

# Error estimates on the approximate finite volume solution of convection diffusion equations with general boundary conditions

Thierry Gallouët<sup>1</sup>, Raphaèle Herbin<sup>2</sup> and Marie Hélène Vignal<sup>3</sup>

**Summary** We study here the convergence of a finite volume scheme for a diffusion-convection equation on an open bounded set of  $\mathbb{R}^d$  ( $d = 2$  or  $3$ ) for which we consider Dirichlet, Neumann or Robin boundary conditions. We consider unstructured meshes which include Voronoï or triangular meshes; we use for the diffusion term an “ $s$  points” (where  $s$  is the number of sides of each cell) finite volume scheme and for the convection term an upstream finite volume scheme. Assuming the exact solution at least in  $H^2$  we prove error estimates in a discrete  $H_0^1$  norm of order the size of the mesh. Discrete Poincaré inequalities then allow to prove error estimates in the  $L^2$  norm.

## 1 Presentation of the problem

Let  $\Omega$  be an open bounded subset of  $\mathbb{R}^d$  ( $d = 2$  or  $3$ ) which is assumed to be polygonal if  $d = 2$  and polyhedral if  $d = 3$ . We denote by  $\partial\Omega$  its boundary and by  $\mathbf{n}$  the unit normal to  $\partial\Omega$  outward to  $\Omega$ .

We consider the following convection diffusion reaction problem:

$$-\Delta u(x) + \operatorname{div}(\mathbf{v}(x) u(x)) + b(x)u(x) = f(x), \quad x \in \Omega, \quad (1)$$

with different boundary conditions and the following hypotheses

**Assumption 1**  $f \in L^2(\Omega)$ ,  $b \in L^\infty(\Omega)$  and  $\mathbf{v} \in C^1(\overline{\Omega}, \mathbb{R}^d)$  such that  $\operatorname{div}\mathbf{v}/2 + b \geq 0$  almost everywhere.

In this paper, we consider three different types of boundary conditions for the previous diffusion convection equation, namely Dirichlet, Neumann or Robin boundary conditions; these conditions are not necessarily homogeneous. This elliptic problem is then discretized with a finite volume scheme: an “ $s$ -points” scheme, where  $s$  is the number of sides of each cell, is used for the diffusion term and an upstream scheme for the convection term.

Let us remark that the analysis which is developed here still holds for equations of the type

$$-\operatorname{div}(k(x)\nabla u(x)) + \operatorname{div}(\mathbf{v}(x) u(x)) + b(x)u(x) = f(x), \quad x \in \Omega, \quad (2)$$

under Assumption 1 with the following hypothesis on  $k$ :

**Assumption 2**  $k$  is a piecewise  $C^1$  function from  $\overline{\Omega}$  to  $\mathbb{R}$  such that there exists  $k_0 \in \mathbb{R}_+^*$  such that  $k(x) \geq k_0$  for a.e.  $x \in \Omega$ .

For the sake of the simplicity of notations we prefer to deal with the Laplace operator here but we shall point out the modifications which take place if the operator  $\operatorname{div}(k\nabla \cdot)$  is considered instead: see remarks 1, 4 and 7 in the case of the Dirichlet boundary conditions. Let us now assume that  $k$  is a tensor satisfying the following hypothesis:

**Assumption 3**  $k$  is a piecewise  $C^1$  function from  $\overline{\Omega}$  to  $\mathbb{R}^{d \times d}$  such that for all  $x \in \overline{\Omega}$ ,  $k(x)$  is a symmetric matrix and such that there exists  $k_0 \in \mathbb{R}_+^*$  such that  $k(x)\xi \cdot \xi \geq k_0$  for a.e.  $x \in \Omega$  and for all  $\xi \in \mathbb{R}^d$ .

---

<sup>1</sup>Université de Provence, CMI, 39 rue Joliot Curie, 13453 Marseille Cedex 13; email gallouet@gyptis.univ-mrs.fr

<sup>2</sup>Université de Provence, CMI, 39 rue Joliot Curie, 13453 Marseille Cedex 13; email herbin@gyptis.univ-mrs.fr

<sup>3</sup>Université P. Sabatier, UFR MIG, MIP, 31062 Toulouse cedex 4; email vignal@mip.ups-tlse.fr

Then one may still write the finite volume scheme and obtain some error estimates, but the assumptions on the mesh have to be modified see Remark 1 and [10], see also the modified scheme of Coudière *et al* [5] for this case. However if the mesh is cartesian and if for all  $x \in \overline{\Omega}$  the matrix  $k(x)$  is diagonal then it is “aligned” with the grid and the analysis is similar to the (non constant) scalar case of Equation (2).

Finite volumes are known to be well adapted to the discretization of conservation equations, particularly in the presence of convection terms. Their theoretical study has recently been undertaken. Two main directions are usually followed in order to obtain convergence properties of finite volume schemes. The first one consists in writing the finite volume as a finite element or mixed finite element method by using some numerical integration, see for instance [1], [2], [18], [19] or [20]; the convergence then follows from the general finite element framework. The second one, see for example [5], [6], [12], [9], [21], [14] or [22], consists in establishing the convergence by using the direct formulation of the finite volume scheme together with some appropriate discrete functional analysis tools. This last approach is considered here. A discrete system is obtained for each type of boundary condition. Existence and uniqueness (sometimes up to a constant like in the continuous case) of the approximate solution is proven. The stability of the scheme is shown in each case by establishing some estimates on the approximate solution which are independent of the mesh size. If the exact solution is assumed to be at least in  $H^2(\Omega)$ , one may then establish the convergence of the scheme by proving error estimates. A first one in a discrete  $H_0^1$  norm is obtained and a second one in  $L^2$  norm follows with the help of discrete Poincaré inequalities. It is also possible to prove error estimates in the  $L^q$  norm, see [4], for all  $q$  such that  $1 \leq q < +\infty$  if  $d = 2$  and such that  $1 \leq q \leq 6$  if  $d = 3$  establishing discrete Sobolev’s imbeddings.

This work is divided in four sections. The first one introduces the admissible meshes which are needed for the discretization of the elliptic problem, and the three following sections correspond to the three types of boundary conditions which we consider here. Homogeneous Dirichlet conditions were studied in e.g. [12], [21], [9], [14], with different assumptions on the data and the mesh; to our knowledge, nonhomogeneous Dirichlet, Neumann and Robin boundary conditions have only been considered up to now in [6] with some simplifying assumptions; the convergence of the method for Neumann and Robin conditions requires some additional work compared to that of the Dirichlet case. In the case of Neumann boundary conditions, a “discrete Poincaré-Wirtinger” inequality needs to be proven in order to obtain an  $L^2$  error estimate. The stability results for both Neumann and Robin boundary conditions are obtained by using a discrete trace inequality which we prove to be true for piecewise constant functions. In the case of the Robin condition, it is interesting to note that an artificial upwinding has to be introduced in the treatment of the boundary condition in order for the scheme to be well defined with no additional condition on the mesh.

## 2 Admissible meshes

**Definition 1 (Admissible meshes)** *A finite volume mesh of  $\Omega$ , denoted by  $\mathcal{T}$ , is given by a family of “control volumes”, which are open polygonal (or polyhedral) convex subsets of  $\Omega$  (with positive measure), a family of subsets of  $\overline{\Omega}$  contained in hyperplanes of  $\mathbb{R}^d$ , denoted by  $\mathcal{E}$  (these are the edges (if  $d = 2$ ) or sides (if  $d = 3$ ) of the control volumes), with strictly positive  $(d - 1)$ -dimensional measure, and a family of points of  $\overline{\Omega}$  denoted by  $\mathcal{P}$ . The finite volume mesh is said to be admissible if the properties (i) to (iv) below are satisfied and restricted admissible if the properties (i) to (v) below are satisfied.*

- (i) *The closure of the union of all the control volumes is  $\overline{\Omega}$ ;*
- (ii) *For any  $K \in \mathcal{T}$ , there exists a subset  $\mathcal{E}_K$  of  $\mathcal{E}$  such that  $\partial K = \overline{K} \setminus K = \cup_{\sigma \in \mathcal{E}_K} \overline{\sigma}$ . Let  $\mathcal{E} = \cup_{K \in \mathcal{T}} \mathcal{E}_K$ .*
- (iii) *For any  $(K, L) \in \mathcal{T}^2$  with  $K \neq L$ , either the  $(d - 1)$ -dimensional Lebesgue measure of  $\overline{K} \cap \overline{L}$  is 0 or  $\overline{K} \cap \overline{L} = \overline{\sigma}$  for some  $\sigma \in \mathcal{E}$ , which will then be denoted by  $K|L$ .*
- (iv) *The family  $\mathcal{P} = (x_K)_{K \in \mathcal{T}}$  is such that  $x_K \in \overline{K}$  (for all  $K \in \mathcal{T}$ ) and, if  $\sigma = K|L$ , it is assumed that  $x_K \neq x_L$ , and that the straight line  $D_{K,L}$  going through  $x_K$  and  $x_L$  is orthogonal to  $K|L$ .*

- (v) For any  $\sigma \in \mathcal{E}$  such that  $\sigma \subset \partial\Omega$ , let  $K$  be the control volume such that  $\sigma \in \mathcal{E}_K$ . If  $x_K \notin \sigma$ , let  $\mathcal{D}_{K,\sigma}$  be the straight line going through  $x_K$  and orthogonal to  $\sigma$ , then the condition  $\mathcal{D}_{K,\sigma} \cap \sigma \neq \emptyset$  is assumed; let  $y_\sigma = \mathcal{D}_{K,\sigma} \cap \sigma$ .

In the sequel, the following notations are used. The mesh size is defined by:  $\text{size}(\mathcal{T}) = \sup\{\text{diam}(K), K \in \mathcal{T}\}$ , where  $\text{diam}(K)$  is the diameter of  $K \in \mathcal{T}$ . For any  $K \in \mathcal{T}$  and  $\sigma \in \mathcal{E}$ ,  $m(K)$  is the  $d$ -dimensional Lebesgue measure of  $K$  (i.e. area if  $d = 2$ , volume if  $d = 3$ ),  $m(\sigma)$  the  $(d - 1)$ -dimensional measure of  $\sigma$ , and  $\mathbf{n}_{K,\sigma}$  denotes the unit normal vector to  $\sigma$  outward to  $K$ . The set of interior (resp. boundary) edges is denoted by  $\mathcal{E}_{\text{int}}$  (resp.  $\mathcal{E}_{\text{ext}}$ ), that is  $\mathcal{E}_{\text{int}} = \{\sigma \in \mathcal{E}; \sigma \not\subset \partial\Omega\}$  (resp.  $\mathcal{E}_{\text{ext}} = \{\sigma \in \mathcal{E}; \sigma \subset \partial\Omega\}$ ). The set of neighbours of  $K$  is denoted by  $\mathcal{N}(K)$ , that is  $\mathcal{N}(K) = \{L \in \mathcal{T}; \exists \sigma \in \mathcal{E}_K, \bar{\sigma} = \overline{K} \cap \overline{L}\}$ . If  $\sigma = K|L$ , we denote by  $d_\sigma$  or  $d_{K|L}$  the Euclidean distance between  $x_K$  and  $x_L$  (which is positive) and by  $d_{K,\sigma}$  the distance from  $x_K$  to  $\sigma$ . If  $\sigma \in \mathcal{E}_K \cap \mathcal{E}_{\text{ext}}$ , let  $d_\sigma$  denote the Euclidean distance between  $x_K$  and  $y_\sigma$  (then,  $d_\sigma = d_{K,\sigma}$ ). For any  $\sigma \in \mathcal{E}$ ; the “transmissibility” through  $\sigma$  is defined by  $\tau_\sigma = m(\sigma)/d_\sigma$  if  $d_\sigma \neq 0$  and  $\tau_\sigma = 0$  if  $d_\sigma = 0$ . In some results and proofs given below, there are summations over  $\sigma \in \mathcal{E}_0$ , with  $\mathcal{E}_0 = \{\sigma \in \mathcal{E}; d_\sigma \neq 0\}$ . For simplicity, (in these results and proofs)  $\mathcal{E} = \mathcal{E}_0$  is assumed.

Admissible (or restricted admissible) meshes include, for instance, meshes made with triangles and rectangles in two space dimensions, and also Voronoï meshes: the latter consists in building a mesh using the orthogonal bisectors from a given family of points (for more details see [7]). Admissible meshes will be used for the Neumann boundary conditions. Property (v) of the restricted admissible meshes is needed for the Dirichlet and Robin boundary conditions.

**Remark 1** In the case of the operator  $\operatorname{div}(k\nabla \cdot)$  which is considered in Equation (2) where  $k$  is a function from  $\overline{\Omega}$  to  $\mathbb{R}$  or  $\mathbb{R}^{d \times d}$  which satisfies Assumption 2 or 3, admissible meshes must satisfy the following additional condition:

- (vi) For any  $K \in \mathcal{T}$ , the restriction  $k|_K$  of the function  $k$  to any given control volume  $K$  belongs to  $C^1(\overline{K})$ .

Furthermore if  $k$  is a piecewise  $C^1$  function from  $\overline{\Omega}$  to  $\mathbb{R}^{d \times d}$ , the orthogonality conditions (iv) and (v) have to be modified into:

- (iv)' For any  $K \in \mathcal{T}$ , let  $k_K$  denote the mean value of  $k$  on  $K$ , that is

$$k_K = \frac{1}{m(K)} \int_K k(x) dx.$$

The set  $\mathcal{T}$  is such that there exists a family of points

$$\mathcal{P} = (x_K)_{K \in \mathcal{T}} \text{ such that } x_K = \cap_{\sigma \in \mathcal{E}_K} \mathcal{D}_{K,\sigma,k} \in \overline{K},$$

where  $\mathcal{D}_{K,\sigma,k}$  is a straigth line perpendicular to  $\sigma$  with respect to the scalar product induced by  $k_K^{-1}$  such that  $\mathcal{D}_{K,\sigma,k} \cap \sigma = \mathcal{D}_{L,\sigma,k} \cap \sigma \neq \emptyset$  if  $\sigma = K|L$ . Furthermore, if  $\sigma = K|L$ , let  $y_\sigma = \mathcal{D}_{K,\sigma,k} \cap \sigma (= \mathcal{D}_{L,\sigma,k} \cap \sigma)$  and assume that  $x_K \neq x_L$ .

- (v)' For any  $\sigma \in \mathcal{E}_{\text{ext}}$ , let  $K$  be the control volume such that  $\sigma \in \mathcal{E}_K$  and let  $\mathcal{D}_{K,\sigma,k}$  be the straight line going through  $x_K$  and orthogonal to  $\sigma$  with respect to the scalar product induced by  $k_K^{-1}$ ; then, there exists  $y_\sigma \in \sigma \cap \mathcal{D}_{K,\sigma,k}$ ; let  $g_\sigma = g(y_\sigma)$ .

### 3 Dirichlet boundary conditions

The first type of boundary condition which we consider is a Dirichlet condition:

$$u(x) = g^D(x), \quad x \in \partial\Omega, \tag{3}$$

where  $g^D \in H^{1/2}(\partial\Omega)$ .

Let us denote by  $\tilde{g}^D$  a function of  $H^1(\Omega)$  such that  $\bar{\gamma}(\tilde{g}^D) = g^D$ , where  $\bar{\gamma}$  denotes the trace operator from  $H^1(\Omega)$  into  $H^{1/2}(\partial\Omega)$ .

Under Assumption 1, there exists a unique variational solution  $u \in H^1(\Omega)$  of (1), (3) by the Lax-Milgram Theorem. That is to say,  $u$  satisfies  $u = \tilde{u} + \tilde{g}^D$  where  $\tilde{u} \in H_0^1(\Omega)$  is the unique solution to

$$\begin{aligned} \int_{\Omega} \left( \nabla \tilde{u}(x) \cdot \nabla \phi(x) + \operatorname{div}(\mathbf{v}(x) \tilde{u}(x)) \phi(x) + b(x) \tilde{u}(x) \phi(x) \right) dx \\ = \int_{\Omega} \left( -\nabla \tilde{g}^D(x) \cdot \nabla \phi(x) - \operatorname{div}(\mathbf{v}(x) \tilde{g}^D(x)) \phi(x) - b(x) \tilde{g}^D(x) \phi(x) + f(x) \phi(x) \right) dx, \end{aligned}$$

for all  $\phi \in H_0^1(\Omega)$ .

In order to obtain an error estimate, we shall need some more regularity on the boundary condition (however, the definition of the finite volume and its convergence require less regularity, see Remark 2):

**Assumption 4**  $g^D \in H^{3/2}(\partial\Omega)$ .

### 3.1 Discretization

The approximate finite volume solution which is sought here is constant on each cell of the mesh. The discrete unknowns are denoted by  $(u_K)_{K \in \mathcal{T}}$ . The principle of classical finite volume schemes is to integrate the equation on each cell of the mesh in order to obtain an equation which is sometimes called the balance equation, for each control volume.

Let  $K \in \mathcal{T}$ , using Green's formula, one has:

$$\begin{aligned} \int_K \left[ -\Delta u(x) + \operatorname{div}(\mathbf{v}(x) u(x)) + b(x) u(x) \right] dx \\ = \sum_{\sigma \in \mathcal{E}_K} \int_{\sigma} \left[ -\nabla u(x) + \mathbf{v}(x) u(x) \right] \cdot \mathbf{n}_{K,\sigma} d\gamma(x) + \int_K b(x) u(x) dx = \int_K f(x) dx, \end{aligned}$$

where  $d\gamma$  is the integration symbol for the  $(d-1)$ -dimensional Lebesgue measure on the considered hyperplane.

For all  $K \in \mathcal{T}$  and all  $\sigma \in \mathcal{E}_K$ , let us denote by  $F_{K,\sigma}$  the approximate diffusion flux (respectively by  $V_{K,\sigma}$ , the approximate convection flux) that is to say an approximation of  $\int_{\sigma} -\nabla u(x) \cdot \mathbf{n}_{K,\sigma} d\gamma(x)$  (respectively of  $\int_{\sigma} \mathbf{v}(x) \cdot \mathbf{n}_{K,\sigma} u(x) d\gamma(x)$ ).

In order to prove the convergence of the scheme, one needs two basic properties. The first one, called conservativity of the scheme, states that the numerical flux through a given edge is conservative, i.e.:

$$F_{K,\sigma} = -F_{L,\sigma} \quad \text{for all } K \in \mathcal{T}, L \in \mathcal{N}(K) \text{ and where } \sigma = K|L. \quad (4)$$

The second one is that  $\frac{1}{m(\sigma)} F_{K,\sigma}$  is a consistent approximation of  $\frac{1}{m(\sigma)} \int_{\sigma} -\nabla u(x) \cdot \mathbf{n}_{K,\sigma} d\gamma(x)$  (for more details see Lemmas 2 and 3). The same properties are required for  $V_{K,\sigma}$ .

The numerical diffusion flux  $F_{K,\sigma}$  is chosen as:

$$F_{K,\sigma} = -m(K|L) \frac{u_L - u_K}{d_{K|L}} \quad \text{if } \sigma = K|L, \quad (5)$$

and

$$F_{K,\sigma} d_{K,\sigma} = -m(\sigma) (u_{\sigma} - u_K) \quad \text{if } \sigma \in \mathcal{E}_{\text{ext}} \cap \mathcal{E}_K. \quad (6)$$

The numerical convective flux  $V_{K,\sigma}$  is obtained with an upstream scheme, that is:

$$V_{K,\sigma} = v_{K,\sigma} u_{\sigma,+} \quad (7)$$

with

$$v_{K,\sigma} = \int_{\sigma} \mathbf{v}(x) \cdot \mathbf{n}_{K,\sigma} d\gamma(x), \quad (8)$$

and

$$\begin{cases} \text{if } \sigma = K|L, & u_{\sigma,+} = \begin{cases} u_K & \text{if } v_{K,\sigma} \geq 0, \\ u_L & \text{otherwise,} \end{cases} \\ \text{if } \sigma \in \mathcal{E}_{\text{ext}} \cap \mathcal{E}_K, & u_{\sigma,+} = \begin{cases} u_K & \text{if } v_{K,\sigma} \geq 0, \\ u_\sigma & \text{otherwise.} \end{cases} \end{cases} \quad (9)$$

and where we set

$$u_\sigma = g^D(y_\sigma), \quad (10)$$

with  $y_\sigma$  defined in Definition 1.

**Remark 2** If Assumption 4 is weakened to  $g^D \in L^2(\partial\Omega)$ ,  $g^D$  is no longer defined pointwise, but (10) may be replaced by  $u_\sigma = \frac{1}{m(\sigma)} \int_{\sigma} g^D(y) d\gamma(y)$ , where  $d\gamma$  stands for the  $(d-1)$ -dimensional Lebesgue measure on  $\sigma$ . In this latter case we do not obtain an error estimate but only a convergence result as in [7].

For all  $K \in \mathcal{T}$ , let  $f_K$  and  $b_K$  denote the mean value of  $f$  and  $b$  on  $K$  that is to say

$$f_K = \frac{1}{m(K)} \int_K f(x) dx \quad \text{and} \quad b_K = \frac{1}{m(K)} \int_K b(x) dx. \quad (11)$$

Then the considered finite volume scheme is defined by the following equations:

$$\sum_{\sigma \in \mathcal{E}_K} (F_{K,\sigma} + v_{K,\sigma} u_{\sigma,+}) + m(K) b_K u_K = m(K) f_K. \quad (12)$$

**Remark 3** The definition of the diffusion flux on a boundary edge (6) allows  $d_{K,\sigma} = 0$ , in this case one has  $u_K = u_\sigma$  and  $F_{K,\sigma}$  becomes an unknown.

