} \in \mathcal{T}_*(E_i)$  the element such that  $e = \partial E_i \cap \partial E_{i,e}$ . Then,

$$\begin{aligned} \|j_{\mathcal{T}_*}(e; v) - j_{\mathcal{T}}(e; v)\|_{0,e} &= \|[A(\Pi_{\mathcal{T}_*}^0 - \Pi_{\mathcal{T}}^0)\nabla v]\|_{0,e} \\ &\leq \|A_E(\Pi_{k-1,E_i}^0 - \Pi_{k-1,E}^0)\nabla v\|_{0,e} + \|A_{\hat{E}_{i,e}}(\Pi_{k-1,E_{i,e}}^0 - \Pi_{k-1,\hat{E}_{i,e}}^0)\nabla v\|_{0,e}, \end{aligned}$$

where  $\hat{E}_{i,e}$  indicates the parent of  $E_{i,e}$ . Using the trace inequality we have

$$\begin{aligned} T_2 &\lesssim \sum_{i=1}^2 \sum_{E' \in \mathcal{T}_*(E_i)} \|(\Pi_{k-1,E'}^0 - \Pi_{k-1,\hat{E}'}^0)\nabla v\|_{0,E'}^2 \\ &\lesssim \sum_{i=1}^2 \sum_{E' \in \mathcal{T}_*(E_i)} (\|\nabla v - \Pi_{k-1,E'}^0 \nabla v\|_{0,E'}^2 + \|\nabla v - \Pi_{k-1,\hat{E}'}^0 \nabla v\|_{0,E'}^2) \end{aligned}$$

Using now the minimality property of  $\Pi_{k-1,E'}^0$  and  $\Pi_{k-1,\hat{E}'}^0$ , we easily get as above

$$T_2 \leq \sum_{E' \in \mathcal{T}(E)} \|\nabla(v - \mathcal{I}_{E'} v)\|_{0,E'}^2 \lesssim \sum_{E' \in \mathcal{T}(E)} S_{E'}(v, v),$$

which, for a sufficiently small  $\epsilon$ , concludes the proof.  $\square$

From this Lemma and the Lipschitz continuity of the residual estimator with respect to the argument  $v$  (whose proof is independent of the used polynomial degree, so we refer to [5, Lemma 5.3]), we immediately deduce the following result.

**Proposition 8.2** (residual estimator reduction on refined elements). *There exist constants  $\mu_r \in (0, 1)$ ,  $c_{er,1} > 0$  and  $c_{er,2} > 0$  independent of  $\mathcal{T}$  such that for any  $v \in \mathbb{V}_{\mathcal{T}}$  and  $w \in \mathbb{V}_{\mathcal{T}_*}$ , and any element  $E \in \mathcal{T}$  which is split into two children  $E_1, E_2 \in \mathcal{T}_*$ , one has*

$$\eta_{\mathcal{T}_*}(E; w, \mathcal{D}) \leq \mu_r \eta_{\mathcal{T}}(E; v, \mathcal{D}) + c_{er,1} S_{\mathcal{T}(E)}^{1/2}(v, v) + c_{er,2} |v - w|_{1,\mathcal{T}(E)}. \quad (8.1)$$

### 8.1.2. The virtual inconsistency estimator

Given  $v \in \mathbb{V}_{\mathcal{T}}$  and  $E \in \mathcal{T}$ , consider the two virtual inconsistency terms  $\Psi_{\mathcal{T},A}(E, v, \mathcal{D})$  and  $\Psi_{\mathcal{T},c}(E, v, \mathcal{D})$  introduced in (6.3). When  $E$  is bisected into  $E_1$  and  $E_2$ , the term  $\Psi_{\mathcal{T},c}(E, v, \mathcal{D})$  is reduced by a factor  $\mu_c < 1$  up to an addend proportional to the stabilization term, i.e., there exists  $c_{vi,c} > 0$  such that

$$\left( \sum_{i=1}^2 \Psi_{\mathcal{T},c}(E_i, v, \mathcal{D}) \right)^{1/2} \leq \mu_c \Psi_{\mathcal{T},c}(E, v, \mathcal{D}) + c_{vi,c} S_E(v, v)^{1/2}. \quad (8.2)$$

This stems from the presence of the factor  $h_E$  in front of the norm  $\|(I - \Pi_k^0)(c_E \Pi_k^0 v)\|_{0,E}$ , with an argument similar to the one used in the proof of Lemma 8.1.

Due to the lack of the factor  $h_E$ , a reduction result similar to (8.2) does not hold for  $\Psi_{\mathcal{T},A}(E, v, \mathcal{D})$ . Indeed, since  $A_E \Pi_{k-1,E}^0 \nabla v \in \mathbb{P}_{2k-2}(E)$ , one may ask whether a constant  $\mu < 1$  exists such that

$$\sum_{i=1}^2 \|(I - \Pi_{k-1,E_i}^0)q\|_{0,E_i}^2 \leq \mu^2 \|(I - \Pi_{k-1,E}^0)q\|_{0,E}^2 \quad \forall q \in \mathbb{P}_{2k-2}(E). \quad (8.3)$$



Unfortunately, the answer is no, as it can be seen numerically, working on the reference element  $\hat{E}$  by affinity and identifying  $\mu^2$  as the largest eigenvalue of a generalized eigenvalue problem. However, the same numerics indicates that if  $\hat{E}$  is split into  $2^m$  triangles of equal area by  $m$  successive levels of uniform bisections, then  $\mu^2$  becomes  $< 1$  for  $m$  large enough, as seen in Table 1.

**Table 1.** Value of  $\mu^2$  in (8.3) for different values of the polynomial degree  $k$  and the level of refinement  $m$ .

	$m = 1$	$m = 2$
$k = 2$	1.0000	0.3153
$k = 3$	1.0000	0.6648

This is indeed predicted by the following result.

