



# An optimization-based numerical method for diffusion problems with sign-changing coefficients

Assyr Abdulle, Martin E. Huber, Simon Lemaire

## ► To cite this version:

Assyr Abdulle, Martin E. Huber, Simon Lemaire. An optimization-based numerical method for diffusion problems with sign-changing coefficients. *Comptes Rendus. Mathématique*, 2017, 355 (4), pp.472-478. 10.1016/j.crma.2017.02.010 . hal-01354092v2

**HAL Id: hal-01354092**

**<https://hal.science/hal-01354092v2>**

Submitted on 12 Dec 2016

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# An optimization-based numerical method for diffusion problems with sign-changing coefficients

Assyr Abdulle<sup>a</sup>, Martin E. Huber<sup>a</sup>, Simon Lemaire<sup>a</sup>

<sup>a</sup>*École Polytechnique Fédérale de Lausanne (EPFL), SB-MATHICSE-ANMC, Station 8, 1015 Lausanne, Switzerland*

---

## Abstract

A new optimization-based numerical method is proposed for the solution of diffusion problems with sign-changing conductivity coefficients. In contrast to existing approaches, our method does not rely on the discretization of a stabilized equation and the convergence of the scheme can be proved without any symmetry assumption on the mesh near the interface where the conductivity sign changes.

## Résumé

**Une méthode d'optimisation pour des problèmes de diffusion avec changement de signe.** Nous proposons une nouvelle méthode, basée sur la résolution d'un problème de minimisation, pour l'approximation de problèmes de diffusion avec changement de signe. Cette approche, qui tire profit d'une reformulation du modèle initial sous la forme d'un problème de transmission, ne repose pas sur la discrétisation d'une équation stabilisée, et la convergence de la méthode est obtenue sans hypothèse de symétrie du maillage dans un voisinage de l'interface où la conductivité change de signe.

---

## Version française abrégée

Dans cette note, nous introduisons une méthode d'optimisation pour l'approximation numérique de problèmes de diffusion dont la conductivité change de signe dans le domaine. La résolution numérique efficace de ce genre de problèmes est importante pour de nombreuses applications (e.g., super-lentilles, invisibilité), mais les méthodes existantes ne sont pour l'instant pas satisfaisantes. Dans [6], les deux approches envisagées reposent (i) sur la discrétisation d'une équation stabilisée, pour laquelle les taux de convergence obtenus sont sous-optimaux, ou (ii) sur des hypothèses de symétrie du maillage autour de l'interface où la conductivité change de signe, exigences pouvant s'avérer très contraignantes pour des interfaces générales (voir [3]) ou en 3D. La méthode numérique introduite ici, qui utilise une reformulation du modèle initial en un problème de transmission, ne repose pas sur l'ajout de dissipation à l'équation, et nous montrons sa convergence pour des problèmes elliptiques présentant un changement de signe sans aucune hypothèse de symétrie sur le maillage. Nous notons que cette méthode numérique a pour la première fois été introduite par Gunzburger et al. [9] (voir aussi [8]), dans un contexte de décomposition de domaine pour des équations elliptiques classiques, sans preuve de convergence. L'application de cet algorithme à des problèmes elliptiques présentant un changement de signe est introduite dans cette note, et la convergence de la méthode est démontrée.

---

*Email addresses:* `assyр.abdulle@epfl.ch` (Assyr Abdulle), `martin.huber@epfl.ch` (Martin E. Huber), `simon.lemaire@epfl.ch` (Simon Lemaire).

## 1. Introduction

Partial differential equations with sign-changing coefficients play an increasingly important role in the modeling of metamaterials, with applications ranging from superlensing to cloaking. In this paper, we consider in an open bounded polytopal domain  $\Omega \subset \mathbb{R}^d$ ,  $d \in \{2, 3\}$  such that  $\bar{\Omega} = \bar{\Omega}_e \cup \bar{\Omega}_i$ , with  $\Omega_e, \Omega_i$  disjoint open polytopal subsets of  $\Omega$  with nonzero measure, the following sign-changing diffusion problem

$$-\operatorname{div}(s \mathbb{K} \nabla \tilde{u}) = f \text{ in } \Omega, \quad \tilde{u} = 0 \text{ on } \partial\Omega, \quad (1)$$

where  $\mathbb{K}$  is a symmetric, uniformly elliptic and bounded matrix-valued field, and  $s : \Omega \rightarrow \mathbb{R}$  is a sign function such that  $s_e := s|_{\Omega_e} = 1$  and  $s_i := s|_{\Omega_i} = -1$ . The subscripts 'e' and 'i' respectively stand for the exterior (for which we assume that  $\partial\Omega_e \cap \partial\Omega$  has nonzero measure) and the interior subdomains, and we denote their interface  $\Gamma = \partial\Omega_e \cap \partial\Omega_i$ . We assume that  $f \in L^2(\Omega)$  and study the following weak formulation of Problem (1): Find  $\tilde{u} \in H_0^1(\Omega)$  such that

$$a(\tilde{u}, \varphi) := (s \mathbb{K} \nabla \tilde{u}, \nabla \varphi)_\Omega = (f, \varphi)_\Omega \quad \forall \varphi \in H_0^1(\Omega). \quad (2)$$