**Remark 4** In the case of a non constant diffusion coefficient as in Equation (2) where  $k$  is a function from  $\Omega$  to  $\mathbb{R}$  satisfying Assumption 2 or from  $\Omega$  to  $\mathbb{R}^{d \times d}$  satisfying Assumption 3, one considers admissible meshes satisfying (vi) of Remark 1 and in the tensor case also (iv)' and (v)' instead of (iv) and (v). For  $K \in \mathcal{T}$  and  $\sigma \in \mathcal{E}_K$ , let

$$k_{K,\sigma} = \left| \frac{1}{m(K)} \int_K k(x) dx \mathbf{n}_{K,\sigma} \right|$$

(where  $|\cdot|$  denotes the Euclidean norm). Note that in the scalar case, this yields in fact  $k_{K,\sigma} = \frac{1}{m(K)} \int_K k(x) dx$ . The exact diffusion fluxes  $k(x) \nabla u \cdot \mathbf{n}_{K,\sigma}$  on an edge  $\sigma$  of the mesh may then be approximated in a consistent way (see [6] and [10]) by replacing the formulae (5) and (6) by:

- internal edges:

$$F_{K,\sigma} = -\tau_\sigma(u_L - u_K), \quad \text{if } \sigma \in \mathcal{E}_{\text{int}}, \quad \sigma = K|L, \quad (13)$$

where

$$\tau_\sigma = m(\sigma) \frac{k_{K,\sigma} k_{L,\sigma}}{k_{K,\sigma} d_{L,\sigma} + k_{L,\sigma} d_{K,\sigma}};$$

- boundary edges:

$$F_{K,\sigma} = -\tau_\sigma(u_\sigma - u_K), \quad \text{if } \sigma \in \mathcal{E}_{\text{ext}} \text{ and } x_K \notin \sigma, \quad (14)$$

where

$$\tau_\sigma = m(\sigma) \frac{k_{K,\sigma}}{d_{K,\sigma}}.$$

### 3.2 Existence, uniqueness and stability of the approximate solution

The proof of existence and uniqueness of the approximate solution may be performed by establishing a discrete maximum principle

**Proposition 1** *Under Assumptions 1 and 4, let  $\mathcal{T}$  be an admissible mesh in the sense of Definition 1. Let  $(u_\sigma)_{\sigma \in \mathcal{E}_{\text{ext}}}$ ,  $(f_K)_{K \in \mathcal{T}}$ ,  $(b_K)_{K \in \mathcal{T}}$  and  $(v_{K,\sigma})_{\sigma \in \mathcal{E}_K, K \in \mathcal{T}}$  be defined by (10), (11) and (8). If  $f_K \geq 0$  for all  $K \in \mathcal{T}$ , and  $u_\sigma \geq 0$ , for all  $\sigma \in \mathcal{E}_{\text{ext}}$ , then if  $(u_K)_{K \in \mathcal{T}}$  is a solution to (12), (5), (6), (7), (9), then  $u_K \geq 0$  for all  $K \in \mathcal{T}$ .*

The proof of this maximum principle may be found in [6] or [9] using a strong formulation and in [7] for a weak formulation. It immediately yields the existence and uniqueness of the solution to the finite volume scheme. We may therefore now define the approximate finite volume solution  $u_{\mathcal{T}}$  by:

$$u_{\mathcal{T}}(x) = u_K \quad \text{for a.e. } x \in K, K \in \mathcal{T}, \quad (15)$$

where  $(u_K)_{K \in \mathcal{T}}$  is the unique solution to (5)-(12).

The following stability estimate on the approximate solution was proven in [7] in the more general case of a semilinear convection diffusion equation.

**Lemma 1** *Under Assumptions 1 and 4, let  $\mathcal{T}$  be an admissible mesh in the sense of Definition 1, and let:*

$$\zeta = \min \left( \min_{K \in \mathcal{T}} \min_{\sigma \in \mathcal{E}_K} \frac{d_{K,\sigma}}{\text{diam}(K)}, \min_{K \in \mathcal{T}} \min_{\sigma \in \mathcal{E}_K} \frac{d_{K,\sigma}}{d_\sigma} \right), \quad (16)$$

*Let  $u_{\mathcal{T}}$  be defined by  $u_{\mathcal{T}}(x) = u_K$  for a.e.  $x \in K$ , and for any  $K \in \mathcal{T}$  where  $(u_K)_{K \in \mathcal{T}}$  is the solution to (5)-(12).*

*Then there exists  $C \in \mathbb{R}_+$ , only depending on  $\Omega$ ,  $g^D$ ,  $\zeta$ ,  $M = \max_{K \in \mathcal{T}} \text{card}(\mathcal{E}_K)$ ,  $b$  and  $f$ , such that:*

$$\|u_{\mathcal{T}}\|_{1,\mathcal{T}} \leq C \quad \text{and} \quad \|u_{\mathcal{T}}\|_{L^2(\Omega)} \leq C. \quad (17)$$

### 3.3 Error estimates

In this section, one proves the convergence of the approximate solution  $u_{\mathcal{T}}$  towards the exact solution  $u$  to (1), (3) assuming  $u \in C^2(\bar{\Omega})$  or  $u \in H^2(\Omega)$  and an additional assumption on the mesh. To do this, we establish error estimates in a discrete  $H_0^1$  norm. Some similar results are also in [5], [6], [9] and [12]. Let us now define the discrete  $H_0^1$  norm of a piecewise constant function from  $\Omega$  to  $\mathbb{R}$ .

**Definition 2 (Discrete  $H_0^1$  norm)** *Let  $\mathcal{T}$  be an admissible finite volume mesh in the sense of Definition 1. For  $u$  which is constant on each control volume of  $\mathcal{T}$ , that is to say  $u(x) = u_K$  for a.e.  $x \in K$ ,  $K \in \mathcal{T}$ , one defines the discrete  $H_0^1$  norm by*

$$\|u\|_{1,\mathcal{T}} = \left( \sum_{\sigma \in \mathcal{E}} \tau_\sigma (D_\sigma u)^2 \right)^{1/2},$$

where  $D_\sigma u = |u_K - u_L|$  if  $\sigma \in \mathcal{E}_{\text{int}}$ ,  $\sigma = K|L$ ,  $D_\sigma u = |u_K|$  if  $\sigma \in \mathcal{E}_{\text{ext}} \cap \mathcal{E}_K$ ,  $\tau_\sigma$  and the sets  $\mathcal{E}$ ,  $\mathcal{E}_{\text{int}}$ ,  $\mathcal{E}_{\text{ext}}$  and  $\mathcal{E}_K$  are defined in Definition 1.

**Remark 5** *Let  $\sigma \in \mathcal{E}_{\text{int}}$  and assume  $\sigma = K|L$  with  $K \in \mathcal{T}$  and  $L \in \mathcal{N}(K)$ . One can see the difference quotient  $D_\sigma u/d_\sigma$  as a discrete normal gradient of  $u$  on  $\sigma$  and therefore also on the diamond shaped dual cell defined by the vertices of  $\sigma$  and  $x_K$  and  $x_L$ , note that the measure of this dual cell is  $d_\sigma m(\sigma)/d$ , where  $d = 2$  or  $3$  is the space dimension. Let  $\sigma \in \mathcal{E}_{\text{ext}}$  and assume  $\sigma \in \mathcal{E}_K$  with  $K \in \mathcal{T}$ . Again in this case, assuming  $u = 0$  on  $\partial\Omega$ , the difference quotient  $D_\sigma u/d_\sigma$  can be seen as a discrete gradient of  $u$  on  $\sigma$  and so on a dual cell defined by the vertices of  $\sigma$  and  $x_K$  of measure  $d_\sigma m(\sigma)/2$ .*

Hence  $\|u\|_{1,\mathcal{T}} = \left( \sum_{\sigma \in \mathcal{E}} d_\sigma m(\sigma) (D_\sigma u/d_\sigma)^2 \right)^{1/2}$  can be seen as a discrete  $H_0^1$  norm of  $u$ .

Let us now state the error estimates under some regularity assumption on the solution.

**Theorem 1 ( $C^2$  regularity)** *Under Assumptions 1 and 4, let  $\mathcal{T}$  be a restricted admissible mesh in the sense of Definition 1 and  $u_{\mathcal{T}}$  be the solution to (5)-(12) and (15). Assume that the unique variational solution  $u$  of Problem (1), (3) satisfies  $u \in C^2(\overline{\Omega})$ . Let  $e_{\mathcal{T}}$  be defined by  $e_{\mathcal{T}}(x) = e_K = u(x_K) - u_K$  if  $x \in K$ ,  $K \in \mathcal{T}$ .*

*Then, there exists  $C > 0$ , only depending on  $u$ ,  $\mathbf{v}$ ,  $b$ ,  $d$  and  $\Omega$ , such that*

$$\|e_{\mathcal{T}}\|_{1,\mathcal{T}} \leq C \text{size}(\mathcal{T}), \quad (18)$$

where  $\|\cdot\|_{1,\mathcal{T}}$  is the discrete  $H_0^1$  norm defined in Definition 2. Furthermore:

$$\|e_{\mathcal{T}}\|_{L^2(\Omega)} \leq C \text{size}(\mathcal{T}). \quad (19)$$

**Theorem 2 ( $H^2$  regularity)** *Under Assumptions 1 and 4, let  $\mathcal{T}$  be a restricted admissible mesh in the sense of Definition 1 and let*

$$\zeta = \min_{K \in \mathcal{T}} \min_{\sigma \in \mathcal{E}_K} \frac{d_{K,\sigma}}{\text{diam}(K)}.$$

*Let  $u_{\mathcal{T}}$  be defined by (5)-(12) and (15). Assume that the unique variational solution  $u$  to (1) and (3) belongs to  $H^2(\Omega)$ . Let  $e_{\mathcal{T}}$  be defined by  $e_{\mathcal{T}}(x) = e_K = u(x_K) - u_K$  if  $x \in K$ ,  $K \in \mathcal{T}$ .*

*Then, there exists  $C$ , only depending on  $u$ ,  $\mathbf{v}$ ,  $b$ ,  $\Omega$ ,  $d$  and  $\zeta$ , such that (18) and (19) hold.*

## Remark 6

1. Inequality (18) (resp. (19)) yields an estimate of order 1 for the discrete  $H_0^1$  norm (resp.  $L^2$  norm) of the error on the solution. Note also that, using  $u \in C^1(\overline{\Omega})$  (or  $u \in H^1(\Omega)$ ), one deduces, from (19), the existence of  $C$ , only depending on  $u$ ,  $b$ ,  $\mathbf{v}$ ,  $d$  and  $\Omega$  (or  $u$ ,  $b$ ,  $\mathbf{v}$ ,  $d$ ,  $\Omega$  and  $\zeta$ ), such that  $\|u - u_{\mathcal{T}}\|_{L^2(\Omega)} \leq C \text{size}(\mathcal{T})$ .
2. In Theorem 2, the function  $e_{\mathcal{T}}$  is still well defined and so is the quantity “ $\nabla u \cdot \mathbf{n}_{\sigma}$ ” on  $\sigma$ , for all  $\sigma \in \mathcal{E}$ . Indeed, since  $u \in H^2(\Omega)$  and  $d \leq 3$ , one has  $u \in C(\overline{\Omega})$ , then  $u(x_K)$  is well defined for all control volume  $K \in \mathcal{T}$  and  $\nabla u \cdot \mathbf{n}_{\sigma}$  belongs to  $L^2(\sigma)$  (for the  $(d-1)$ -Lebesgue measure on  $\sigma$ ) for all  $\sigma \in \mathcal{E}$ .
3. Note that, under Assumptions 1 and 4 with  $b = \mathbf{v} = g^D = 0$  the (unique) variational solution of (1), (3) is necessarily in  $H^2(\Omega)$  provided that  $\Omega$  is convex.
4. Thanks to (iv) of Definition 1,  $\zeta$ , in Theorem 2, is always defined. The important fact is that the constant  $C$  in the error estimates (18) and (19) depends on  $\zeta$ .
5. In Theorem 2 if one considers rectangular meshes for which the points  $x_K$  are located at the centers of the cells, then  $\zeta$  depends on the ratio of the length and the width of rectangles; if one considers meshes which are made out of triangles for which the points  $x_K$  are located at the intersection of the orthogonal bisectors of the triangles, then all angles of triangles must be lower than  $\pi/2 + \eta$  where  $\eta$  is a positive constant and  $\zeta$  depends only on  $\eta$ .
6. In the case of the pure diffusion operator, one may approximate the diffusive fluxes through the mesh interfaces up to the second order with the same difference quotient by using appropriate meshes such as rectangles or Voronoï meshes. Indeed if  $y_{\sigma}$  is the mid-point of  $[x_K, x_L]$  for any  $\sigma = K|L \in \mathcal{E}_{\text{int}}$ , the proof of Lemma 3 may be adapted to show that the consistency error on the flux  $R_{K,\sigma}$  defined in (24) below is of order 2 with respect to the mesh size  $\text{size}(\mathcal{T})$  instead of 1 as in (26) below. This is in good agreement with numerical results which are presented in the recent paper [3] for a related co-volume scheme. In the case of the diffusion-convection operator however the consistency error of the upwind approximation of the convection flux remains of order one for any mesh. Hence the order one estimate which we obtain here seems to be sharp for the upwind scheme. Order one is also obtained in the numerical results of [3] for the convective case.

### Proof of Theorems 1 and 2

One proceeds in two steps. In the first step, one proves that the approximation of the fluxes is consistent. In the second step, one establishes error estimates using this property and the conservativity of the scheme (see (4)).

#### Step 1 (Consistency)

Let  $K$  be a control volume and  $\sigma \in \mathcal{E}_K$ . We define the exact diffusion flux  $\overline{F}_{K,\sigma}$  and the exact convection flux  $\overline{V}_{K,\sigma}$  by:

$$\overline{F}_{K,\sigma} = - \int_{\sigma} \nabla u(x) \cdot \mathbf{n}_{K,\sigma} d\gamma(x) \quad \text{and} \quad \overline{V}_{K,\sigma} = \int_{\sigma} u(x) \mathbf{v}(x) \cdot \mathbf{n}_{K,\sigma} d\gamma(x). \quad (20)$$

Next for all  $K \in \mathcal{T}$  and all  $\sigma \in \mathcal{E}_K$ , let  $F_{K,\sigma}^*$  and  $V_{K,\sigma}^*$  be defined by:

$$F_{K,\sigma}^* = -m(K|L) \frac{u(x_L) - u(x_K)}{d_{K|L}}, \quad \text{if } \sigma = K|L \in \mathcal{E}_K \cap \mathcal{E}_{\text{int}}, \quad (21)$$

$$d_{K,\sigma} F_{K,\sigma}^* = -m(\sigma) (u(y_\sigma) - u(x_K)), \quad \text{if } \sigma \in \mathcal{E}_{\text{ext}} \cap \mathcal{E}_K, \quad (22)$$

$$V_{K,\sigma}^* = v_{K,\sigma} u(x_{\sigma,+}), \quad (23)$$

where  $x_{\sigma,+} = x_K$  (resp.  $x_L$ ) if  $\sigma \in \mathcal{E}_{\text{int}}$ ,  $\sigma = K|L$  and  $v_{K,\sigma} \geq 0$  (resp.  $v_{K,\sigma} < 0$ ) and  $x_{\sigma,+} = x_K$  (resp.  $y_\sigma$ ) if  $\sigma \in \mathcal{E}_K \cap \mathcal{E}_{\text{ext}}$  and  $v_{K,\sigma} \geq 0$  (resp.  $v_{K,\sigma} < 0$ ). Then, the consistency error on the diffusion and convection fluxes may be defined as

$$R_{K,\sigma} = \frac{1}{m(\sigma)} (\overline{F}_{K,\sigma} - F_{K,\sigma}^*) \quad \text{and} \quad r_{K,\sigma} = \frac{1}{m(\sigma)} (\overline{V}_{K,\sigma} - V_{K,\sigma}^*). \quad (24)$$

Moreover, we define

$$\rho_K = \frac{1}{m(K)} \int_K b(x) (u(x) - u(x_K)) dx. \quad (25)$$

Thanks to the regularity of  $u$  and  $\mathbf{v}$ , a Taylor expansion immediately yields the following lemma which gives the consistency of scheme in a finite volume sense when  $u \in C^2(\overline{\Omega})$ .

**Lemma 2** *Under Assumptions 1 and 4, let  $\mathcal{T}$  be a restricted admissible mesh in the sense of Definition 1. Assume that the unique variational solution  $u$  of Problem (1), (3) satisfies  $u \in C^2(\overline{\Omega})$ . Then there exists  $C > 0$ , only depending on  $u$ ,  $b$  and  $\mathbf{v}$ , such that*

$$|R_{K,\sigma}| + |r_{K,\sigma}| + |\rho_K| \leq C \text{size}(\mathcal{T}),$$

for any  $K \in \mathcal{T}$  and  $\sigma \in \mathcal{E}_K$ , where  $R_{K,\sigma}$ ,  $r_{K,\sigma}$  and  $\rho_K$  are defined by (24) and (25).

We prove a similar lemma when  $u$  only belongs to  $H^2(\Omega)$ .

**Lemma 3** *Under Assumptions 1 and 4, let  $\mathcal{T}$  be a restricted admissible mesh in the sense of Definition 1*

$$\zeta = \min_{K \in \mathcal{T}} \min_{\sigma \in \mathcal{E}_K} \frac{d_{K,\sigma}}{\text{diam}(K)}.$$

For  $K \in \mathcal{T}$ , let  $\mathcal{V}_{K,\sigma} = \{tx_K + (1-t)x, x \in \sigma, t \in [0,1]\}$ . For  $\sigma \in \mathcal{E}_{\text{int}}$ , let  $\mathcal{V}_\sigma = \mathcal{V}_{K,\sigma} \cup \mathcal{V}_{L,\sigma}$  where  $K$  and  $L$  are the control volumes such that  $\sigma = K|L$ . For  $\sigma \in \mathcal{E}_{\text{ext}} \cap \mathcal{E}_K$ , let  $\mathcal{V}_\sigma = \mathcal{V}_{K,\sigma}$ .

Assume that the unique variational solution  $u$  to (1), (3) belongs to  $H^2(\Omega)$ . Then there exists  $C_1$ , only depending on  $d$  and  $\zeta$ , and  $C_2$ , only depending on  $d$ ,  $\mathbf{v}$ ,  $\zeta$  and  $p$  such that for all  $K \in \mathcal{T}$  and all  $\sigma \in \mathcal{E}_K$ ,

$$|R_{K,\sigma}| \leq C_1 \text{size}(\mathcal{T}) (m(\sigma)d_\sigma)^{-1/2} \|u\|_{H^2(\mathcal{V}_\sigma)}, \quad (26)$$

$$|r_{K,\sigma}| \leq C_2 \text{size}(\mathcal{T}) (m(\sigma)d_\sigma)^{-1/p} \|u\|_{W^{1,p}(\mathcal{V}_\sigma)}, \quad (27)$$

and

$$|\rho_K| \leq \|b\|_\infty \operatorname{size}(\mathcal{T}) m(K)^{-1/p} \|u\|_{W^{1,p}(K)}. \quad (28)$$

for all  $p > d$  and such that  $p < +\infty$  if  $d = 2$  and  $p \leq 6$  if  $d = 3$ ,

where, for all  $p$  such that  $1 \leq p < +\infty$ :

$$\|u\|_{W^{1,p}(\mathcal{V}_\sigma)}^p = \|u\|_{L^p(\mathcal{V}_\sigma)}^p + \sum_{i=1}^d \|D_i u\|_{L^p(\mathcal{V}_\sigma)}^p,$$

$D_i$  is the (weak) derivative with respect to the component  $z_i$  of  $z = (z_1, \dots, z_d)^t \in \mathbb{R}^d$ .

### Proof of Lemma 3

First note that thanks to Sobolev's imbeddings, if  $u \in H^2(\Omega)$  then  $u \in W^{1,p}(\Omega)$  for all  $p$  such that  $1 \leq p < +\infty$  if  $d = 2$  and such that  $1 \leq p \leq 6$  if  $d = 3$ . Then (27) and (28) are well defined.

Let  $\sigma \in \mathcal{E}$ . Since  $u \in H^2(\Omega)$ , the restriction of  $u$  to  $\mathcal{V}_\sigma$  belongs to  $H^2(\mathcal{V}_\sigma)$ . The space  $C^2(\overline{\mathcal{V}_\sigma})$  is dense in  $H^2(\mathcal{V}_\sigma)$  (see, for instance, [17], this can be proved quite easily by a regularization technique). Then, using a density argument, one needs only to prove (26), (27) and (28) for  $u \in C^2(\overline{\mathcal{V}_\sigma})$ . Therefore let us first assume that  $u \in C^2(\overline{\mathcal{V}_\sigma})$ . The density argument will be proven for (26) in the sequel. It is straightforward for (27) and (28).

First, one proves (26) if  $\sigma \in \mathcal{E}_{\text{int}}$ . Let  $K$  and  $L$  be the two control volumes such that  $\sigma = K|L$ . It is possible to assume, for simplicity of notations and without loss of generality, that  $\sigma = 0 \times \tilde{\sigma}$ , with some  $\tilde{\sigma} \subset \mathbb{R}^{d-1}$ , and  $x_K = (-d_{K,\sigma}, 0)^t$ ,  $x_L = (d_{L,\sigma}, 0)^t$ .

A Taylor expansion, using  $u \in C^2(\overline{\mathcal{V}_\sigma})$  gives, for a.e. (for the  $(d-1)$ -Lebesgue measure)  $x = (0, \tilde{x})^t \in \sigma$ ,

$$u(x_L) - u(x) = \nabla u(x) \cdot (x_L - x) + \int_0^1 H(u)(tx + (1-t)x_L)(x_L - x) \cdot (x_L - x) t dt,$$

where  $H(u)(z)$  denotes the Hessian matrix of  $u$  at point  $z$ , and

$$u(x_K) - u(x) = \nabla u(x) \cdot (x_K - x) + \int_0^1 H(u)(tx + (1-t)x_K)(x_K - x) \cdot (x_K - x) t dt.$$

Remark that  $x_L - x_K = \mathbf{n}_{K,\sigma} d_\sigma$ ; subtracting one equation off the other and integrating over  $\sigma$  yields  $|R_{K,\sigma}| \leq B_{K,\sigma} + B_{L,\sigma}$ , with, for some  $C_3$  only depending on  $d$ ,

$$B_{K,\sigma} = \frac{C_3}{m(\sigma)d_\sigma} \int_\sigma \int_0^1 |H(u)(tx + (1-t)x_K)| |x_K - x|^2 t dt d\gamma(x),$$

where  $|H(u)(z)|^2 = \sum_{i,j=1}^d |D_i D_j u(z)|^2$ .