**Lemma 8.3.** *Let  $E \in \mathcal{T}$ . For any polynomial degree  $k \geq 1$  there exists a minimal  $m \in \mathbb{N}$  and a constant  $\mu = \mu_m < 1$  independent of  $E$  such that, if  $E$  is partitioned into  $2^m$  elements  $E_i$  of equal area by  $m$  levels of uniform newest vertex bisection, it holds*

$$\sum_{i=1}^{2^m} \|(I - \Pi_{k-1, E_i}^0)q\|_{0, E_i}^2 \leq \mu^2 \|(I - \Pi_{k-1, E}^0)q\|_{0, E}^2 \quad \forall q \in \mathbb{P}_{2k-2}(E). \quad (8.4)$$

*Proof.* Since by construction  $h_{E_i} = 2^{-m/2}h_E$ , classical approximation results give

$$\sum_{i=1}^{2^m} \|(I - \Pi_{k-1, E_i}^0)q\|_{0, E_i}^2 \leq C_k 2^{-m} h_E^2 |q|_{1, E}^2$$

for some constant  $C_k$  depending on  $k$ . Replacing  $q$  by  $q - \Pi_{k-1, E}^0 q$  leaves the left-hand side unchanged, whereas on the right-hand side an inverse inequality yields

$$\sum_{i=1}^{2^m} \|(I - \Pi_{k-1, E_i}^0)q\|_{0, E_i}^2 \leq C_k C_{\text{inv}, k} 2^{-m} \|q - \Pi_{k-1, E}^0 q\|_{0, E}^2.$$

One concludes taking as  $m$  the smallest integer such that  $\mu_m^2 := C_k C_{\text{inv}, k} 2^{-m} < 1$ .  $\square$

Based on these results, let  $\mathcal{T}_*^m$  be a refinement of  $\mathcal{T}$  in which the element  $E$  has undergone  $m$  levels of uniform refinements by newest vertex bisection, and has been replaced by  $2^m$  subelements  $E_i$ . Given  $v \in \mathbb{V}_{\mathcal{T}}$ , let us set

$$\Psi_{\mathcal{T}_*^m, A}^2(E; v, \mathcal{D}) = \sum_{i=1}^{2^m} \|(I - \Pi_{E_i, k-1}^0)(A_E \Pi_{E_i, k-1}^0 \nabla v)\|_{0, E_i}^2.$$

**Lemma 8.4.** *There exist constants  $\rho_A < 1$  and  $c_{vi, A} > 0$  such that for any  $v \in \mathbb{V}_{E, k}$*

$$\Psi_{\mathcal{T}_*^m, A}(E; v, \mathcal{D}) \leq \rho_A \Psi_{\mathcal{T}, A}(E; v, \mathcal{D}) + c_{vi, A} S_E^{1/2}(v, v).$$

*Proof.* Write

$$\begin{aligned} \|(I - \Pi_{E_i, k-1}^0)(A_E \Pi_{E_i, k-1}^0 \nabla v)\|_{0, E_i} &\leq \|(I - \Pi_{E_i, k-1}^0)(A_E \Pi_{E, k-1}^0 \nabla v)\|_{0, E_i} \\ &\quad + \|A_E(\Pi_{E_i, k-1}^0 \nabla v - \Pi_{E, k-1}^0 \nabla v)\|_{0, E_i}, \end{aligned}$$

sum over  $i$ , and conclude using (8.4) and the usual arguments based on the minimality of the  $L^2$ -orthogonal projections.  $\square$

Let us set

$$\Psi_{\mathcal{T}_*^m}^2(E, v, \mathcal{D}) := \Psi_{\mathcal{T}_*^m, A}^2(E; v, \mathcal{D}) + \Psi_{\mathcal{T}_*^m, c}^2(E; v, \mathcal{D})$$

with

$$\Psi_{\mathcal{T}_*^m, c}^2(E; v, \mathcal{D}) = \sum_{i=1}^{2^m} h_{E_i}^2 \|(I - \Pi_{E_i, k}^0)(c_E \Pi_{E_i, k}^0 v)\|_{0, E_i}^2.$$

Applying a bound similar to (8.2) to the successive level of refinements, we arrive at the following result.

**Lemma 8.5.** *There exist constants  $\mu_{vi} < 1$  and  $c_{vi,1} > 0$  such that for any  $v \in \mathbb{V}_{E,k}$*

$$\Psi_{\mathcal{T}_*^m}(E; v, \mathcal{D}) \leq \mu_{vi} \Psi_{\mathcal{T}}(E; v, \mathcal{D}) + c_{vi,1} S_E^{1/2}(v, v).$$

Combining this estimate with the Lipschitz continuity property of the virtual inconsistency estimator, we obtain the following result.

**Proposition 8.6** (virtual inconsistency estimator reduction on refined elements). *There exist constants  $\mu_{vi} \in (0, 1)$ ,  $c_{vi,1} > 0$  and  $c_{vi,2} > 0$  independent of  $\mathcal{T}$  such that for any  $v \in \mathbb{V}_{\mathcal{T}}$  and  $w \in \mathbb{V}_{\mathcal{T}_*^m}$ , and any element  $E \in \mathcal{T}$  which is split into  $2^m$  children  $E_i \in \mathbb{V}_{\mathcal{T}_*^m}$ , one has*

$$\Psi_{\mathcal{T}_*^m}(E; w, \mathcal{D}) \leq \mu_{vi} \Psi_{\mathcal{T}}(E; v, \mathcal{D}) + c_{vi,1} S_E^{1/2}(v, v) + c_{vi,2} |v - w|_{1, E}. \quad (8.5)$$

## 8.2. Quasi-orthogonality property

Let  $u_{\mathcal{T}_*} \in \mathbb{V}_{\mathcal{T}_*}$  be the solution of Problem (3.4) on the refined mesh  $\mathcal{T}_*$ . Hereafter we establish relations between the two energy errors  $\|u - u_{\mathcal{T}}\|$  and  $\|u - u_{\mathcal{T}_*}\|$ . The first result follows from Proposition 5.5 and Lemma 3.1; the proof is independent of the used polynomial degree, so we refer to [5, Proposition 5.7].