Several approaches have been developed to study the well-posedness of Problem (2). We first mention the T-coercivity theory of Bonnet-Ben Dhia et al. [4]. In this context, well-posedness (in the classical Hadamard sense or in the Fredholm sense) of Problem (2) is equivalent to finding a linear, bounded and bijective operator  $T : H_0^1(\Omega) \rightarrow H_0^1(\Omega)$  such that the bilinear form  $a(\cdot, T\cdot)$  is coercive (for classical well-posedness) or weakly coercive (for the Fredholm case). The operator  $T$  plays the role of an explicit inf – sup operator. However, proving T-coercivity can be difficult for complex geometries, especially in 3D. More recently, another viewpoint has come out that consists in studying the “limit” as  $\delta \rightarrow 0^+$  of the well-posed stabilized problem: Find  $\tilde{u}_\delta \in H_0^1(\Omega)$  such that

$$a_\delta(\tilde{u}_\delta, \varphi) := a(\tilde{u}_\delta, \varphi) - i\delta(\mathbb{K} \nabla \tilde{u}_\delta, \nabla \varphi)_{\Omega_i} = (f, \varphi)_\Omega \quad \forall \varphi \in H_0^1(\Omega), \quad (3)$$

where  $H_0^1(\Omega)$  is complex-valued. This approach has been developed in the context of Helmholtz equations by Nguyen [11,12], and makes use of transport operators called reflections.

Different strategies have been studied by Chesnel and Ciarlet Jr. in [6] for the numerical approximation of Problem (2) (assuming that it is well-posed in the classical Hadamard sense). The first approach is based on a simplicial discretization  $\mathcal{T}_h$  of  $\Omega$ , that respects the interface  $\Gamma$  and the construction of a conforming finite element space  $V_0(\mathcal{T}_h)$  that is stable by the operator  $T$ . Well-posedness and optimal convergence rates can be shown for this approach. In practice, T-stability is achieved by means of symmetric meshes near the interface, whose construction is a nontrivial task for complicated interfaces (see [3]) or 3D problems. On general meshes, two main approaches have been investigated by Chesnel and Ciarlet Jr. The first one consists in building a (mesh-dependent) operator  $T_h$ , such that the bilinear form  $a$  is  $T_h$ -coercive on  $V_0(\mathcal{T}_h)$  (see, e.g., [13]). This kind of approach is limited by the fact that, in general, well-posedness cannot be proved for the whole range of admissible coefficients. The second idea consists in discretizing the stabilized equation (3), and in scaling the dissipation  $\delta$  as a function of  $h$ . However, and as studied in [6], this approach leads to suboptimal convergence rates.

In this note, we aim at proposing a new discretization approach, that is based on a reformulation of Problem (1) as a transmission problem, whose solution is obtained numerically from an optimization procedure. The approach we propose is proved to converge, and does not rely on any symmetry assumption on the mesh. As opposed to T-coercivity theory, we do not need to assume well-posedness of Problem (2) in the Hadamard sense (i.e. for any  $f \in H^{-1}(\Omega)$ ) to prove convergence of our numerical method. We only need to assume that, for the given  $f \in L^2(\Omega)$  we consider, the solution to Problem (2) exists and is unique. This allows us to treat cases that cannot be analyzed using the T-coercivity theory (cf. Remark 3). The numerical method we consider here has been introduced by Gunzburger et al. [9] (see also [8]), in the context of domain decomposition for classical elliptic equations. However, up to now, the convergence

of this method has never been proved. We provide the first proof of convergence of such a numerical method, with application to sign-changing elliptic equations, for which this approach seems particularly promising.

## 2. Transmission problem and numerical method

We provide in this section, under suitable regularity assumptions, an alternative characterization of the solution to Problem (2), that will be our starting point for the design of the numerical method.

### 2.1. Transmission problem

Recall that  $f \in L^2(\Omega)$ . Hence,  $s\mathbb{K}\nabla\tilde{u} \in \mathbf{H}(\text{div}; \Omega)$ . Denoting, for  $\alpha \in \{e, i\}$ , by  $\mathbf{n}_\alpha$  the unit normal vector to  $\Gamma$  pointing out of  $\Omega_\alpha$ , and introducing  $\tilde{g}$  such that

$$\tilde{g} := \left[ (\mathbb{K}\nabla\tilde{u})_{|\Omega_e} \cdot \mathbf{n}_e \right]_{|\Gamma}, \quad (4)$$

we have that  $\tilde{g} \in H^{-1/2}(\Gamma)$ . Here,  $H^{-1/2}(\Gamma)$  denotes the dual of  $H^{1/2}(\partial\Omega_i)$  when  $\Gamma = \partial\Omega_i$  or of  $H_{00}^{1/2}(\Gamma)$  otherwise. For  $\alpha \in \{e, i\}$ , and for any  $g \in H^{-1/2}(\Gamma)$ , we consider in the subdomain  $\Omega_\alpha$  the problem

$$s_\alpha(\mathbb{K}\nabla u_\alpha(g), \nabla\varphi_\alpha)_{\Omega_\alpha} = (f, \varphi_\alpha)_{\Omega_\alpha} + s_\alpha\langle g, \varphi_\alpha \rangle_\Gamma \quad \forall \varphi_\alpha \in H_{0\setminus\Gamma}^1(\Omega_\alpha), \quad (5)$$

