



# First-order least-squares method for the obstacle problem

Thomas Führer<sup>1</sup>

Received: 29 October 2018 / Revised: 14 May 2019 / Published online: 22 October 2019  
© Springer-Verlag GmbH Germany, part of Springer Nature 2019

## Abstract

We define and analyse a least-squares finite element method for a first-order reformulation of the obstacle problem. Moreover, we derive variational inequalities that are based on similar but non-symmetric bilinear forms. A priori error estimates including the case of non-conforming convex sets are given and optimal convergence rates are shown for the lowest-order case. We provide a posteriori bounds that can be used as error indicators in an adaptive algorithm. Numerical studies are presented.

**Mathematics Subject Classification** 65N30 · 65N12 · 49J40

## 1 Introduction

Many physical problems are of obstacle type, or more generally, described by variational inequalities [25,29]. In this article we consider, as a model problem, the classical obstacle problem where one seeks the equilibrium position of an elastic membrane constrained to lie over an obstacle. Another important example of an elliptic obstacle problem is the bending of a plate over an obstacle.

There exists already a long history of numerical methods, in particular finite element methods, see e.g., the books [16,17] for an overview on the topic. However, the literature on least-squares methods for obstacle problems is scarce. In fact, until the writing of this paper only [9] was available for the classical obstacle problem where the idea goes back to a Nitsche-based method for contact problems introduced and analyzed in [11]. An analysis of first-order least-squares finite element methods for Signorini problems can be found in [1] and more recently [26]. Let us also mention the pioneering work [14] for the a priori analysis of a classical finite element scheme. Newer articles include [18–20] where mixed and stabilized methods are considered.

Least-squares finite element methods are a widespread class of numerical schemes and their basic idea is to approximate the solution by minimizing the residual in some

---

✉ Thomas Führer  
tofuhrer@mat.uc.cl

<sup>1</sup> Facultad de Matemáticas, Pontificia Universidad Católica de Chile, Santiago, Chile

given norm. Let us recall some important properties of least-squares finite element methods, a detailed list is given in the introduction of the overview article [5], see also the book [6].

- *Unconstrained stability* One feature of least-squares schemes is that the methods are stable for all pairings of discrete spaces.
- *Adaptivity* Another feature is that a posteriori error bounds are obtained by simply evaluating the least-squares functional. For instance, standard least-squares methods for the Poisson problem [6] are based on minimizing residuals in  $L^2$  norms, which can be localized and, then, be used as error indicators in an adaptive algorithm.

The main purpose of this paper is to close the gap in the literature and define least-squares based methods for the obstacle problem. In particular, we want to study if the aforementioned properties transfer to the case of obstacle problems. Let us shortly describe the functional our method is based on. For simplicity assume a zero obstacle (the remainder of the paper deals with general non-zero obstacles). Then, the problem reads

$$-\Delta u \geq f, \quad u \geq 0, \quad (-\Delta u - f)u = 0$$

in some domain  $\Omega$  and  $u|_{\partial\Omega} = 0$ . Introducing the Lagrange multiplier (or reaction force)  $\lambda = -\Delta u - f$  and  $\sigma = \nabla u$ , we rewrite the problem as a first-order system, see also [2,3,9,18],

$$-\operatorname{div} \sigma - \lambda = f, \quad \sigma - \nabla u = 0, \quad u \geq 0, \quad \lambda \geq 0, \quad \lambda u = 0.$$

Note that  $f \in L^2(\Omega)$  does not imply more regularity for  $u$  so that  $\lambda \in H^{-1}(\Omega)$  lives only in the dual space in general. However, observe that  $\operatorname{div} \sigma + \lambda = -f \in L^2(\Omega)$  and therefore the functional

$$J((u, \sigma, \lambda); f) := \|\operatorname{div} \sigma + \lambda + f\|^2 + \|\nabla u - \sigma\|^2 + \langle \lambda, u \rangle,$$

where  $\langle \cdot, \cdot \rangle$  denotes a duality pairing, is well-defined for  $\operatorname{div} \sigma + \lambda \in L^2(\Omega)$ . We will show that minimizing  $J$  over a convex set with the additional constraints  $u \geq 0$ ,  $\lambda \geq 0$  is equivalent to solving the obstacle problem. We will consider the variational inequality associated to this problem with corresponding bilinear form  $a(\cdot, \cdot)$ . An issue that arises is that  $a(\cdot, \cdot)$  is not necessarily coercive. However, as it turns out, a simple scaling of the first term in the functional ensures coercivity on the whole space. In view of the aforementioned properties, this means that our method is *unconstrained stable*. The recent work [18] based on a Lagrange formulation (without reformulation to a first-order system) considers augmenting the trial spaces with bubble functions (mixed method) resp. adding residual terms (stabilized method) to obtain stability. The authors extended their work also to plate-bending problems, see [20].

Another motivation of the proposed first-order reformulation is that it allows to simultaneously approximate displacements and stresses. In many problems of structural engineering the stress is usually the primary quantity of interest. For the present

problem of an elastic membrane the stress is directly related to the gradient and for the problem of bending a plate over an obstacle the physical quantities of interest are the bending moments.

Furthermore, we will see that the functional  $J$  evaluated at some discrete approximation  $(u_h, \sigma_h, \lambda_h)$  with  $u_h, \lambda_h \geq 0$  is an upper bound for the error. Note that for  $\lambda_h \in L^2(\Omega)$  the duality  $\langle \lambda_h, u_h \rangle$  reduces to the  $L^2$  inner product. Thus, all the terms in the functional can be localized and used as error indicators.

Additionally, we will derive and analyse other variational inequalities that are also based on the first-order reformulation. The resulting methods are quite similar to the least-squares scheme since they share the same residual terms. The only difference is that the complementary condition  $\lambda u = 0$  is incorporated in a different, non-symmetric, way. We will present a uniform analysis that covers the least-squares formulation and the novel variational inequalities of the obstacle problem.

Finally, we point out that the use of adaptive schemes for obstacle problems is quite natural. First, the solutions may suffer from singularities stemming from the geometry, and second, the free boundary is a priori unknown. There exists plenty of literature on a posteriori estimators resp. adaptivity for finite elements methods for the obstacle problem, see e.g., [4,7,10,27,28,31,32] to name a few. Many of the estimators are based on the use of a discrete Lagrange multiplier which is obtained in a postprocessing step. In contrast, our proposed methods simultaneously approximate the Lagrange multiplier. This allows for a simple analysis of reliable a posteriori bounds.

## 1.1 Outline

The remainder of the paper is organized as follows. In Sect. 2 we describe the model problem, introduce the corresponding first-order system and based on that reformulation define our least-squares method. Then, Sect. 3 deals with the definition and analysis of different variational inequalities. In Sect. 4 we provide an a posteriori analysis and numerical studies are presented in Sect. 5. The appendix contains an example, which shows that  $a(\cdot, \cdot)$  is not coercive in general, and proofs of some auxiliary results.

## 2 Least-squares method

In Sects. 2.1 and 2.2 we describe the model problem and introduce the reader to our notation. Then, Sect. 2.3 is devoted to the definition and analysis of a least-squares functional.

### 2.1 Model problem

Let  $\Omega \subset \mathbb{R}^n$ ,  $n = 2, 3$  denote a polygonal Lipschitz domain with boundary  $\Gamma = \partial\Omega$ . For given  $f \in L^2(\Omega)$  and  $g \in H^1(\Omega)$  with  $g|_{\Gamma} \leq 0$  we consider the classical obstacle problem: find a solution  $u$  to

$$-\Delta u \geq f \quad \text{in } \Omega, \quad (1a)$$

$$u \geq g \quad \text{in } \Omega, \quad (1b)$$

$$(u - g)(-\Delta u - f) = 0 \quad \text{in } \Omega, \quad (1c)$$

$$u = 0 \quad \text{on } \Gamma. \quad (1d)$$

It is well-known that this problem admits a unique solution  $u \in H_0^1(\Omega)$ , and it can be equivalently characterized by the variational inequality: find  $u \in H_0^1(\Omega)$ ,  $u \geq g$  such that

$$\int_{\Omega} \nabla u \cdot \nabla(v - u) \, dx \geq \int_{\Omega} f(v - u) \, dx \quad \text{for all } v \in H_0^1(\Omega), v \geq g, \quad (2)$$

see [25]. For a more detailed description of the involved function spaces we refer to Sect. 2.2 below.

## 2.2 Notation and function spaces

We use the common notation for Sobolev spaces  $H_0^1(\Omega)$ ,  $H^s(\Omega)$  ( $s > 0$ ). Let  $(\cdot, \cdot)$  denote the  $L^2(\Omega)$  inner product, which induces the norm  $\|\cdot\|$ . The dual of  $H_0^1(\Omega)$  is denoted by  $H^{-1}(\Omega) := (H_0^1(\Omega))^*$ , where duality  $\langle \cdot, \cdot \rangle$  is understood with respect to the extended  $L^2(\Omega)$  inner product. We equip  $H^{-1}(\Omega)$  with the dual norm

$$\|\lambda\|_{-1} := \sup_{0 \neq v \in H_0^1(\Omega)} \frac{\langle \lambda, v \rangle}{\|\nabla v\|}.$$

Recall Friedrichs' inequality

$$\|u\| \leq C_F \|\nabla v\| \quad \text{for } v \in H_0^1(\Omega),$$

where  $0 < C_F = C_F(\Omega) \leq \text{diam}(\Omega)$ . Thus, by definition we have  $\|\lambda\|_{-1} \leq C_F \|\lambda\|$  for  $\lambda \in L^2(\Omega)$ .

Let  $\text{div} : L^2(\Omega) := L^2(\Omega)^n \rightarrow H^{-1}(\Omega)$  denote the generalized divergence operator, i.e.,  $\langle \text{div } \sigma, u \rangle := -(\sigma, \nabla u)$  for all  $\sigma \in L^2(\Omega)$ ,  $u \in H_0^1(\Omega)$ . This operator is bounded,

$$\|\text{div } \sigma\|_{-1} = \sup_{0 \neq v \in H_0^1(\Omega)} \frac{\langle \text{div } \sigma, v \rangle}{\|\nabla v\|} = \sup_{0 \neq v \in H_0^1(\Omega)} \frac{-(\sigma, \nabla v)}{\|\nabla v\|} \leq \|\sigma\|.$$

Let  $v \in H^1(\Omega)$ . We say  $v \geq 0$  if  $v \geq 0$  a.e. in  $\Omega$ . Moreover,  $\lambda \geq 0$  for  $\lambda \in H^{-1}(\Omega)$  means that  $\langle \lambda, v \rangle \geq 0$  for all  $v \in H_0^1(\Omega)$  with  $v \geq 0$ . We also interpret  $v \geq w$  as  $v - w \geq 0$  for  $v, w \in H^1(\Omega)$ .