**Proposition 8.7** (comparison of the energy error under refinement). *For any  $\delta \in (0, 1]$  there exists a constant  $C_E > 0$  independent of  $\mathcal{T}$  and  $\delta$  such that*

$$\|u - u_{\mathcal{T}_*}\|^2 \leq (1 + \delta) \|u - u_{\mathcal{T}}\|^2 - \|u_{\mathcal{T}_*} - u_{\mathcal{T}}\|^2 + C_E \left(1 + \frac{1}{\delta}\right) (S_{\mathcal{T}}(u_{\mathcal{T}}, u_{\mathcal{T}}) + S_{\mathcal{T}_*}(u_{\mathcal{T}_*}, u_{\mathcal{T}_*})).$$

Next result extends Corollary 5.8 in [5].

**Proposition 8.8** (quasi-orthogonality of energy errors without stabilization). *Given any  $\delta \in (0, \frac{1}{4}]$ , there exists  $\gamma_{\delta} > 0$  such that for any  $\gamma > \gamma_{\delta}$ , it holds*

$$\|u - u_{\mathcal{T}_*}\|^2 \leq (1 + 4\delta) \|u - u_{\mathcal{T}}\|^2 - \|u_{\mathcal{T}_*} - u_{\mathcal{T}}\|^2 + 2\delta (\Psi_{\mathcal{T}}^2(u_{\mathcal{T}}, \mathcal{D}) + \Psi_{\mathcal{T}_*}^2(u_{\mathcal{T}_*}, \mathcal{D})).$$

*Proof.* Let  $e := \|u - u_{\mathcal{T}}\|$ ,  $e_* := \|u - u_{\mathcal{T}_*}\|$ ,  $S := S_{\mathcal{T}}(u_{\mathcal{T}}, u_{\mathcal{T}})$ ,  $S_* := S_{\mathcal{T}_*}(u_{\mathcal{T}_*}, u_{\mathcal{T}_*})$ ,  $\eta := \eta_{\mathcal{T}}(u_{\mathcal{T}}, \mathcal{D})$ ,  $\Psi := \Psi_{\mathcal{T}}(u_{\mathcal{T}}, \mathcal{D})$ ,  $\Psi_* := \Psi_{\mathcal{T}_*}(u_{\mathcal{T}_*}, \mathcal{D})$  and  $E := \|u_{\mathcal{T}} - u_{\mathcal{T}_*}\|$ . From Corollary 7.4 and (2.3), we get  $\eta^2 \leq \frac{S}{c_{apost}} + \frac{e^2}{c_{apost} c_{\mathcal{B}}}$ , while, from Proposition 7.5,  $S \leq \frac{C_B}{\gamma^2} (\eta^2 + \Psi^2)$ . Combining them, we have

$$\left(1 - \frac{C_B}{\gamma^2 c_{apost}}\right) S \leq \frac{C_B}{\gamma^2} \left(\frac{e^2}{c_{apost} c_{\mathcal{B}}} + \Psi^2\right).$$

Doing the same on  $\mathcal{T}_*$  and defining

$$\bar{C} := \left(1 - \frac{C_B}{c_{apost}}\right)^{-1} C_B \max \left\{1, \frac{1}{c_{apost} c_B}\right\} \leq \left(1 - \frac{C_B}{\gamma^2 c_{apost}}\right)^{-1} C_B \max \left\{1, \frac{1}{c_{apost} c_B}\right\}$$

provided  $\gamma^2 \geq 1$ , we get

$$S \leq \frac{\bar{C}}{\gamma^2} (e^2 + \Psi^2), \quad S_* \leq \frac{\bar{C}}{\gamma^2} (e_*^2 + \Psi_*^2).$$

Employing Proposition 8.7, we obtain

$$e_*^2 \leq (1 + \delta)e^2 - E^2 + C_E \left(1 + \frac{1}{\delta}\right) \frac{\bar{C}}{\gamma^2} (e^2 + e_*^2 + \Psi^2 + \Psi_*^2).$$

If we define  $D := C_E \left(1 + \frac{1}{\delta}\right) \bar{C}$ ,

$$\left(1 - \frac{D}{\gamma^2}\right) e_*^2 \leq \left(1 + \delta + \frac{D}{\gamma^2}\right) e^2 - E^2 + \frac{D}{\gamma^2} (\Psi^2 + \Psi_*^2).$$

By choosing  $\gamma$  such that

$$\frac{1}{\gamma^2} \leq \frac{\delta}{D}, \tag{8.6}$$

we get

$$(1 - \delta)e_*^2 \leq (1 + 2\delta)e^2 - E^2 + \delta(\Psi^2 + \Psi_*^2),$$

which concludes the proof by observing that  $\frac{1+2\delta}{1-\delta} \leq 1 + 4\delta$  and  $\frac{\delta}{1-\delta} \leq 2\delta$ , when  $\delta \leq \frac{1}{4}$ .  $\square$

## 9. The module GALERKIN

Let us consider a  $\Lambda$ -admissible input mesh  $\mathcal{T}_0$ , a set of approximated data  $\mathcal{D}$  which consist of piecewise polynomials of degree  $k - 1$  on  $\mathcal{T}_0$ , and a tolerance  $\epsilon > 0$ . The call

$$[\mathcal{T}, u_{\mathcal{T}}] = \text{GALERKIN}(\mathcal{T}_0, \mathcal{D}, \epsilon)$$

produces a  $\Lambda$ -admissible refined mesh  $\mathcal{T}$  and the Galerkin approximation  $u_{\mathcal{T}} \in \mathbb{V}_{\mathcal{T}}$ , such as

$$\|u - u_{\mathcal{T}}\| \leq C_G \epsilon,$$

where  $u$  is the solution of Problem (2.2) and  $C_G = \sqrt{c^{\mathcal{B}} \max\{C_{U_1}, C_{U_2}\}}$ , with  $c^{\mathcal{B}}$  is defined in (2.3) and  $C_{U_1}, C_{U_2}$  in Corollary 7.6. We obtain it by iterating the sequence

$$\text{SOLVE} \rightarrow \text{ESTIMATE} \rightarrow \text{MARK} \rightarrow \text{REFINE}.$$

At each step, a  $\Lambda$ -admissible mesh  $\mathcal{T}_j$  and the associated solution  $u_j$  of the discrete Problem (3.4) are produced. The process stops when the condition  $\eta_{\mathcal{T}_j}^2(u_j, \mathcal{D}) + \Psi_{\mathcal{T}_j}^2(u_j, \mathcal{D}) \leq \epsilon^2$  is reached.