where  $H_{0\setminus\Gamma}^1(\Omega_\alpha)$  is the space of functions in  $H^1(\Omega_\alpha)$  that vanish on  $\partial\Omega_\alpha \setminus \Gamma$ . In  $\Omega_e$ , Problem (5) always admits a unique solution  $u_e(g) \in H_{0\setminus\Gamma}^1(\Omega_e)$ , as the measure of  $\partial\Omega_e \cap \partial\Omega$  is nonzero. In  $\Omega_i$ , Problem (5) also admits a unique solution  $u_i(g) \in H_{0\setminus\Gamma}^1(\Omega_i)$  if the measure of  $\partial\Omega_i \cap \partial\Omega$  is nonzero. Otherwise, the problem in  $\Omega_i$  is purely Neumann and we assume that the flux  $g \in H^{-1/2}(\Gamma)$  satisfies  $(f, 1)_{\Omega_i} - \langle g, 1 \rangle_\Gamma = 0$  to ensure that Problem (5) admits a solution, that is unique up to an additive constant. We fix the constant by imposing  $(u_i(g), 1)_\Gamma = (u_e(g), 1)_\Gamma$ . Finally, for  $g \in H^{-1/2}(\Gamma)$ , we denote by  $u(g)$  the function such that  $u(g)_{|\Omega_\alpha} := u_\alpha(g)$ ,  $\alpha \in \{e, i\}$ .

**Proposition 2.1 (Characterization of the solution to Problem (2))** *We assume that (2) admits a unique solution  $\tilde{u} \in H_0^1(\Omega)$ , and that there exists  $\tilde{s} > \frac{1}{2}$  such that  $\tilde{u}_{|\Omega_\alpha} \in H^{1+\tilde{s}}(\Omega_\alpha)$  for  $\alpha \in \{e, i\}$ . Then,*

- the flux  $\tilde{g}$  defined in (4) belongs to  $L^2(\Gamma)$ , and satisfies  $\tilde{g} = -\left[ (-\mathbb{K}\nabla\tilde{u})_{|\Omega_i} \cdot \mathbf{n}_i \right]_{|\Gamma}$ ;
- $u(\tilde{g}) = \tilde{u} \in H_0^1(\Omega)$ ;
- almost everywhere on the interface  $\Gamma$ ,  $u_e(\tilde{g}) = u_i(\tilde{g})$ ;
- the problem  $\inf_{g \in L^2(\Gamma)} \|u_e(g) - u_i(g)\|_{0,\Gamma}^2$  admits  $\tilde{g}$  as its unique solution.

*Remark 1 (Assumptions on  $\tilde{u}$ )* The existence (and uniqueness) of a solution  $\tilde{u} \in H_0^1(\Omega)$  to Problem (2) is satisfied in practice in a large variety of situations (cf., e.g., [12]). The assumption  $\tilde{u}_{|\Omega_\alpha} \in H^{1+\tilde{s}}(\Omega_\alpha)$  for  $\alpha \in \{e, i\}$ ,  $\tilde{s} > \frac{1}{2}$ , made in Proposition 2.1 is convenient for the analysis, as it enables to work in  $L^2(\Gamma)$  instead of working in  $H^{-1/2}(\Gamma)$ . This theoretical assumption is quite strict, but not mandatory in practice for the method to be applicable. In a forthcoming work [1], we will consider a test-case for which this regularity assumption is violated, and show that numerical convergence can still be observed.

## 2.2. Minimization problem and numerical method

We consider a family of simplicial conformal discretizations  $\{\mathcal{T}_h\}_h$  of  $\Omega$ , that respects the interface  $\Gamma$  and is shape-regular in the sense of Ciarlet [7]. The subscript  $h$  stands for the meshsize, i.e. the maximum diameter of all the simplices in  $\mathcal{T}_h$ . We denote by  $\Gamma_h$  the set of faces of the mesh  $\mathcal{T}_h$  belonging to  $\Gamma$ . For  $k \in \mathbb{N}^*$  and  $\alpha \in \{e, i\}$ , we introduce the space  $V^k(\mathcal{T}_{h,\alpha}) := \{v \in H^1(\Omega_\alpha) \mid v|_T \in \mathbb{P}_d^k(T), \forall T \in \mathcal{T}_{h,\alpha}\}$ , where  $\mathbb{P}_d^k(T)$  is the space of  $d$ -variate polynomial functions of total degree less or equal to  $k$  in  $T \in \mathcal{T}_h$ . We also introduce the subspace  $V_{0\Gamma}^k(\mathcal{T}_{h,\alpha}) := V^k(\mathcal{T}_{h,\alpha}) \cap H_{0\Gamma}^1(\Omega_\alpha)$ .