Define the space

$$V := H_0^1(\Omega) \times L^2(\Omega) \times H^{-1}(\Omega)$$

with norm

$$\|v\|_V^2 := \|\nabla v\|^2 + \|\tau\|^2 + \|\mu\|_{-1}^2 \quad \text{for } v = (v, \tau, \mu) \in V$$

and the space

$$U := \{(u, \sigma, \lambda) \in V : \operatorname{div} \sigma + \lambda \in L^2(\Omega)\}$$

with norm

$$\|u\|_U^2 := \|\nabla u\|^2 + \|\sigma\|^2 + \|\operatorname{div} \sigma + \lambda\|^2 \quad \text{for } u = (u, \sigma, \lambda) \in U.$$

Observe that  $\|\cdot\|_U$  is a stronger norm than  $\|\cdot\|_V$ , i.e.,

$$\begin{aligned} \|\nabla u\|^2 + \|\sigma\|^2 + \|\lambda\|_{-1}^2 &\leq \|\nabla u\|^2 + \|\sigma\|^2 + 2\|\operatorname{div} \sigma + \lambda\|_{-1}^2 + 2\|\operatorname{div} \sigma\|_{-1}^2 \\ &\leq \|\nabla u\|^2 + 3\|\sigma\|^2 + 2C_F^2\|\operatorname{div} \sigma + \lambda\|^2. \end{aligned}$$

Our first least-squares formulation will be based on the minimization over the non-empty, convex and closed subset

$$K^s := \{(u, \sigma, \lambda) \in U : u - g \geq 0, \lambda \geq 0\},$$

where  $g$  is the given obstacle function. We will also derive and analyse variational inequalities based on non-symmetric bilinear forms that utilize the sets

$$\begin{aligned} K^0 &:= \{(u, \sigma, \lambda) \in U : u - g \geq 0\}, \\ K^1 &:= \{(u, \sigma, \lambda) \in U : \lambda \geq 0\}. \end{aligned}$$

Clearly,  $K^s \subset K^j$  for  $j = 0, 1$ .

We write  $A \lesssim B$  if there exists a constant  $C > 0$ , independent of quantities of interest, such that  $A \leq CB$ . Analogously we define  $A \gtrsim B$ . If  $A \lesssim B$  and  $B \lesssim A$  hold then we write  $A \simeq B$ .

### 2.3 Least-squares functional

Let  $u \in H_0^1(\Omega)$  denote the unique solution of the obstacle problem (1). Define  $\lambda := -\Delta u - f \in H^{-1}(\Omega)$  and  $\sigma := \nabla u$ . Problem (1) can equivalently be written as the first-order problem

$$-\operatorname{div} \sigma - \lambda = f \quad \text{in } \Omega, \tag{3a}$$

$$\sigma - \nabla u = 0 \quad \text{in } \Omega, \tag{3b}$$

$$u \geq g \quad \text{in } \Omega, \tag{3c}$$

$$\lambda \geq 0 \quad \text{in } \Omega, \tag{3d}$$

$$(u - g)\lambda = 0 \quad \text{in } \Omega, \tag{3e}$$

$$u = 0 \quad \text{on } \Gamma. \quad (3f)$$

Observe that  $\operatorname{div} \boldsymbol{\sigma} + \lambda \in L^2(\Omega)$  and that the unique solution  $\mathbf{u} = (u, \boldsymbol{\sigma}, \lambda) \in U$  satisfies  $\mathbf{u} \in K^s$ . We consider the functional

$$J(\mathbf{u}; f, g) := \|\operatorname{div} \boldsymbol{\sigma} + \lambda + f\|^2 + \|\nabla u - \boldsymbol{\sigma}\|^2 + \langle \lambda, u - g \rangle$$

for  $\mathbf{u} = (u, \boldsymbol{\sigma}, \lambda) \in U$ ,  $f \in L^2(\Omega)$ ,  $g \in H_0^1(\Omega)$  and the minimization problem: find  $\mathbf{u} \in K^s$  with

$$J(\mathbf{u}; f, g) = \min_{\mathbf{v} \in K^s} J(\mathbf{v}; f, g). \quad (4)$$

Note that the definition of the functional only makes sense if  $g \in H_0^1(\Omega)$ .

**Theorem 1** *If  $f \in L^2(\Omega)$ ,  $g \in H_0^1(\Omega)$ , then problems (3) and (4) are equivalent. In particular, there exists a unique solution  $\mathbf{u} \in K^s$  of (4) and it holds that*

$$J(\mathbf{v}; f, g) \geq C_J \|\mathbf{v} - \mathbf{u}\|_U^2 \quad \text{for all } \mathbf{v} \in K^s. \quad (5)$$

The constant  $C_J > 0$  depends only on  $\Omega$ .

**Proof** Let  $\mathbf{u} := (u, \boldsymbol{\sigma}, \lambda) = (u, \nabla u, -\Delta u - f) \in K^s$  denote the unique solution of (3). Observe that  $J(\mathbf{v}; f, g) \geq 0$  for all  $\mathbf{v} \in K^s$  and  $J(\mathbf{u}; f, g) = 0$ , thus,  $\mathbf{u}$  minimizes the functional. Suppose (5) holds and that  $\mathbf{u}^* \in K^s$  is another minimizer. Then, (5) proves that  $\mathbf{u} = \mathbf{u}^*$ . It only remains to show (5). Let  $\mathbf{v} = (v, \boldsymbol{\tau}, \mu) \in K^s$ . Note that all terms in  $J(\mathbf{v}; f, g)$  are non-negative. Since  $f = -\operatorname{div} \boldsymbol{\sigma} - \lambda$  and  $\nabla u - \boldsymbol{\sigma} = 0$  we have with the constant  $C_F > 0$  that

$$\begin{aligned} J(\mathbf{v}; f, g) &= \|\operatorname{div}(\boldsymbol{\tau} - \boldsymbol{\sigma}) + (\mu - \lambda)\|^2 + \|\nabla(v - u) - (\boldsymbol{\tau} - \boldsymbol{\sigma})\|^2 + \langle \mu, v - g \rangle \\ &= \frac{1}{1 + C_F^2} \left( (1 + C_F^2) \|\operatorname{div}(\boldsymbol{\tau} - \boldsymbol{\sigma}) + (\mu - \lambda)\|^2 \right. \\ &\quad \left. + (1 + C_F^2) \|\nabla(v - u) - (\boldsymbol{\tau} - \boldsymbol{\sigma})\|^2 \right. \\ &\quad \left. + (1 + C_F^2) \langle \mu, v - g \rangle \right) \\ &\geq \frac{1}{1 + C_F^2} \left( (1 + C_F^2) \|\operatorname{div}(\boldsymbol{\tau} - \boldsymbol{\sigma}) + (\mu - \lambda)\|^2 \right. \\ &\quad \left. + \|\nabla(v - u) - (\boldsymbol{\tau} - \boldsymbol{\sigma})\|^2 + \langle \mu, v - g \rangle \right). \end{aligned}$$

Moreover,  $\langle \lambda, u - g \rangle = 0$  and  $\langle \lambda, v - g \rangle \geq 0$ ,  $\langle \mu, u - g \rangle \geq 0$ . We estimate

$$\begin{aligned} \langle \mu, v - g \rangle &= \langle \mu, v - u \rangle + \langle \mu, u - g \rangle + \langle \lambda, u - g \rangle \\ &\geq \langle \mu, v - u \rangle + \langle \lambda, u - g \rangle + \langle \lambda, g - v \rangle \\ &= \langle \mu, v - u \rangle + \langle \lambda, u - v \rangle = \langle \mu - \lambda, v - u \rangle. \end{aligned}$$

Define  $\mathbf{w} := (w, \boldsymbol{\chi}, v) := \mathbf{v} - \mathbf{u}$ . Then, the Cauchy–Schwarz inequality, Young’s inequality and the definition of the divergence operator yield

$$\begin{aligned} J(\mathbf{v}; f, g) &\gtrsim (1 + C_F^2) \|\operatorname{div}(\boldsymbol{\tau} - \boldsymbol{\sigma}) + (\mu - \lambda)\|^2 \\ &\quad + \|\nabla(v - u) - (\boldsymbol{\tau} - \boldsymbol{\sigma})\|^2 + \langle \mu, v - g \rangle \\ &\geq (1 + C_F^2) \|\operatorname{div} \boldsymbol{\chi} + v\|^2 + \|\nabla w - \boldsymbol{\chi}\|^2 + \langle v, w \rangle \\ &= (1 + C_F^2) \|\operatorname{div} \boldsymbol{\chi} + v\|^2 + \|\nabla w\|^2 \\ &\quad + \|\boldsymbol{\chi}\|^2 - \langle \nabla w, \boldsymbol{\chi} \rangle + \langle \operatorname{div} \boldsymbol{\chi}, w \rangle + \langle v, w \rangle \\ &\geq (1 + C_F^2) \|\operatorname{div} \boldsymbol{\chi} + v\|^2 + \frac{1}{2} \|\nabla w\|^2 + \frac{1}{2} \|\boldsymbol{\chi}\|^2 + \langle \operatorname{div} \boldsymbol{\chi} + v, w \rangle. \end{aligned}$$

Application of the Cauchy–Schwarz inequality, Friedrichs’ inequality and Young’s inequality gives us for the last term and  $\delta > 0$

$$|\langle \operatorname{div} \boldsymbol{\chi} + v, w \rangle| \leq C_F \|\operatorname{div} \boldsymbol{\chi} + v\| \|\nabla w\| \leq C_F^2 \frac{\delta^{-1}}{2} \|\operatorname{div} \boldsymbol{\chi} + v\|^2 + \frac{\delta}{2} \|\nabla w\|^2.$$

Putting altogether and choosing  $\delta = \frac{1}{2}$  we end up with

$$\begin{aligned} J(\mathbf{v}; f, g) &\gtrsim (1 + C_F^2) \|\operatorname{div} \boldsymbol{\chi} + v\|^2 + \|\nabla w - \boldsymbol{\chi}\|^2 + \langle \mu, v - g \rangle \\ &\geq (1 + C_F^2) \|\operatorname{div} \boldsymbol{\chi} + v\|^2 + \|\nabla w - \boldsymbol{\chi}\|^2 + \langle v, w \rangle \\ &\geq \|\operatorname{div} \boldsymbol{\chi} + v\|^2 + \frac{1}{4} \|\nabla w\|^2 + \frac{1}{2} \|\boldsymbol{\chi}\|^2 \simeq \|\mathbf{w}\|_U^2 = \|\mathbf{v} - \mathbf{u}\|_U^2, \end{aligned}$$

which finishes the proof. □

**Remark 2** Note that (5) measures the error of any function  $\mathbf{v} \in K^s$ , in particular, it can be used as a posteriori error estimator when  $\mathbf{v} \in K_h^s \subset K^s$  is a discrete approximation. However, in practice the condition  $K_h^s \subset K^s$  is often hard to realize. Below we introduce a simple scaling of the first term in the least-squares functional that allows us to prove coercivity of the associated bilinear form on the whole space  $U$ .