The quantity  $B_{L,\sigma}$  is obtained from  $B_{K,\sigma}$  by changing  $K$  in  $L$ . One uses a change of variables in  $B_{K,\sigma}$ . Indeed, one sets  $z = tx + (1-t)x_K$ . Since  $|x_K - x| \leq \operatorname{diam}(K)$  and  $dz = t^{d-1} d_{K,\sigma} dt d\gamma(x)$ , one obtains, using  $z_1 = (t-1) d_{K,\sigma}$ ,  $z = (z_1, \bar{z})^t$  with  $\bar{z} \in \mathbb{R}^{d-1}$ ,

$$B_{K,\sigma} \leq \frac{C_3 (\operatorname{diam}(K))^2}{m(\sigma) d_\sigma} \int_{\mathcal{V}_{K,\sigma}} |H(u)(z)| \frac{(d_{K,\sigma})^{d-2}}{(d_{K,\sigma} (z_1 + d_{K,\sigma}))^{d-2}} dz.$$

This gives with the Cauchy-Schwarz inequality,

$$B_{K,\sigma} \leq \frac{C_3 (d_{K,\sigma})^{d-3} (\operatorname{diam}(K))^2}{m(\sigma) d_\sigma} \left( \int_{\mathcal{V}_{K,\sigma}} |H(u)(z)|^2 dz \right)^{1/2} \left( \int_{\mathcal{V}_{K,\sigma}} \frac{1}{(z_1 + d_{K,\sigma})^{(d-2)2}} dz \right)^{1/2}. \quad (29)$$

For  $d = 2$ , remarking that  $m(\mathcal{V}_{K,\sigma}) = (d_{K,\sigma} m(\sigma))/2$ , (29) gives

$$B_{K,\sigma} \leq \frac{C_3 (\operatorname{diam}(K))^2}{\sqrt{2} (m(\sigma) d_\sigma)^{1/2} (d_\sigma d_{K,\sigma})^{1/2}} \left( \int_{\mathcal{V}_{K,\sigma}} |H(u)(z)|^2 dz \right)^{1/2}.$$

A similar estimate holds on  $B_{L,\sigma}$  by changing  $K$  in  $L$  and  $d_{K,\sigma}$  in  $d_{L,\sigma}$ . Since  $d_{K,\sigma}, d_{L,\sigma} \geq \zeta \text{diam}(K)$  and  $d_\sigma = d_{K,\sigma} + d_{L,\sigma} \geq 2\zeta \text{diam}(K)$ , these estimates on  $B_{K,\sigma}$  and  $B_{L,\sigma}$  yield (26) for some  $C$  only depending on  $d$  and  $\zeta$ .

For  $d = 3$ ,

$$B_{K,\sigma} \leq \frac{C_3 (\text{diam}(K))^2}{(\text{m}(\sigma) d_\sigma^2 d_{K,\sigma})^{1/2}} \left( \int_{\mathcal{V}_{K,\sigma}} |H(u)(z)|^2 dz \right)^{1/2} \leq \frac{C_3 \text{size}(\mathcal{T})}{\sqrt{2} \zeta (\text{m}(\sigma) d_\sigma)^{1/2}} \|H(u)\|_{L^2(\mathcal{V}_{K,\sigma})}.$$

With a similar estimate on  $B_{L,\sigma}$ , this yields (26) for some  $C$  only depending on  $d$  and  $\zeta$ .

Now, one proves (26) if  $\sigma \in \mathcal{E}_{\text{ext}}$ . Let  $K$  be the control volume such that  $\sigma \in \mathcal{E}_K$ . One can assume, without loss of generality, that  $x_K = 0$  and  $\sigma = d_{K,\sigma} \times \tilde{\sigma}$  with  $\tilde{\sigma} \subset \mathbb{R}^{d-1}$ . The above proof gives (see Definition 1 for the definition of  $y_\sigma$ ), with some  $C_4$  only depending on  $d$  and  $\zeta$ ,

$$\left| \frac{u(y_\sigma) - u(x_K)}{d_{K,\sigma}} - \frac{1}{\text{m}(\tilde{\sigma})} \int_{\tilde{\sigma}} \nabla u(x) \cdot \mathbf{n}_{K,\sigma} d\gamma(x) \right|^2 \leq C_4 \frac{(\text{size}(\mathcal{T}))^2}{\text{m}(\sigma) d_\sigma} \int_{\mathcal{V}_{\tilde{\sigma}}} |H(u)(z)|^2 dz, \quad (30)$$

with  $\hat{\sigma} = \{(d_{K,\sigma}/2, \tilde{x}/2), \tilde{x} \in \tilde{\sigma}\}$ , and  $\mathcal{V}_{\hat{\sigma}} = \{ty_\sigma + (1-t)x, x \in \hat{\sigma}, t \in [0, 1]\} \cup \{tx_K + (1-t)x, x \in \hat{\sigma}, t \in [0, 1]\}$ . Note that  $\text{m}(\hat{\sigma}) = \text{m}(\sigma)/2^{d-1}$  and that  $\mathcal{V}_{\hat{\sigma}} \subset \mathcal{V}_\sigma$ .

One has now to compare  $I_\sigma = \frac{1}{\text{m}(\sigma)} \int_\sigma \nabla u(x) \cdot \mathbf{n}_{K,\sigma} d\gamma(x)$  with  $I_{\hat{\sigma}} = \frac{1}{\text{m}(\hat{\sigma})} \int_{\hat{\sigma}} \nabla u(x) \cdot \mathbf{n}_{K,\sigma} d\gamma(x)$ . A Taylor expansion gives

$$I_\sigma - I_{\hat{\sigma}} = \frac{1}{\text{m}(\sigma)} \int_\sigma \int_{\frac{1}{2}}^1 H(u)(x_K + t(x - x_K))(x - x_K) \cdot \mathbf{n}_{K,\sigma} dt d\gamma(x).$$

The change of variables in this last integral  $z = x_K + t(x - x_K)$ , which gives  $dz = 2d_{K,\sigma}t^{d-1}dt d\gamma(x)$ , yields, with  $E_\sigma = \{tx + (1-t)x_K, x \in \sigma, t \in [\frac{1}{2}, 1]\}$  and some  $C_5$  only depending on  $d$  (note that  $t \geq \frac{1}{2}$ ),

$$|I_\sigma - I_{\hat{\sigma}}| \leq \frac{C_5}{\text{m}(\sigma) d_{K,\sigma}} \int_{E_\sigma} |H(u)(z)| |x - x_K| dz.$$

Then, using once more the Cauchy-Schwarz inequality and  $|x - x_K| \leq \text{diam}(K)$ ,

$$|I_\sigma - I_{\hat{\sigma}}|^2 \leq \frac{C_6 (\text{diam}(K))^2}{\text{m}(\sigma) d_\sigma} \int_{E_\sigma} |H(u)(z)|^2 dz \leq \frac{C_6 (\text{size}(\mathcal{T}))^2}{\text{m}(\sigma) d_\sigma} \int_{\mathcal{V}_\sigma} |H(u)(z)|^2 dz, \quad (31)$$

with some  $C_6$  only depending on  $d$ .

Inequalities (30) and (31) yield (26) for some  $C$  only depending on  $d$  and  $\zeta$  for  $u \in C^2(\overline{\mathcal{V}}_\sigma)$ . Taking  $C$  convenient for  $\sigma \in \mathcal{E}_{\text{int}}$  and  $\mathcal{E}_{\text{ext}}$  gives (26) for all  $\sigma \in \mathcal{E}$ .

Now for the density argument, let  $u \in H^2(\mathcal{V}_\sigma)$  and let  $(u_n)_{n \in \mathbb{N}} \subset C^2(\overline{\mathcal{V}}_\sigma)$  be a sequence which converges to  $u$  in the  $H^2(\mathcal{V}_\sigma)$  norm. Thanks to the previous result, one has

$$\left| \frac{u_n(x_L) - u_n(x_K)}{d_{K,\sigma}} - \frac{1}{\text{m}(\sigma)} \int_\sigma \nabla u_n(x) \cdot \mathbf{n}_{K,\sigma} d\gamma(x) \right| \leq C_1 \text{size}(\mathcal{T}) (\text{m}(\sigma) d_\sigma)^{-1/2} \|u_n\|_{H^2(\mathcal{V}_\sigma)}$$

Thanks to Sobolev imbeddings the sequence  $(u_n)_{n \in \mathbb{N}} \subset C^2(\overline{\mathcal{V}}_\sigma)$  converges to  $u \in H^2(\mathcal{V}_\sigma)$  uniformly and the sequence  $(\nabla u_n \cdot \mathbf{n}_{K,\sigma} \subset L^2(\sigma))$  converges to  $\nabla u \cdot \mathbf{n}_{K,\sigma}$  in  $L^2(\sigma)$  and therefore in  $L^1(\sigma)$ . Hence one pass to the limit in the left hand side term and obviously in the right hand side too. This gives (26) for some  $C$  only depending on  $d$  and  $\zeta$  for  $u \in H^2(\mathcal{V}_\sigma)$ .

Let us now prove (27) in the case  $\sigma \in \mathcal{E}_{\text{int}}$ ; let  $\sigma = K|L$  with  $K \in \mathcal{T}$  and  $L \in \mathcal{N}(K)$ . One assumes  $v_{K,\sigma} \geq 0$  (the case  $v_{K,\sigma} < 0$  works in the same way) so

$$|r_{K,\sigma}| = \left| \frac{1}{\text{m}(\sigma)} \int_\sigma \mathbf{v}(x) \cdot \mathbf{n}_{K,\sigma} (u(x) - u(x_K)) d\gamma(x) \right|.$$

It is possible to assume, for simplicity of notations and without loss of generality, that  $\sigma = 0 \times \tilde{\sigma}$ , with some  $\tilde{\sigma} \subset \mathbb{R}^{d-1}$ , and  $x_K = (-d_{K,\sigma}, 0)^t$ . A Taylor expansion, using  $u \in C^1(\overline{\mathcal{V}_\sigma})$  gives with  $x = (0, \tilde{x})^t \in \sigma$

$$|r_{K,\sigma}| \leq \sup_{x \in \Omega} |\mathbf{v}(x)| \frac{\text{size}(\mathcal{T})}{m(\sigma)} \int_{\tilde{\sigma}} \int_0^1 \left| \nabla u((t-1)d_{K,\sigma}, t\tilde{x}) \right| dt d\tilde{x}.$$

Let  $p > d$  be such that  $p < +\infty$  if  $d = 2$  and  $p \leq 6$  if  $d = 3$ , let  $p'$  be its conjugate exponent that is  $\frac{1}{p} + \frac{1}{p'} = 1$ . Thanks to Hölder's inequality:

$$\begin{aligned} |r_{K,\sigma}| &\leq \sup_{x \in \overline{\Omega}} |\mathbf{v}(x)| \frac{\text{size}(\mathcal{T})}{m(\sigma)} \left( \int_{\tilde{\sigma}} \int_0^1 \left| \nabla u((t-1)d_{K,\sigma}, t\tilde{x}) \right| t^{d-1} d_{K,\sigma} dt d\tilde{x} \right)^{1/p} \\ &\quad \times \left( \int_{\tilde{\sigma}} \int_0^1 \frac{1}{(t^{d-1} d_{K,\sigma})^{p'/p}} dt d\tilde{x} \right)^{1/p'}. \end{aligned}$$

Using a change of variables such that  $(\tilde{x}, t) \mapsto z = ((t-1)d_{K,\sigma}, t\tilde{x})$  and remarking that  $\frac{p'}{p}(d-1) = (p'-1)(d-1) < 1$  since  $p > d$ , one obtains

$$\begin{aligned} |r_{K,\sigma}| &\leq \sup_{x \in \overline{\Omega}} |\mathbf{v}(x)| \|u\|_{W^{1,p}(\mathcal{V}_{K,\sigma})} \text{size}(\mathcal{T}) (m(\sigma) d_{K,\sigma})^{-1/p} \left( \int_0^1 \frac{1}{t^{(p'-1)(d-1)}} dt \right)^{1/p'} \\ &= \frac{\sup_{x \in \overline{\Omega}} |\mathbf{v}(x)|}{(1 - (p'-1)(d-1))^{1/p'}} \|u\|_{W^{1,p}(\mathcal{V}_{K,\sigma})} \text{size}(\mathcal{T}) (m(\sigma) d_{K,\sigma})^{-1/p}. \quad (32) \end{aligned}$$

Remarking that  $d_\sigma = d_{K,\sigma} + d_{L,\sigma} \geq 2\zeta \text{diam}(K) \geq 2\zeta d_{K,\sigma}$  one obtains (27) for some  $C$  only depending on  $\mathbf{v}$ ,  $\zeta$  and  $p$ .

Now let us prove (27) for  $\sigma \in \mathcal{E}_{\text{ext}} \cap \mathcal{E}_K$ ,  $K \in \mathcal{T}$ . If  $v_{K,\sigma} \geq 0$ , the proof of (27) is identical to the case  $\sigma \in \mathcal{E}_{\text{int}}$ , so one assumes  $v_{K,\sigma} < 0$ ; hence:

$$|r_{K,\sigma}| = \left| \frac{1}{m(\sigma)} \int_{\sigma} \mathbf{v}(x) \cdot \mathbf{n}_{K,\sigma} (u(x) - u(y_\sigma)) d\gamma(x) \right|.$$

One can assume, without loss of generality, that  $x_K = 0$  and  $\sigma = d_{K,\sigma} \times \tilde{\sigma}$  with  $\tilde{\sigma} \subset \mathbb{R}^{d-1}$ . We introduce  $\hat{\sigma} = \{(\frac{d_{K,\sigma}}{2}, \frac{x}{2}), x \in \tilde{\sigma}\}$ . Note that  $m(\hat{\sigma}) = \frac{m(\sigma)}{2^{d-1}}$ , then:

$$|r_{K,\sigma}| \leq \sup_{x \in \overline{\Omega}} |\mathbf{v}(x)| \left( \frac{1}{m(\sigma) m(\hat{\sigma})} \int_{\sigma} \int_{\hat{\sigma}} |u(x) - u(y)| d\gamma(x) d\gamma(y) + \frac{1}{m(\hat{\sigma})} \int_{\hat{\sigma}} |u(y) - u(y_\sigma)| d\gamma(y) \right).$$

Then using a Taylor expansion, a change of variables and Hölder's inequality (for more details see the proof of (32)), one has:

$$|r_{K,\sigma}| \leq C \|u\|_{W^{1,p}(\mathcal{V}_\sigma)} \text{size}(\mathcal{T}) (m(\sigma) d_\sigma)^{-1/p},$$

for any  $p > d$  such that  $p < +\infty$  if  $d = 2$  and  $p \leq 6$  if  $d = 3$  and where  $C$  only depends on  $\mathbf{v}$ ,  $d$  and  $p$ . Finally let us prove (28). Using a Taylor expansion, one obtains

$$|\rho_K| \leq \frac{\|b\|_\infty \text{size}(\mathcal{T})}{m(K)} \int_K \int_0^1 |\nabla u(t x + (1-t)x_K)| dt dx.$$

Using the change of variables such that  $x \mapsto z = t x + (1-t)x_K$  and denoting by  $K_t$  the image of  $K$  by this change of variables, one obtains:

$$|\rho_K| \leq \frac{\|b\|_\infty \text{size}(\mathcal{T})}{m(K)} \int_0^1 \int_{K_t} \chi_{K_t}(z) \frac{|\nabla u(z)|}{t^d} dz dt,$$

where  $\chi_{K_t}$  is the characteristic function of  $K_t$ .

Thanks to Hölder's inequality and using  $m(K_t) \leq t^d m(K)$ , one has:

$$|\rho_K| \leq \frac{\|b\|_\infty \text{size}(\mathcal{T})}{m(K)} \int_0^1 \left( \int_K |\nabla u(z)|^p dz \right)^{1/p} \frac{(m(K_t))^{1/p'}}{t^d} dt \leq \frac{\|b\|_\infty \text{size}(\mathcal{T})}{m(K)^{1/p}} \|u\|_{W^{1,p}(K)} \int_0^1 \frac{1}{t^{d/p}} dt,$$

for all  $p > d$  such that  $p < +\infty$  if  $d = 2$  and  $p \leq 6$  if  $d = 3$ . As  $p > d$  we obtain (28). This concludes the proof of Lemma 3 and also step 1.

### Step 2 (Error estimates)

Now, one proves Estimates (18) and (19).

As  $u$  is the exact solution to (1), (3), for all  $K \in \mathcal{T}$ , one has:

$$\sum_{\sigma \in \mathcal{E}_K} (\bar{F}_{K,\sigma} + \bar{V}_{K,\sigma}) + \int_K b(x) u(x) dx = \int_K f(x) dx. \quad (33)$$

Substracting (12) off the previous equation, using (24) and the regularity of  $u$  yields

$$\sum_{\sigma \in \mathcal{E}_K} ((F_{K,\sigma}^* - F_{K,\sigma}) + (V_{K,\sigma}^* - V_{K,\sigma})) + b_K m(K) e_K = -m(K) \rho_K - \sum_{\sigma \in \mathcal{E}_K} m(\sigma) (R_{K,\sigma} + r_{K,\sigma}),$$

Multiplying the result by  $e_K$ , summing for  $K \in \mathcal{T}$ , and noting that

$$\sum_{K \in \mathcal{T}} \sum_{\sigma \in \mathcal{E}_K} (F_{K,\sigma}^* - F_{K,\sigma}) e_K = \sum_{\sigma \in \mathcal{E}} |D_\sigma e|^2 \tau_\sigma = \|e\|_{1,\mathcal{T}}^2,$$

yields

$$\begin{aligned} \|e_{\mathcal{T}}\|_{1,\mathcal{T}}^2 + \sum_{K \in \mathcal{T}} \sum_{\sigma \in \mathcal{E}_K} v_{K,\sigma} e_{\sigma,+} e_K + \int_{\Omega} b(x) (e_{\mathcal{T}}(x))^2 dx \\ \leq - \sum_{K \in \mathcal{T}} m(K) \rho_K e_K - \sum_{K \in \mathcal{T}} \sum_{\sigma \in \mathcal{E}_K} m(\sigma) (R_{K,\sigma} + r_{K,\sigma}) e_K, \end{aligned} \quad (34)$$

where  $|D_\sigma e|$  is defined in Definition 2 and  $e_{\sigma,+} = u(x_{\sigma,+}) - u_{\sigma,+}$ .

Reordering the summation over the set of edges, one has

$$\sum_{K \in \mathcal{T}} \sum_{\sigma \in \mathcal{E}_K} v_{K,\sigma} e_{\sigma,+} e_K = \sum_{\sigma \in \mathcal{E}} v_\sigma (e_{\sigma,+} - e_{\sigma,-}) e_{\sigma,+} = \frac{1}{2} \sum_{\sigma \in \mathcal{E}} v_\sigma ((e_{\sigma,+} - e_{\sigma,-})^2 + (e_{\sigma,+}^2 - e_{\sigma,-}^2)),$$

where  $v_\sigma = |\int_\sigma \mathbf{v}(x) \cdot \mathbf{n}_\sigma d\gamma(x)|$ ,  $\mathbf{n}_\sigma$  being a unit normal vector to  $\sigma$ , and  $e_{\sigma,-}$  is the downstream value to  $\sigma$  with respect to  $\mathbf{v}$ , that is to say if  $\sigma = K|L$ , then  $e_{\sigma,-} = e_K$  if  $v_{K,\sigma} \leq 0$ , and  $e_{\sigma,-} = e_L$  otherwise; if  $\sigma \in \mathcal{E}_K \cap \mathcal{E}_{\text{ext}}$ , then  $e_{\sigma,-} = u_K$  if  $v_{K,\sigma} \leq 0$  and  $e_{\sigma,-} = 0$  otherwise.

Now note that:

$$\sum_{\sigma \in \mathcal{E}} v_\sigma (e_{\sigma,+}^2 - e_{\sigma,-}^2) = \sum_{K \in \mathcal{T}} \left( \int_{\partial K} \mathbf{v}(x) \cdot \mathbf{n}_K d\gamma(x) \right) e_K^2 = \int_{\Omega} (\text{div } \mathbf{v}(x)) e_{\mathcal{T}}^2(x) dx.$$

Then, one obtains

$$\sum_{K \in \mathcal{T}} \sum_{\sigma \in \mathcal{E}_K} v_{K,\sigma} e_{\sigma,+} e_K \geq \frac{1}{2} \int_{\Omega} (\text{div } \mathbf{v}(x)) e_{\mathcal{T}}^2(x) dx,$$

and so, using this result in (34),

$$\|e_{\mathcal{T}}\|_{1,\mathcal{T}}^2 + \int_{\Omega} \left( \frac{\text{div } \mathbf{v}(x)}{2} + b(x) \right) e_{\mathcal{T}}^2(x) dx \leq - \sum_{K \in \mathcal{T}} m(K) \rho_K e_K - \sum_{K \in \mathcal{T}} \sum_{\sigma \in \mathcal{E}_K} m(\sigma) (R_{K,\sigma} + r_{K,\sigma}) e_K. \quad (35)$$

Let us now deal with the consistency error terms: By Young's inequality, for all  $\delta > 0$ , one has

$$-\sum_{K \in \mathcal{T}} m(K) \rho_K e_K \leq \frac{\delta}{2} \|e_{\mathcal{T}}\|_{L^2(\Omega)}^2 + \frac{1}{2\delta} \sum_{K \in \mathcal{T}} m(K) \rho_K^2.$$

Hence if  $u \in C^2(\overline{\Omega})$ , using Lemma 2, one obtains, for all  $\delta > 0$ :

$$-\sum_{K \in \mathcal{T}} m(K) \rho_K e_K \leq \frac{\delta}{2} \|e_{\mathcal{T}}\|_{L^2(\Omega)}^2 + \frac{C}{\delta} (\text{size}(\mathcal{T}))^2, \quad (36)$$

where  $C$  only depends on  $u$ ,  $b$  and  $\Omega$ .

If  $u$  is only in  $H^2(\Omega)$ , thanks to Lemma 3 and to Hölder's inequality one has, for all  $\delta > 0$  and all  $p > d$  such that  $p < +\infty$  if  $d = 2$  and  $p \leq 6$  if  $d = 3$ :

$$\begin{aligned} -\sum_{K \in \mathcal{T}} m(K) \rho_K e_K &\leq \frac{\delta}{2} \|e_{\mathcal{T}}\|_{L^2(\Omega)}^2 + \frac{(\|b\|_\infty \text{size}(\mathcal{T}))^2}{2\delta} \sum_{K \in \mathcal{T}} m(K)^{1-2/p} \|u\|_{W^{1,p}(K)}^2 \\ &\leq \frac{\delta}{2} \|e_{\mathcal{T}}\|_{L^2(\Omega)}^2 + \frac{(\|b\|_\infty \text{size}(\mathcal{T}))^2}{2\delta} \|u\|_{W^{1,p}(\Omega)}^2 m(\Omega)^{(p-2)/p}, \end{aligned}$$

choosing  $p = 4$ , one obtains (36) for all  $\delta > 0$ , where  $C$  only depends on  $b$ ,  $u$  and  $\Omega$ .