For  $\alpha \in \{e, i\}$ , and for any  $g \in H^{-1/2}(\Gamma)$ , we consider, in the subdomain  $\Omega_\alpha$ , the following conforming approximation of Problem (5):

$$s_\alpha(\mathbb{K} \nabla u_{h,\alpha}(g), \nabla \varphi_{h,\alpha})_{\Omega_\alpha} = (f, \varphi_{h,\alpha})_{\Omega_\alpha} + s_\alpha \langle g, \varphi_{h,\alpha} \rangle_\Gamma \quad \forall \varphi_{h,\alpha} \in V_{0\Gamma}^k(\mathcal{T}_{h,\alpha}). \quad (6)$$

As in the continuous case, Problem (6) always admits a unique solution  $u_{h,e}(g) \in V_{0\Gamma}^k(\mathcal{T}_{h,e})$  in  $\Omega_e$ , and  $u_{h,i}(g) \in V_{0\Gamma}^k(\mathcal{T}_{h,i})$  in  $\Omega_i$  when the measure of  $\partial\Omega_i \cap \partial\Omega$  is nonzero. For the case of a pure Neumann problem in  $\Omega_i$ , we assume that  $(f, 1)_{\Omega_i} - \langle g, 1 \rangle_\Gamma = 0$ , ensuring existence of a solution and uniqueness up to an additive constant, that we fix imposing  $(u_{h,i}(g), 1)_\Gamma = (u_{h,e}(g), 1)_\Gamma$ . For  $g \in H^{-1/2}(\Gamma)$ , we introduce  $u_h(g)$  such that  $u_h(g)|_{\Omega_\alpha} := u_{h,\alpha}(g)$ ,  $\alpha \in \{e, i\}$ .

To define the minimization problem, we introduce  $F^k(\Gamma_h) := \{q \in L^2(\Gamma) \mid q|_F \in \mathbb{P}_{d-1}^k(F), \forall F \in \Gamma_h\}$  when the measure of  $\partial\Omega_i \cap \partial\Omega$  is nonzero, and we add the constraint  $(q, 1)_\Gamma = (f, 1)_{\Omega_i}$  otherwise. We then define the functional  $J_h : F^k(\Gamma_h) \rightarrow \mathbb{R}_+$  such that  $J_h(g_h) := \|u_{h,e}(g_h) - u_{h,i}(g_h)\|_{0,\Gamma}^2 + \lambda(h) \|g_h\|_{0,\Gamma}^2$ , where  $\lambda : \mathbb{R}_+^* \rightarrow \mathbb{R}_+^*$  is a function such that  $\lim_{h \rightarrow 0} \lambda(h) = 0$ , and we consider the minimization problem

$$\inf_{g_h \in F^k(\Gamma_h)} J_h(g_h). \quad (7)$$

The continuous function  $J_h$  is coercive on  $F^k(\Gamma_h)$ , and hence admits at least one minimizer: there exists  $\tilde{g}_h \in F^k(\Gamma_h)$  such that

$$J_h(\tilde{g}_h) \leq J_h(g_h) \quad \text{for any } g_h \in F^k(\Gamma_h). \quad (8)$$

The approximation of  $\tilde{u}$  we then consider is  $u_h(\tilde{g}_h)$ . The following result is a direct consequence of (8) and of the properties of the  $L^2$ -orthogonal projector from  $L^2(\Gamma)$  onto  $F^k(\Gamma_h)$ .

**Lemma 2.1 (Estimate on  $J_h(\tilde{g}_h)$ )** *Assume that  $\tilde{g}$  defined in (4) is in  $L^2(\Gamma)$ . Then, the following holds:*

$$J_h(\tilde{g}_h) \leq \|u_{h,e}(\tilde{g}) - u_{h,i}(\tilde{g})\|_{0,\Gamma}^2 + \lambda(h) \|\tilde{g}\|_{0,\Gamma}^2. \quad (9)$$

## 3. Convergence

From now on, we write  $A \lesssim B$  when there exists a constant  $c > 0$ , possibly depending on  $\mathbb{K}$ ,  $k$ , and on the geometry, but independent of  $h$ , and of  $\tilde{u}$ ,  $\tilde{g}$ , such that  $A \leq cB$ .

**Lemma 3.1 (Approximation properties)** *Suppose that the assumptions of Proposition 2.1 are fulfilled. Then, letting  $p = \min(k, \tilde{s})$ , there exists  $\delta \in (0, \frac{1}{2}]$  such that*

$$\begin{aligned} \|\nabla_h(u_h(\tilde{g}) - \tilde{u})\|_{0,\Omega} &\lesssim h^p \left( |\tilde{u}|_{1+p,\Omega_e} + |\tilde{u}|_{1+p,\Omega_i} \right), \\ \|u_{h,e}(\tilde{g}) - u_{h,i}(\tilde{g})\|_{0,\Gamma} &\leq \sum_{\alpha \in \{e, i\}} \|u_{h,\alpha}(\tilde{g}) - \tilde{u}\|_{0,\Gamma} \lesssim h^{p+\delta} \left( |\tilde{u}|_{1+p,\Omega_e} + |\tilde{u}|_{1+p,\Omega_i} \right). \end{aligned}$$