In particular, the modules are defined as follows:

- $[u_{\mathcal{T}}] = \text{SOLVE}(\mathcal{T}, \mathcal{D})$  produces the solution of Problem (3.4) with data  $\mathcal{D}$ ;
- $[\{\eta_{\mathcal{T}}(\cdot; u_{\mathcal{T}}, \mathcal{D})\}, \{\Psi_{\mathcal{T}}(\cdot; u_{\mathcal{T}}, \mathcal{D})\}] = \text{ESTIMATE}(\mathcal{T}, u_{\mathcal{T}})$  computes the local estimators on  $\mathcal{T}$ ;
- $[\mathcal{M}] = \text{MARK}(\mathcal{T}, \{\eta_{\mathcal{T}}(\cdot; u_{\mathcal{T}}, \mathcal{D})\}, \{\Psi_{\mathcal{T}}(\cdot; u_{\mathcal{T}}, \mathcal{D})\}, \theta)$  implements the Dörfler criterion [15] and finds an almost minimal set  $\mathcal{M}$  of elements in  $\mathcal{T}$  such that

$$\theta \left( \eta_{\mathcal{T}}^2(u_{\mathcal{T}}, \mathcal{D}) + \Psi_{\mathcal{T}}^2(u_{\mathcal{T}}, \mathcal{D}) \right) \leq \sum_{E \in \mathcal{M}} \left( \eta_{\mathcal{T}}^2(E; u_{\mathcal{T}}, \mathcal{D}) + \Psi_{\mathcal{T}}^2(E; u_{\mathcal{T}}, \mathcal{D}) \right), \quad (9.1)$$

for a given parameter  $\theta \in (0, 1)$ ;

- $[\mathcal{T}_*] = \text{REFINE}(\mathcal{T}, \mathcal{M}, \Lambda)$  returns a  $\Lambda$ -admissible refined mesh obtained from  $\mathcal{T}$  by suitable newest-vertex bisections of the elements in  $\mathcal{M}$ , and possibly of other elements to fulfil the  $\Lambda$ -admissibility condition.

It is worth adding some details about the procedure REFINE. Let  $E \in \mathcal{M}$  be an element marked for refinement. For simplicity, hereafter let us set  $\eta := \eta_{\mathcal{T}}(E; u_{\mathcal{T}}, \mathcal{D})$  and  $\Psi := \Psi_{\mathcal{T}}(E; u_{\mathcal{T}}, \mathcal{D})$ . The refinement of  $E$  is performed as follows:

- If  $\eta \geq \Psi$ , then  $E$  is bisected once;
- If  $\eta < \Psi$ , then  $E$  is bisected  $m$ -times, where  $m$  has been introduced in Section 8.1.2 (see Lemma 8.3).

Denote by  $\mathcal{P}(E)$  the partition of  $E$  so obtained, and set  $\eta_*^2 := \sum_{E' \in \mathcal{P}(E)} \eta_{\mathcal{P}(E)}^2(E'; u_{\mathcal{T}}, \mathcal{D})$  and  $\Psi_*^2 := \sum_{E' \in \mathcal{P}(E)} \Psi_{\mathcal{P}(E)}^2(E'; u_{\mathcal{T}}, \mathcal{D})$ . Then, recalling Lemmas 8.1 and 8.5, one gets when  $\eta \geq \Psi$

$$\eta_* + \Psi_* \leq \frac{\mu_r + 1}{2}(\eta + \Psi) + c S_{\mathcal{T}(E)}^{1/2}(u_{\mathcal{T}}, u_{\mathcal{T}}).$$

Indeed,  $\Psi$  can be written as  $\Psi = \lambda\eta$  for a certain  $\lambda \in [0, 1]$  and

$$\begin{aligned} \eta_* + \Psi_* &\leq \mu_r \eta + \lambda \eta + c S_{\mathcal{T}(E)}^{1/2}(u_{\mathcal{T}}, u_{\mathcal{T}}) \\ &= \frac{\mu_r + \lambda}{1 + \lambda} (1 + \lambda) \eta + c S_{\mathcal{T}(E)}^{1/2}(u_{\mathcal{T}}, u_{\mathcal{T}}) \\ &= \frac{(\mu_r + \lambda)}{1 + \lambda} (\eta + \Psi) + c S_{\mathcal{T}(E)}^{1/2}(u_{\mathcal{T}}, u_{\mathcal{T}}). \end{aligned}$$

In the case  $\eta < \Psi$ ,

$$\eta_* + \Psi_* \leq \max(\mu_r^m, \mu_{vi})(\eta + \Psi) + c S_{\mathcal{T}(E)}^{1/2}(u_{\mathcal{T}}, u_{\mathcal{T}}).$$

In all cases, it holds

$$\eta_* + \Psi_* \leq \max\left(\frac{\mu_r + 1}{2}, \mu_{vi}\right)(\eta + \Psi) + c S_{\mathcal{T}(E)}^{1/2}(u_{\mathcal{T}}, u_{\mathcal{T}}), \quad (9.2)$$

which shows that in each marked element the sum of the two estimators is reduced under refinement, up to the stabilization term. Note that for values  $k = 2$  or  $3$  of the polynomial degree of practical use, two bisections ( $m = 2$ ) are enough when  $\eta < \Psi$ .