Furthermore, thanks to the conservativity property of the scheme (see (4)), one has  $R_{K,\sigma} = -R_{L,\sigma}$  and  $r_{K,\sigma} = -r_{L,\sigma}$  for  $\sigma \in \mathcal{E}_{\text{int}}$  such that  $\sigma = K|L$ . Let  $R_\sigma = |R_{K,\sigma}|$  and  $r_\sigma = |r_{K,\sigma}|$  if  $\sigma \in \mathcal{E}_K$ . Reordering the summation over the edges and using Young's inequality, one obtains

$$\begin{aligned} \left| \sum_{K \in \mathcal{T}} \sum_{\sigma \in \mathcal{E}_K} m(\sigma) (R_{K,\sigma} + r_{K,\sigma}) e_K \right| &\leq \sum_{\sigma \in \mathcal{E}} m(\sigma) (D_\sigma e)(R_\sigma + r_\sigma) \\ &\leq \frac{1}{2} \sum_{\sigma \in \mathcal{E}} \frac{m(\sigma)}{d_\sigma} (D_\sigma e)^2 + \frac{1}{2} \sum_{\sigma \in \mathcal{E}} m(\sigma) d_\sigma (R_\sigma + r_\sigma)^2. \end{aligned} \quad (37)$$

Now, using Lemma 2, if  $u \in C^2(\overline{\Omega})$ , or Lemma 3 (with  $p = 4$ ) and Hölder's inequality, if  $u$  is only in  $H^2(\Omega)$  (for more details see the proof of inequality (36)), and remarking that  $\sum_{\sigma \in \mathcal{E}} m(\sigma) d_\sigma = d m(\Omega)$ , (37) yields the existence of  $C$ , only depending on  $u$ ,  $\mathbf{v}$ ,  $d$  and  $\Omega$  if  $u \in C^2(\overline{\Omega})$  and on  $u$ ,  $\mathbf{v}$ ,  $d$ ,  $\zeta$  and  $\Omega$  if  $u$  is only in  $H^2(\Omega)$ , such that

$$\left| \sum_{K \in \mathcal{T}} \sum_{\sigma \in \mathcal{E}_K} m(\sigma) (R_{K,\sigma} + r_{K,\sigma}) e_K \right| \leq \frac{1}{2} \|e_{\mathcal{T}}\|_{1,\mathcal{T}}^2 + C (\text{size}(\mathcal{T}))^2.$$

Hence, (35), (36) and the previous inequality yield for all  $\delta > 0$

$$\frac{1}{2} \|e_{\mathcal{T}}\|_{1,\mathcal{T}}^2 + \int_{\Omega} \left( \frac{\text{div} \mathbf{v}(x)}{2} + b(x) \right) e_{\mathcal{T}}^2(x) dx \leq \frac{\delta}{2} \|e_{\mathcal{T}}\|_{L^2(\Omega)}^2 + C \left( 1 + \frac{1}{\delta} \right) (\text{size}(\mathcal{T}))^2, \quad (38)$$

where  $C$  depends only on  $b$ ,  $u$ ,  $\mathbf{v}$ ,  $d$  and  $\Omega$  if  $u \in C^2(\overline{\Omega})$  and on  $b$ ,  $u$ ,  $\mathbf{v}$ ,  $d$ ,  $\zeta$  and  $\Omega$  if  $u$  is only in  $H^2(\Omega)$ .

If there exists  $\delta > 0$  such that  $\text{div} \mathbf{v}/2 + b \geq \delta/2$ , this inequality yields Estimate (18) and Estimate (19). Otherwise, to obtain Estimate (18), one uses the inequality (38) with  $\delta = \frac{1}{2(\text{diam}(\Omega))^2}$ , the positivity of  $\text{div} \mathbf{v}/2 + b$  and a discrete Poincaré inequality which is proved in [9] or [7] and which we recall here:

**Lemma 4 (Discrete Poincaré inequality)** *Let  $\mathcal{T}$  be an admissible finite volume mesh in the sense of Definition 1 and  $u$  be a function which is constant on each cell of  $\mathcal{T}$ , that is  $u(x) = u_K$  if  $x \in K$ ,  $K \in \mathcal{T}$ . Then*

$$\|u\|_{L^2(\Omega)} \leq \text{diam}(\Omega) \|u\|_{1,\mathcal{T}},$$

where  $\|\cdot\|_{1,\mathcal{T}}$  is the discrete  $H_0^1$  norm defined in Definition 2.

Using the above lemma once more and Estimate (18) gives Estimate (19).

This concludes the proofs of Theorems 1 and 2 and shows that the numerical solution converges towards the exact solution to (1), (3).  $\blacksquare$

**Remark 7** *The error estimates of Theorems 1 and 2 still hold if a non constant piecewise  $C^1$  diffusion scalar coefficient is considered i.e. if  $k$  satisfies Assumption 2 and if Equation (2) is discretized by the scheme (7)-(15).*

*In the general tensor case however, some more restrictive assumptions are needed on the mesh in order to obtain an error estimate: see [6] and [10]. More precisely, in the error estimate (18) and (19), the real number  $C$  now depends on  $\zeta_1$  and  $\zeta_2 \in \mathbb{R}_+$  such that*

$$\begin{aligned}\zeta_1 (\text{size}(\mathcal{T}))^2 &\leq m(K) \leq \zeta_2 (\text{size}(\mathcal{T}))^2, \\ \zeta_1 \text{size}(\mathcal{T}) &\leq m(\sigma) \leq \zeta_2 \text{size}(\mathcal{T}), \\ \zeta_1 \text{size}(\mathcal{T}) &\leq d_\sigma \leq \zeta_2 \text{size}(\mathcal{T}).\end{aligned}$$

## 4 Neumann boundary conditions

The second boundary condition we consider is a Neumann condition:

$$\nabla u(x) \cdot \mathbf{n}(x) = g^N(x), \quad x \in \partial\Omega, \quad (39)$$

where

**Assumption 5**  *$b = 0$  a.e. on  $\Omega$ ,  $\text{div } \mathbf{v} = 0$  on  $\Omega$ ,  $\mathbf{v} \cdot \mathbf{n} = 0$  on  $\partial\Omega$ , and  $g^N \in H^{1/2}(\partial\Omega)$  satisfies the following compatibility relation:  $\int_{\partial\Omega} g^N(x) d\gamma(x) + \int_{\Omega} f(x) dx = 0$ .*

Then under Assumptions 1 and 5, by Lax-Milgram Theorem, there exists a unique variational solution  $u \in H^1(\Omega)$  such that  $\int_{\Omega} u(x) dx = 0$ , of (1), (39). That is to say  $u \in H^1(\Omega)$  such that  $\int_{\Omega} u(x) dx = 0$  satisfies for all  $\phi \in H^1(\Omega)$

$$\int_{\Omega} (\nabla u(x) \cdot \nabla \phi(x) + \text{div}(\mathbf{v}(x) u(x)) \phi(x)) dx = \int_{\partial\Omega} g^N(x) \bar{\gamma}(\phi)(x) d\gamma(x) + \int_{\Omega} f(x) \phi(x) dx,$$

where  $\bar{\gamma}$  denotes the trace operator from  $H^1(\Omega)$  to  $H^{1/2}(\partial\Omega)$  and  $d\gamma$  is the integration symbol for the  $(d-1)$ -dimensional Lebesgue measure.

**Remark 8** *The assumptions  $b = 0$ ,  $\text{div } \mathbf{v} = 0$  and  $\mathbf{v} \cdot \mathbf{n} = 0$  are sufficient to prove the coercivity of the bilinear form of the variational formulation. However, if the hypotheses on  $\mathbf{v}$  and  $b$  of Assumption 5 are not satisfied, we do not need a compatibility relation. This latter case is therefore treated in section 5 which deals with Robin boundary conditions.*

### 4.1 Discretization

We use the same notations as in the previous section. Let  $\mathcal{T}$  be an admissible mesh in the sense of Definition 1 and  $u_K$  be the discrete unknown associated with the control volume  $K$  for all  $K \in \mathcal{T}$ . Let us integrate Equation (1) on each cell of the mesh; the diffusion flux is discretized on interior edges only since it is known on the boundary of  $\Omega$ ; an upstream scheme is used for the convection term and one obtains:

$$\sum_{\sigma \in \mathcal{E}_K \cap \mathcal{E}_{\text{int}}} (F_{K,\sigma} + v_{K,\sigma} u_{\sigma,+}) = m(K) f_K + \sum_{\sigma \in \mathcal{E}_K \cap \mathcal{E}_{\text{ext}}} g_{\sigma}^N, \quad (40)$$

where  $F_{K,\sigma}$  is defined by (5) if  $\sigma \in \mathcal{E}_{\text{int}}$ ,  $\sigma = K|L$ ,

$$g_{\sigma}^N = \int_{\sigma} g^N(x) d\gamma(x) \quad \text{if } \sigma \in \mathcal{E}_{\text{ext}} \quad (41)$$

and where  $v_{K,\sigma}$  and  $f_K$  are defined by (8) and (11) and  $u_{\sigma,+}$  is defined by:

$$\text{if } \sigma = K|L, \quad u_{\sigma,+} = \begin{cases} u_K & \text{if } v_{K,\sigma} \geq 0, \\ u_L & \text{otherwise.} \end{cases} \quad (42)$$

## 4.2 Existence, uniqueness and stability of the approximate solution

One proves the following proposition which gives the existence of the approximate solution and the uniqueness up to a constant like in the continuous case.

**Proposition 2** *Under Assumptions 1 and 5, let  $\mathcal{T}$  be an admissible mesh in the sense of Definition 1. Then, there exists a solution  $u_{\mathcal{T}}$  to (15), (40), (41), (42), (5), (8) and (11). This solution is unique up to a constant.*

**Proof of Proposition 2**

Let us first study the kernel of the linear operator defined by the left hand side of (40). For all  $K \in \mathcal{T}$ , suppose that  $\sum_{\sigma \in \mathcal{E}_{\text{ext}} \cap \mathcal{E}_K} g_{\sigma}^N + m(K) f_K = 0$ . Let us denote by  $K_0$  a cell of  $\mathcal{T}$  such that  $u_{K_0} = \min_{K \in \mathcal{T}} u_K$ . Since  $\operatorname{div} \mathbf{v} = 0$  on  $\Omega$  and  $\mathbf{v} \cdot \mathbf{n} = 0$  on  $\partial\Omega$ , one has  $\sum_{\sigma \in \mathcal{E}_{K_0} \cap \mathcal{E}_{\text{int}}} \int_{\sigma} \mathbf{v}(x) \cdot \mathbf{n}_{K_0,\sigma} d\gamma(x) = 0$ ; then, using the numerical scheme (40), one gets:

$$\sum_{L \in \mathcal{N}(K_0)} -m(K_0|L) \frac{u_L - u_{K_0}}{d_{K_0|L}} + \sum_{\sigma \in \mathcal{E}_{K_0} \cap \mathcal{E}_{\text{int}}} v_{K_0,\sigma} (u_{\sigma,+} - u_{K_0}) = 0.$$

Now remarking that  $u_{K_0} \leq u_L$  for all  $L \in \mathcal{N}(K_0)$  and  $\sum_{\sigma \in \mathcal{E}_{K_0} \cap \mathcal{E}_{\text{int}}} v_{K_0,\sigma} (u_{\sigma,+} - u_{K_0}) \leq 0$ , one obtains  $u_L = u_{K_0}$  for any neighbour  $L$  of  $K_0$ . Since  $\Omega$  is connected, one has  $u_L = u_K$  for all  $(K, L) \in \mathcal{T}^2$ . So the dimension of the kernel of the linear operator defined by the left hand side of (40) is 1. Let us now study its image. First note that its dimension is  $\operatorname{card}(\mathcal{T}) - 1$  where  $\operatorname{card}(\mathcal{T})$  is the number of control volumes of the mesh.

Summing Equation (40) over  $K \in \mathcal{T}$ , remarking that  $\mathbf{v} \cdot \mathbf{n} = 0$  on  $\partial\Omega$  and  $v_{K,\sigma} = -v_{L,\sigma}$  for all  $K \in \mathcal{T}$  and all  $L \in \mathcal{N}(K)$ , one obtains:

$$\sum_{K \in \mathcal{T}} \left( \sum_{\sigma \in \mathcal{E}_{\text{ext}} \cap \mathcal{E}_K} g_{\sigma}^N + m(K) f_K \right) = \int_{\partial\Omega} g^N(x) d\gamma(x) + \int_{\Omega} f(x) dx = 0.$$

So assuming  $\int_{\partial\Omega} g^N(x) d\gamma(x) + \int_{\Omega} f(x) dx = 0$  there exists a solution to (15), (40), (41), (42), (5), (8) and (11) and this solution is unique up to a constant.

Let us introduce, like in the Dirichlet case, the discrete  $H^1$  semi-norm of a function from  $\Omega$  to  $\mathbb{R}$  which is constant on each control volume (or cell) of  $\mathcal{T}$ .

**Definition 3 (Discrete  $H^1$  semi-norm)** *Let  $\mathcal{T}$  be an admissible finite volume mesh in the sense of Definition 1. Let  $u$  be a function which is constant on each control volume of  $\mathcal{T}$ , that is  $u(x) = u_K$  if  $x \in K$ ,  $K \in \mathcal{T}$ , one defines the discrete  $H^1$  semi-norm by*

$$|u|_{1,\mathcal{T}} = \left( \sum_{\sigma \in \mathcal{E}_{\text{int}}} \frac{m(\sigma)}{d_{\sigma}} (D_{\sigma} u)^2 \right)^{1/2},$$

where  $D_{\sigma} u = |u_K - u_L|$  if  $\sigma = K|L$  and the set  $\mathcal{E}_{\text{int}}$  is defined in Definition 1.

Let us now give an estimate on the approximate solution.

**Lemma 5** Under Assumptions 1 and 5, let  $\mathcal{T}$  be an admissible mesh in the sense of Definition 1. Let  $u_{\mathcal{T}}$  be the solution to (15), (40), (41), (42), (5), (8) and (11) such that  $\int_{\Omega} u_{\mathcal{T}}(x) dx = \alpha$ . Then there exists  $C \in \mathbb{R}_+$  depending only on  $\Omega$  such that

$$|u_{\mathcal{T}}|_{1,\mathcal{T}} \leq C \left( \|g^N\|_{L^2(\partial\Omega)} + \|f\|_{L^2(\Omega)} + \alpha \right). \quad (43)$$

where  $|\cdot|_{1,\mathcal{T}}$  is defined in Definition 3

**Proof** of Lemma 5

Let  $K \in \mathcal{T}$ , we multiply (40) by  $u_K$  and we sum the result over  $K \in \mathcal{T}$ , we obtain:

$$\sum_{K \in \mathcal{T}} \sum_{\substack{\sigma \in \mathcal{E}_K \cap \mathcal{E}_{\text{int}} \\ \sigma = K|L}} \left( \frac{u_K - u_L}{d_{K|L}} u_K m(\sigma) + v_{K,\sigma} u_{\sigma,+} u_K \right) = \sum_{K \in \mathcal{T}} m(K) f_K u_K + \sum_{K \in \mathcal{T}} \sum_{\substack{\sigma \in \mathcal{E}_K \cap \mathcal{E}_{\text{ext}}}} g_{\sigma}^N u_K.$$

Now let us note that:

$$\sum_{K \in \mathcal{T}} \sum_{\substack{\sigma \in \mathcal{E}_K \cap \mathcal{E}_{\text{int}} \\ \sigma = K|L}} \frac{u_K - u_L}{d_{K|L}} u_K m(\sigma) = \sum_{\sigma \in \mathcal{E}_{\text{int}}} \frac{m(\sigma)}{d_{\sigma}} (D_{\sigma} u)^2 \quad (44)$$

Furthermore, for  $\sigma \in \mathcal{E}_{\text{int}}$ ,  $\sigma = K|L$ , let  $u_{\sigma,-} = u_K$  if  $v_{K,\sigma} < 0$  and  $u_{\sigma,-} = u_L$  otherwise and  $v_{\sigma} = |v_{K,\sigma}| = |v_{L,\sigma}|$ . Then, one has

$$\begin{aligned} \sum_{K \in \mathcal{T}} \sum_{\substack{\sigma \in \mathcal{E}_K \cap \mathcal{E}_{\text{int}} \\ \sigma = K|L}} v_{K,\sigma} u_{\sigma,+} u_K &= \sum_{\sigma \in \mathcal{E}_{\text{int}}} v_{\sigma} \left( (u_{\sigma,+})^2 - u_{\sigma,+} u_{\sigma,-} \right) \\ &= \sum_{\sigma \in \mathcal{E}_{\text{int}}} v_{\sigma} \left( \frac{(u_{\sigma,+} - u_{\sigma,-})^2}{2} + \frac{(u_{\sigma,+})^2}{2} - \frac{(u_{\sigma,-})^2}{2} \right) \geq \frac{1}{2} \sum_{K \in \mathcal{T}} \sum_{\sigma \in \mathcal{E}_{\text{int}}} v_{K,\sigma} (u_K)^2 \\ &= \frac{1}{2} \int_{\Omega} \operatorname{div}(\mathbf{v}(x)) (u_{\mathcal{T}}(x))^2 dx - \frac{1}{2} \int_{\partial\Omega} \mathbf{v}(x) \cdot \mathbf{n}(x) (u_{\mathcal{T}})^2 d\gamma(x). \end{aligned} \quad (45)$$

Therefore by Assumption 5 and Young's inequality, we get for all  $\delta > 0$ :

$$|u_{\mathcal{T}}|_{1,\mathcal{T}}^2 \leq \delta \left( \|g^N\|_{L^2(\partial\Omega)}^2 + \|f\|_{L^2(\Omega)}^2 \right) + \frac{1}{4\delta} \left( \|u_{\mathcal{T}}\|_{L^2(\partial\Omega)}^2 + \|u_{\mathcal{T}}\|_{L^2(\Omega)}^2 \right),$$

where  $u_{\mathcal{T}}(x) = u_K$  for almost every  $x \in \sigma$ ,  $\sigma \in \mathcal{E}_{\text{ext}} \cap \mathcal{E}_K$ .

In order to conclude, one uses Lemmas 6 and 7, which are stated and proved below and obtains

$$|u_{\mathcal{T}}|_{1,\mathcal{T}}^2 \leq \delta \left( \|g^N\|_{L^2(\partial\Omega)}^2 + \|f\|_{L^2(\Omega)}^2 \right) + \frac{C}{\delta} \left( |u_{\mathcal{T}}|_{1,\mathcal{T}}^2 + \frac{2}{m(\Omega)} \alpha^2 \right)$$

for all  $\delta > 0$  and where  $C$  only depends on  $\Omega$ .

Choosing  $\delta = 2C$  gives (43).

**Lemma 6 (Discrete Poincaré-Wirtinger inequality)** Let  $\mathcal{T}$  be an admissible mesh in the sense of Definition 1. Let  $u$  be a function which is constant on each cell of  $\mathcal{T}$ , that is  $u(x) = u_K$  if  $x \in K$ ,  $K \in \mathcal{T}$ . Then

$$\|u\|_{L^2(\Omega)}^2 \leq C |u|_{1,\mathcal{T}}^2 + 2(m(\Omega))^{-1} \left( \int_{\Omega} u(x) dx \right)^2,$$

where  $C$  only depends on  $\Omega$  and  $|\cdot|_{1,\mathcal{T}}$  is defined in Definition 3.

### Proof of Lemma 6

Let  $\mathcal{T}$  be an admissible mesh and  $u$  be a function which is constant on each cell of  $\mathcal{T}$ . Let  $m_\Omega(u)$  be the mean value of  $u$  over  $\Omega$ , that is

$$m_\Omega(u) = \frac{1}{m(\Omega)} \int_{\Omega} u(x) dx.$$

Note that

$$\|u\|_{L^2(\Omega)}^2 \leq 2\|u - m_\Omega(u)\|_{L^2(\Omega)}^2 + 2(m_\Omega(u))^2 m(\Omega),$$

and therefore the proof of Lemma 6 is reduced to the proof of the existence of  $D \geq 0$ , only depending on  $\Omega$ , such that

$$\|u - m_\Omega(u)\|_{L^2(\Omega)}^2 \leq D|u|_{1,\mathcal{T}}^2. \quad (46)$$

The proof of (46) may be decomposed into three steps (indeed, if  $\Omega$  is convex, the first step is sufficient).

#### Step 1 (Estimate on a convex part of $\Omega$ )

Let  $\omega$  be an open convex subset of  $\Omega$ ,  $\omega \neq \emptyset$  and  $m_\omega(u)$  be the mean value of  $u$  on  $\omega$ . In this step, one proves that there exists  $C_0$ , depending only on  $\omega$ , such that

$$\|u(x) - m_\omega(u)\|_{L^2(\omega)}^2 \leq \frac{1}{m(\omega)} C_0 |u|_{1,\mathcal{T}}^2. \quad (47)$$

(Taking  $\omega = \Omega$ , this proves (46) and Lemma 6 in the case where  $\omega$  is convex.)