*Proof* We first assume that we have, both in  $\Omega_e$  and  $\Omega_i$ , mixed Dirichlet–Neumann boundary conditions. Then, classical arguments (particularly, Aubin–Nitsche duality argument, and a multiplicative continuous trace inequality) yield

$$\|u_{h,\alpha}(\tilde{g}) - \tilde{u}\|_{0,\Omega_\alpha} + h^{\frac{\delta_\alpha^m}{2}} \|u_{h,\alpha}(\tilde{g}) - \tilde{u}\|_{0,\Gamma} + h^{\delta_\alpha^m} \|\nabla(u_{h,\alpha}(\tilde{g}) - \tilde{u})\|_{0,\Omega_\alpha} \lesssim h^{p+\delta_\alpha^m} |\tilde{u}|_{1+p,\Omega_\alpha}, \quad (10)$$

for  $\alpha \in \{e, i\}$  and with  $\delta_\alpha^m \in (0, 1]$  only depending on the geometry of  $\Omega_\alpha$ . By setting  $\delta = \min\left(\frac{\delta_e^m}{2}, \frac{\delta_i^m}{2}\right)$ , we obtain the stated result. We next treat the case of a pure Neumann problem in  $\Omega_i$ . Due to the choice  $(u_{h,i}(\tilde{g}), 1)_\Gamma = (u_{h,e}(\tilde{g}), 1)_\Gamma$  to fix the constant, a straightforward Aubin–Nitsche duality argument fails. We therefore consider an auxiliary function  $\bar{u}_{h,i}(\tilde{g})$ , defined as the unique discrete solution of the same Neumann problem in  $\Omega_i$ , but for which we fix the constant as in the continuous problem, namely  $(\bar{u}_{h,i}(\tilde{g}), 1)_{\Omega_i} = (\tilde{u}, 1)_{\Omega_i}$ . Then, there exists  $\delta_i^n \in (0, 1]$ , only depending on the geometry of  $\Omega_i$ , such that

$$\|\bar{u}_{h,i}(\tilde{g}) - \tilde{u}\|_{0,\Omega_i} + h^{\frac{\delta_i^n}{2}} \|\bar{u}_{h,i}(\tilde{g}) - \tilde{u}\|_{0,\Gamma} + h^{\delta_i^n} \|\nabla(\bar{u}_{h,i}(\tilde{g}) - \tilde{u})\|_{0,\Omega_i} \lesssim h^{p+\delta_i^n} |\tilde{u}|_{1+p,\Omega_i}. \quad (11)$$

Note that, a priori,  $\delta_i^n \neq \delta_i^m$  since regularity results for pure Neumann or mixed problems are different in general. Next, we observe that  $u_{h,i}(\tilde{g}) - \bar{u}_{h,i}(\tilde{g}) = c_h$ , for a constant  $c_h \in \mathbb{R}$ . We then write

$$|\Gamma| c_h = (u_{h,i}(\tilde{g}) - \bar{u}_{h,i}(\tilde{g}), 1)_\Gamma = (u_{h,e}(\tilde{g}) - \bar{u}_{h,i}(\tilde{g}), 1)_\Gamma = (u_{h,e}(\tilde{g}) - \tilde{u}, 1)_\Gamma + (\tilde{u} - \bar{u}_{h,i}(\tilde{g}), 1)_\Gamma,$$

which yields, using (10) and (11),

$$|c_h| \leq |\Gamma|^{-1/2} \left( \|u_{h,e}(\tilde{g}) - \tilde{u}\|_{0,\Gamma} + \|\bar{u}_{h,i}(\tilde{g}) - \tilde{u}\|_{0,\Gamma} \right) \lesssim h^{p+\frac{\delta_e^m}{2}} |\tilde{u}|_{1+p,\Omega_e} + h^{p+\frac{\delta_i^n}{2}} |\tilde{u}|_{1+p,\Omega_i}. \quad (12)$$

We finally obtain from (12)

$$\|u_{h,i}(\tilde{g}) - \tilde{u}\|_{0,\Gamma} \leq \|u_{h,e}(\tilde{g}) - \tilde{u}\|_{0,\Gamma} + 2\|\bar{u}_{h,i}(\tilde{g}) - \tilde{u}\|_{0,\Gamma} \lesssim h^{p+\frac{\delta_e^m}{2}} |\tilde{u}|_{1+p,\Omega_e} + h^{p+\frac{\delta_i^n}{2}} |\tilde{u}|_{1+p,\Omega_i},$$

which enables to conclude the proof, setting in that case  $\delta = \min\left(\frac{\delta_e^m}{2}, \frac{\delta_i^n}{2}\right)$ .  $\square$

We note that, as a consequence of Lemma 3.1, it is always possible to choose  $\lambda(h)$  in (7) such that  $\|u_{h,e}(\tilde{g}) - u_{h,i}(\tilde{g})\|_{0,\Gamma}^2 = O(\lambda(h))$  or  $\|u_{h,e}(\tilde{g}) - u_{h,i}(\tilde{g})\|_{0,\Gamma}^2 = o(\lambda(h))$ . We then have the

**Theorem 3.2 (Convergence of the method)** *Suppose that the assumptions of Proposition 2.1 are fulfilled. Denote  $u_h(\tilde{g}_h)$  by  $\tilde{u}_h$ . Then, up to a choice of  $\lambda(h)$  in (7) such that  $\|u_{h,e}(\tilde{g}) - u_{h,i}(\tilde{g})\|_{0,\Gamma}^2 = O(\lambda(h))$ , there holds as  $h \rightarrow 0$ :*

$$\tilde{g}_h \rightarrow \tilde{g} \quad \text{in } L^2(\Gamma), \quad \nabla_h \tilde{u}_h \rightarrow \nabla \tilde{u} \quad \text{in } L^2(\Omega)^d, \quad \tilde{u}_h \rightarrow \tilde{u} \quad \text{in } L^2(\Omega). \quad (13)$$