This refinement may create non-admissible hanging nodes, i.e., hanging nodes with global index larger than  $\Lambda$ . To remove them and guarantee  $\Lambda$ -admissibility of  $\mathcal{T}_*$ , further refinements should be applied. For the realization of this technical part, we refer to Section 11.1 in [6].

The following section proves the convergence of the GALERKIN algorithm.

## 10. Convergence property of GALERKIN

**Proposition 10.1** (global estimators reduction). *Let  $u_{\mathcal{T}} \in \mathbb{V}_{\mathcal{T}}$  be the solution of the discrete variational Problem (3.4). There exist constants  $\rho \in (0, 1)$  and  $C_{ger,1}, C_{ger,2} > 0$  independent of  $\mathcal{T}$  such that, if  $\mathcal{T}_*$  is the refinement of  $\mathcal{T}$  obtained by applying REFINE, one has for any  $w \in \mathbb{V}_{\mathcal{T}_*}$*

$$\begin{aligned} & \eta_{\mathcal{T}_*}^2(w, \mathcal{D}) + \Psi_{\mathcal{T}_*}^2(w, \mathcal{D}) \\ & \leq \rho \left( \eta_{\mathcal{T}}^2(u_{\mathcal{T}}, \mathcal{D}) + \Psi_{\mathcal{T}}^2(u_{\mathcal{T}}, \mathcal{D}) \right) + C_{ger,1} S_{\mathcal{T}}(u_{\mathcal{T}}, u_{\mathcal{T}}) + C_{ger,2} \|u_{\mathcal{T}} - w\|_{1,\Omega}^2. \end{aligned} \quad (10.1)$$

*Proof.* One can reach the conclusion e.g., as in [5, proof of Proposition 5.5], using the bound (9.2) in each element  $E$  marked for refinement.  $\square$

**Theorem 10.2** (contraction property of GALERKIN). *Let  $\mathcal{M} \subset \mathcal{T}$  be the set of the marked elements relative to the solution  $u_{\mathcal{T}} \in \mathbb{V}_{\mathcal{T}}$  of the discrete variational Problem (3.4). If  $\mathcal{T}_*$  is the refinement of  $\mathcal{T}$  obtained by applying REFINE, then for  $\gamma$  sufficiently large there exist  $\alpha \in (0, 1)$  and  $\beta > 0$ ,  $\zeta > 0$  such that*

$$\| \|u - u_{\mathcal{T}_*}\|^2 + \beta \eta_{\mathcal{T}_*}^2(u_{\mathcal{T}_*}, \mathcal{D}) + \zeta \Psi_{\mathcal{T}_*}^2(u_{\mathcal{T}_*}, \mathcal{D}) \leq \alpha \left( \|u - u_{\mathcal{T}}\|^2 + \beta \eta_{\mathcal{T}}^2(u_{\mathcal{T}}, \mathcal{D}) + \zeta \Psi_{\mathcal{T}}^2(u_{\mathcal{T}}, \mathcal{D}) \right).$$

*Proof.* To simplify notation, we set again  $e = \|u - u_{\mathcal{T}}\|$ ,  $e_* = \|u - u_{\mathcal{T}_*}\|$ ,  $S = S_{\mathcal{T}}(u_{\mathcal{T}}, u_{\mathcal{T}})$ ,  $S_* = S_{\mathcal{T}_*}(u_{\mathcal{T}_*}, u_{\mathcal{T}_*})$ ,  $\eta = \eta_{\mathcal{T}}(u_{\mathcal{T}}, \mathcal{D})$ ,  $\eta_* = \eta_{\mathcal{T}_*}(u_{\mathcal{T}_*}, \mathcal{D})$ ,  $\Psi = \Psi_{\mathcal{T}}(u_{\mathcal{T}}, \mathcal{D})$ ,  $\Psi_* = \Psi_{\mathcal{T}_*}(u_{\mathcal{T}_*}, \mathcal{D})$  and  $E = \| \|u_{\mathcal{T}} - u_{\mathcal{T}_*}\| \|$ . From Proposition 8.8,

$$e_*^2 \leq (1 + 4\delta)e^2 - E^2 + 2\delta(\Psi + \Psi_*),$$

whereas using Proposition 10.1 and Proposition 7.5, we get

$$\eta_*^2 + \Psi_*^2 \leq \rho(\eta^2 + \Psi^2) + C_{ger,1} S + \frac{C_{ger,2}}{c_{\mathcal{B}}} E^2 \leq \left( \rho + \frac{C_{ger,1} C_{\mathcal{B}}}{\gamma^2} \right) (\eta^2 + \Psi^2) + \frac{C_{ger,2}}{c_{\mathcal{B}}} E^2.$$

Combining them, we get

$$\begin{aligned} e_*^2 + \beta \eta_*^2 + (\beta - 2\delta) \Psi_*^2 & \leq (1 + 4\delta)e^2 + \left( \frac{\beta C_{ger,2}}{c_{\mathcal{B}}} - 1 \right) E^2 \\ & + \beta \left( \rho + \frac{C_{ger,1} C_{\mathcal{B}}}{\gamma^2} \right) \eta^2 + \beta \left( \rho + \frac{C_{ger,1} C_{\mathcal{B}}}{\gamma^2} + \frac{2\delta}{\beta} \right) \Psi^2, \end{aligned}$$

which suggests choosing  $\beta$  such that

$$\frac{\beta C_{ger,2}}{c_{\mathcal{B}}} = 1. \quad (10.2)$$

Next, we write

$$\begin{aligned} e_*^2 + \beta \eta_*^2 + (\beta - 2\delta) \Psi_*^2 & \leq (1 - \delta)e^2 + 5\delta e^2 \\ & + \beta \left( \rho + \frac{C_{ger,1} C_{\mathcal{B}}}{\gamma^2} \right) \eta^2 + \beta \left( \rho + \frac{C_{ger,1} C_{\mathcal{B}}}{\gamma^2} + \frac{2\delta}{\beta} \right) \Psi^2, \end{aligned}$$