For given  $f \in L^2(\Omega)$ ,  $g \in H_0^1(\Omega)$ , and fixed parameter  $\beta > 0$  define the bilinear form  $a_\beta: U \times U \rightarrow \mathbb{R}$  and the functional  $F_\beta: U \rightarrow \mathbb{R}$  by

$$a_\beta(\mathbf{u}, \mathbf{v}) := \beta(\operatorname{div} \boldsymbol{\sigma} + \lambda, \operatorname{div} \boldsymbol{\tau} + \mu) + \langle \nabla u - \boldsymbol{\sigma}, \nabla v - \boldsymbol{\tau} \rangle + \frac{1}{2}(\langle \mu, u \rangle + \langle \lambda, v \rangle), \tag{6}$$

$$F_\beta(\mathbf{v}) := -\beta(f, \operatorname{div} \boldsymbol{\tau} + \mu) + \frac{1}{2} \langle \mu, g \rangle \tag{7}$$

for all  $\mathbf{u} = (u, \boldsymbol{\sigma}, \lambda)$ ,  $\mathbf{v} = (v, \boldsymbol{\tau}, \mu) \in U$ . We stress that  $a_1(\cdot, \cdot)$  and  $F_1(\cdot)$  induce the functional  $J(\cdot; \cdot)$ , i.e.,

$$J(\mathbf{u}; f, g) = a_1(\mathbf{u}, \mathbf{u}) - 2F_1(\mathbf{u}) + (f, f).$$

Since  $J$  is differentiable it is well-known that the solution  $\mathbf{u} \in K^s$  of (4) satisfies the variational inequality

$$a_1(\mathbf{u}, \mathbf{v} - \mathbf{u}) \geq F_1(\mathbf{v} - \mathbf{u}) \quad \text{for all } \mathbf{v} \in K^s. \quad (8)$$

Conversely, if  $J$  is also convex in  $K^s$ , then any solution of (8) solves (4). However,  $J$  is convex on  $K^s$  iff  $a_1(\mathbf{v} - \mathbf{w}, \mathbf{v} - \mathbf{w}) \geq 0$  for all  $\mathbf{v}, \mathbf{w} \in K^s$ , which is not true in general, see the example from ‘‘Appendix A’’. In Sect. 3 below we will show that for sufficiently large  $\beta > 1$  the bilinear form  $a_\beta(\cdot, \cdot)$  is coercive, even on the whole space  $U$ . This has the advantage that we can prove unique solvability of the continuous problem and its discretization simultaneously. More important, in practice this allows the use of non-conforming subsets  $K_h^s \not\subseteq K^s$ .

### 3 Variational inequalities

In this section we introduce and analyse different variational inequalities. The idea of including the complementary condition in different ways has also been used in [15] to derive DPG methods for contact problems.

We define the bilinear forms  $b_\beta, c_\beta: U \times U \rightarrow \mathbb{R}$  and functionals  $G_\beta, H_\beta$  by

$$\begin{aligned} b_\beta(\mathbf{u}, \mathbf{v}) &:= \beta(\operatorname{div} \boldsymbol{\sigma} + \lambda, \operatorname{div} \boldsymbol{\tau} + \mu) + (\nabla u - \boldsymbol{\sigma}, \nabla v - \boldsymbol{\tau}) + \langle \lambda, v \rangle, \\ c_\beta(\mathbf{u}, \mathbf{v}) &:= \beta(\operatorname{div} \boldsymbol{\sigma} + \lambda, \operatorname{div} \boldsymbol{\tau} + \mu) + (\nabla u - \boldsymbol{\sigma}, \nabla v - \boldsymbol{\tau}) + \langle \mu, u \rangle, \\ G_\beta(\mathbf{v}) &:= -\beta(f, \operatorname{div} \boldsymbol{\tau} + \mu) \\ H_\beta(\mathbf{v}) &:= -\beta(f, \operatorname{div} \boldsymbol{\tau} + \mu) + \langle \mu, g \rangle. \end{aligned}$$

Let  $\mathbf{u} = (u, \boldsymbol{\sigma}, \lambda) \in K^s \subset K^j$  ( $j = 0, 1$ ) denote the unique solution of (3) with  $f \in L^2(\Omega)$ ,  $g \in H_0^1(\Omega)$ . Recall that  $\operatorname{div} \boldsymbol{\sigma} + \lambda = -f$ . Testing this identity with  $\operatorname{div} \boldsymbol{\tau} + \mu$ , multiplying with  $(\beta - 1)$  and adding it to (8) we see that the solution  $\mathbf{u} \in K^s$  satisfies the variational inequality

$$a_\beta(\mathbf{u}, \mathbf{v} - \mathbf{u}) \geq F_\beta(\mathbf{v} - \mathbf{u}) \quad \text{for all } \mathbf{v} \in K^s. \quad (\text{VIa})$$

For the derivation of our second variational inequality let  $\mathbf{u} = (u, \boldsymbol{\sigma}, \lambda) \in K^0$  denote the unique solution of (3) with  $f \in L^2(\Omega)$ ,  $g \in H^1(\Omega)$ ,  $g|_\Gamma \leq 0$ . Recall that  $\lambda = -\Delta u - f$ . By (2) we have that

$$\langle \lambda, v - u \rangle = (\nabla u, \nabla(v - u)) - (f, v - u) \geq 0$$

for all  $v \in H_0^1(\Omega)$ ,  $v \geq g$ . Thus,  $\mathbf{u} \in K^0$  satisfies the variational inequality

$$b_\beta(\mathbf{u}, \mathbf{v} - \mathbf{u}) \geq G_\beta(\mathbf{v} - \mathbf{u}) \quad \text{for all } \mathbf{v} \in K^0. \quad (\text{VIb})$$

Our final method is based on the observation that for  $\mu \geq 0$ , we have that  $\langle \mu, u - g \rangle \geq 0$  for  $u \geq g \in H_0^1(\Omega)$ . Together with  $\langle \lambda, u - g \rangle = 0$  we conclude  $\langle \mu - \lambda, u - g \rangle \geq 0$ . Thus,  $\mathbf{u} \in K^1$  satisfies the variational inequality



$$c_\beta(\mathbf{u}, \mathbf{v} - \mathbf{u}) \geq H_\beta(\mathbf{v} - \mathbf{u}) \quad \text{for all } \mathbf{v} \in K^1. \tag{VIc}$$

Note that  $a_\beta$  is symmetric, whereas  $b_\beta, c_\beta$  are not.

### 3.1 Solvability

In what follows we analyse the (unique) solvability of the variational inequalities (VIa)–(VIc) in a uniform manner (including discretizations).

**Lemma 3** *Suppose  $\beta > 0$ . Let  $A \in \{a_\beta, b_\beta, c_\beta\}$ . There exists  $C_\beta > 0$  depending only on  $\beta > 0$  and  $\Omega$  such that*

$$|A(\mathbf{u}, \mathbf{v})| \leq C_\beta \|\mathbf{u}\|_U \|\mathbf{v}\|_U \quad \text{for all } \mathbf{u}, \mathbf{v} \in U.$$

If  $\beta \geq 1 + C_F^2$ , then  $A$  is coercive, i.e.,

$$C \|\mathbf{u}\|_U^2 \leq A(\mathbf{u}, \mathbf{u}) \quad \text{for all } \mathbf{u} \in U.$$

The constant  $C > 0$  is independent of  $\beta$  and  $\Omega$ .

**Proof** We prove boundedness of  $A = a_\beta$ . Let  $\mathbf{u} = (u, \boldsymbol{\sigma}, \lambda), \mathbf{v} = (v, \boldsymbol{\tau}, \mu) \in U$  be given. The Cauchy–Schwarz inequality together with the Friedrichs’ inequality and boundedness of the divergence operator yields

$$\begin{aligned} |a_\beta(\mathbf{u}, \mathbf{v})| &\leq \beta \|\operatorname{div} \boldsymbol{\sigma} + \lambda\| \|\operatorname{div} \boldsymbol{\tau} + \mu\| + \|\nabla u - \boldsymbol{\sigma}\| \|\nabla v - \boldsymbol{\tau}\| \\ &\quad + \frac{1}{2} (\langle \operatorname{div} \boldsymbol{\tau} + \mu, u \rangle - \langle \operatorname{div} \boldsymbol{\tau}, u \rangle + \langle \operatorname{div} \boldsymbol{\sigma} + \lambda, v \rangle - \langle \operatorname{div} \boldsymbol{\sigma}, v \rangle) \\ &\leq \beta \|\operatorname{div} \boldsymbol{\sigma} + \lambda\| \|\operatorname{div} \boldsymbol{\tau} + \mu\| + \|\nabla u - \boldsymbol{\sigma}\| \|\nabla v - \boldsymbol{\tau}\| \\ &\quad + \frac{1}{2} ((C_F \|\operatorname{div} \boldsymbol{\tau} + \mu\| + \|\boldsymbol{\tau}\|) \|\nabla u\| + (C_F \|\operatorname{div} \boldsymbol{\sigma} + \lambda\| + \|\boldsymbol{\sigma}\|) \|\nabla v\|). \end{aligned}$$

This shows boundedness of  $a_\beta(\cdot, \cdot)$ . Similarly, one concludes boundedness of  $b_\beta(\cdot, \cdot)$  and  $c_\beta(\cdot, \cdot)$ .