Noting that

$$\int_{\omega} (u(x) - m_\omega(u))^2 dx \leq \frac{1}{m(\omega)} \int_{\omega} \int_{\omega} (u(x) - u(y))^2 dy dx,$$

(47) is proved provided that there exists  $C_0 \in \mathbb{R}_+$ , only depending on  $\omega$ , such that

$$\int_{\omega} \int_{\omega} (u(x) - u(y))^2 dx dy \leq C_0 |u|_{1,\mathcal{T}}^2. \quad (48)$$

For  $\sigma \in \mathcal{E}_{\text{int}}$ , let the function  $\chi_\sigma$  from  $\mathbb{R}^d \times \mathbb{R}^d$  to  $\{0, 1\}$  be defined by

$$\begin{aligned} \chi_\sigma(x, y) &= 1, \text{ if } x, y \in \overline{\omega}, [x, y] \cap \sigma \neq \emptyset, \\ \chi_\sigma(x, y) &= 0, \text{ if } x \notin \overline{\omega} \text{ or } y \notin \overline{\omega} \text{ or } [x, y] \cap \sigma = \emptyset. \end{aligned}$$

(Recall that  $[x, y] = \{tx + (1-t)y, t \in [0, 1]\}$ .) For a.e.  $x, y \in \omega$ , one has, with  $D_\sigma u = |u_K - u_L|$  if  $\sigma \in \mathcal{E}_{\text{int}}$ ,  $\sigma = K|L$ ,

$$(u(x) - u(y))^2 \leq \left( \sum_{\sigma \in \mathcal{E}_{\text{int}}} |D_\sigma u| \chi_\sigma(x, y) \right)^2,$$

(note that the convexity of  $\omega$  is used here) which yields, thanks to the Cauchy-Schwarz inequality,

$$(u(x) - u(y))^2 \leq \sum_{\sigma \in \mathcal{E}_{\text{int}}} \frac{|D_\sigma u|^2}{d_\sigma c_{\sigma,y-x}} \chi_\sigma(x, y) \sum_{\sigma \in \mathcal{E}_{\text{int}}} d_\sigma c_{\sigma,y-x} \chi_\sigma(x, y), \quad (49)$$

with

$$c_{\sigma,y-x} = \left| \frac{y-x}{|y-x|} \cdot \mathbf{n}_\sigma \right|,$$

recall that  $\mathbf{n}_\sigma$  is a unit normal vector to  $\sigma$ , and that  $x_K - x_L = \pm d_\sigma \mathbf{n}_\sigma$  if  $\sigma \in \mathcal{E}_{\text{int}}$ ,  $\sigma = K|L$ . For a.e.  $x, y \in \omega$ , one has

$$\sum_{\sigma \in \mathcal{E}_{\text{int}}} d_\sigma c_{\sigma,y-x} \chi_\sigma(x, y) = \left| (x_K - x_L) \cdot \frac{y-x}{|y-x|} \right|,$$

for some convenient control volumes  $K$  and  $L$ , depending on  $x$  and  $y$  (the convexity of  $\omega$  is used again here). Therefore,

$$\sum_{\sigma \in \mathcal{E}_{\text{int}}} d_\sigma c_{\sigma,y-x} \chi_\sigma(x, y) \leq \text{diam}(\omega).$$

Thus, integrating (49) with respect to  $x$  and  $y$  in  $\omega$ ,

$$\int_{\omega} \int_{\omega} (u(x) - u(y))^2 dx dy \leq \text{diam}(\omega) \int_{\omega} \int_{\omega} \sum_{\sigma \in \mathcal{E}_{\text{int}}} \frac{|D_\sigma u|^2}{d_\sigma c_{\sigma,y-x}} \chi_\sigma(x, y) dx dy,$$

which gives, by a change of variables,

$$\int_{\omega} \int_{\omega} (u(x) - u(y))^2 dx dy \leq \text{diam}(\omega) \int_{\mathbb{R}^d} \left( \sum_{\sigma \in \mathcal{E}_{\text{int}}} \frac{|D_\sigma u|^2}{d_\sigma c_{\sigma,z}} \int_{\omega} \chi_\sigma(x, x+z) dx \right) dz. \quad (50)$$

Noting that, if  $|z| > \text{diam}(\omega)$ ,  $\chi_\sigma(x, x+z) = 0$ , for a.e.  $x \in \omega$ , and

$$\int_{\omega} \chi_\sigma(x, x+z) dx \leq m(\sigma) |z \cdot n_\sigma| = m(\sigma) |z| c_{\sigma,z} \text{ for a.e. } z \in \mathbb{R}^d,$$

therefore, with (50):

$$\int_{\omega} \int_{\omega} (u(x) - u(y))^2 dx dy \leq (\text{diam}(\omega))^2 m(B_\omega) \sum_{\sigma \in \mathcal{E}_{\text{int}}} \frac{m(\sigma) |D_\sigma u|^2}{d_\sigma},$$

where  $B_\omega$  denotes the ball of  $\mathbb{R}^d$  of center 0 and radius  $\text{diam}(\omega)$ .

This inequality proves (48) and then (47) with  $C_0 = (\text{diam}(\omega))^2 m(B_\omega)$  (which only depends on  $\omega$ ). Taking  $\omega = \Omega$ , it concludes the proof of Lemma 6 in the case where  $\Omega$  is convex.

### Step 2 (Estimate with respect to the mean value on a part of the boundary)

In this step, one proves the same inequality than (47) but with the mean value of  $u$  on a (arbitrary) part  $I$  of the boundary of  $\omega$  instead of  $m_\omega(u)$  and with a convenient  $C_1$  depending on  $I$ ,  $\Omega$  and  $\omega$  instead of  $C_0$ .

More precisely, let  $\omega$  be a polygonal open convex subset of  $\Omega$  and let  $I \subset \partial\omega$ , with  $m(I) > 0$  ( $m(I)$  is the  $(d-1)$ -Lebesgue measure of  $I$ ). Assume that  $I$  is included in a hyperplane of  $\mathbb{R}^d$ . Let  $\gamma u$  be the “trace” of  $u$  on the boundary of  $\omega$ , that is  $\gamma u(x) = u_K$  if  $x \in \partial\omega \cap \overline{K}$ , for  $K \in \mathcal{T}$  (if  $x \in \overline{K} \cap \overline{L}$ , the choice of  $\gamma u(x)$  between  $u_K$  and  $u_L$  does not matter). Let  $m_I(u)$  be the mean value of  $\gamma u$  on  $I$ . This step is devoted to the proof of the existence of  $C_1$ , only depending on  $\Omega$ ,  $\omega$  and  $I$ , such that

$$\|u(x) - m_I(u)\|_{L^2(\omega)}^2 \leq C_1 |u|_{1,\mathcal{T}}^2. \quad (51)$$

For the sake of simplicity, only the case  $d = 2$  is considered here. Since  $I$  is included in a hyperplane, it may be assumed, without loss of generality, that  $I = \{0\} \times J$ , with  $J \subset \mathbb{R}$  and  $\omega \subset \mathbb{R}_+ \times \mathbb{R}$  (one uses here the convexity of  $\omega$ ).

Let  $\alpha = \max\{x_1, x = (x_1, x_2)^t \in \overline{\omega}\}$  and  $a = (\alpha, \beta)^t \in \overline{\omega}$ . In the following,  $a$  is fixed. For a.e.  $x = (x_1, x_2)^t \in \omega$  and for a.e. (for the 1-Lebesgue measure)  $y = (0, \bar{y})^t \in I$  (with  $\bar{y} \in J$ ), one sets  $z(x, y) = ta + (1-t)y$  with  $t = x_1/\alpha$ . Note that, thanks to the convexity of  $\omega$ ,  $z(x, y) = (z_1, z_2)^t \in \overline{\omega}$ , with  $z_1 = x_1$ . The following inequality holds:

$$\pm(u(x) - \gamma u(y)) \leq |u(x) - u(z(x, y))| + |u(z(x, y)) - \gamma u(y)|.$$

In the following, the notation  $C_i$ ,  $i \in \mathbb{N}^*$ , will be used for quantities only depending on  $\Omega$ ,  $\omega$  and  $I$ .

Let us integrate the above inequality over  $y \in I$ , take the power 2, from the Cauchy-Schwarz inequality, an integration over  $x \in \omega$  leads to

$$\begin{aligned} \int_{\omega} (u(x) - m_I(u))^2 dx &\leq \frac{2}{m(I)} \int_{\omega} \int_I (u(x) - u(z(x, y)))^2 d\gamma(y) dx \\ &+ \frac{2}{m(I)} \int_{\omega} \int_I (u(z(x, y)) - u(y))^2 d\gamma(y) dx. \end{aligned}$$

Then,

$$\int_{\omega} (u(x) - m_I(u))^2 dx \leq \frac{2}{m(I)} (A + B),$$

with, since  $\omega$  is convex,

$$A = \int_{\omega} \int_I \left( \sum_{\sigma \in \mathcal{E}_{\text{int}}} |D_{\sigma} u| \chi_{\sigma}(x, z(x, y)) \right)^2 d\gamma(y) dx,$$

and

$$B = \int_{\omega} \int_I \left( \sum_{\sigma \in \mathcal{E}_{\text{int}}} |D_{\sigma} u| \chi_{\sigma}(z(x, y), y) \right)^2 d\gamma(y) dx.$$

Recall that, for  $\xi, \eta \in \overline{\omega}$ ,  $\chi_{\sigma}(\xi, \eta) = 1$  if  $[\xi, \eta] \cap \sigma \neq \emptyset$  and  $\chi_{\sigma}(\xi, \eta) = 0$  if  $[\xi, \eta] \cap \sigma = \emptyset$ . Let us now look for some bounds of  $A$  and  $B$  of the form  $C|u|_{1,\mathcal{T}}^2$ .

The bound for  $A$  is easy. Using the Cauchy-Schwarz inequality and the fact that

$$\sum_{\sigma \in \mathcal{E}_{\text{int}}} c_{\sigma, x-z(x, y)} d_{\sigma} \chi_{\sigma}(x, z(x, y)) \leq \text{diam}(\omega)$$

(recall that  $c_{\sigma, \eta} = |\frac{\eta}{|\eta|} \cdot \mathbf{n}_{\sigma}|$  (for  $\eta \in \mathbb{R}^2 \setminus 0$ ) gives

$$A \leq C_2 \int_{\omega} \int_I \sum_{\sigma \in \mathcal{E}_{\text{int}}} \frac{|D_{\sigma} u|^2 \chi_{\sigma}(x, z(x, y))}{c_{\sigma, x-z(x, y)} d_{\sigma}} dx d\gamma(y).$$

Since  $z_1 = x_1$ , one has  $c_{\sigma, x-z(x, y)} = c_{\sigma, e}$ , with  $e = (0, 1)^t$ . Let us perform the integration of the right hand side of the previous inequality, with respect to the first component of  $x$ , denoted by  $x_1$ , first. The result of the integration with respect to  $x_1$  is bounded by  $|u|_{1,\mathcal{T}}^2$ . Then, integrating with respect to  $x_2$  and  $y \in I$  gives  $A \leq C_3 |u|_{1,\mathcal{T}}^2$ .

In order to obtain a bound  $B$ , one remarks, as for  $A$ , that

$$B \leq C_4 \int_{\omega} \int_I \sum_{\sigma \in \mathcal{E}_{\text{int}}} \frac{|D_{\sigma} u|^2 \chi_{\sigma}(z(x, y), y)}{c_{\sigma, y-z(x, y)} d_{\sigma}} dx d\gamma(y).$$

In the right hand side of this inequality, the integration with respect to  $y \in I$  is transformed into an integration with respect to  $\xi = (\xi_1, \xi_2)^t \in \sigma$ , this yields (note that  $c_{\sigma, y-z(x, y)} = c_{\sigma, a-y}$ )

$$B \leq C_4 \sum_{\sigma \in \mathcal{E}_{\text{int}}} \frac{|D_{\sigma} u|^2}{d_{\sigma}} \int_{\omega} \int_{\sigma} \frac{\psi_{\sigma}(x, \xi)}{c_{I, a-y(\xi)}} \frac{|a-y(\xi)|}{|a-\xi|} dx d\gamma(\xi),$$

where  $y(\xi) = s\xi + (1-s)a$ , with  $s\xi_1 + (1-s)\alpha = 0$ , and where  $\psi_{\sigma}$  is defined by

$$\begin{aligned} \psi_{\sigma}(x, \xi) &= 1, \text{ if } y(\xi) \in I \text{ and } \xi_1 \leq x_1 \\ \psi_{\sigma}(x, \xi) &= 0, \text{ if } y(\xi) \notin I \text{ or } \xi_1 > x_1. \end{aligned}$$

Noting that  $c_{I, a-y(\xi)} \geq C_5 > 0$ , one deduces that

$$B \leq C_6 \sum_{\sigma \in \mathcal{E}_{\text{int}}} \frac{|D_{\sigma} u|^2}{d_{\sigma}} \int_{\sigma} \left( \int_{\omega} \psi_{\sigma}(x, \xi) \frac{|a-y(\xi)|}{|a-\xi|} dx \right) d\gamma(\xi) \leq C_7 |u|_{1,\mathcal{T}}^2,$$

with, for instance,  $C_7 = C_6(\text{diam}(\omega))^2$ . The bounds on  $A$  and  $B$  yield (51).

**Step 3 (proof of (46))**

Let us now prove that there exists  $D \in \mathbb{R}_+$ , only depending on  $\Omega$  such that (46) hold. Since  $\Omega$  is a polygonal set ( $d = 2$  or  $3$ ), there exists a finite number of disjoint convex polygonal sets, denoted by  $\{\Omega_1, \dots, \Omega_n\}$ , such that  $\overline{\Omega} = \cup_{i=1}^n \overline{\Omega_i}$ . Let  $I_{i,j} = \overline{\Omega_i} \cap \overline{\Omega_j}$ , and  $B$  be the set of couples  $(i, j) \in \{1, \dots, n\}^2$  such that  $i \neq j$  and the  $(d-1)$ -dimensional Lebesgue measure of  $I_{i,j}$ , denoted by  $m(I_{i,j})$ , is positive.

Let  $m_i$  denote the mean value of  $u$  on  $\Omega_i$ ,  $i \in \{1, \dots, n\}$ , and  $m_{i,j}$  denote the mean value of  $u$  on  $I_{i,j}$ ,  $(i, j) \in B$ . (For  $\sigma \in \mathcal{E}_{\text{int}}$ , in order that  $u$  be defined on  $\sigma$ , a.e. for the  $(d-1)$ -dimensional Lebesgue measure, let  $K \in \mathcal{T}$  be a control volume such that  $\sigma \in \mathcal{E}_K$ , one sets  $u = u_K$  on  $\sigma$ .) By definition,  $m_{i,j} = m_{j,i}$  for all  $(i, j) \in B$ .

Step 1 gives the existence of  $C_i$ ,  $i \in \{1, \dots, n\}$ , only depending on  $\Omega$  (since the  $\Omega_i$  only depend on  $\Omega$ ), such that

$$\|u - m_i\|_{L^2(\Omega_i)}^2 \leq C_i |u|_{1,\mathcal{T}}^2, \quad \forall i \in \{1, \dots, n\}, \quad (52)$$

Step 2 gives the existence of  $C_{i,j}$ ,  $i, j \in B$ , only depending on  $\Omega$ , such that

$$\|u - m_{i,j}\|_{L^2(\Omega_i)}^2 \leq C_{i,j} |u|_{1,\mathcal{T}}^2, \quad \forall (i, j) \in B.$$

Then, one has  $(m_i - m_{i,j})^2 m(\Omega_i) \leq 2(C_i + C_{i,j}) |u|_{1,\mathcal{T}}^2$ , for all  $(i, j) \in B$ . Since  $\Omega$  is connected, the above inequality yields the existence of  $M$ , only depending on  $\Omega$ , such that  $|m_i - m_j| \leq M |u|_{1,\mathcal{T}}$  for all  $(i, j) \in \{1, \dots, n\}^2$ , and therefore  $|m_\Omega(u) - m_i| \leq M |u|_{1,\mathcal{T}}$  for all  $i \in \{1, \dots, n\}$ . Then, (52) yields the existence of  $D$ , only depending on  $\Omega$ , such that (46) holds. This completes the proof of Lemma 6.  $\blacksquare$

**Lemma 7** *Let  $\mathcal{T}$  be an admissible finite volume mesh in the sense of Definition 1 and  $u$  be a function which is constant on each cell of  $\mathcal{T}$  and each edge of  $\mathcal{E}_{\text{ext}}$ , that is  $u(x) = u_K$  if  $x \in K$ ,  $K \in \mathcal{T}$  and  $u(x) = u_\sigma$  if  $x \in \sigma$ ,  $\sigma \in \mathcal{E}_{\text{ext}}$ . Let  $\Gamma \subset \partial\Omega$  such that its  $(d-1)$ -dimensional measure  $m(\Gamma) \neq 0$  and  $\mathcal{O} \subset \Omega$  such that its  $d$ -dimensional measure  $m(\mathcal{O}) \neq 0$ . Then there exists  $C$ , only depending on  $\Omega$ , such that*

$$\|u\|_{L^2(\Omega)}^2 \leq C \left( |u|_{1,\mathcal{T}}^2 + \|u\|_{L^2(\Gamma)}^2 \right), \quad (53)$$

and

$$\|u\|_{L^2(\partial\Omega)}^2 \leq C \left( |u|_{1,\mathcal{T}}^2 + \|u\|_{L^2(\mathcal{O})}^2 \right), \quad (54)$$

where  $|\cdot|_{1,\mathcal{T}}$  is the discrete  $H_0^1$  norm defined in Definition 4.

**Proof of Lemma 7**

We proceed in two steps. The first two steps deal with the proof of (53) while the third step deals with (54). The first step consists in proving (53) on a part of  $\Omega$  with a boundary containing  $\Gamma$ . In the second step we use a discrete trace inequality which is stated in Lemma 8 to conclude the proof of the announced result on  $\Omega$ .

**Step 1**

We can assume without loss of generality that  $\Gamma$  is included in a hyperplane of  $\mathbb{R}^d$ , indeed if it is not we can split  $\Gamma$  in several parts included in hyperplanes of  $\mathbb{R}^d$  since  $\Omega$  is polygonal if  $d = 2$  or polyhedral if  $d = 3$ . For  $x, y \in \mathbb{R}^d$ , one defines  $[x, y] = \{tx + (1-t)y ; t \in [0, 1]\}$ . Let us define

$$\mathcal{O}(\Gamma) = \left\{ x \in \Omega ; \exists y \in \Gamma \text{ such that } (x - y) \cdot y = 0 \text{ and } [x, y] \subset \Omega \right\}. \quad (55)$$

Then we choose a coordinate system such that a point  $y \in \Gamma$  has for coordinate  $(0, \tilde{y})$  with  $\tilde{y} \in I \subset \mathbb{R}^{d-1}$  and such that if we consider a point  $x \in \mathcal{O}(\Gamma)$  with  $x = (x_1, \tilde{x})$ ,  $\tilde{x} \in I$ , then  $x_1 > 0$ .

Let us denote by  $\eta$  the first unit vector of the coordinate system, so  $\eta = (1, 0)$  if  $d = 2$  and  $\eta = (1, 0, 0)$  if  $d = 3$ . For  $\sigma \in \mathcal{E}$ , we define  $\chi_\sigma$  from  $\mathbb{R}^d \times \mathbb{R}^d$  to  $\{0, 1\}$  by  $\chi_\sigma(x, y) = 1$  if  $\sigma \cap [x, y] \neq \emptyset$  and  $\chi_\sigma(x, y) = 0$  otherwise.

For  $\tilde{y} \in I$ , let denote by  $\mathcal{D}_{\tilde{y}, \eta}$  the semi-line defined by its origin  $(0, \tilde{y})$  and the vector  $\eta$ , and by  $\alpha(\tilde{y})$  the real such that  $(\alpha(\tilde{y}), \tilde{y}) \in \mathcal{D}_{\tilde{y}, \eta} \cap \partial\Omega$  and  $[(0, \tilde{y}), (\alpha(\tilde{y}), \tilde{y})] \subset \overline{\Omega}$ .

Let  $y = (0, \tilde{y}) \in \Gamma$  and  $x_1 \in ]0, \alpha(\tilde{y})[$ , then

$$|u(x_1, \tilde{y})| \leq |u(0, \tilde{y})| + \sum_{\sigma \in \mathcal{E}} (D_\sigma u) \chi_\sigma((\alpha(\tilde{y}), \tilde{y}), y),$$

where  $D_\sigma u$  is defined in Definition 4.

Using the Cauchy-Schwarz inequality, and setting  $c_\sigma = |\mathbf{n}_\sigma \cdot \eta|$  where  $\mathbf{n}_\sigma$  is a unit normal vector to  $\sigma$ , one gets

$$|u(x_1, \tilde{y})|^2 \leq 2|u(0, \tilde{y})|^2 + 2 \left( \sum_{\sigma \in \mathcal{E}} \frac{(D_\sigma u)^2}{d_\sigma c_\sigma} \chi_\sigma((\alpha(\tilde{y}), \tilde{y}), y) \right) \left( \sum_{\sigma \in \mathcal{E}} d_\sigma c_\sigma \chi_\sigma((\alpha(\tilde{y}), \tilde{y}), y) \right).$$

Remarking that  $\sum_{\sigma \in \mathcal{E}} d_\sigma c_\sigma \chi_\sigma((\alpha(\tilde{y}), \tilde{y}), y) \leq \text{diam}(\Omega)$  and integrating with respect to  $x_1$  and  $\tilde{y}$ , one obtains

$$\|u\|_{L^2(\mathcal{O}(\Gamma))}^2 \leq 2\|u\|_{L^2(\Gamma)}^2 + 2(\text{diam}(\Omega))^2 \int_I \sum_{\sigma \in \mathcal{E}} \frac{(D_\sigma u)^2}{d_\sigma c_\sigma} \chi_\sigma((\alpha(\tilde{y}), \tilde{y}), y) d\tilde{y}.$$

Since  $\int_I \chi_\sigma((\alpha(\tilde{y}), \tilde{y}), y) d\tilde{y} \leq m(\sigma) c_\sigma$ , one has

$$\|u\|_{L^2(\mathcal{O}(\Gamma))}^2 \leq C \left( \|u\|_{L^2(\Gamma)}^2 + \|u\|_{1,\mathcal{T}}^2 \right),$$

where  $C$  only depends on  $\Omega$ . This concludes the first step.

## Step 2 Proof of (53)

By compactness of the boundary of  $\partial\mathcal{O}(\Gamma)$  (where  $\mathcal{O}(\Gamma)$  is defined by (55) and  $\partial\mathcal{O}(\Gamma)$  denotes its boundary), there exists a finite number of hyperplanes of  $\mathbb{R}^d$ ,  $\{\Gamma_i, i = 1, \dots, N\}$ , such that  $\partial\mathcal{O}(\Gamma) \subset \cup_{i=1}^N \Gamma_i$  and  $\overline{\Gamma}_i \cap \overline{\Gamma}_j \subset \mathbb{R}^{d-2}$  for  $i, j \in \{1, \dots, N\}$ ,  $i \neq j$ .

Let  $j \in \{1, \dots, N\}$ , then, thanks to Lemma 8 which is stated and proved below, one has:

$$\|\gamma u\|_{L^2(\Gamma_j \cap \Omega)}^2 \leq C_1 \left( \|u\|_{L^2(\mathcal{O}(\Gamma))}^2 + \|u\|_{1,\mathcal{T}}^2 \right), \quad (56)$$

where  $\gamma u$  denotes the “discrete trace” of  $u$ , that is  $\gamma u = u_K$  for all  $x \in \sigma$  such that  $\sigma \in \mathcal{E}_{\text{ext}} \cap \mathcal{E}_K$  and  $C_1$  only depends on  $\Omega$ .

Let us define

$$\mathcal{O}(\Gamma_j \cap \Omega) = \left\{ x \in \Omega ; \exists y \in \Gamma_j \cap \Omega \text{ such that } (x - y) \cdot y = 0 \text{ and } [x, y] \subset \Omega \setminus \mathcal{O}(\Gamma) \right\}$$

Then applying the first step to  $\Gamma_j \cap \Omega$  instead of  $\Gamma$ , one gets

$$\|u\|_{L^2(\mathcal{O}(\Gamma_j \cap \Omega))}^2 \leq C_2 \left( \|u\|_{L^2(\Gamma_j \cap \Omega)}^2 + \|u\|_{1,\mathcal{T}}^2 \right),$$

where  $C_2$  only depends on  $\Omega$ .

Then using (56)

$$\|u\|_{L^2(\mathcal{O}(\Gamma_j \cap \Omega))}^2 \leq (C_2 + C_1 C_2) \left( \|u\|_{L^2(\mathcal{O}(\Gamma))}^2 + \|u\|_{1,\mathcal{T}}^2 \right),$$

and thanks to Step 1

$$\|u\|_{L^2(\mathcal{O}(\Gamma_j \cap \Omega))}^2 \leq C \left( \|u\|_{L^2(\Gamma)}^2 + \|u\|_{1,\mathcal{T}}^2 \right),$$

where  $C$  only depends on  $\Omega$ .