and we invoke Corollary 7.6 to write

$$e^2 \leq c^{\mathcal{B}} C_{apost} \left( 1 + \frac{C_{\mathcal{B}}}{\gamma^2} \right) \eta^2 + c^{\mathcal{B}} C_{apost} \frac{C_{\mathcal{B}}}{\gamma^2} \Psi^2,$$

which gives

$$e_*^2 + \beta \eta_*^2 + (\beta - 2\delta) \Psi_*^2 \leq (1 - \delta)e^2 + \beta \left( \rho + \frac{C_{ger,1} C_B}{\gamma^2} + \frac{5\delta}{\beta} c^{\mathcal{B}} C_{apost} \left( 1 + \frac{C_B}{\gamma^2} \right) \right) \eta^2 \\ + \beta \left( \rho + \frac{C_{ger,1} C_B}{\gamma^2} + \frac{2\delta}{\beta} + \frac{5\delta}{\beta} c^{\mathcal{B}} C_{apost} \frac{C_B}{\gamma^2} \right) \Psi^2.$$

We now choose  $\gamma$  and  $\delta$  such that

$$\rho + \frac{C_{ger,1} C_B}{\gamma^2} + \frac{5\delta}{\beta} c^{\mathcal{B}} C_{apost} \left( 1 + \frac{C_B}{\gamma^2} \right) \leq \frac{1 + \rho}{2}$$

which holds true if

$$\frac{C_{ger,1} C_B}{\gamma^2} \leq \frac{1 - \rho}{4} \quad \text{and} \quad \frac{5\delta}{\beta} c^{\mathcal{B}} C_{apost} (1 + C_B) \leq \frac{1 - \rho}{4} \quad (10.3)$$

(recall that we already assumed  $\gamma^2 \geq 1$ ). Similarly, we choose  $\gamma$  and  $\delta$  such that

$$\beta \left( \rho + \frac{C_{ger,1} C_B}{\gamma^2} + \frac{2\delta}{\beta} + \frac{5\delta}{\beta} c^{\mathcal{B}} C_{apost} \frac{C_B}{\gamma^2} \right) \leq (\beta - 2\delta) \frac{1 + \rho}{2},$$

which holds true if  $\gamma$  satisfies the first condition in (10.3), whereas  $\delta$  satisfies

$$\left( 2 + 5c^{\mathcal{B}} C_{apost} C_B + \frac{1 + \rho}{\beta} \right) \delta \leq \frac{1 - \rho}{4}. \quad (10.4)$$

This proves the result, if we define  $\zeta := \beta - 2\delta$ , with  $\beta$  defined by (10.2) and  $\delta < \frac{\beta}{2}$ , and

$$\alpha := \min \left( 1 - \delta, \frac{1 + \rho}{2} \right). \quad (10.5)$$

The conditions on  $\gamma$  and  $\delta$  which lead to the desired estimate are given in (8.6), (10.3) and (10.4).  $\square$

## 11. Numerical experiments

The aim of this numerical test is to confirm the convergence of our GALERKIN algorithm. We consider a classical test with an  $L$ -shaped domain  $\Omega = (-1, 1)^2 \setminus (-1, 0)^2$  and the reaction-diffusion problem (2.1), with polynomial coefficients of order one for the case  $k = 2$ , i.e.,

$$\mathbf{A} = \begin{bmatrix} 2 + y & 0 \\ 0 & 2 + x \end{bmatrix}, \quad c = x + y + 3,$$

and polynomials of order two for the case  $k = 3$ ,

$$\mathbf{A} = \begin{bmatrix} 2 + y^2 & 0 \\ 0 & 2 + x^2 \end{bmatrix}, \quad c = x^2 + y^2.$$

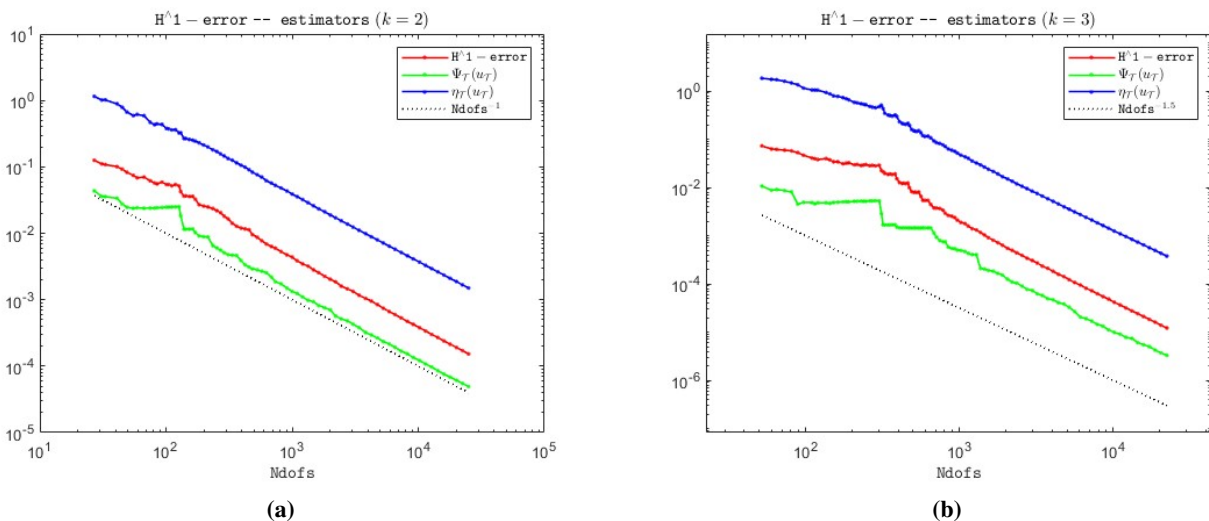
The forcing term  $f$  and the Dirichlet boundary conditions are chosen so that the solution of the problem results

$$u_{\text{ex}}(r, \beta) = r^{2/3} \sin\left(\frac{2}{3}\left(\beta + \frac{\pi}{2}\right)\right), \quad (11.1)$$