For the proof of coercivity, observe that  $a_\beta(\mathbf{w}, \mathbf{w}) = b_\beta(\mathbf{w}, \mathbf{w}) = c_\beta(\mathbf{w}, \mathbf{w})$  for all  $\mathbf{w} \in U$ . We stress that coercivity directly follows from the arguments given in the proof of Theorem 1. Note that the choice of  $\beta$  yields

$$A(\mathbf{w}, \mathbf{w}) \geq (1 + C_F^2) \|\operatorname{div} \boldsymbol{\chi} + v\|^2 + \|\nabla w - \boldsymbol{\chi}\|^2 + \langle v, w \rangle$$

for  $\mathbf{w} = (w, \boldsymbol{\chi}, v) \in U$ . The right-hand side can be further estimated following the argumentation as in the proof of Theorem 1 which gives us

$$(1 + C_F^2) \|\operatorname{div} \boldsymbol{\chi} + v\|^2 + \|\nabla w - \boldsymbol{\chi}\|^2 + \langle v, w \rangle \gtrsim \|\mathbf{w}\|_U^2.$$

This finishes the proof. □

**Remark 4** Recall that  $C_F \leq \text{diam}(\Omega)$ . Therefore, we can always choose  $\beta = 1 + \text{diam}(\Omega)^2$  to ensure coercivity of our bilinear forms. We stress that a choice of  $\beta$  of order  $\text{diam}(\Omega)$  is not only sufficient to ensure coercivity but also necessary in general as the example from ‘‘Appendix A’’ shows. Another possibility is to rescale  $\Omega$  such that  $\text{diam}(\Omega) \leq 1$  which implies that we can choose  $\beta = 2$ . Furthermore, observe that a scaling of  $\Omega$  transforms (1) to an equivalent obstacle problem (with appropriate redefined functions  $f, g$ ). To be more precise, define  $\tilde{u}(x) := u(dx)$  with  $d := \text{diam}(\Omega) > 0$  and  $u \in H_0^1(\Omega)$  the solution of (1). Moreover, set  $\tilde{f}(x) = d^2 f(dx)$ ,  $\tilde{g}(x) := g(dx)$ . Then,  $\tilde{u}$  solves (1) in  $\tilde{\Omega} := \{x/d: x \in \Omega\}$  with  $f, g$  replaced by  $\tilde{f}, \tilde{g}$ .

The variational inequalities (VIa)–(VIc) are of the first kind and we use a standard framework for the analysis (Lions–Stampacchia theorem), see [16,17,25].

**Theorem 5** Suppose  $\beta \geq 1 + C_F^2$ . Let  $A \in \{a_\beta, b_\beta, c_\beta\}$  and let  $F : U \rightarrow \mathbb{R}$  denote a bounded linear functional. If  $K \subseteq U$  is a non-empty convex and closed subset, then the variational inequality

$$\text{Find } \mathbf{u} \in K \text{ s.t. } A(\mathbf{u}, \mathbf{v} - \mathbf{u}) \geq F(\mathbf{v} - \mathbf{u}) \text{ for all } \mathbf{v} \in K \tag{9}$$

admits a unique solution.

In particular, for  $f \in L^2(\Omega), g \in H_0^1(\Omega)$  each of the problems (VIa)–(VIc) has a unique solution and the problems are equivalent to (3).

**Proof** With the assumption on  $\beta$ , Lemma 3 proves that the bilinear forms are coercive and bounded. Then, unique solvability of (9) follows from the Lions–Stampacchia theorem, see e.g., [16,17,25].

Unique solvability of (VIa)–(VIc) follows since the functionals  $F_\beta, G_\beta, H_\beta$  are linear and bounded: For example, boundedness of  $F_\beta$  can be seen from

$$\begin{aligned} |F_\beta(\mathbf{v})| &= | -\beta(f, \text{div } \boldsymbol{\tau} + \mu) + \frac{1}{2}(\text{div } \boldsymbol{\tau} + \mu, g) - \frac{1}{2}\langle \text{div } \boldsymbol{\tau}, g \rangle | \\ &\leq \beta \|f\| \|\text{div } \boldsymbol{\tau} + \mu\| + \frac{1}{2} \|\text{div } \boldsymbol{\tau} + \mu\| \|g\| + \frac{1}{2} \|\boldsymbol{\tau}\| \|\nabla g\| \\ &\lesssim (\|f\| + \|\nabla g\|) \|\mathbf{v}\|_U. \end{aligned}$$

The same arguments prove that  $G_\beta$  and  $H_\beta$  are bounded.

Finally, equivalence to (3) follows since all problems admit unique solutions and by construction the solution of (3) also solves each of the problems (VIa)–(VIc).  $\square$

**Remark 6** We stress that the assumption  $g \in H_0^1(\Omega)$  is necessary. If  $g \in H^1(\Omega)$  then the term  $\langle \mu, g \rangle$  in  $F_\beta, H_\beta$  is not well-defined. However, this term does not appear in  $G_\beta$  and therefore the variational inequality in (VIb) admits a unique solution if we only assume  $g \in H^1(\Omega)$  with  $g|_\Gamma \leq 0$ .

**Remark 7** The variational inequality (VIa) corresponds to a least-squares finite element method with convex functional

$$J_\beta(\mathbf{u}; f, g) := a_\beta(\mathbf{u}, \mathbf{u}) - 2F_\beta(\mathbf{u}) + \beta(f, f).$$

Then, Theorem 5 proves that the problem

$$J_\beta(\mathbf{u}; f, g) = \min_{\mathbf{v} \in K} J_\beta(\mathbf{v}; f, g)$$

admits a unique solution for all non-empty convex and closed sets  $K \subseteq U$ . Moreover,  $J_\beta(\mathbf{u}; f, g) \simeq J(\mathbf{u}; f, g)$  for  $\mathbf{u} \in K^s$ , so that this problem is equivalent to (4) for  $K = K^s$ .

### 3.2 A priori analysis

The following three results provide general bounds on the approximation error. The proofs are based on standard arguments, see e.g., [14]. We give details for the proof of the first result, the others follow the same lines of argumentation and are left to the reader.

**Theorem 8** *Suppose  $\beta \geq 1 + C_F^2$ . Let  $\mathbf{u} \in K^s$  denote the solution of (VIa), where  $f \in L^2(\Omega)$ ,  $g \in H_0^1(\Omega)$ . Let  $K_h \subset U$  denote a non-empty convex and closed subset and let  $\mathbf{u}_h \in K_h$  denote the solution of (9) with  $A = a_\beta$ ,  $F = F_\beta$  and  $K = K_h$ . It holds that*

$$\begin{aligned} \|\mathbf{u} - \mathbf{u}_h\|_U^2 &\leq C_{\text{opt}} \left( \inf_{\mathbf{v}_h \in K_h} (\|\mathbf{u} - \mathbf{v}_h\|_U^2 + |\langle \lambda, \mathbf{v}_h - \mathbf{u} \rangle + \langle \mu_h - \lambda, \mathbf{u} - g \rangle|) \right. \\ &\quad \left. + \inf_{\mathbf{v} \in K^s} |\langle \lambda, \mathbf{v} - \mathbf{u}_h \rangle + \langle \mu - \lambda_h, \mathbf{u} - g \rangle| \right). \end{aligned}$$

The constant  $C_{\text{opt}} > 0$  depends only on  $\beta$  and  $\Omega$ .

**Proof** Throughout let  $\mathbf{v} = (v, \boldsymbol{\tau}, \mu) \in K^s$ ,  $\mathbf{v}_h = (v_h, \boldsymbol{\tau}_h, \mu_h) \in K_h$  and let  $\mathbf{u} = (u, \boldsymbol{\sigma}, \lambda) \in K^s$  denote the exact solution of (VIa). Thus,  $\text{div } \boldsymbol{\sigma} + \lambda + f = 0$  and  $\nabla u - \boldsymbol{\sigma} = 0$ . For arbitrary  $\mathbf{w} = (w, \boldsymbol{\chi}, \nu) \in U$  it holds that

$$\begin{aligned} a_\beta(\mathbf{u}, \mathbf{w}) &= \beta(\text{div } \boldsymbol{\sigma} + \lambda, \text{div } \boldsymbol{\chi} + \nu) + (\nabla u - \boldsymbol{\sigma}, \nabla w - \boldsymbol{\chi}) + \frac{1}{2}(\langle \lambda, w \rangle + \langle \nu, u \rangle) \\ &= -\beta(f, \text{div } \boldsymbol{\chi} + \nu) + \frac{1}{2}\langle \nu, g \rangle + \frac{1}{2}(\langle \lambda, w \rangle + \langle \nu, u - g \rangle) \\ &= F_\beta(\mathbf{w}) + \frac{1}{2}(\langle \lambda, w \rangle + \langle \nu, u - g \rangle). \end{aligned} \tag{10}$$

Using coercivity of  $a_\beta(\cdot, \cdot)$ , identity (10) and the fact that  $\mathbf{u}_h$  solves the discretized variational inequality (on  $K_h$ ) shows that

$$\begin{aligned} \|\mathbf{u} - \mathbf{u}_h\|_U^2 &\lesssim a_\beta(\mathbf{u} - \mathbf{u}_h, \mathbf{u} - \mathbf{u}_h) \\ &= a_\beta(\mathbf{u}, \mathbf{u} - \mathbf{u}_h) - a_\beta(\mathbf{u}_h, \mathbf{u} - \mathbf{v}_h) - a_\beta(\mathbf{u}_h, \mathbf{v}_h - \mathbf{u}_h) \\ &\leq F_\beta(\mathbf{u} - \mathbf{u}_h) + \frac{1}{2}(\langle \lambda, u - u_h \rangle + \langle \lambda - \lambda_h, u - g \rangle) \\ &\quad - a_\beta(\mathbf{u}_h, \mathbf{u} - \mathbf{v}_h) - F_\beta(\mathbf{v}_h - \mathbf{u}_h) \\ &= F_\beta(\mathbf{u} - \mathbf{v}_h) + \frac{1}{2}(\langle \lambda, u - u_h \rangle + \langle \lambda - \lambda_h, u - g \rangle) - a_\beta(\mathbf{u}_h, \mathbf{u} - \mathbf{v}_h). \end{aligned}$$