*If we further choose  $\lambda(h)$  such that  $\|u_{h,e}(\tilde{g}) - u_{h,i}(\tilde{g})\|_{0,\Gamma}^2 = o(\lambda(h))$ , then we have*

$$\tilde{g}_h \rightarrow \tilde{g} \quad \text{in } L^2(\Gamma), \quad \nabla_h \tilde{u}_h \rightarrow \nabla \tilde{u} \quad \text{in } L^2(\Omega)^d, \quad \tilde{u}_h \rightarrow \tilde{u} \quad \text{in } L^2(\Omega). \quad (14)$$

*Proof* We begin by remarking that, for  $\alpha \in \{e, i\}$ ,

$$u_{h,\alpha}(\tilde{g}_h) = u_{h,\alpha}(\tilde{g}) + v_{h,\alpha}, \quad (15)$$

where  $v_{h,\alpha} \in V_{0\backslash\Gamma}^k(\mathcal{T}_{h,\alpha})$  satisfies

$$(\mathbb{K} \nabla v_{h,\alpha}, \nabla \varphi_{h,\alpha})_{\Omega_\alpha} = ((\tilde{g}_h - \tilde{g}), \varphi_{h,\alpha})_\Gamma \quad \forall \varphi_{h,\alpha} \in V_{0\backslash\Gamma}^k(\mathcal{T}_{h,\alpha}). \quad (16)$$

From (16), we readily infer

$$\|\nabla v_{h,\alpha}\|_{0,\Omega_\alpha}^2 \lesssim \|\tilde{g}_h - \tilde{g}\|_{0,\Gamma} \|v_{h,\alpha}\|_{0,\Gamma}. \quad (17)$$

Besides, the following estimate can be obtained from (9):

$$\|\tilde{g}_h\|_{0,\Gamma}^2 \leq \frac{\|u_{h,e}(\tilde{g}) - u_{h,i}(\tilde{g})\|_{0,\Gamma}^2}{\lambda(h)} + \|\tilde{g}\|_{0,\Gamma}^2. \quad (18)$$

As soon as  $\|u_{h,e}(\tilde{g}) - u_{h,i}(\tilde{g})\|_{0,\Gamma}^2/\lambda(h) \leq c$ , the combination of (18), (17), and of a continuous trace inequality (as well as of an appropriate Poincaré inequality [5, Equation (1.1)] when the measure of  $\partial\Omega_i \cap \partial\Omega$  is zero) enables to infer that there exist  $\tilde{g}_0 \in L^2(\Gamma)$ , and  $v_\alpha \in H_0^1(\Omega_\alpha)$ ,  $\alpha \in \{e, i\}$ , so that the following convergences hold as  $h \rightarrow 0$ , up to a subsequence (retaining the same notation):

$$\begin{aligned} \tilde{g}_h &\rightharpoonup \tilde{g}_0 && \text{in } L^2(\Gamma), && \nabla v_{h,\alpha} &\rightharpoonup \nabla v_\alpha && \text{in } L^2(\Omega_\alpha)^d, \\ v_{h,\alpha} &\rightarrow v_\alpha && \text{in } L^2(\Omega_\alpha), && v_{h,\alpha} &\rightarrow v_\alpha && \text{in } L^2(\Gamma). \end{aligned} \quad (19)$$

Now, from (9) we infer

$$\|u_{h,e}(\tilde{g}_h) - u_{h,i}(\tilde{g}_h)\|_{0,\Gamma}^2 \leq \|u_{h,e}(\tilde{g}) - u_{h,i}(\tilde{g})\|_{0,\Gamma}^2 + \lambda(h)\|\tilde{g}\|_{0,\Gamma}^2 \leq (c + \|\tilde{g}\|_{0,\Gamma}^2)\lambda(h). \quad (20)$$

Owing to the fact that  $\lambda(h)$  tends to zero as  $h$  vanishes, we deduce from (20) that  $\|u_{h,e}(\tilde{g}_h) - u_{h,i}(\tilde{g}_h)\|_{0,\Gamma} \rightarrow 0$  as  $h \rightarrow 0$ . Combining this last result with (15) and Lemma 3.1, we get that  $\|v_{h,e} - v_{h,i}\|_{0,\Gamma} \rightarrow 0$  as  $h \rightarrow 0$  and hence, owing to (19), that  $v_e = v_i$  almost everywhere on  $\Gamma$ . Passing to the limit in (16) using (19) and a strongly converging interpolant for test functions, it can be shown that  $v \in H_0^1(\Omega)$  such that  $v|_{\Omega_\alpha} := v_\alpha$  for  $\alpha \in \{e, i\}$  satisfies  $(s\mathbb{K}\nabla v, \nabla\varphi)_\Omega = 0$  for all  $\varphi \in H_0^1(\Omega)$ , which implies, from the well-posedness of Problem (2), that  $v \equiv 0$  and  $\tilde{g}_0 \equiv \tilde{g}$  (this last result is inferred considering the limit equation in a subdomain, and using a density argument along with the fact that  $\tilde{g}_0, \tilde{g} \in L^2(\Gamma)$ ). The uniqueness of the limits implies that the whole sequences converge in (19). Collecting these last results, (15), and Lemma 3.1, we prove (13). If we further choose  $\lambda(h)$  such that  $\|u_{h,e}(\tilde{g}) - u_{h,i}(\tilde{g})\|_{0,\Gamma}^2/\lambda(h) \rightarrow 0$  as  $h \rightarrow 0$ , we get from (18) and from the weak convergence of  $\tilde{g}_h$  towards  $\tilde{g}$  that