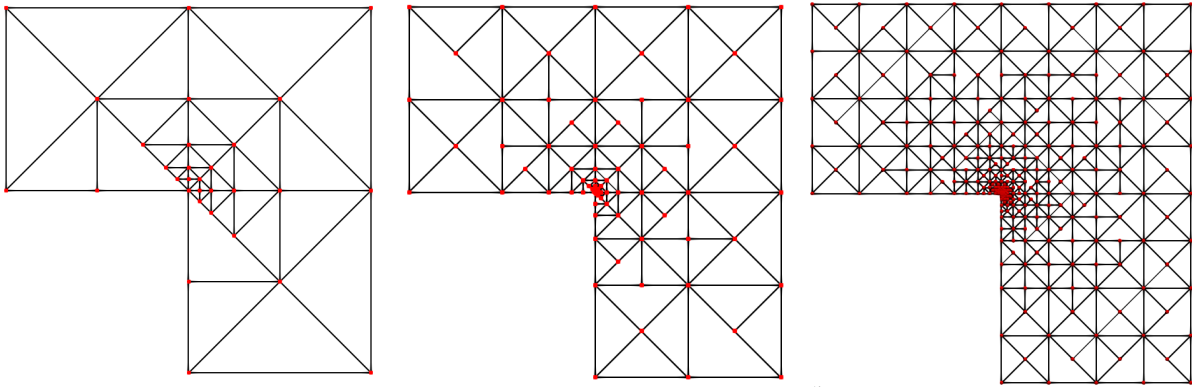
where  $r$  and  $\beta$  are the polar coordinates. It is possible to prove that there exists a  $p$  with  $p \in \left(\frac{2}{k+1}, \frac{2}{k+2/3}\right)$  such that  $u_{\text{ex}} \in W_p^k(\Omega)$  when  $p \geq 1$ , and  $u_{\text{ex}} \in L_p^k(\Omega)$  when  $p \in (0, 1)$ , where  $W_p^k(\Omega)$  and  $L_p^k(\Omega)$  indicate respectively Sobolev and Lipschitz spaces. Then, according to the theory of approximation classes [16, 17], we expect the maximal rate of convergence, i.e.,  $\text{Ndofs}^{-k/2}$ , where  $\text{Ndofs}$  is the number of the degrees of freedom. We apply the adaptive algorithm as described in Section 9 and for the marking strategy (9.1) we consider  $\theta = 0.5$ . In order to compute the VEM error, we consider the computable quantity:

$$\text{H}^1 - \text{error} := \frac{\left(\sum_{E \in \mathcal{T}} \|\nabla(u_{\text{ex}} - \Pi_E^\nabla u_{\mathcal{T}})\|_{0,E}^2\right)^{1/2}}{\|\nabla u_{\text{ex}}\|_{0,\Omega}}.$$

In Figure 6, we represent the evolution of  $\text{H}^1 - \text{error}$  and the estimator terms  $\eta_{\mathcal{T}}(u_{\mathcal{T}}, \mathcal{D})$  and  $\Psi_{\mathcal{T}}(u_{\mathcal{T}}, \mathcal{D})$ , which confirms the results of Corollary 7.6. Furthermore, we notice that after a transient phase, the error and the estimator terms decays reach asymptotically the theoretical optimal rate  $\text{Ndofs}^{-1.0}$  (for the case  $k = 2$ ) and  $\text{Ndofs}^{-1.5}$  (for the case  $k = 3$ ). In Figure 7, we depict the meshes after 20, 35 and 50 loops of the adaptive algorithm in the case  $k = 2$ . We highlight the presence of hanging nodes in the different meshes loops.



**Figure 6.**  $\text{H}^1 - \text{error}$  (red), the residual type term  $\eta_{\mathcal{T}}(u_{\mathcal{T}}, \mathcal{D})$  (blue), the inconsistency term  $\Psi_{\mathcal{T}}(u_{\mathcal{T}}, \mathcal{D})$  (green), and the expected optimal decay (dashed) in the case  $k = 2$  (a) and  $k = 3$  (b).



**Figure 7.** The partition of the domain after 20 loops (first), 35 loops (second), and 50 loops (third) of the adaptive algorithm of order  $k = 2$ .

## 12. Conclusions

In this paper, we presented an adaptive VEM of order  $k \geq 2$  on nonconforming triangular meshes. In the analysis, the space  $\mathbb{V}_{\mathcal{T}}^0$  of continuous, piecewise polynomials functions of degree  $k$  on the triangulation  $\mathcal{T}$  plays a fundamental role. Indeed, it is contained in the global VEM space,  $\mathbb{V}_{\mathcal{T}}^0 \subseteq \mathbb{V}_{\mathcal{T}}$ , and guarantees a quasi-orthogonality property for any refinement  $\mathcal{T}_*$  of  $\mathcal{T}$ , since  $\mathbb{V}_{\mathcal{T}}^0 \subseteq \mathbb{V}_{\mathcal{T}_*}^0$ . By pivoting on this space, we proved an a posteriori error estimate which does not contain the stabilization term appearing in the VEM discrete formulation. Consequently, we established the convergence of the adaptive VEM algorithm, by a contraction argument.

Extensions of our work include:

- The complexity and optimality analysis of the two step algorithm AVEM mentioned in the Introduction to account for non-polynomial data;
- The study of a variant of the adaptive algorithm in which the polynomial degree  $k$  may take large values, in the spirit of a  $p$ -version;
- The treatment of more general polygonal meshes which, as remarked in [5], seems non-trivial. The main difficulties lay in the choice of a suitable refinement strategy in replacement of the newest-vertex bisection used here, and in the lack of a conforming space  $\mathbb{V}_{\mathcal{T}}^0$  for general polygonal meshes.