Iterating this process so long as a part of  $\Omega$  has not been reached, we obtain (53) where  $C$  only depends on  $\Omega$  which concludes the proof of (53).

### Step 3 Proof of (54)

Thanks to Lemma 8

$$\|u\|_{L^2(\partial\Omega)}^2 \leq C \left( \|u\|_{L^2(\Omega \setminus \mathcal{O})}^2 + \|u\|_{1,\mathcal{T}}^2 \right),$$

where  $C$  only depends on  $\Omega$ .

Let us denote by  $\partial\mathcal{O}$  the boundary of  $\mathcal{O}$ . We denote by  $\gamma u$  the discrete trace of  $u$  on  $\partial\mathcal{O}$ , that is for all  $x \in \partial\mathcal{O}$ , if  $x \in K$ ,  $K \in \mathcal{T}$  then  $\gamma u(x) = u_K$ , if  $x \in \sigma$ ,  $\sigma \in \mathcal{E}_{\text{ext}}$  then  $\gamma u(x) = u_\sigma$ , if  $x \in \sigma$ ,  $\sigma \in \mathcal{E}_{\text{int}}$ ,  $\sigma = K|L \subset \partial\mathcal{O}$  and  $K \subset \mathcal{O}$  then  $\gamma u(x) = u_L$ , finally if  $x \in \sigma$ ,  $\sigma \in \mathcal{E}_{\text{int}}$  and  $\sigma = K|L \not\subset \partial\mathcal{O}$  then  $\gamma u(x) = u_L$  or  $u_K$ . Using (53), one gets

$$\|u\|_{L^2(\Omega \setminus \mathcal{O})}^2 \leq C \left( \|\gamma u\|_{L^2(\partial\mathcal{O})}^2 + \|u\|_{1,\mathcal{T}}^2 \right),$$

where  $C$  only depends on  $\Omega$ .

Using Lemma 8 once more, one obtains

$$\|\gamma u\|_{L^2(\partial\mathcal{O})}^2 \leq C \left( \|u\|_{L^2(\mathcal{O})}^2 + \|u\|_{1,\mathcal{T}}^2 \right),$$

where  $C$  only depends on  $\Omega$ .

These three results yield (54). This concludes the proof of Lemma 7. ■

**Lemma 8 (Trace inequality)** *Let  $\Omega$  be an open bounded polygonal subset of  $\mathbb{R}^d$ ,  $d = 2$  or  $3$ . Let  $\mathcal{T}$  be an admissible mesh in the sense of Definition 1, and  $u$  be a function from  $\Omega$  to  $\mathbb{R}$  which is constant on each control volume of the mesh. Let  $u_K$  be the value of  $u$  in the control volume  $K$ . Let  $\gamma u$  be defined by  $\gamma u = u_K$  a.e. (for the  $(d-1)$ -dimensional Lebesgue measure) on  $\sigma$ , if  $\sigma \in \mathcal{E}_{\text{ext}}$  and  $\sigma \in \mathcal{E}_K$ . Then, there exists  $C$ , only depending on  $\Omega$ , such that*

$$\|\gamma u\|_{L^2(\partial\Omega)} \leq C(|u|_{1,\mathcal{T}} + \|u\|_{L^2(\Omega)}). \quad (57)$$

### Proof of Lemma 8

By compactness of the boundary of  $\partial\Omega$ , there exists a finite number of open hyper-rectangles ( $d = 2$  or  $3$ ),  $\{R_i, i = 1, \dots, N\}$ , and normalized vectors of  $\mathbb{R}^d$ ,  $\{\eta_i, i = 1, \dots, N\}$ , such that

$$\begin{cases} \partial\Omega \subset \cup_{i=1}^N R_i, \\ \eta_i \cdot \mathbf{n}(x) \geq \alpha > 0 \text{ for all } x \in R_i \cap \partial\Omega, i \in \{1, \dots, N\}, \\ \{x + t\eta_i, x \in R_i \cap \partial\Omega, t \in \mathbb{R}_+\} \cap R_i \subset \Omega, \end{cases}$$

where  $\alpha$  is some positive number and  $\mathbf{n}(x)$  is the normal vector to  $\partial\Omega$  at  $x$ , inward to  $\Omega$ . Let  $\{\alpha_i, i = 1, \dots, N\}$  be a family of functions such that  $\sum_{i=1}^N \alpha_i(x) = 1$ , for all  $x \in \partial\Omega$ ,  $\alpha_i \in C_c^\infty(\mathbb{R}^d, \mathbb{R}_+)$  and  $\alpha_i = 0$  outside of  $R_i$ , for all  $i = 1, \dots, N$ . Let  $\Gamma_i = R_i \cap \partial\Omega$ ; let us prove that there exists  $C_i$  only depending on  $\alpha$  and  $\alpha_i$  such that

$$\|\alpha_i \gamma u\|_{L^2(\Gamma_i)} \leq C_i (|u|_{1,\mathcal{T}} + \|u\|_{L^2(\Omega)}). \quad (58)$$

The existence of  $C$ , only depending on  $\Omega$ , such that (57) holds, follows easily (taking  $C = \sum_{i=1}^N C_i$ , and using  $\sum_{i=1}^N \alpha_i(x) = 1$ , note that  $\alpha$  and  $\alpha_i$  depend only on  $\Omega$ ). It remains to prove (58).

Let us introduce some notations. For  $\sigma \in \mathcal{E}$  and  $K \in \mathcal{T}$ , define  $\chi_\sigma$  and  $\chi_K$  from  $\mathbb{R}^d \times \mathbb{R}^d$  to  $\{0, 1\}$  by  $\chi_\sigma(x, y) = 1$ , if  $[x, y] \cap \sigma \neq \emptyset$ ,  $\chi_\sigma(x, y) = 0$ , if  $[x, y] \cap \sigma = \emptyset$ , and  $\chi_K(x, y) = 1$ , if  $[x, y] \cap K \neq \emptyset$ ,  $\chi_K(x, y) = 0$ , if  $[x, y] \cap K = \emptyset$ .

Let  $i \in \{1, \dots, N\}$  and let  $x \in \Gamma_i$ . There exists a unique  $t > 0$  such that  $x + t\eta_i \in \partial R_i$ , let  $y(x) = x + t\eta_i$ . For  $\sigma \in \mathcal{E}$ , let  $z_\sigma(x) = [x, y(x)] \cap \sigma$  if  $[x, y(x)] \cap \sigma \neq \emptyset$  and is reduced to one point. For  $K \in \mathcal{T}$ , let  $\xi_K(x), \eta_K(x)$  be such that  $[x, y(x)] \cap K = [\xi_K(x), \eta_K(x)]$  if  $[x, y(x)] \cap K \neq \emptyset$ .

One has, for a.e. (for the  $(d-1)$ -dimensional Lebesgue measure)  $x \in \Gamma_i$ ,

$$|\alpha_i \gamma u(x)| \leq \sum_{\sigma=K|L \in \mathcal{E}_{\text{int}}} |\alpha_i(z_\sigma(x))(u_K - u_L)| \chi_\sigma(x, y(x)) + \sum_{K \in \mathcal{T}} |(\alpha_i(\xi_K(x)) - \alpha_i(\eta_K(x)))u_K| \chi_K(x, y(x)),$$

that is,

$$|\alpha_i \gamma u(x)|^2 \leq A(x) + B(x) \quad (59)$$

with

$$A(x) = 2 \left( \sum_{\sigma=K|L \in \mathcal{E}_{\text{int}}} |\alpha_i(z_\sigma(x))(u_K - u_L)| \chi_\sigma(x, y(x)) \right)^2,$$

$$B(x) = 2 \left( \sum_{K \in \mathcal{T}} |(\alpha_i(\xi_K(x)) - \alpha_i(\eta_K(x)))u_K| \chi_K(x, y(x)) \right)^2.$$

A bound on  $A(x)$  is obtained for a.e.  $x \in \Gamma_i$ , by remarking that, from the Cauchy-Schwarz inequality:

$$A(x) \leq D_1 \sum_{\sigma \in \mathcal{E}_{\text{int}}} \frac{|D_\sigma u|^2}{d_\sigma c_\sigma} \chi_\sigma(x, y(x)) \sum_{\sigma \in \mathcal{E}_{\text{int}}} d_\sigma c_\sigma \chi_\sigma(x, y(x)),$$

where  $D_1$  only depends on  $\alpha_i$  and  $c_\sigma = |\eta_i \cdot \mathbf{n}_\sigma|$ . (Recall that  $D_\sigma u = |u_K - u_L|$ .) Since

$$\sum_{\sigma \in \mathcal{E}_{\text{int}}} d_\sigma c_\sigma \chi_\sigma(x, y(x)) \leq \text{diam}(\Omega),$$

this yields:

$$A(x) \leq \text{diam}(\Omega) D_1 \sum_{\sigma \in \mathcal{E}_{\text{int}}} \frac{|D_\sigma u|^2}{d_\sigma c_\sigma} \chi_\sigma(x, y(x)).$$

Then, since

$$\int_{\Gamma_i} \chi_\sigma(x, y(x)) d\gamma(x) \leq \frac{1}{\alpha} c_\sigma m(\sigma),$$

there exists  $D_2$ , only depending on  $\Omega$ , such that

$$A = \int_{\Gamma_i} A(x) d\gamma(x) \leq D_2 |u|_{1,\mathcal{T}}^2.$$

A bound  $B(x)$  for a.e.  $x \in \Gamma_i$  is obtained with the Cauchy-Schwarz inequality:

$$B(x) \leq D_3 \sum_{K \in \mathcal{T}} u_K^2 \chi_K(x, y(x)) |\xi_K(x) - \eta_K(x)| \sum_{K \in \mathcal{T}} |\xi_K(x) - \eta_K(x)| \chi_K(x, y(x)),$$

where  $D_3$  only depends on  $\alpha_i$ . Since

$$\sum_{K \in \mathcal{T}} |\xi_K(x) - \eta_K(x)| \chi_K(x, y(x)) \leq \text{diam}(\Omega) \text{ and } \int_{\Gamma_i} \chi_K(x, y(x)) |\xi_K(x) - \eta_K(x)| d\gamma(x) \leq \frac{1}{\alpha} m(K),$$

there exists  $D_4$ , only depending on  $\Omega$ , such that

$$B = \int_{\Gamma_i} B(x) d\gamma(x) \leq D_4 \|u\|_{L^2(\Omega)}^2.$$

Integrating (59) over  $\Gamma_i$ , the bounds on  $A$  and  $B$  lead (58) for some convenient  $C_i$  and it concludes the proof of Lemma 8.  $\blacksquare$

### 4.3 Error estimates

**Theorem 3 ( $C^2$  regularity)** Under Assumptions 1 and 5, let  $\mathcal{T}$  be an admissible mesh in the sense of Definition 1. One assumes that the unique variational solution  $u \in H^1(\Omega)$ , such that  $\int_{\Omega} u(x) dx = 0$ , of Problem (1), (39) satisfies  $u \in C^2(\bar{\Omega})$ . Let  $u_{\mathcal{T}}$  be the solution to (15), (40), (41), (42), (5), (8) and (11), such that  $\sum_{K \in \mathcal{T}} m(K) u_K = \sum_{K \in \mathcal{T}} m(K) u(x_K)$ , where  $x_K$  is defined in Definition 1. Let  $e_{\mathcal{T}}$  be defined by  $e_{\mathcal{T}}(x) = e_K = u(x_K) - u_K$  if  $x \in K$ ,  $K \in \mathcal{T}$ .

Then, there exists  $C > 0$  only depending on  $u$ ,  $\mathbf{v}$ ,  $d$  and  $\Omega$  such that

$$|e_{\mathcal{T}}|_{1,\mathcal{T}} \leq C \text{size}(\mathcal{T}), \quad (60)$$

where  $|\cdot|_{1,\mathcal{T}}$  is the discrete  $H^1$  semi-norm defined in Definition 3. Furthermore

$$\|e_{\mathcal{T}}\|_{L^2(\Omega)} \leq C \text{size}(\mathcal{T}). \quad (61)$$

**Theorem 4 ( $H^2$  regularity)** Under Assumptions 1 and 5, let  $\mathcal{T}$  be an admissible mesh in the sense of Definition 1 and let

$$\zeta = \min_{K \in \mathcal{T}} \min_{\sigma \in \mathcal{E}_K} \frac{d_{K,\sigma}}{\text{diam}(K)}.$$

One assumes that the unique variational solution  $u \in H^1(\Omega)$ , such that  $\int_{\Omega} u(x) dx = 0$ , of Problem (1), (39) satisfies  $u \in H^2(\Omega)$ . Let  $u_{\mathcal{T}}$  be the solution to (15), (40), (41), (42), (5), (8) and (11), such that  $\sum_{K \in \mathcal{T}} m(K) u_K = \sum_{K \in \mathcal{T}} m(K) u(x_K)$ , where  $x_K$  is defined in Definition 1. Let  $e_{\mathcal{T}}$  be defined by  $e_{\mathcal{T}}(x) = e_K = u(x_K) - u_K$  if  $x \in K$ ,  $K \in \mathcal{T}$ .

Then, there exists  $C$ , only depending on  $u$ ,  $\mathbf{v}$ ,  $d$ ,  $\Omega$  and  $\zeta$ , such that (60) and (61) hold.

**Proof** of Theorems 3 and 4

As  $u$  is the exact solution to (1), (39), one has:

$$\sum_{\sigma \in \mathcal{E}_K \cap \mathcal{E}_{\text{int}}} (\bar{F}_{K,\sigma} + \bar{V}_{K,\sigma}) = \int_K f(x) dx + \sum_{\sigma \in \mathcal{E}_K \cap \mathcal{E}_{\text{ext}}} g_{\sigma}^N,$$

where  $\bar{F}_{K,\sigma}$  and  $\bar{V}_{K,\sigma}$  are defined by (20).

Substracting (40) off the previous equation yields

$$\sum_{\sigma=K|L \in \mathcal{E}_{\text{int}}} (F_{K,\sigma}^* - F_{K,\sigma}) + \sum_{\sigma \in \mathcal{E}_K \cap \mathcal{E}_{\text{int}}} (V_{K,\sigma}^* - V_{K,\sigma}) = - \sum_{\sigma=K|L \in \mathcal{E}_{\text{int}}} m(\sigma) R_{K,\sigma} - \sum_{\sigma \in \mathcal{E}_K \cap \mathcal{E}_{\text{int}}} m(\sigma) r_{K,\sigma}, \quad (62)$$

where  $F_{K,\sigma}^*$  is defined by (21) and  $V_{K,\sigma}^* = v_{K,\sigma} u(x_{\sigma,+})$ ,  $\forall \sigma \in \mathcal{E}_K \cap \mathcal{E}_{\text{int}}$ ,  $\forall K \in \mathcal{T}$ , where  $x_{\sigma,+} = x_K$  (resp.  $x_L$ ) if  $\sigma \in \mathcal{E}_{\text{int}}$ ,  $\sigma = K|L$  and  $v_{K,\sigma} \geq 0$  (resp.  $v_{K,\sigma} \leq 0$ ), finally  $R_{K,\sigma}$  and  $r_{K,\sigma}$  are defined by (24).

Multiplying (62) by  $e_K$ , summing for  $K \in \mathcal{T}$  and noting that

$$\sum_{K \in \mathcal{T}} \sum_{\sigma=K|L \in \mathcal{E}_{\text{int}}} (F_{K,\sigma}^* - F_{K,\sigma}) e_K = \sum_{\sigma \in \mathcal{E}_{\text{int}}} |D_{\sigma} e|^2 \frac{m(\sigma)}{d_{\sigma}} = |e|_{1,\mathcal{T}}^2,$$

where  $|\cdot|_{1,\mathcal{T}}^2$  is defined in Definition 3, yield

$$|e_{\mathcal{T}}|_{1,\mathcal{T}}^2 + \sum_{K \in \mathcal{T}} \sum_{\sigma=K|L \in \mathcal{E}_{\text{int}}} v_{K,\sigma} e_{\sigma,+} e_K \leq - \sum_{K \in \mathcal{T}} \sum_{\sigma=K|L \in \mathcal{E}_{\text{int}}} m(\sigma) (R_{K,\sigma} + r_{K,\sigma}) e_K, \quad (63)$$

where  $e_{\sigma,+} = u(x_{\sigma,+}) - u_{\sigma,+}$ .

Reordering the summation over the set of edges, one has

$$\sum_{K \in \mathcal{T}} \sum_{\sigma=K|L \in \mathcal{E}_{\text{int}}} v_{K,\sigma} e_{\sigma,+} e_K = \sum_{\sigma \in \mathcal{E}_{\text{int}}} v_{\sigma} (e_{\sigma,+} - e_{\sigma,-}) e_{\sigma,+} = \frac{1}{2} \sum_{\sigma \in \mathcal{E}_{\text{int}}} v_{\sigma} ((e_{\sigma,+} - e_{\sigma,-})^2 + (e_{\sigma,+}^2 - e_{\sigma,-}^2)),$$

where, for all  $\sigma \in \mathcal{E}_{\text{int}}$ ,  $v_\sigma = |\int_\sigma \mathbf{v}(x) \cdot \mathbf{n} d\gamma(x)|$ ,  $\mathbf{n}$  being a unit normal vector to  $\sigma$ , and  $e_{\sigma,-}$  is the downstream value to  $\sigma$  with respect to  $\mathbf{v}$ , i.e. if  $\sigma = K|L$ , then  $e_{\sigma,-} = e_K$  if  $v_{K,\sigma} \leq 0$ , and  $e_{\sigma,-} = e_L$  otherwise. Thanks to the assumptions  $\operatorname{div} \mathbf{v} = 0$  on  $\Omega$  and  $\mathbf{v} \cdot \mathbf{n} = 0$  on  $\partial\Omega$ , one obtains

$$\begin{aligned} \frac{1}{2} \sum_{\sigma \in \mathcal{E}_{\text{int}}} v_\sigma (e_{\sigma,+}^2 - e_{\sigma,-}^2) &= \frac{1}{2} \sum_{\sigma \in \mathcal{E}_{\text{int}}} v_\sigma (e_{\sigma,+}^2 - e_{\sigma,-}^2) + \frac{1}{2} \sum_{K \in \mathcal{T}} \sum_{\sigma \in \mathcal{E}_{\text{ext}}} v_{K,\sigma} |e_K|^2 \\ &= \frac{1}{2} \sum_{K \in \mathcal{T}} \left( \int_{\partial K} \mathbf{v}(x) \cdot \mathbf{n}_K d\gamma(x) \right) e_K^2 = \frac{1}{2} \int_{\Omega} (\operatorname{div} \mathbf{v}(x)) e_{\mathcal{T}}^2(x) dx = 0. \end{aligned}$$

Hence, (63) yields

$$|e_{\mathcal{T}}|_{1,\mathcal{T}}^2 \leq - \sum_{K \in \mathcal{T}} \sum_{\sigma \in \mathcal{E}_K} m(\sigma) (R_{K,\sigma} + r_{K,\sigma}) e_K. \quad (64)$$

Thanks to the conservativity property of the scheme (see (4)), one has  $R_{K,\sigma} = -R_{L,\sigma}$  and  $r_{K,\sigma} = -r_{L,\sigma}$  for  $\sigma \in \mathcal{E}_{\text{int}}$  such that  $\sigma = K|L$ . Let  $R_\sigma = |R_{K,\sigma}|$  and  $r_\sigma = |r_{K,\sigma}|$  if  $\sigma \in \mathcal{E}_K$ . Reordering the summation over the edges and using Young's inequality, one then obtains

$$\begin{aligned} \left| \sum_{K \in \mathcal{T}} \sum_{\sigma=K|L \in \mathcal{E}_{\text{int}}} m(\sigma) (R_{K,\sigma} + r_{K,\sigma}) e_K \right| &\leq \sum_{\sigma \in \mathcal{E}_{\text{int}}} m(\sigma) (D_\sigma e) (R_\sigma + r_\sigma) \leq \\ &\leq \frac{1}{2} \sum_{\sigma \in \mathcal{E}_{\text{int}}} \frac{m(\sigma)}{d_\sigma} (D_\sigma e)^2 + \frac{1}{2} \sum_{\sigma \in \mathcal{E}_{\text{int}}} m(\sigma) d_\sigma (R_\sigma + r_\sigma)^2. \quad (65) \end{aligned}$$

Using Lemma 2, if  $u \in C^2(\overline{\Omega})$ , or Hölder's inequality and Lemma 3 (with  $p = 4$ ), if  $u$  is only in  $H^2(\Omega)$  (for more details see the proof of inequality (36)), and remarking that  $\sum_{\sigma \in \mathcal{E}} m(\sigma) d_\sigma = d m(\Omega)$ , (64) and (65) yield the existence of  $C$ , only depending on  $u$ ,  $\mathbf{v}$ ,  $d$  and  $\Omega$  if  $u \in C^2(\overline{\Omega})$  and on  $u$ ,  $\mathbf{v}$ ,  $d$ ,  $\Omega$  and  $\zeta$  if  $u$  is only in  $H^2(\Omega)$ , such that

$$|e_{\mathcal{T}}|_{1,\mathcal{T}}^2 \leq C (\operatorname{size}(\mathcal{T}))^2.$$

This estimate gives (60). In order to obtain (61), we use a discrete Poincaré-Wirtinger inequality which is given in Lemma 6. This concludes the proofs of Theorems 3 and 4. ■

## 5 Robin boundary condition

The last type of boundary condition we consider is a Robin condition:

$$\nabla u(x) \cdot \mathbf{n}(x) + \lambda(x) u(x) = g^F(x), \quad x \in \partial\Omega, \quad (66)$$

with

**Assumption 6**  $g^F \in H^{1/2}(\partial\Omega)$ ,  $\lambda \in L^\infty(\partial\Omega)$  such that  $\mathbf{v} \cdot \mathbf{n}/2 + \lambda \geq 0$  a.e. on  $\partial\Omega$ . Furthermore, if  $\mathbf{v}(x) \cdot \mathbf{n}(x)/2 + \lambda(x) = 0$  for almost every  $x \in \partial\Omega$  then one assumes the existence of  $\mathcal{O} \subset \overline{\Omega}$  such that its  $d$ -dimensional measure  $m(\mathcal{O}) \neq 0$  and such that  $\operatorname{div}(\mathbf{v})/2 + b \neq 0$  a.e. on  $\mathcal{O}$ .