$$\|\tilde{g}\|_{0,\Gamma} \leq \liminf_{h \rightarrow 0} \|\tilde{g}_h\|_{0,\Gamma} = \limsup_{h \rightarrow 0} \|\tilde{g}_h\|_{0,\Gamma} \leq \|\tilde{g}\|_{0,\Gamma}, \quad (21)$$

which states the strong convergence of  $\tilde{g}_h$  towards  $\tilde{g}$  in  $L^2(\Gamma)$ . Testing (16) with  $v_{h,\alpha}$ , and using (21) combined with the weak convergence result for  $v_{h,\alpha}$  in  $L^2(\Gamma)$  of (19), we finally get the strong convergence of  $\nabla v_{h,\alpha}$  towards 0 in  $L^2(\Omega_\alpha)^d$  for  $\alpha \in \{e, i\}$ , proving (14) and concluding the proof.  $\square$

#### 4. Numerical validation

We consider the 2D symmetric cavity test-case of [6], for which  $\Omega := (-1, 1) \times (0, 1)$  and  $\Omega_e := (-1, 0) \times (0, 1)$  (it follows that  $\Omega_i = (0, 1) \times (0, 1)$  and  $\Gamma = \{0\} \times [0, 1]$ ). For this test-case, both the exterior and the interior problems feature mixed Dirichlet–Neumann boundary conditions. For the particular geometry considered here, elliptic regularity holds in both subdomains (see, e.g., [2, Remark I.3.6]), so that  $\delta = \frac{1}{2}$  in Lemma 3.1. The tensor  $\mathbb{K}$  is taken isotropic, with constant value in each subdomain ( $k_{|\Omega_e} = k_e > 0$  and  $k_{|\Omega_i} = k_i > 0$ ). For this particular setting, it is known [6] that Problem (2) is well-posed in the classical Hadamard sense if and only if the contrast  $\nu := -\frac{k_i}{k_e}$  is different from  $-1$ . The exact solution  $\tilde{u} \in H_0^1(\Omega)$  we consider is

$$\tilde{u}(x, y) := \begin{cases} \left( (x+1)^2 - \frac{2k_e - k_i}{k_e - k_i}(x+1) \right) \sin(\pi y) & \text{in } \Omega_e, \\ \frac{k_e}{k_e - k_i}(x-1) \sin(\pi y) & \text{in } \Omega_i, \end{cases}$$

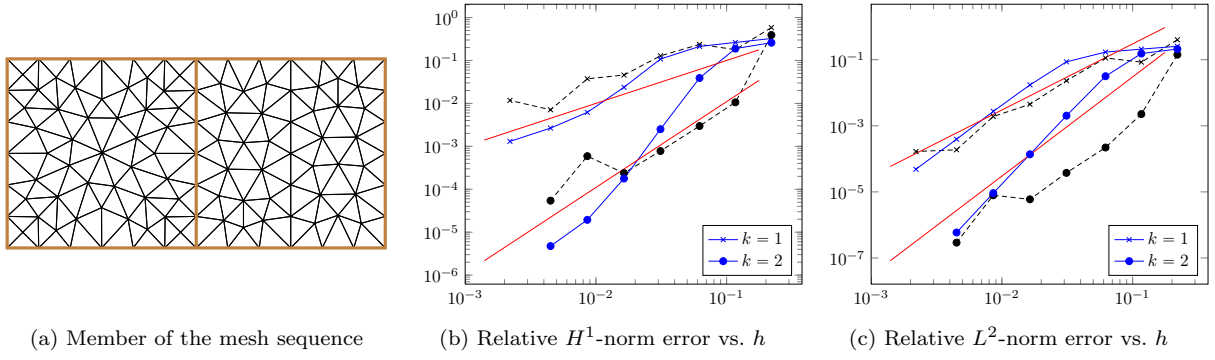


Figure 1. Relative errors on  $\Omega$  vs. meshsize for the symmetric cavity test-case with  $\nu = -1.001$ .

and is associated to the right-hand side  $f = -\operatorname{div}(s k \nabla \tilde{u}) \in L^2(\Omega)$ . It can be easily seen that, for  $\alpha \in \{e, i\}$ ,  $\tilde{u}|_{\Omega_\alpha} \in H^{1+l}(\Omega_\alpha)$  for any  $l > 0$ , meaning that  $p = k$  in Lemma 3.1.