Note that  $0 = \langle \lambda, u - g \rangle \leq \langle \lambda, v - g \rangle$  and  $\langle \lambda, u - g \rangle \leq \langle \mu, u - g \rangle$ . Hence,

$$\begin{aligned} \langle \lambda, u - u_h \rangle + \langle \lambda - \lambda_h, u - g \rangle &= \langle \lambda, u - g + g - u_h \rangle + \langle \lambda - \lambda_h, u - g \rangle \\ &\leq \langle \lambda, v - g + g - u_h \rangle + \langle \mu - \lambda_h, u - g \rangle. \end{aligned}$$

This and identity (10) with  $w = u - v_h$  imply that

$$\begin{aligned} F_\beta(u - v_h) - a_\beta(u_h, u - v_h) + \frac{1}{2}(\langle \lambda, u - u_h \rangle + \langle \lambda - \lambda_h, u - g \rangle) \\ \leq a_\beta(u - u_h, u - v_h) - \frac{1}{2}(\langle \lambda, u - v_h \rangle + \langle \lambda - \mu_h, u - g \rangle) \\ + \frac{1}{2}(\langle \lambda, v - u_h \rangle + \langle \mu - \lambda_h, u - g \rangle). \end{aligned}$$

Putting altogether, boundedness of  $a_\beta(\cdot, \cdot)$  and an application of Young's inequality with parameter  $\delta > 0$  show that

$$\begin{aligned} \|u - u_h\|_U^2 &\lesssim \frac{\delta}{2} \|u - u_h\|_U^2 + \frac{\delta^{-1}}{2} \|u - v_h\|_U^2 + |\langle \lambda, v_h - u \rangle + \langle \mu_h - \lambda, u - g \rangle| \\ &\quad + |\langle \lambda, v - u_h \rangle + \langle \mu - \lambda_h, u - g \rangle|. \end{aligned}$$

Subtracting the term  $\delta/2 \|u - u_h\|_U^2$  for some sufficiently small  $\delta > 0$  finishes the proof since  $v \in K^s$ ,  $v_h \in K_h$  are arbitrary.  $\square$

**Theorem 9** Suppose  $\beta \geq 1 + C_F^2$ . Let  $u \in K^0$  denote the solution of (VIb), where  $f \in L^2(\Omega)$ ,  $g \in H^1(\Omega)$  with  $g|_\Gamma \leq 0$ . Let  $K_h \subset U$  denote a non-empty convex and closed subset and let  $u_h \in K_h$  denote the solution of (9) with  $A = b_\beta$ ,  $F = G_\beta$ , and  $K = K_h$ . It holds that

$$\|u - u_h\|_U^2 \leq C_{\text{opt}} \left( \inf_{v_h \in K_h} (\|u - v_h\|_U^2 + |\langle \lambda, v_h - u \rangle|) + \inf_{v \in K^0} |\langle \lambda, v - u_h \rangle| \right).$$

The constant  $C_{\text{opt}} > 0$  depends only on  $\beta$  and  $\Omega$ .

**Theorem 10** Suppose  $\beta \geq 1 + C_F^2$ . Let  $u \in K^1$  denote the solution of (VIc), where  $f \in L^2(\Omega)$ ,  $g \in H_0^1(\Omega)$ . Let  $K_h \subset U$  denote a non-empty convex and closed subset and let  $u_h \in K_h$  denote the solution of (9) with  $A = c_\beta$ ,  $F = H_\beta$ , and  $K = K_h$ . It holds that

$$\|u - u_h\|_U^2 \leq C_{\text{opt}} \left( \inf_{v_h \in K_h} (\|u - v_h\|_U^2 + |\langle \mu_h - \lambda, u - g \rangle|) + \inf_{v \in K^1} |\langle \mu - \lambda_h, u - g \rangle| \right).$$

The constant  $C_{\text{opt}} > 0$  depends only on  $\beta$  and  $\Omega$ .

### 3.3 Discretization

Let  $\mathcal{T}$  denote a regular triangulation of  $\Omega$ ,  $\bigcup_{T \in \mathcal{T}} \bar{T} = \bar{\Omega}$ . We assume that  $\mathcal{T}$  is  $\kappa$ -shape regular, i.e.,

$$\sup_{T \in \mathcal{T}} \frac{\text{diam}(T)^n}{|T|} \leq \kappa < \infty.$$

Moreover, let  $\mathcal{V}$  denote the vertices of the mesh  $\mathcal{T}$  and  $\mathcal{V}_0 := \mathcal{V} \setminus \Gamma$ . Let  $h_{\mathcal{T}} \in L^\infty(\Omega)$  denote the mesh-size function,  $h_{\mathcal{T}}|_T := h_T := \text{diam}(T)$  for  $T \in \mathcal{T}$ . Set  $h := \max_{T \in \mathcal{T}} \text{diam}(T)$ . We use standard finite element spaces for the discretization. Let  $\mathcal{P}^p(\mathcal{T})$  denote the space of  $\mathcal{T}$ -elementwise polynomials of degree less or equal than  $p \in \mathbb{N}_0$ . Let  $\mathcal{RT}^p(\mathcal{T})$  denote the Raviart–Thomas space of degree  $p \in \mathbb{N}_0$ ,  $\mathcal{S}_0^{p+1}(\mathcal{T}) := \mathcal{P}^{p+1}(\mathcal{T}) \cap H_0^1(\Omega)$ , and

$$U_{hp} := \mathcal{S}_0^{p+1}(\mathcal{T}) \times \mathcal{RT}^p(\mathcal{T}) \times \mathcal{P}^p(\mathcal{T}).$$

Clearly,  $U_{hp} \subset U$ . We stress that the polynomial degree is chosen, so that the best approximation in the norm  $\|\cdot\|_U$  is of order  $h^{p+1}$ .

To define admissible convex sets for the discrete variational inequalities we need to put constraints on functions from the space  $\mathcal{S}_0^{p+1}(\mathcal{T})$  or from  $\mathcal{P}^p(\mathcal{T})$  or both. Let us remark that for a polynomial degree  $\geq 2$  such constraints are not straightforward to implement. One possibility would be to impose such constraints pointwise and then analyse the consistency error. We comment on the case  $p = 1$  and  $n = 2$  below. For some  $hp$ -FEM method for elliptic obstacle problems we refer to [2,3]. In order to avoid such quite technical treatments and for a simpler representation of the basic ideas we consider from now on the lowest-order case only, where the linear constraints can easily be built in. To that end define the non-empty convex subsets

$$K_h^s := \{(v_h, \boldsymbol{\tau}_h, \mu_h) \in U_{h0} : \mu_h \geq 0, v_h(x) \geq g(x) \text{ for all } x \in \mathcal{V}_0\}, \tag{11a}$$

$$K_h^0 := \{(v_h, \boldsymbol{\tau}_h, \mu_h) \in U_{h0} : v_h(x) \geq g(x) \text{ for all } x \in \mathcal{V}_0\}, \tag{11b}$$

$$K_h^1 := \{(v_h, \boldsymbol{\tau}_h, \mu_h) \in U_{h0} : \mu_h \geq 0\}. \tag{11c}$$

In the definition of  $K_h^s, K_h^0$  we assume  $g \in H^1(\Omega) \cap C^0(\bar{\Omega})$  so that the point evaluation is well-defined.

Let us shortly comment on how to incorporate the constraints for the higher-order space  $U_{h1}$  and  $n = 2$ . Let  $\mathcal{V}_m$  denote the midpoints of interior edges of the triangulation  $\mathcal{T}$ . Then, a choice for the discrete convex set is

$$K_{h1}^s := \{(v_h, \boldsymbol{\tau}_h, \mu_h) \in U_{h1} : \mu_h \geq 0, v_h(z) \geq g(z) \text{ for all } z \in \mathcal{V}_0 \cup \mathcal{V}_m\}.$$

In the same manner one defines  $K_{h1}^0$  resp.  $K_{h1}^1$ .

### 3.4 Auxiliary results

For the analysis of the convergence rates we use the nodal interpolation operator  $I_h: H^2(\Omega) \rightarrow \mathcal{S}^1(\mathcal{T}) := \mathcal{P}^1(\mathcal{T}) \cap C^0(\overline{\Omega})$ , the Raviart–Thomas projector  $\Pi_h^{\text{div}}: H^1(\Omega)^n \rightarrow \mathcal{RT}^0(\mathcal{T})$ , and the  $L^2(\Omega)$  projector  $\Pi_h: L^2(\Omega) \rightarrow \mathcal{P}^0(\mathcal{T})$ . Observe that with  $v \geq 0, \mu \geq 0$  we have (with sufficient regularity) that  $I_h v \geq 0, \Pi_h \mu \geq 0$ . Moreover, recall the commutativity property  $\text{div } \Pi_h^{\text{div}} = \Pi_h \text{div}$ , as well as the approximation properties

$$\|v - I_h v\| + h \|\nabla(v - I_h v)\| \lesssim h^2 \|D^2 v\|, \tag{12}$$

$$\|\boldsymbol{\tau} - \Pi_h^{\text{div}} \boldsymbol{\tau}\| \lesssim h \|\nabla \boldsymbol{\tau}\|, \tag{13}$$

$$\|\mu - \Pi_h \mu\| \lesssim \|h_{\mathcal{T}} \nabla_{\mathcal{T}} \mu\|. \tag{14}$$

Here,  $\nabla \boldsymbol{\tau}$  is understood componentwise,  $\nabla_{\mathcal{T}} \mu$  denotes the  $\mathcal{T}$ -elementwise gradient of  $\mu \in H^1(\mathcal{T}) := \{v \in L^2(\Omega): v|_T \in H^1(T), T \in \mathcal{T}\}$ . Set  $\|v\|_{H^1(\mathcal{T})}^2 := \|v\|^2 + \|\nabla_{\mathcal{T}} v\|^2$ . The involved constants depend only on the  $\kappa$ -shape regularity of  $\mathcal{T}$  but are otherwise independent of  $\mathcal{T}$ . Furthermore, for  $\mu \in L^2(\Omega)$ , it also holds that

$$\|\mu - \Pi_h \mu\|_{-1} \lesssim \|h_{\mathcal{T}}(\mu - \Pi_h \mu)\|,$$

which follows from the definition of the dual norm, the projection and approximation property of  $\Pi_h$ .