Then, under Assumptions 1 and 6 the Lax-Milgram theorem ensures the existence of a unique variational solution  $u \in H^1(\Omega)$  of (1), (66). That is to say  $u \in H^1(\Omega)$  satisfies

$$\begin{aligned} \int_{\Omega} [\nabla u(x) \cdot \nabla \phi(x) + \operatorname{div}(\mathbf{v}(x) u(x)) \phi(x) + b(x) u(x) \phi(x)] dx + \int_{\partial\Omega} \lambda(x) \bar{\gamma}(u(x)) \bar{\gamma}(\phi)(x) d\gamma(x) \\ = \int_{\partial\Omega} g^F(x) \bar{\gamma}(\phi)(x) d\gamma(x) + \int_{\Omega} f(x) \phi(x) dx, \quad \text{for all } \phi \in H^1(\Omega), \end{aligned}$$

where  $\bar{\gamma}$  denotes the trace operator from  $H^1(\Omega)$  into  $H^{1/2}(\partial\Omega)$  and  $d\gamma$  is the integration symbol for the  $(d-1)$ -dimensional Lebesgue measure on  $\partial\Omega$ .

**Remark 9** Assumptions 1 and 6 give the coercivity of the elliptic operator associated to the above variational equality. It does not need a compatibility relation, even in the case  $\lambda = 0$  a.e.; in this last case, even though the boundary condition looks like a Neumann condition, the solution behaves as if the problem were a Robin condition and the proof of the error estimate is the same as for a Robin condition. The case  $\lambda = 0$  under assumptions 1 and 6 is therefore treated in this section.

## 5.1 Discretization

Let  $\mathcal{T}$  be an admissible mesh in the sense of Definition 1. The discretization of the diffusion-convection equation (1) with a Robin boundary condition is performed with the help of some auxiliary unknowns which are defined on the edges of the boundary. These may be eliminated when solving the linear system. We shall however keep them throughout our study because they simplify several expressions in the error estimate. Hence in this section the discrete unknowns are  $(u_K)_{K \in \mathcal{T}} \cup (u_\sigma)_{\sigma \in \mathcal{E}_{\text{ext}}}$ . In order to obtain the discretized equation let us as usual integrate (1) on each cell of the mesh. Using a “four points” finite volume scheme for the diffusion terms and an upstream scheme for the convection terms, one gets, for all  $K \in \mathcal{T}$ ,

$$\sum_{\sigma \in \mathcal{E}_K} \left[ F_{K,\sigma} + v_{K,\sigma} u_{\sigma,+} \right] + b_K m(K) u_K = m(K) f_K, \quad (67)$$

where, for all  $K \in \mathcal{T}$  and all  $\sigma \in \mathcal{E}_K$ ,  $v_{K,\sigma}$ ,  $f_K$  and  $b_K$  are defined by (8) and (11),  $u_{\sigma,+}$  is defined by (9). Furthermore  $F_{K,\sigma}$  is defined by (5) if  $\sigma \in \mathcal{E}_{\text{int}}$ ,  $\sigma = K|L$ , and by (6) if  $\sigma \in \mathcal{E}_{\text{ext}}$ , and we set for all  $\sigma \in \mathcal{E}_{\text{ext}}$

$$g_\sigma^F = \int_\sigma g^F(x) d\gamma(x) \quad \text{and} \quad \lambda_\sigma = \frac{1}{m(\sigma)} \int_\sigma \lambda(x) d\gamma(x). \quad (68)$$

There remains to give the equations associated with the boundary unknowns  $(u_\sigma)_{\sigma \in \mathcal{E}_{\text{ext}}}$ . These are obtained by discretizing (66). The discretization which we choose involves the upstream value  $u_{\sigma,+}$  in order for the scheme to be well defined with no additional condition on the mesh (see remarks 10 and 11). It writes:

$$-F_{K,\sigma} + (m(\sigma) \lambda_\sigma + v_{K,\sigma}) u_\sigma - v_{K,\sigma} u_{\sigma,+} = g_\sigma^F, \quad \text{for all } \sigma \in \mathcal{E}_K \cap \mathcal{E}_{\text{ext}}, \quad (69)$$

### Remark 10

- Using (6) and (69), one can eliminate  $u_\sigma$  for all  $\sigma \in \mathcal{E}_{\text{ext}}$  in (67) and obtain

$$u_\sigma = \frac{\left( (v_{K,\sigma} \top 0) d_{K,\sigma} + m(\sigma) \right) u_K + d_{K,\sigma} g_\sigma^F}{m(\sigma) + \left( m(\sigma) \lambda_\sigma + v_{K,\sigma} - (v_{K,\sigma} \perp 0) \right) d_{K,\sigma}},$$

where for all  $a, b \in \mathbb{R}$ ,  $a \top b = \max(a, b)$  and  $a \perp b = \min(a, b)$ . Again, the numerical unknowns are  $(u_K)_{K \in \mathcal{T}}$ .

- In order to discretize the boundary condition on an edge  $\sigma \in \mathcal{E}_{\text{ext}}$  of  $K \in \mathcal{T}$ , we use a non centered scheme summing and subtracting  $v_{K,\sigma}$ . This choice is performed, even though to our knowledge there is no physical background to this choice, in order to prove existence, uniqueness and convergence towards the exact solution, with no restriction on the mesh (see Remark 11), for  $\lambda$  such that there exists a subset of  $\partial\Omega$  with a non zero  $(d-1)$ -dimensional measure and such that  $\lambda < 0$  on this subset. In fact, it would be more natural to discretize the boundary condition as follows:

$$-F_{K,\sigma} + \lambda_\sigma m(\sigma) u_\sigma = g_\sigma^F, \quad \forall \sigma \in \mathcal{E}_K \cap \mathcal{E}_{\text{ext}}, \quad \forall K \in \mathcal{T}. \quad (70)$$

We shall give the idea of the proof for this scheme in Remarks 11 and 12 and see that for negative values of  $\lambda$  the convergence proof requires further assumptions on the mesh. Hence, (69) will be preferred for the discretization of the boundary condition so as to be able to handle negative values of  $\lambda$  with no additional condition on the mesh.

## 5.2 Existence, uniqueness and stability of the approximate solution

Let us first introduce as in the previous sections the discrete  $H^1$  norm of a function which is constant on each cell of the mesh and each edge on the boundary.

**Definition 4 (Discrete  $H^1$  norm)** Let  $\mathcal{T}$  be an admissible finite volume mesh in the sense of Definition 1. Let  $u$  be a function which is constant on each control volume of  $\mathcal{T}$  and on each edge on the boundary with  $u(x) = u_K$  if  $x \in K$ ,  $K \in \mathcal{T}$ , and  $u(x) = u_\sigma$  if  $x \in \sigma$ ,  $\sigma \in \mathcal{E}_{\text{ext}}$ , one defines the discrete  $H^1$  semi-norm by

$$|u|_{1,\mathcal{T}} = \left( \sum_{\sigma \in \mathcal{E}} \frac{m(\sigma)}{d_\sigma} (D_\sigma u)^2 \right)^{1/2},$$

where  $D_\sigma u = |u_K - u_L|$  if  $\sigma = K|L$  and  $D_\sigma u = |u_K - u_\sigma|$  if  $\sigma \in \mathcal{E}_{\text{ext}} \cap \mathcal{E}_K$ ,  $K \in \mathcal{T}$ .

**Proposition 3** Under Assumptions 1 and 6, let  $\mathcal{T}$  be an admissible mesh in the sense of Definition 1. Then there exists a unique solution  $(u_K)_{K \in \mathcal{T}} \cup (u_\sigma)_{\sigma \in \mathcal{E}_{\text{ext}}}$  to (67), (69), (68), (5), (6), (8), (9) and (11). Furthermore let  $u_\mathcal{T}$  be defined a.e. from  $\bar{\Omega}$  to  $\mathbb{R}$  by  $u_\mathcal{T}(x) = u_K$  if  $x \in K$ ,  $K \in \mathcal{T}$  and  $u_\mathcal{T}(x) = u_\sigma$  if  $x \in \sigma$ ,  $\sigma \in \mathcal{E}_{\text{ext}}$ ; then there exists  $C \in \mathbb{R}_+$  depending only on  $\Omega$  such that

$$|u_\mathcal{T}|_{1,\mathcal{T}} \leq C \left( \|g^F\|_{L^2(\partial\Omega)} + \|f\|_{L^2(\Omega)} \right), \quad (71)$$

and

$$\|u_\mathcal{T}\|_{L^2(\partial\Omega)} + \|u_\mathcal{T}\|_{L^2(\Omega)} \leq C \left( \|g^F\|_{L^2(\partial\Omega)} + \|f\|_{L^2(\Omega)} \right). \quad (72)$$

### Proof of Proposition 3

Let  $K \in \mathcal{T}$  and  $\sigma \in \mathcal{E}_{\text{ext}} \cap \mathcal{E}_K$ , then we multiply (67) by  $u_K$  and (69) by  $u_\sigma$ ; summing the results, we get

$$\begin{aligned} & \sum_{\substack{\sigma \in \mathcal{E}_K \cap \mathcal{E}_{\text{int}} \\ \sigma = K|L}} \frac{u_K - u_L}{d_{K|L}} u_K m(\sigma) + \sum_{\sigma \in \mathcal{E}_K} v_{K,\sigma} u_{\sigma,+} u_K + b_K m(K) (u_K)^2 \\ & + \sum_{\sigma \in \mathcal{E}_K \cap \mathcal{E}_{\text{ext}}} m(\sigma) \left( \frac{u_K - u_\sigma}{d_{K,\sigma}} u_K + \frac{u_\sigma - u_K}{d_{K,\sigma}} u_\sigma + \left( \lambda_\sigma + \frac{1}{m(\sigma)} v_{K,\sigma} \right) (u_\sigma)^2 - \frac{1}{m(\sigma)} v_{K,\sigma} u_{\sigma,+} u_\sigma \right) \\ & = \sum_{K \in \mathcal{T}} m(K) u_K f_K + \sum_{\sigma \in \mathcal{E}_{\text{ext}}} u_\sigma g_\sigma^F. \end{aligned}$$

Summing the result over  $K \in \mathcal{T}$ , using (44), (45) and Young's inequality, one gets for all  $\delta > 0$  and all  $\varepsilon > 0$

$$\begin{aligned} & |u_\mathcal{T}|_{1,\mathcal{T}}^2 + \int_\Omega \left( \frac{\text{div}(\mathbf{v}(x))}{2} + b(x) \right) (u_\mathcal{T}(x))^2 dx \\ & + \sum_{K \in \mathcal{T}} \sum_{\sigma \in \mathcal{E}_K \cap \mathcal{E}_{\text{ext}}} \left[ v_{K,\sigma} \left( u_{\sigma,+} u_K - \frac{(u_K)^2}{2} - u_{\sigma,+} u_\sigma + \frac{(u_\sigma)^2}{2} \right) + \left( \lambda_\sigma m(\sigma) + \frac{v_{K,\sigma}}{2} \right) (u_\sigma)^2 \right] \\ & \leq \frac{2}{\delta} \|f\|_{L^2(\Omega)}^2 + \frac{2}{\varepsilon} \|g^F\|_{L^2(\partial\Omega)}^2 + \frac{\delta}{2} \|u_\mathcal{T}\|_{L^2(\Omega)}^2 + \frac{\varepsilon}{2} \|u_\mathcal{T}\|_{L^2(\partial\Omega)}^2. \end{aligned}$$

Remarking that for all  $K \in \mathcal{T}$  and all  $\sigma \in \mathcal{E}_{\text{ext}} \cap \mathcal{E}_K$ ,  $v_{K,\sigma} \left( u_{\sigma,+} u_K - \frac{(u_K)^2}{2} - u_{\sigma,+} u_\sigma + \frac{(u_\sigma)^2}{2} \right) \geq 0$ , hence:

$$\begin{aligned} & |u_\mathcal{T}|_{1,\mathcal{T}}^2 + \int_\Omega \left( \frac{\text{div}(\mathbf{v}(x))}{2} + b(x) \right) (u_\mathcal{T}(x))^2 dx + \int_{\partial\Omega} \left( \lambda(x) + \frac{\mathbf{v}(x) \cdot \mathbf{n}(x)}{2} \right) (u_\mathcal{T}(x))^2 dx \\ & \leq \frac{2}{\delta} \|f\|_{L^2(\Omega)}^2 + \frac{2}{\varepsilon} \|g^F\|_{L^2(\partial\Omega)}^2 + \frac{\delta}{2} \|u_\mathcal{T}\|_{L^2(\Omega)}^2 + \frac{\varepsilon}{2} \|u_\mathcal{T}\|_{L^2(\partial\Omega)}^2, \end{aligned} \quad (73)$$

for all  $\delta > 0$  and all  $\varepsilon > 0$ .

Hence if there exists  $\delta > 0$  and  $\varepsilon > 0$  such that  $\operatorname{div}(\mathbf{v})/2 + b > \delta/2$  a.e. on  $\Omega$  and  $\mathbf{v} \cdot \mathbf{n}/2 + \lambda > \varepsilon/2$  a.e. on  $\partial\Omega$ , (73) gives (71). Otherwise, thanks to Assumption 6 there are two cases.

The first one is when there exists  $\Gamma \subset \partial\Omega$  such that its  $(d-1)$ -dimensional measure  $m(\Gamma) \neq 0$  and such that  $i_b = \inf_{x \in \Gamma} (\mathbf{v}(x) \cdot \mathbf{n}(x)/2 + \lambda(x)) \neq 0$ .

The second one is when there exists  $\mathcal{O} \subset \Omega$  such that its  $d$ -dimensional measure  $m(\mathcal{O}) \neq 0$  and such that  $i_i = \inf_{x \in \mathcal{O}} (\operatorname{div}(\mathbf{v}(x))/2 + b(x)) \neq 0$ .

In both cases one uses Lemma 7 in (73).

In the first case, one obtains

$$|u_\mathcal{T}|_{1,\mathcal{T}}^2 + i_b \|u_\mathcal{T}\|_{L^2(\Gamma)}^2 \leq \frac{2}{\delta} \|f\|_{L^2(\Omega)}^2 + \frac{2}{\varepsilon} \|g^F\|_{L^2(\partial\Omega)}^2 + C_\Omega \left( \frac{\delta}{2} + \frac{\varepsilon}{2} (1 + C_\Omega) \right) \left( \|u_\mathcal{T}\|_{L^2(\Gamma)}^2 + |u_\mathcal{T}|_{1,\mathcal{T}}^2 \right),$$

for all  $\delta > 0$ , all  $\varepsilon > 0$  and where  $C_\Omega$  only depends on  $\Omega$ .

In the second case, one gets

$$|u_\mathcal{T}|_{1,\mathcal{T}}^2 + i_i \|u_\mathcal{T}\|_{L^2(\mathcal{O})}^2 \leq \frac{2}{\delta} \|f\|_{L^2(\Omega)}^2 + \frac{2}{\varepsilon} \|g^F\|_{L^2(\partial\Omega)}^2 + C_\Omega \left( \frac{\varepsilon}{2} + \frac{\delta}{2} (1 + C_\Omega) \right) \left( \|u_\mathcal{T}\|_{L^2(\mathcal{O})}^2 + |u_\mathcal{T}|_{1,\mathcal{T}}^2 \right),$$

for all  $\delta > 0$ , all  $\varepsilon > 0$  and where  $C_\Omega$  only depends on  $\Omega$ .

Then a well adapted choice of  $\delta$  and  $\varepsilon$  gives in both cases (71). And using once more Lemma 7 gives (72).

Now let us assume  $f = 0$  on  $\Omega$  and  $g = 0$  on  $\partial\Omega$  then, thanks to (72),  $u_K = 0$  for all  $K \in \mathcal{T}$  and  $u_\sigma = 0$  for all  $\sigma \in \mathcal{E}_{\text{ext}}$ . This proves uniqueness and therefore existence since the dimension of the space is finite (equal to the number of discrete unknowns). ■

**Remark 11** If the discretization (70) is used instead of (69), remarking that

$$\begin{aligned} \sum_{K \in \mathcal{T}} \sum_{\sigma \in \mathcal{E}_K \cap \mathcal{E}_{\text{ext}}} & \left( v_{K,\sigma} u_{\sigma,+} - v_{K,\sigma} \frac{(u_K)^2}{2} + m(\sigma) \lambda_\sigma (u_\sigma)^2 \right) \\ &= \sum_{K \in \mathcal{T}} \left[ \sum_{\sigma \in \mathcal{E}_K \cap \mathcal{E}_{\text{ext}}} \left[ \frac{(u_K - u_\sigma)^2}{2} |v_{K,\sigma}| + \left( m(\sigma) \lambda_\sigma + \frac{v_{K,\sigma}}{2} \right) (u_\sigma)^2 \right] + \sum_{\substack{\sigma \in \mathcal{E}_K \cap \mathcal{E}_{\text{ext}} \\ v_{K,\sigma} \geq 0}} u_\sigma (u_K - u_\sigma) v_{K,\sigma} \right], \end{aligned}$$

computations similar to those of the above proof yield:

$$\begin{aligned} |u_\mathcal{T}|_{1,\mathcal{T}}^2 + \int_{\Omega} \left( \frac{\operatorname{div}(\mathbf{v}(x))}{2} + b \right) (u_\mathcal{T}(x))^2 dx \\ + \sum_{K \in \mathcal{T}} \sum_{\sigma \in \mathcal{E}_K \cap \mathcal{E}_{\text{ext}}} \left[ \frac{(u_K - u_\sigma)^2}{2} |v_{K,\sigma}| + \left( m(\sigma) \lambda_\sigma + \frac{v_{K,\sigma}}{2} \right) (u_\sigma)^2 \right] + \sum_{\substack{\sigma \in \mathcal{E}_K \cap \mathcal{E}_{\text{ext}} \\ v_{K,\sigma} \geq 0}} \lambda_\sigma d_{K,\sigma} m(\sigma) (u_\sigma)^2 \\ \leq \frac{2}{\delta} \|f\|_{L^2(\Omega)}^2 + \frac{2}{\varepsilon} \|g^F\|_{L^2(\partial\Omega)}^2 + \frac{\delta}{2} \|u_\mathcal{T}\|_{L^2(\Omega)}^2 + \frac{\varepsilon}{2} \|u_\mathcal{T}\|_{L^2(\partial\Omega)}^2, \end{aligned}$$

for all  $\delta > 0$  and all  $\varepsilon > 0$ . So if  $\lambda \geq 0$  a.e. on  $\partial\Omega$ , this inequality gives Proposition 3, otherwise one must assume some more restrictive assumption on the mesh as already mentionned in Remark 10; for instance one might assume  $m(\sigma) \lambda_\sigma + \frac{1}{2} v_{K,\sigma} + \lambda_\sigma d_{K,\sigma} m(\sigma) \geq 0$  if  $v_{K,\sigma} \geq 0$ .

We may now define the approximate solution by

$$\begin{cases} u_\mathcal{T}(x) = u_K & \text{if } x \in K, K \in \mathcal{T}, \\ u_\mathcal{T}(x) = u_\sigma & \text{if } x \in \sigma, \sigma \in \mathcal{E}_{\text{ext}}. \end{cases} \quad (74)$$

### 5.3 Error estimate

We prove in this section an error estimate in a discrete  $H^1$  semi-norm assuming  $u \in C^2(\overline{\Omega})$  or  $u \in H^2(\Omega)$  (with more restrictive assumptions on the mesh in the latter case).

**Theorem 5 ( $C^2$  regularity)** *Under Assumptions 1 and 6, let  $\mathcal{T}$  be a restricted admissible mesh in the sense of Definition 1. Let  $u_{\mathcal{T}}$  be the solution to (74), (67), (69), (68), (5), (6), (8) and (11).*

*Assume that the unique variational solution  $u$  of Problem (1), (66) satisfies  $u \in C^2(\overline{\Omega})$ . Let  $e_{\mathcal{T}}$  be defined by  $e_{\mathcal{T}}(x) = e_K = u(x_K) - u_K$  if  $x \in K$ ,  $K \in \mathcal{T}$  and  $e_{\mathcal{T}}(x) = e_{\sigma} = u(y_{\sigma}) - u_{\sigma}$  if  $x \in \sigma$ ,  $\sigma \in \mathcal{E}_{\text{ext}}$ .*

*Then, there exists  $C$  only depending on  $d$ ,  $u$ ,  $\mathbf{v}$ ,  $b$ ,  $\lambda$  and  $\Omega$  such that*

$$|e_{\mathcal{T}}|_{1,\mathcal{T}} \leq C \text{size}(\mathcal{T}), \quad (75)$$

where  $|\cdot|_{1,\mathcal{T}}$  is the discrete  $H_0^1$  norm defined in Definition 4.

Furthermore

$$\|e_{\mathcal{T}}\|_{L^2(\Omega)} + \|e_{\mathcal{T}}\|_{L^2(\partial\Omega)} \leq C \text{size}(\mathcal{T}). \quad (76)$$

One proves a similar result when  $u$  is only in  $H^2(\Omega)$ , assuming more restrictive hypotheses on  $\mathcal{T}$ .

**Theorem 6 ( $H^2$  regularity)** *Under Assumptions 1 and 6, let  $\mathcal{T}$  be a restricted admissible mesh in the sense of Definition 1 and let*

$$\zeta = \min_{K \in \mathcal{T}} \min_{\sigma \in \mathcal{E}_K} \frac{d_{K,\sigma}}{\text{diam}(K)}.$$

*Let  $u_{\mathcal{T}}$  be the solution to (74), (67), (69), (68), (5), (6), (8) and (11). Assume that the unique variational solution  $u$  of Problem (1), (66) satisfies  $u$  belongs to  $H^2(\Omega)$ . Let  $e_{\mathcal{T}}$  be defined by  $e_{\mathcal{T}}(x) = e_K = u(x_K) - u_K$  if  $x \in K$ ,  $K \in \mathcal{T}$  and  $e_{\mathcal{T}}(x) = e_{\sigma} = u(y_{\sigma}) - u_{\sigma}$  if  $x \in \sigma$ ,  $\sigma \in \mathcal{E}_{\text{ext}}$ .*

*Then, there exists  $C$ , only depending on  $u$ ,  $\mathbf{v}$ ,  $b$ ,  $\lambda$ ,  $\Omega$  and  $\zeta$ , such that (75) and (76) hold.*

#### Proof of Theorems 5 and 6

One proceeds, like in the Dirichlet case, in two steps. In the first one, one proves the consistency of the scheme, in a finite volume sense. Then in the second step, using this result and the conservativity of the scheme (see (4)), one proves error estimates.

#### Step 1

For all  $K \in \mathcal{T}$  and all  $\sigma \in \mathcal{E}_K$  let

$$\rho_K = \frac{1}{m(K)} \int_K b(x) (u(x) - u(x_K)) dx \quad \text{and} \quad r_{K,\sigma} = \frac{1}{m(\sigma)} \int_{\sigma} \mathbf{v}(x) \cdot \mathbf{n}_{K,\sigma} (u(x) - u(x_{\sigma,+})) d\gamma(x),$$

with  $x_K$  defined in Definition 1 and  $x_{\sigma,+} = x_K$  if  $v_{K,\sigma} \geq 0$ ,  $x_{\sigma,+} = x_L$  if  $v_{K,\sigma} < 0$  and  $\sigma \in \mathcal{E}_K \cap \mathcal{E}_{\text{int}}$ ,  $\sigma = K|L$ , finally  $x_{\sigma,+} = x_{\sigma}$  if  $v_{K,\sigma} < 0$  and  $\sigma \in \mathcal{E}_K \cap \mathcal{E}_{\text{ext}}$ .