As in [6], we choose  $k_e = 1$  and  $k_i = 1.001$ , so that  $\nu = -1.001$ . We run the computations on a sequence of non-symmetric (with respect to the interface  $\Gamma$ ), unstructured meshes. Results are depicted on Figure 1 for linear ( $k = 1$ ) and second-order ( $k = 2$ ) FEM, for our approach (solid blue), and for classical FEM (dashed black) applied to (2). For our approach, the parameter  $\lambda(h)$  is chosen as  $O(h^\beta)$  with  $\beta \leq (2k+1)$ , as required by Theorem 3.2 and Lemma 3.1. Here we choose  $\lambda(h) = 0.002 h^3$  for  $k = 1$  and  $\lambda(h) = 0.002 h^{4.2}$  for  $k = 2$ . We first observe that all convergence plots for our approach are strictly monotone, as opposed to classical FEM. In the  $L^2$ -norm, for  $h$  sufficiently small, we observe a slightly super-convergent behavior for our approach for both  $k = 1$  and  $k = 2$ . In the  $H^1$ -norm, for both orders, our method seems to reach the expected convergence rates for  $h$  sufficiently small, and clearly outperforms classical FEM. The choice of a small multiplicative coefficient 0.002 in the function  $\lambda$  is guided by the fact that the norm of  $\tilde{g}$  is big for such a contrast. Hence, for coarse meshes for which convergence is far from being reached, one has to give less weight to the second term of the functional  $J_h$  in the minimization process, to have a chance to recover a correct approximation of  $\tilde{g}$ . This weight is useless when  $h$  is sufficiently small.

*Remark 2 (Stabilization parameter)* For coarse meshes, one has to tune the parameter  $\lambda(h)$  in a nontrivial, contrast-depending way. This has to do with the regularization we use. If the  $L^2$ -orthogonal projection  $\Pi_h^k \tilde{g}$  of  $\tilde{g}$  onto  $F^k(\Gamma_h)$  was known, then we would stabilize the functional  $J_h$  with  $\mu(h) \|g_h - \Pi_h^k \tilde{g}\|_{0,\Gamma}^2$  instead of  $\lambda(h) \|g_h\|_{0,\Gamma}^2$ , yielding, in view of (18), an optimal estimate on  $\|\tilde{g}_h - \Pi_h^k \tilde{g}\|_{0,\Gamma}$ .

*Remark 3* It has to be noted that the symmetric cavity test-case can be analyzed by means of  $T$ -coercivity. However, it cannot be on general meshes, which is the main outcome of our approach. Furthermore, there are also many cases (even in 2D) that cannot be analyzed, at the continuous level, using the  $T$ -coercivity approach. This is for example the case of the cloaking device studied in [10, Section 4]. In that case, the operator associated to the problem, viewed as an operator from  $H_0^1(\Omega)$  to  $H^{-1}(\Omega)$  is not Fredholm, as its range is not closed.  $T$ -coercivity is hence inapplicable. However, it is known that, for compatible right-hand sides  $f \in L^2(\Omega)$ , the solution exists and is unique in  $H_0^1(\Omega)$  [10]. The convergence of our numerical method for compatible loadings as well as numerical experiments for such problems will be investigated in a forthcoming work [1].

## References

- [1] A. Abdulle and S. Lemaire. An optimization-based method for sign-changing PDEs. In preparation.



- [2] C. Bernardi, Y. Maday, and F. Rapetti. *Discrétisations variationnelles de problèmes aux limites elliptiques*, volume 45 of *Mathématiques & Applications*. Springer-Verlag, Berlin, 2004.
- [3] A.-S. Bonnet-Ben Dhia, C. Carvalho, and P. Ciarlet Jr. Mesh requirements for the finite element approximation of problems with sign-changing coefficients. Submitted (2016). Preprint hal-01335153.
- [4] A.-S. Bonnet-Ben Dhia, L. Chesnel, and P. Ciarlet Jr. T-coercivity for scalar interface problems between dielectrics and metamaterials. *ESAIM: Math. Model. Numer. Anal. (M2AN)*, 46:1363–1387, 2012.
- [5] S. C. Brenner. Poincaré–Friedrichs inequalities for piecewise  $H^1$  functions. *SIAM J. Numer. Anal.*, 41(1):306–324, 2003.
- [6] L. Chesnel and P. Ciarlet Jr. T-coercivity and continuous Galerkin methods: application to transmission problems with sign-changing coefficients. *Numer. Math.*, 124:1–29, 2013.
- [7] P. G. Ciarlet. *The finite element method for elliptic problems*, volume 40 of *Classics in Applied Mathematics*. Society for Industrial and Applied Mathematics (SIAM), Philadelphia, PA, 2002. Reprint of the 1978 original [North-Holland, Amsterdam].
- [8] M. D. Gunzburger, M. Heinkenschloss, and H. K. Lee. Solution of elliptic partial differential equations by an optimization-based domain decomposition method. *Appl. Math. Comput.*, 113(2-3):111–139, 2000.
- [9] M. D. Gunzburger, J. S. Peterson, and H. K. Lee. An optimization-based domain decomposition method for partial differential equations. *Comput. Math. Appl.*, 37(10):77–93, 1999.
- [10] H.-M. Nguyen. Negative index materials and their applications: recent mathematics progress. *Chinese Annals of Mathematics*. To appear.
- [11] H.-M. Nguyen. Asymptotic behavior of solutions to the Helmholtz equations with sign-changing coefficients. *Transactions of the American Mathematical Society*, 367:6581–6595, 2015.
- [12] H.-M. Nguyen. Limiting absorption principle and well-posedness for the Helmholtz equation with sign-changing coefficients. *Journal de Mathématiques Pures et Appliquées*, 106(2):342–374, 2016.
- [13] S. Nicaise and J. Venel. A posteriori error estimates for a finite element approximation of transmission problems with sign-changing coefficients. *J. Comput. Appl. Math.*, 235:4272–4282, 2011.