The proof of optimal a priori convergence rates will also rely on the following two results. Scaling arguments and the continuous embedding  $H^2(T_{\text{ref}}) \hookrightarrow C^0(\overline{T_{\text{ref}}})$  show the next result. Here,  $\overline{T_{\text{ref}}}$  denotes some reference element.

**Lemma 11** *There exists a constant  $C > 0$  depending only on  $T_{\text{ref}}$  and  $\kappa$ -shape regularity of the triangulation such that*

$$\|v\|_{L^\infty(T)} \leq C |T|^{-1/2} (\|v\|_T + h_T \|\nabla v\|_T + h_T^2 \|D^2 v\|_T) \text{ for all } v \in H^2(T). \tag{15}$$

The next result is proven along the lines of [12, Lemma 7]. For completeness we present the proof of the nontrivial result adapted to our situation in ‘‘Appendix B’’. For each element  $T \in \mathcal{T}$  and  $v \in H^2(T)$  we define the level set

$$T_C(v) := \{x \in T: v(x) = 0\} \text{ as well as the set } T_{\text{NC}}(v) := \{x \in T: v(x) \neq 0\}.$$

Note that  $v \in H^2(T)$  implies that these sets are measurable. Moreover,  $|T_C(v)| + |T_{\text{NC}}(v)| = |T|$ .

**Lemma 12** *Let  $v \in H^2(T)$ . Assume  $|T_C(v)| > 0$ . Then,*

$$\|v\|_T \leq Ch_T \frac{|T|^{1/2}}{|T_C(v)|^{1/2}} \|\nabla v\|_T$$

and in particular

$$\|\nabla v\|_T \leq Ch_T \frac{|T|^{1/2}}{|T_C(v)|^{1/2}} \|D^2 v\|_T.$$

Here,  $C = \sqrt{n}$  for  $n = 2, 3$ .

### 3.5 Optimal a priori convergence rates

**Theorem 13** Suppose  $\beta \geq 1 + C_F^2$ . Let  $\mathbf{u} \in K^s$  denote the solution of (VIa) with data  $f \in L^2(\Omega)$ ,  $g \in H_0^1(\Omega)$ . Let  $K_h^s$  denote the set defined in (11a) and let  $\mathbf{u}_h \in K_h^s$  denote the solution of (9) with  $A = a_\beta$ ,  $F = F_\beta$ , and  $K = K_h^s$ . If  $u \in H^2(\Omega)$ ,  $g \in H^2(\Omega)$  and  $f \in H^1(\mathcal{T})$ , then

$$\|\mathbf{u} - \mathbf{u}_h\|_U \leq C_{\text{app}} h (\|u\|_{H^2(\Omega)} + \|\nabla_{\mathcal{T}} f\| + \|\lambda\| + \|g\|_{H^2(\Omega)}).$$

The constant  $C_{\text{app}} > 0$  depends only on  $\beta$ ,  $\Omega$ , and  $\kappa$ -shape regularity of  $\mathcal{T}$ .

**Proof** Choose  $\mathbf{v}_h = (I_h u, \Pi_h^{\text{div}} \boldsymbol{\sigma}, \Pi_h \lambda) \in K_h^s$ . The commutativity property of  $\Pi_h^{\text{div}}$  shows that

$$\text{div}(\boldsymbol{\sigma} - \Pi_h^{\text{div}} \boldsymbol{\sigma}) + \lambda - \Pi_h \lambda = (1 - \Pi_h)(\text{div} \boldsymbol{\sigma} + \lambda) = (1 - \Pi_h) f.$$

Therefore, using the approximation properties of the involved operators proves

$$\begin{aligned} \|\mathbf{u} - \mathbf{v}_h\|_U &\leq \|(1 - \Pi_h) f\| + \|\boldsymbol{\sigma} - \Pi_h^{\text{div}} \boldsymbol{\sigma}\| + \|\nabla(u - I_h u)\| \\ &\lesssim h \|\nabla_{\mathcal{T}} f\| + h \|u\|_{H^2(\Omega)}. \end{aligned}$$

Moreover,

$$|\langle \lambda, I_h u - u \rangle| \leq \|\lambda\| h^2 \|D^2 u\| \lesssim h^2 (\|u\|_{H^2(\Omega)}^2 + \|\lambda\|^2).$$

We have to estimate the term

$$|\langle \Pi_h \lambda - \lambda, u - g \rangle|.$$

Define  $T_C := T_C(u - g)$  and  $T_{\text{NC}} := T_{\text{NC}}(u - g)$ . Note that these two sets are measurable and we have that  $|T_C| + |T_{\text{NC}}| = |T|$ . We consider three cases: First, assume that  $|T_C| = 0$ . This implies that  $u - g > 0$  a.e. in  $T$  but since  $(u - g)\lambda = 0$  we infer that  $\lambda = 0$  a.e. in  $T$ . Therefore,  $\Pi_h \lambda|_T = 0$  and we have that  $\langle \Pi_h \lambda - \lambda, u - g \rangle_T = 0$ . Second, assume that  $|T_{\text{NC}}| = 0$ . But then,  $u - g = 0$  a.e. in  $T$  and we have again  $\langle \Pi_h \lambda - \lambda, u - g \rangle_T = 0$ . The final case to be considered is  $|T_{\text{NC}}| > 0, |T_C| > 0$ : We have that

$$|\langle \Pi_h \lambda - \lambda, u - g \rangle_T| = |\langle \lambda, (\Pi_h - 1)(u - g) \rangle_T|$$

$$\leq \|\lambda\|_{L^1(T)} \|(1 - \Pi_h)(u - g)\|_{L^\infty(T)}.$$

Note that  $\lambda|_{T_{\text{NC}}} = 0$ . Thus,  $\|\lambda\|_{L^1(T)} = \|\lambda\|_{L^1(T_C)} \leq |T_C|^{1/2} \|\lambda\|_T$ . For the second term we apply Lemma 11 with  $v = (1 - \Pi_h)(u - g)$  and together with the approximation property of  $\Pi_h$  we get the estimate

$$\begin{aligned} \|(1 - \Pi_h)(u - g)\|_{L^\infty(T)} &\lesssim |T|^{-1/2} (\|(1 - \Pi_h)(u - g)\|_T \\ &\quad + h_T \|\nabla(u - g)\|_T + h_T^2 \|D^2(u - g)\|_T) \\ &\lesssim |T|^{-1/2} (h_T \|\nabla(u - g)\|_T + h_T^2 \|D^2(u - g)\|_T). \end{aligned}$$

We can estimate the gradient term by applying the second inequality of Lemma 12 which gives us

$$\|\nabla(u - g)\|_T \lesssim \frac{h_T |T|^{1/2}}{|T_C|^{1/2}} \|D^2(u - g)\|_T.$$

Clearly  $|T_C|^{1/2} \leq |T|^{1/2}$ , thus  $|T|^{-1/2} \leq |T_C|^{-1/2}$  and we conclude that

$$\|(1 - \Pi_h)(u - g)\|_{L^\infty(T)} \lesssim \frac{h_T^2}{|T_C|^{1/2}} \|D^2(u - g)\|_T.$$

Using  $\|\lambda\|_{L^1(T)} \leq |T_C|^{1/2} \|\lambda\|_T$  then yields that

$$\begin{aligned} |\langle \Pi_h \lambda - \lambda, u - g \rangle_T| &\lesssim |T_C|^{1/2} \|\lambda\|_T \frac{h_T^2}{|T_C|^{1/2}} \|D^2(u - g)\|_T \\ &\leq h_T^2 \left( \|\lambda\|_T^2 + \|u\|_{H^2(T)}^2 + \|g\|_{H^2(T)}^2 \right). \end{aligned}$$

Summing up we have that

$$\begin{aligned} \inf_{v_h \in K_h^s} & \left( \|u - v_h\|_U^2 + |\langle \lambda, v_h - u \rangle + \langle \mu_h - \lambda, u - g \rangle| \right) \\ & \lesssim h^2 \left( \|u\|_{H^2(\Omega)}^2 + \|\nabla_T f\|^2 + \|\lambda\|^2 + \|g\|_{H^2(\Omega)}^2 \right). \end{aligned}$$

Therefore, in view of Theorem 8 it only remains to estimate the consistency error

$$\inf_{v \in K^s} |\langle \lambda, v - u_h \rangle + \langle \mu - \lambda_h, u - g \rangle|.$$

Define  $v := (v, \chi, \mu) := (v, 0, \lambda_h) \in U$  with  $v := \sup\{u_h, g\}$  and observe that  $v \in K^s$ . This directly leads to  $\langle \mu - \lambda_h, u - g \rangle = 0$ . For the remaining term we follow the seminal work [14] of Falk. The same lines as in the proof of [14, Lemma 4] show that

$$|\langle \lambda, v - u_h \rangle| \leq \|\lambda\| \|v - u_h\| \leq \|\lambda\| \|g - I_h g\| \lesssim h^2 \|g\|_{H^2(\Omega)} \|\lambda\|.$$



This finishes the proof. □

The proof of the following result can be obtained in the same fashion as the previous one and is therefore omitted.