Furthermore, if  $\sigma \in \mathcal{E}_K \cap \mathcal{E}_{\text{int}}$ ,  $\sigma = K|L$ , one has

$$R_{K,\sigma} = -\frac{1}{m(\sigma)} \int_{\sigma} \left( \nabla u(x) \cdot \mathbf{n}_{K,\sigma} - \frac{u(x_L) - u(x_K)}{d_{K,\sigma}} \right) d\gamma(x),$$

and, if  $\sigma \in \mathcal{E}_K \cap \mathcal{E}_{\text{ext}}$   $x_L$  is replaced by  $y_{\sigma}$  where  $y_{\sigma}$  is defined in Definition 1.

In a same way, one uses, for all  $\sigma \in \mathcal{E}_K \cap \mathcal{E}_{\text{ext}}$

$$\tilde{R}_{K,\sigma} = \frac{1}{m(\sigma)} \int_{\sigma} (\lambda(x) + \mathbf{v}(x) \cdot \mathbf{n}_{K,\sigma}) (u(x) - u(y_{\sigma})) d\gamma(x), \quad (77)$$

One recalls that Lemmas 2 and 3 hold. Moreover, using Taylor expansions, one proves the following result:

**Lemma 9** Under Assumptions 1 and 6, let  $\mathcal{T}$  be an admissible mesh in the sense of Definition 1. Assume that the unique variational solution  $u$  of Problem (1), (66) satisfies  $u \in C^2(\bar{\Omega})$ . Then there exists  $C > 0$ , only depending on  $u$ ,  $\mathbf{v}$  and  $\lambda$ , such that

$$|\tilde{R}_{K,\sigma}| \leq C \operatorname{size}(\mathcal{T}),$$

for any  $K \in \mathcal{T}$  and  $\sigma \in \mathcal{E}_K$ .

With the same technique as was used for  $r_{K,\sigma}$  in Lemma 3 we prove a similar lemma when  $u$  is only in  $H^2(\Omega)$ :

**Lemma 10** Under Assumptions 1 and 6, let  $\mathcal{T}$  be an admissible mesh in the sense of Definition 1 and let

$$\zeta = \min_{K \in \mathcal{T}} \min_{\sigma \in \mathcal{E}_K} \frac{d_{K,\sigma}}{\operatorname{diam}(K)}.$$

Assume that the unique variational solution,  $u$ , to (1), (66) belongs to  $H^2(\Omega)$ . Then there exists  $C$ , only depending on  $\lambda$ ,  $d$ ,  $\mathbf{v}$ ,  $\zeta$  and  $p$  such that, for all  $K \in \mathcal{T}$  and all  $\sigma \in \mathcal{E}_K$ ,

$$|\tilde{R}_{K,\sigma}| \leq C \operatorname{size}(\mathcal{T}) (\mathbf{m}(\sigma) d_\sigma)^{-1/p} \|u\|_{W^{1,p}(\mathcal{V}_\sigma)}, \quad (78)$$

for all  $p > d$  and such that  $p < +\infty$  if  $d = 2$  and  $p \leq 6$  if  $d = 3$ , where  $\mathcal{V}_\sigma$  is defined in Lemma 3.

This concludes the proof of the scheme consistency in a finite volume sense, i.e. step 1.

## Step 2

Let  $K \in \mathcal{T}$ , since  $u$  is the exact solution to (1), (66), one has:

$$\sum_{\sigma \in \mathcal{E}_K} \int_\sigma (-\nabla u(x) \cdot \mathbf{n}_{K,\sigma} + \mathbf{v}(x) \cdot \mathbf{n}_{K,\sigma} u(x)) d\gamma(x) + \int_K b(x) u(x) dx = \int_K f(x) dx,$$

Substracting (67) off the previous equation, one gets for all  $K \in \mathcal{T}$

$$\begin{aligned} \sum_{\substack{\sigma \in \mathcal{E}_K \cap \mathcal{E}_{\text{int}} \\ \sigma = K|L}} -\frac{e_L - e_K}{d_{K|L}} \mathbf{m}(\sigma) + \sum_{\substack{\sigma \in \mathcal{E}_K \cap \mathcal{E}_{\text{ext}} \\ \sigma = K|L}} -\frac{e_\sigma - e_K}{d_{K,\sigma}} \mathbf{m}(\sigma) + \sum_{\sigma \in \mathcal{E}_K} v_{K,\sigma} e_{\sigma,+} + b_K m(K) e_K \\ = -m(K) \rho_K - \sum_{\sigma \in \mathcal{E}_K} \mathbf{m}(\sigma) (R_{K,\sigma} + r_{K,\sigma}), \end{aligned} \quad (79)$$

where, for all  $\sigma \in \mathcal{E}_K$ ,  $e_{\sigma,+} = e_K$  if  $v_{K,\sigma} \geq 0$ ,  $e_{\sigma,+} = e_L$  if  $v_{K,\sigma} < 0$  and  $\sigma \in \mathcal{E}_K \cap \mathcal{E}_{\text{int}}$ ,  $\sigma = K|L$ , finally  $e_{\sigma,+} = e_\sigma$  if  $v_{K,\sigma} < 0$  and  $\sigma \in \mathcal{E}_K \cap \mathcal{E}_{\text{ext}}$ .

In a similar way, using (66) and (69), one has for all  $K \in \mathcal{T}$

$$\mathbf{m}(\sigma) \frac{e_\sigma - e_K}{d_{K,\sigma}} + (\mathbf{m}(\sigma) \lambda_\sigma + v_{K,\sigma}) e_\sigma - v_{K,\sigma} e_{\sigma,+} = \mathbf{m}(\sigma) (R_{K,\sigma} - \tilde{R}_{K,\sigma} + r_{K,\sigma}). \quad (80)$$

We then multiply (80) by  $e_\sigma$ , we sum the result over  $\sigma \in \mathcal{E}_K$ , we multiply (79) by  $e_K$ , we sum these two equalities and we finally sum the result over  $K \in \mathcal{T}$ . Using for the left hand side term the same technique as the one used in the proof of Proposition 3, one obtains

$$\begin{aligned} |e_{\mathcal{T}}|_{1,\mathcal{T}}^2 + \int_\Omega \left( \frac{\operatorname{div}(\mathbf{v}(x))}{2} + b(x) \right) (e_{\mathcal{T}}(x))^2 dx + \int_{\partial\Omega} \left( \lambda(x) + \frac{\mathbf{v}(x) \cdot \mathbf{n}(x)}{2} \right) (e_{\mathcal{T}}(x))^2 dx \\ \leq - \sum_{K \in \mathcal{T}} m(K) \rho_K e_K - \sum_{K \in \mathcal{T}} \sum_{\sigma \in \mathcal{E}_K} \mathbf{m}(\sigma) (R_{K,\sigma} + r_{K,\sigma}) e_K \\ + \sum_{K \in \mathcal{T}} \sum_{\substack{\sigma \in \mathcal{E}_K \cap \mathcal{E}_{\text{ext}}}} \mathbf{m}(\sigma) (R_{K,\sigma} - \tilde{R}_{K,\sigma} + r_{K,\sigma}) e_\sigma. \end{aligned}$$

Using Young's inequality, Lemma 2 if  $u \in C^2(\bar{\Omega})$  and Lemma 3 (with  $p = 4$ ) and Hölder's inequality if  $u$  is only in  $H^2(\Omega)$  (for more details see the proof of Theorems 1 and 2), one gets for all  $\delta > 0$

$$\begin{aligned} |e_{\mathcal{T}}|_{1,\mathcal{T}}^2 + \int_{\Omega} \left( \frac{\operatorname{div}(\mathbf{v}(x))}{2} + b(x) \right) (e_{\mathcal{T}}(x))^2 dx + \int_{\partial\Omega} \left( \lambda(x) + \frac{\mathbf{v}(x) \cdot \mathbf{n}(x)}{2} \right) (e_{\mathcal{T}}(x))^2 dx \\ \leq \frac{C}{2\delta} (\operatorname{size}(\mathcal{T}))^2 + \frac{\delta}{2} \|e_{\mathcal{T}}\|_{L^2(\Omega)}^2 - \sum_{K \in \mathcal{T}} \sum_{\sigma \in \mathcal{E}_K} m(\sigma) (R_{K,\sigma} + r_{K,\sigma}) e_K \\ + \sum_{K \in \mathcal{T}} \sum_{\sigma \in \mathcal{E}_K \cap \mathcal{E}_{\text{ext}}} m(\sigma) (R_{K,\sigma} - \tilde{R}_{K,\sigma} + r_{K,\sigma}) e_{\sigma}, \end{aligned}$$

where  $C$  only depends on  $u$ ,  $b$  and  $\Omega$  if  $u \in C^2(\bar{\Omega})$  and on  $u$ ,  $b$ ,  $\Omega$  and  $\zeta$  if  $u$  is only in  $H^2(\Omega)$ .

Thanks to the conservativity property of the scheme (see (4)), one has, for all  $\sigma \in \mathcal{E}_{\text{int}}$ ,  $\sigma = K|L$ ,  $R_{K,\sigma} = -R_{L,\sigma}$  and  $r_{K,\sigma} = -r_{L,\sigma}$ . Then using this result, Young's inequality, Lemma 2 if  $u \in C^2(\bar{\Omega})$  and Lemma 3 (with  $p = 4$ ) and Hölder's inequality if  $u$  is only in  $H^2(\Omega)$  (for more details see the proof of Theorems 1 and 2)), one obtains

$$- \sum_{K \in \mathcal{T}} \sum_{\sigma \in \mathcal{E}_K} m(\sigma) (R_{K,\sigma} + r_{K,\sigma}) e_K + \sum_{K \in \mathcal{T}} \sum_{\sigma \in \mathcal{E}_K \cap \mathcal{E}_{\text{ext}}} m(\sigma) R_{K,\sigma} e_{\sigma} \leq \frac{1}{2} |e_{\mathcal{T}}|_{1,\mathcal{T}}^2 + C (\operatorname{size}(\mathcal{T}))^2,$$

where  $C$  only depends on  $\mathbf{v}$ ,  $d$ ,  $u$  and  $\Omega$  if  $u \in C^2(\bar{\Omega})$  and on  $\mathbf{v}$ ,  $d$ ,  $u$ ,  $\Omega$  and  $\zeta$  if  $u$  is only in  $H^2(\Omega)$ .

The two previous inequalities yield for all  $\delta > 0$

$$\begin{aligned} \frac{|e_{\mathcal{T}}|_{1,\mathcal{T}}^2}{2} + \int_{\Omega} \left( \frac{\operatorname{div}(\mathbf{v}(x))}{2} + b(x) \right) (e_{\mathcal{T}}(x))^2 dx + \int_{\partial\Omega} \left( \lambda(x) + \frac{\mathbf{v}(x) \cdot \mathbf{n}(x)}{2} \right) (e_{\mathcal{T}}(x))^2 dx \\ \leq C \left( 1 + \frac{1}{\delta} \right) (\operatorname{size}(\mathcal{T}))^2 + \frac{\delta}{2} \|e_{\mathcal{T}}\|_{L^2(\Omega)}^2 - \sum_{K \in \mathcal{T}} \sum_{\sigma \in \mathcal{E}_K \cap \mathcal{E}_{\text{ext}}} m(\sigma) \tilde{R}_{K,\sigma} e_{\sigma}, \end{aligned}$$

where  $C$  only depends on  $d$ ,  $u$ ,  $b$ ,  $\mathbf{v}$ , and  $\Omega$  if  $u \in C^2(\bar{\Omega})$  and on  $d$ ,  $u$ ,  $b$ ,  $\mathbf{v}$ ,  $\Omega$  and  $\zeta$  if  $u$  is only in  $H^2(\Omega)$ . Finally, using Young's inequality, Lemma 9 if  $u \in C^2(\bar{\Omega})$  and Lemma 10 (with  $p = 4$ ) and Hölder's inequality if  $u$  is only in  $H^2(\Omega)$  (for more details see the proof of Theorems 1 and 2)), one obtains for all  $\varepsilon > 0$  and all  $\delta > 0$

$$\begin{aligned} \frac{|e_{\mathcal{T}}|_{1,\mathcal{T}}^2}{2} + \int_{\Omega} \left( \frac{\operatorname{div}(\mathbf{v}(x))}{2} + b(x) \right) (e_{\mathcal{T}}(x))^2 dx + \int_{\partial\Omega} \left( \lambda(x) + \frac{\mathbf{v}(x) \cdot \mathbf{n}(x)}{2} \right) (e_{\mathcal{T}}(x))^2 dx \\ \leq C \left( 1 + \frac{1}{\delta} + \frac{1}{\varepsilon} \right) (\operatorname{size}(\mathcal{T}))^2 + \frac{\delta}{2} \|e_{\mathcal{T}}\|_{L^2(\Omega)}^2 + \frac{\varepsilon}{2} \|e_{\mathcal{T}}\|_{L^2(\partial\Omega)}^2, \end{aligned} \tag{81}$$

where  $C$  only depends on  $d$ ,  $u$ ,  $b$ ,  $\mathbf{v}$ ,  $\lambda$  and  $\Omega$  if  $u \in C^2(\bar{\Omega})$  and on  $d$ ,  $u$ ,  $b$ ,  $\mathbf{v}$ ,  $\lambda$ ,  $\Omega$  and  $\zeta$  if  $u$  is only in  $H^2(\Omega)$ .

**Remark 12** If  $\lambda \geq 0$  a.e. on  $\partial\Omega$  and if one uses (70) instead of (69) in order to discretize the boundary condition. One proves (81), using Remark 11 for the left hand side. For the right hand side, one introduces

$$\tilde{R}_{K,\sigma} = \frac{1}{m(\sigma)} \int_{\sigma} \lambda(x) (u(x) - u(y_{\sigma})) d\gamma(x),$$

then using a technique similar to the one used in the above proof, one gets (81).

Hence if there exists  $\delta > 0$  and  $\varepsilon > 0$  such that  $\operatorname{div}(\mathbf{v})/2 + b > \delta/2$  a.e. on  $\Omega$  and  $\mathbf{v} \cdot \mathbf{n}/2 + \lambda > \varepsilon/2$  a.e. on  $\partial\Omega$ , (81) gives (75) and (76). Otherwise, thanks to Assumption 6 there are two cases.

The first one is when there exists  $\Gamma \subset \partial\Omega$  such that its  $(d-1)$ -dimensional measure  $m(\Gamma) \neq 0$  and such that  $i_b = \inf_{x \in \Gamma} (\mathbf{v}(x) \cdot \mathbf{n}(x)/2 + \lambda(x)) \neq 0$ .

The second one is when there exists  $\mathcal{O} \subset \Omega$  such that its  $d$ -dimensional measure  $m(\mathcal{O}) \neq 0$  and such that  $i_i = \inf_{x \in \mathcal{O}} (\operatorname{div}(\mathbf{v}(x))/2 + b(x)) \neq 0$ .

In both cases one uses Lemma 7 in (81).

In the first case, one obtains

$$\begin{aligned} \frac{|e_{\mathcal{T}}|_{1,\mathcal{T}}^2}{2} + i_b \|e_{\mathcal{T}}\|_{L^2(\Gamma)}^2 &\leq C \left( 1 + \frac{1}{\delta} + \frac{1}{\varepsilon} \right) (\operatorname{size}(\mathcal{T}))^2 \\ &\quad + \frac{C_\Omega}{2} (\delta + \varepsilon C_\Omega) \|e_{\mathcal{T}}\|_{L^2(\Gamma)}^2 + \frac{C_\Omega}{2} (\delta + \varepsilon + \varepsilon C_\Omega) |e_{\mathcal{T}}|_{1,\mathcal{T}}^2, \end{aligned}$$

for all  $\delta > 0$ , all  $\varepsilon > 0$ , where  $C$  only depends on  $\lambda, d, u, b, \mathbf{v}$ , and  $\Omega$  if  $u \in C^2(\bar{\Omega})$  and on  $\lambda, d, u, b, \mathbf{v}, \Omega$  and  $\zeta$  if  $u$  is only in  $H^2(\Omega)$ , and where  $C_\Omega$  only depends on  $\Omega$ .

In the second case, one gets

$$\begin{aligned} \frac{|e_{\mathcal{T}}|_{1,\mathcal{T}}^2}{2} + i_i \|e_{\mathcal{T}}\|_{L^2(\mathcal{O})}^2 &\leq C \left( 1 + \frac{1}{\delta} + \frac{1}{\varepsilon} \right) (\operatorname{size}(\mathcal{T}))^2 \\ &\quad + \frac{C_\Omega}{2} (\varepsilon + \delta C_\Omega) \|e_{\mathcal{T}}\|_{L^2(\mathcal{O})}^2 + \frac{C_\Omega}{2} (\varepsilon + \delta + \delta C_\Omega) |e_{\mathcal{T}}|_{1,\mathcal{T}}^2, \end{aligned}$$

for all  $\delta > 0$ , all  $\varepsilon > 0$ , where  $C$  only depends on  $\lambda, d, u, b, \mathbf{v}$ , and  $\Omega$  if  $u \in C^2(\bar{\Omega})$  and on  $\lambda, d, u, b, \mathbf{v}, \Omega$  and  $\zeta$  if  $u$  is only in  $H^2(\Omega)$ , and where  $C_\Omega$  only depends on  $\Omega$ .

Then a well adapted choice of  $\delta$  and  $\varepsilon$  gives in both cases (75) and using once more Lemma 7 yields (76). This concludes the proof of Theorems 5 and 6. ■

**Acknowledgement** The authors are grateful to the referees for their thorough reading of the article and helpful remarks which helped to improve this article.

## References

- [1] A. AGOUZAL, J. BARANGER, J.-F. MAITRE and F. OUDIN (1995), Connection between finite volume and mixed finite element methods for a diffusion problem with non constant coefficients, with application to Convection Diffusion, *East-West Journal on Numerical Mathematics*. **3**, 4, 237-254.
- [2] J. BARANGER, J.-F. MAITRE and F. OUDIN (1996), Connection between finite volume and mixed finite element methods, *Modél. Math. Anal. Numér.* **30**, 3, 4, 444-465.
- [3] S. H. CHOU and P. S. VASSILEVSKI (1999), A general mixed covolume framework for constructing conservative schemes for elliptic problems, *Math. Comp.* **68**, 991-1011.
- [4] Y. COUDIÈRE, T. GALLOUËT and R. HERBIN, ,  $L^p$  and  $L^\infty$  error estimates for the approximate solution of convection diffusion equations by finite volume schemes on Voronoï meshes, in preparation.
- [5] Y. COUDIÈRE, J.P. VILA and P. VILLEDIEU (1999), Convergence rate of a finite volume scheme for a two-dimensional convection diffusion problem, to appear in M2AN.
- [6] R. EYMARD , T. GALLOUËT and R. HERBIN , The finite volume method, *to appear in Handbook of Numerical Analysis*, J.L. Lions and P.G. Ciarlet eds.
- [7] R. EYMARD, T. GALLOUËT and R. HERBIN , Convergence of finite volume schemes for semilinear convection diffusion equations, accepted for publication in *Numer. Math.* (1999).
- [8] B. HEINRICH (1987), Finite Difference Methods on Irregular Networks, International Series of Numerical Mathematics, **Vol. 82**, *Birkhäuser Verlag, Basel*, A generalized approach to second order elliptic problems.

- [9] R. HERBIN (1995), An error estimate for a finite volume scheme for a diffusion-convection problem on a triangular mesh, *Num. Meth. P.D.E.* **11**, 165-173.
- [10] R. HERBIN (1996), Finite volume methods for diffusion convection equations on general meshes, in *Finite volumes for complex applications, Problems and Perspectives*, F. Benkhaldoun and R. Vilsmeier eds, Hermes, 153-160.
- [11] R.D. LAZAROV and I.D. MISHEV (1996), Finite volume methods for reaction diffusion problems in: F. Benkhaldoun and R. Vilsmeier eds, *Finite volumes for complex applications, Problems and Perspectives* (Hermes, Paris), 233-240.
- [12] R.D. LAZAROV, I.D. MISHEV and P.S. VASSILEVSKI (1996), Finite volume methods for convection diffusion problems, *SIAM J. Numer. Anal.* **33**, 31-55.
- [13] T.A. MANTEUFEL, and A.B. WHITE (1986), The numerical solution of second order boundary value problem on non uniform meshes, *Math. Comput.* **47**, 511-536.
- [14] I. D. MISHEV (1998), Finite volume methods on Voronoï meshes, *Num. Meth. P.D.E.* **14**, 2, 193-212.
- [15] K.W. MORTON and E. SÜLI (1991), Finite volume methods and their analysis, *IMA J. Numer. Anal.* **11**, 241-260.
- [16] K.W. MORTON (1996), *Numerical Solutions of Convection-Diffusion problems* (Chapman and Hall, London).
- [17] J. NEČAS (1967), Les méthodes directes en théorie des équations elliptiques, (Masson, Paris).
- [18] J.-M. THOMAS and D. TRUJILLO (1994), Finite volume variational formulation. Applications to domain decomposition methods, *Contemporary Mathematics* **157**, 127–132.
- [19] J.-M. THOMAS and D. TRUJILLO (1995), Analysis of finite volumes methods, in: world Scientific, A. Bourgeat, C. Carasso, S. Luckaus, A. Michelic ed., *Proceedings of Mathematical Modeling through Porous Media* (St Etienne).
- [20] VANSELOW R. (1996), Relations between FEM and FVM, in: F. Benkhaldoun and R. Vilsmeier eds, *Finite volumes for complex applications, Problems and Perspectives* (Hermes, Paris), 217-223.
- [21] P.S. VASSILEVSKI, S.I. PETROVA and R.D. LAZAROV (1992), Finite difference schemes on triangular cell-centered grids with local refinement, *SIAM J. Sci. Stat. Comput.* **13**, 6, 1287–1313.
- [22] M.H. VIGNAL (1996), Convergence of a finite volume scheme for a system of an elliptic equation and a hyperbolic equation *Modél. Math. Anal. Numér.* **30**, 7, 841-872.
- [23] M.H. VIGNAL (1997), Schémas Volumes Finis pour des équations elliptiques ou hyperboliques avec conditions aux limites, convergence et estimations d'erreur *Thèse de Doctorat*, Ecole Normale Supérieure de Lyon.