## Use of AI tools declaration

The authors declare they have not used Artificial Intelligence (AI) tools in the creation of this article.

## Acknowledgments

The authors performed this research in the framework of the Italian MIUR Award “Dipartimenti di Eccellenza 2018-2022” granted to the Department of Mathematical Sciences, Politecnico di Torino (CUP: E11G18000350001). CC was partially supported by the Italian MIUR through the PRIN grant 201752HKH8; DF thanks the INdAM-GNCS project “Metodi numerici per lo studio di strutture



geometriche parametriche complesse” (CUP: E53C22001930001). The authors are members of the Italian INdAM-GNCS research group.

### Conflict of interest

The authors declare no conflicts of interest.

### References

1. B. Ahmad, A. Alsaedi, F. Brezzi, L. D. Marini, A. Russo, Equivalent projectors for virtual element methods, *Comput. Math. Appl.*, **66** (2013), 376–391. <http://dx.doi.org/10.1016/j.camwa.2013.05.015>
2. P. F. Antonietti, F. Dassi, E. Manuzzi, Machine learning based refinement strategies for polyhedra, *J. Comput. Phys.*, **469** (2022), 111531. <http://dx.doi.org/10.1016/j.jcp.2022.111531>
3. L. Beirão da Veiga, F. Brezzi, A. Cangiani, G. Manzini, L. D. Marini, A. Russo, Basic principles of virtual element methods, *Math. Mod. Meth. Appl. Sci.*, **23** (2013), 199–2014. <http://dx.doi.org/10.1142/S0218202512500492>
4. L. Beirão da Veiga, F. Brezzi, L. D. Marini, A. Russo, The hitchhiker’s guide to the virtual element method, *Math. Mod. Meth. Appl. Sci.*, **24** (2014), 1541–1573. <http://dx.doi.org/10.1142/S021820251440003X>
5. L. Beirão da Veiga, C. Canuto, R. H. Nochetto, G. Vacca, M. Verani, Adaptive VEM: stabilization-free a posteriori error analysis and contraction property, *SIAM J. Numer. Anal.*, **61** (2023), 457–494. <http://dx.doi.org/10.1137/21M1458740>
6. L. Beirão da Veiga, C. Canuto, R. H. Nochetto, G. Vacca, M. Verani, Adaptive VEM for variable data: convergence and optimality, *IMA J. Numer. Anal.*, 2023. <http://dx.doi.org/10.1093/imanum/drad085>
7. L. Beirão da Veiga, C. Lovandina, A. Russo, Stability analysis for the virtual element method, *Math. Mod. Meth. Appl. Sci.*, **27** (2017), 2557–2594. <http://dx.doi.org/10.1142/S021820251750052X>
8. L. Beirão da Veiga, G. Manzini, Residual a posteriori error estimation for the virtual element method for elliptic problems, *ESAIM: M2AN*, **49** (2015), 577–599. <http://dx.doi.org/10.1051/m2an/2014047>
9. S. Berrone, A. Borio, A. D’Auria, Refinement strategies for polygonal meshes applied to adaptive VEM discretization, *Finite Elem. Anal. Des.*, **186** (2021), 103502. <http://dx.doi.org/10.1016/j.finel.2020.103502>
10. S. Berrone, A. D’Auria, A new quality preserving polygonal mesh refinement algorithm for polygonal element methods, *Finite Elem. Anal. Des.*, **207** (2022), 103770. <http://dx.doi.org/10.1016/j.finel.2022.103770>
11. P. Binev, W. Dahmen, R. DeVore, Adaptive finite element methods with convergence rates, *Numer. Math.*, **97** (2004), 219–268. <http://dx.doi.org/10.1007/s00211-003-0492-7>

12. A. Cangiani, E. H. Georgoulis, T. Pryer, O. J. Sutton, A posteriori error estimates for the virtual element method, *Numer. Math.*, **137** (2017), 857–893. <http://dx.doi.org/10.1007/s00211-017-0891-9>
13. C. Carstensen, M. Feischl, M. Page, D. Praetorius, Axioms of adaptivity, *Comput. Math. Appl.*, **67** (2014), 1195–1253. <http://dx.doi.org/10.1016/j.camwa.2013.12.003>
14. J. M. Cascon, C. Kreuzer, R. H. Nochetto, K. G. Siebert, Quasi-optimal convergence rate for an adaptive finite element method, *SIAM J. Numer. Anal.*, **46** (2008), 2524–2550. <http://dx.doi.org/10.1137/07069047X>
15. W. Dörfler, A convergent adaptive algorithm for Poisson’s equation, *SIAM J. Numer. Anal.*, **33** (1996), 1106–1124. <http://dx.doi.org/10.1137/0733054>
16. F. D. Gaspoz, P. Morin, Approximation classes for adaptive higher order finite element approximation, *Math. Comp.*, **83** (2014), 2127–2160. <http://dx.doi.org/10.1090/s0025-5718-2013-02777-9>
17. F. D. Gaspoz, P. Morin, Errata to “Approximation classes for adaptive higher order finite element approximation”, *Math. Comp.*, **86** (2017), 1525–1526. <http://dx.doi.org/10.1090/mcom/3243>
18. R. H. Nochetto, A. Veerer, Primer of adaptive finite element methods, In: S. Bertoluzza, R. H. Nochetto, A. Quarteroni, K. G. Siebert, A. Veerer, *Multiscale and adaptivity: modeling, numerics and applications*, Lecture Notes in Mathematics, Berlin, Heidelberg: Springer, **2040** (2012), 125–225. <http://dx.doi.org/10.1007/978-3-642-24079-9>



AIMS Press

©2023 the Author(s), licensee AIMS Press. This is an open access article distributed under the terms of the Creative Commons Attribution License (<http://creativecommons.org/licenses/by/4.0>)