**Theorem 14** *Suppose  $\beta \geq 1 + C_F^2$ . Let  $\mathbf{u} \in K^0$  denote the solution of (VIb) with data  $f \in L^2(\Omega)$ ,  $g \in H^1(\Omega)$ ,  $g|_\Gamma \leq 0$ . Let  $\mathbf{u}_h \in K_h$  denote the solution of (9) with  $A = b_\beta$ ,  $F = G_\beta$ , and  $K = K_h$ , where either  $K_h = K_h^s$  or  $K_h = K_h^0$ . If  $u \in H^2(\Omega)$ ,  $g \in H^2(\Omega)$  and  $f \in H^1(\mathcal{T})$ , then*

$$\|\mathbf{u} - \mathbf{u}_h\|_U \leq C_{\text{app}} h (\|\mathbf{u}\|_{H^2(\Omega)} + \|\nabla_{\mathcal{T}} f\| + \|\lambda\| + \|g\|_{H^2(\Omega)}).$$

The constant  $C_{\text{app}} > 0$  depends only on  $\beta$ ,  $\Omega$ , and  $\kappa$ -shape regularity of  $\mathcal{T}$ .

Finally, we show convergence rates for problem (VIc) and its approximation. Note that for the sets  $K_h^1, K_h^s$  defined in (11c), (11a) it holds that  $K_h^s \subset K_h^1 \subset K^1$  and thus the consistency error, see Theorem 10, vanishes. The proof is similar to the one of Theorem 13 and is therefore left to the reader.

**Theorem 15** *Suppose  $\beta \geq 1 + C_F^2$ . Let  $\mathbf{u} \in K^1$  denote the solution of (VIc) with data  $f \in L^2(\Omega)$ ,  $g \in H_0^1(\Omega)$ . Let  $\mathbf{u}_h \in K_h$  denote the solution of (9) with  $A = c_\beta$ ,  $F = H_\beta$ , and  $K = K_h$ , where either  $K_h = K_h^s$  or  $K_h = K_h^1$ . If  $u \in H^2(\Omega)$ ,  $g \in H^2(\Omega)$  and  $f \in H^1(\mathcal{T})$ , then*

$$\|\mathbf{u} - \mathbf{u}_h\|_U \leq C_{\text{app}} h (\|\mathbf{u}\|_{H^2(\Omega)} + \|\nabla_{\mathcal{T}} f\| + \|\lambda\| + \|g\|_{H^2(\Omega)}).$$

The constant  $C_{\text{app}} > 0$  depends only on  $\beta$ ,  $\Omega$ , and  $\kappa$ -shape regularity of  $\mathcal{T}$ .

### 4 A posteriori analysis

In this section we derive reliable error bounds that can be used as an a posteriori estimator. We define

$$\text{osc} := \text{osc}(f) := \|(1 - \Pi_h)f\|.$$

The estimator below includes the residual term

$$\eta^2 := \eta(\mathbf{u}_h, f)^2 := \|\text{div } \boldsymbol{\sigma}_h + \lambda_h + \Pi_h f\|^2 + \|\nabla u_h - \boldsymbol{\sigma}_h\|^2,$$

which can be localized. The derivation of our estimators is quite simple and is based on the following observation. Let  $\mathbf{u} \in K^s \subset K^j$  denote the unique solution of (3) and let  $\mathbf{u}_h \in U_{h0}$  be arbitrary. Take  $\beta = 1 + C_F^2$  and recall that by Lemma 3 it holds that  $a_\beta(\mathbf{v}, \mathbf{v}) = b_\beta(\mathbf{v}, \mathbf{v}) = c_\beta(\mathbf{v}, \mathbf{v}) \gtrsim \|\mathbf{v}\|_U^2$  for all  $\mathbf{v} \in U$ . Then, together with the Pythagoras theorem  $\|\mu\|^2 = \|(1 - \Pi_h)\mu\|^2 + \|\Pi_h\mu\|^2$  for  $\mu \in L^2(\Omega)$  and using  $\text{div } \boldsymbol{\sigma} + \lambda + f = 0$ ,  $\nabla u = \boldsymbol{\sigma}$ ,  $\text{div } \boldsymbol{\sigma}_h + \lambda_h \in \mathcal{P}^0(\mathcal{T})$ , it follows that

$$\|\mathbf{u} - \mathbf{u}_h\|_U^2 \lesssim \beta \|\text{div } \boldsymbol{\sigma}_h + \lambda_h + f\|^2 + \|\nabla u_h - \boldsymbol{\sigma}_h\|^2 + \langle \lambda_h - \lambda, u_h - u \rangle$$

$$\begin{aligned}
&= \beta \|\operatorname{div} \boldsymbol{\sigma}_h + \lambda_h + \Pi_h f\|^2 + \beta \operatorname{osc}^2 \\
&\quad + \|\nabla u_h - \boldsymbol{\sigma}_h\|^2 + \langle \lambda_h - \lambda, u_h - u \rangle \\
&\leq \beta(\eta^2 + \operatorname{osc}^2) + \langle \lambda_h - \lambda, u_h - u \rangle.
\end{aligned} \tag{16}$$

The remaining results in this section are proved by estimating the duality term  $\langle \lambda_h - \lambda, u_h - u \rangle$  from (16). In particular, the proof of the next result employs only  $\lambda_h \geq 0$ . We will need the positive resp. negative part of a function  $v: \Omega \rightarrow \mathbb{R}$ ,

$$v_+ := \max\{0, v\}, \quad v_- := -\min\{0, v\}.$$

This definition implies that  $v = v_+ - v_-$ . The ideas of estimating the duality term are similar as in [18,31] and references therein, see also [15] for a related estimate for Signorini-type problems. Note that we do not need to assume  $g \in H_0^1(\Omega)$ .

**Theorem 16** *Let  $\mathbf{u} \in K^s$  denote the solution of (3). Let  $\mathbf{u}_h \in K_h$ , where  $K_h \in \{K_h^s, K_h^1\}$ , be arbitrary. The error satisfies*

$$\|\mathbf{u} - \mathbf{u}_h\|_U^2 \leq C_{\text{rel}}(\eta^2 + \rho^2 + \operatorname{osc}^2),$$

where the estimator contribution  $\rho$  is given by

$$\rho^2 := \langle \lambda_h, (u_h - g)_+ \rangle + \|\nabla(g - u_h)_+\|^2.$$

The constant  $C_{\text{rel}} > 0$  depends only on  $\Omega$ .

**Proof** In view of estimate (16) we only have to tackle the term  $\langle \lambda_h - \lambda, u_h - u \rangle$ . Define  $v_h := \max\{u_h, g\}$ . Clearly,  $v_h \geq g$  and  $v_h \in H_0^1(\Omega)$ . Note that  $\lambda = -\Delta u - f \in H^{-1}(\Omega)$ . Therefore,  $\langle \lambda, v \rangle = (\nabla u, \nabla v) - (f, v)$  for all  $v \in H_0^1(\Omega)$  and using the variational inequality for the exact solution (2) yields

$$\begin{aligned}
-\langle \lambda, u_h - u \rangle &= -\langle \lambda, u_h - v_h \rangle - \langle \lambda, v_h - u \rangle \leq -\langle \lambda, u_h - v_h \rangle \\
&= \langle \lambda, (u_h - g)_- \rangle = \langle \lambda - \lambda_h, (u_h - g)_- \rangle + \langle \lambda_h, (u_h - g)_- \rangle \\
&\leq \frac{\delta}{2} \|\lambda - \lambda_h\|_{-1}^2 + \frac{\delta^{-1}}{2} \|\nabla(u_h - g)_-\|^2 + \langle \lambda_h, (u_h - g)_- \rangle
\end{aligned}$$

for all  $\delta > 0$ . Employing  $\lambda_h \geq 0$ ,  $g - u \leq 0$ , and  $v + v_- = v_+$  we further infer that

$$\begin{aligned}
\langle \lambda_h - \lambda, u_h - u \rangle &\leq \langle \lambda_h, u_h - g + (u_h - g)_- \rangle + \langle \lambda_h, g - u \rangle \\
&\quad + \frac{\delta}{2} \|\lambda - \lambda_h\|_{-1}^2 + \frac{\delta^{-1}}{2} \|\nabla(u_h - g)_-\|^2 \\
&\leq \langle \lambda_h, (u_h - g)_+ \rangle + \frac{\delta}{2} \|\lambda - \lambda_h\|_{-1}^2 + \frac{\delta^{-1}}{2} \|\nabla(u_h - g)_-\|^2.
\end{aligned}$$

Recall that  $\|\lambda - \lambda_h\|_{-1} \leq \|\mathbf{u} - \mathbf{u}_h\|_V \lesssim \|\mathbf{u} - \mathbf{u}_h\|_U$ , where the involved constant depends only on  $\Omega$ . Thus, choosing  $\delta > 0$  sufficiently small the proof is concluded with (16).  $\square$

We could derive a similar estimate if  $u_h \in K_h^0$  by changing the role of  $u_h$  and  $\lambda_h$  resp.  $u$  and  $\lambda$  in the proof. However, this leads to an estimator with a non-local term. To see this, suppose  $g = 0$ . Then, following the last proof we get

$$\langle \lambda_h - \lambda, u_h - u \rangle \leq \langle (\lambda_h)_+, u_h \rangle + \frac{\delta}{2} \|\nabla(u - u_h)\|^2 + \frac{\delta^{-1}}{2} \|(\lambda_h)_-\|_{-1}^2$$

for  $\delta > 0$ . For the total error this would yield

$$\|u - u_h\|_U^2 \lesssim \eta^2 + \text{osc}^2 + \langle (\lambda_h)_+, u_h \rangle + \|(\lambda_h)_-\|_{-1}^2.$$

The last term is not localizable and therefore it is not feasible to use this estimate as an a posteriori error estimator in an adaptive algorithm.

**Remark 17** The derived estimator is efficient up to the term  $\rho$ , i.e.,

$$\eta^2 + \text{osc}^2 \lesssim \|u - u_h\|_U^2.$$

To see this, we employ the Pythagoras theorem to obtain

$$\eta^2 + \text{osc}^2 = \|\text{div } \sigma_h + \lambda_h + f\|^2 + \|\nabla u_h - \sigma_h\|^2.$$

Then,  $\text{div } \sigma + \lambda = -f$ ,  $\nabla u = \sigma$  and the triangle inequality prove the asserted estimate. The proof of the efficiency estimate  $\rho \lesssim \|u - u_h\|_U$  (up to possible data resp. obstacle oscillations) is an open problem, see also the related works [1, 18].

### 5 Examples

In this section we present numerical studies that demonstrate the performance of our proposed methods in different situations:

- In Sect. 5.3 we consider a problem on the unit square with smooth obstacle and known smooth solution.
- In Sect. 5.4 we consider the example from [4, Section 5.2] where the solution is known and exhibits a singularity.
- In Sect. 5.5 we consider a problem on an L-shaped domain with a pyramid-like obstacle and unknown solution.

Before we come to a detailed discussion on the numerical studies some remarks are in order. In all examples we choose  $\beta = 1 + \text{diam}(\Omega)^2$  to ensure coercivity of the bilinear forms (Lemma 3). This also implies that the Galerkin matrices associated to the bilinear forms  $a_\beta$ ,  $b_\beta$ , and  $c_\beta$  are positive definite. Choosing standard basis functions for  $\mathcal{S}_0^1(\mathcal{T})$  (nodal basis),  $\mathcal{RT}^0(\mathcal{T})$  (lowest-order Raviart–Thomas basis) and  $\mathcal{P}^0(\mathcal{T})$  (characteristic functions), the constraints in the discrete convex sets  $K_h^\star$ , where  $\star = 0, \star = 1$  or  $\star = s$ , are straightforward to impose. The resulting discrete variational inequalities are then solved using a (primal-dual) active set strategy, see e.g., [21–23].

## 5.1 Active set method and discrete variational inequalities

In this section we first define and collect results on the (*primal-dual*) *active set method*. Then, we recall the variational inequalities (VIa)–(VIc) and write down their discrete variants.

### 5.1.1 Active set method

Let  $\mathcal{N} = \{1, \dots, N\}$ ,  $N \in \mathbb{N}$ , and let  $\mathcal{N}_\gamma \subseteq \mathcal{N}$  be a non-empty subset. We set  $\mathcal{N}_\omega := \mathcal{N} \setminus \mathcal{N}_\gamma$ . For a vector  $\mathbf{x} \in \mathbb{R}^N$  we write  $\mathbf{x} = 0$  if all components are equal to 0. Similarly,  $\mathbf{x} \geq 0$  means that all components of  $\mathbf{x}$  are  $\geq 0$ . For a subset  $\mathcal{I} \subseteq \mathcal{N}$ ,  $\mathbf{x}_\mathcal{I} = 0$  means  $x_i = 0$  for all  $i \in \mathcal{I}$ . We also use the notation  $\mathbf{x}_\mathcal{I} \geq 0$ , which means  $x_i \geq 0$  for all  $i \in \mathcal{I}$  and  $\mathbf{x}_\mathcal{I} \geq \mathbf{y}_\mathcal{I}$  stands for  $\mathbf{x}_\mathcal{I} - \mathbf{y}_\mathcal{I} \geq 0$ .

For  $\mathbf{g} \in \mathbb{R}^N$  we consider the convex set

$$K := K_{\mathbf{g}} := \{\mathbf{x} \in \mathbb{R}^N : \mathbf{x}_{\mathcal{N}_\gamma} \geq \mathbf{g}_{\mathcal{N}_\gamma}\}.$$

Let  $\mathbf{S} \in \mathbb{R}^{N \times N}$  denote a positive definite (but possibly non-symmetric) matrix, and  $\mathbf{b} \in \mathbb{R}^N$  some arbitrary vector. We consider the variational inequality: find  $\mathbf{x} \in K$ , such that

$$\langle \mathbf{S}\mathbf{x}, \mathbf{y} - \mathbf{x} \rangle_2 \geq \langle \mathbf{b}, \mathbf{y} - \mathbf{x} \rangle_2 \quad \text{for all } \mathbf{y} \in K, \quad (17)$$

where  $\langle \cdot, \cdot \rangle_2$  denotes the Euclidean inner product on  $\mathbb{R}^N$ . Since  $\mathbf{S}$  is positive definite this problem admits a unique solution. It is well-known that problem (17) can be rewritten as follows: find  $(\mathbf{x}, \boldsymbol{\lambda}) \in \mathbb{R}^N \times \mathbb{R}^N$  such that

$$\mathbf{S}\mathbf{x} - \boldsymbol{\lambda} = \mathbf{b}, \quad (18a)$$

$$\boldsymbol{\lambda}_{\mathcal{N}_\omega} = 0, \quad (18b)$$

$$\boldsymbol{\lambda}_{\mathcal{N}_\gamma} = \max\{0, \boldsymbol{\lambda}_{\mathcal{N}_\gamma} - C(\mathbf{x}_{\mathcal{N}_\gamma} - \mathbf{g}_{\mathcal{N}_\gamma})\}, \quad (18c)$$

where  $\max\{\cdot, \cdot\}$  denotes the componentwise maximum and  $C > 0$  is some constant. Note that the solution is independent of  $C$ . Now following the seminal work [21] one defines a (semi-smooth) Newton method for solving (18). The same lines of argumentation as in [21] show that the method can be written as an active set strategy. The algorithm adapted to our situation is given in Algorithm 1.

The solution of the linear system in Line 8 of Algorithm 1 can be written (with  $\mathcal{I} = \mathcal{I}^k$ ,  $\mathcal{J} = \mathcal{J}^k$ ) as

$$\begin{pmatrix} \mathbf{S}_{\mathcal{I}\mathcal{I}} & \mathbf{S}_{\mathcal{I}\mathcal{J}} \\ \mathbf{S}_{\mathcal{J}\mathcal{I}} & \mathbf{S}_{\mathcal{J}\mathcal{J}} \end{pmatrix} \begin{pmatrix} \mathbf{x}_\mathcal{I} \\ \mathbf{x}_\mathcal{J} \end{pmatrix} - \begin{pmatrix} \boldsymbol{\lambda}_\mathcal{I} \\ \boldsymbol{\lambda}_\mathcal{J} \end{pmatrix} = \begin{pmatrix} \mathbf{b}_\mathcal{I} \\ \mathbf{b}_\mathcal{J} \end{pmatrix}.$$

With the constraints  $\mathbf{x}_\mathcal{J} = \mathbf{g}_\mathcal{J}$  and  $\boldsymbol{\lambda}_\mathcal{I} = 0$  this reduces to the solution of the system

$$\mathbf{S}_{\mathcal{I}\mathcal{I}}\mathbf{x}_\mathcal{I} = \mathbf{b}_\mathcal{I} - \mathbf{S}_{\mathcal{I}\mathcal{J}}\mathbf{g}_\mathcal{J}$$

---

**Algorithm 1** Active Set Method for solving (17)

---

**Input:**  $\mathbf{x}^0, \boldsymbol{\lambda}^0$  and  $C > 0$

- 1: Set  $k = 0$
- 2: **while** TRUE **do**
- 3:  $\widehat{\boldsymbol{\lambda}}_{\mathcal{N}_\gamma}^k \leftarrow \max\{0, \boldsymbol{\lambda}_{\mathcal{N}_\gamma}^k - C(\mathbf{x}_{\mathcal{N}_\gamma}^k - \mathbf{g}_{\mathcal{N}_\gamma})\}$  and  $\widehat{\boldsymbol{\lambda}}_{\mathcal{N}_\omega}^k \leftarrow 0$
- 4: Determine set of *active* ( $\mathcal{J}$ ) and *inactive* ( $\mathcal{I}$ ) degrees of freedom
 
$$\mathcal{J}^k \leftarrow \{j \in \mathcal{N}_\gamma : \widehat{\boldsymbol{\lambda}}_j^k > 0\},$$

$$\mathcal{I}^k \leftarrow \mathcal{N} \setminus \mathcal{J}^k.$$
- 5: **if**  $k \geq 1$  &  $\mathcal{J}^k = \mathcal{J}^{k-1}$  **then**
- 6:     **return**  $\mathbf{x} = \mathbf{x}^k$
- 7: **end if**
- 8: Solve
 
$$\begin{aligned} \mathbf{S}\mathbf{x} - \boldsymbol{\lambda} &= \mathbf{b}, \\ \boldsymbol{\lambda}_{\mathcal{I}^k} &= 0, \\ \mathbf{x}_{\mathcal{J}^k} &= \mathbf{g}_{\mathcal{J}^k}. \end{aligned}$$
- 9: Set  $\mathbf{x}^{k+1} := \mathbf{x}$  and  $\boldsymbol{\lambda}^{k+1} := \boldsymbol{\lambda}$
- 10: Increase counter  $k \leftarrow k + 1$
- 11: **end while**

---

and the definition

$$\boldsymbol{\lambda}_{\mathcal{J}} := \mathbf{S}_{\mathcal{J}\mathcal{I}}\mathbf{x}_{\mathcal{I}} + \mathbf{S}_{\mathcal{J}\mathcal{J}}\mathbf{g}_{\mathcal{J}} - \mathbf{b}_{\mathcal{J}}.$$

Since  $\mathbf{S}$  is positive definite the subblock  $\mathbf{S}_{\mathcal{I}\mathcal{I}}$  is as well and thus  $\mathbf{S}_{\mathcal{I}\mathcal{I}}$  is invertible.

Some remarks are in order. We can follow the analysis of [21] to see that the basic (local) convergence result holds true in our case as well.

**Proposition 18** [21, Theorem 3.1] *If the initial guess  $(\mathbf{x}^0, \boldsymbol{\lambda}^0)$  is sufficient close to the exact solution  $(\mathbf{x}, \boldsymbol{\lambda})$  of (18) then the iterates  $(\mathbf{x}^k, \boldsymbol{\lambda}^k)$  in Algorithm 1 converge superlinearly to  $(\mathbf{x}, \boldsymbol{\lambda})$ .*

The stopping criterion in Line 5 can be replaced by other criterions. Here, we choose  $\mathcal{J}^k = \mathcal{J}^{k-1}$  because then we know that we have hit the exact solution of (17). The proof of the following result is a slight modification of the proof of [23, Lemma 3.1] and the interested reader can find it in ‘‘Appendix C’’.

**Lemma 19** *If the stopping criterion in Line 5 of Algorithm 1 is satisfied, then  $\mathbf{x} = \mathbf{x}^k$  is the solution of (17).*

### 5.1.2 Discrete variational inequalities

In this section we recall the discrete versions of the variational inequalities (VIa)–(VIc) and present them in matrix-vector form. They fit into the abstract framework given in Sect. 5.1.1.

Let us recall the discrete space from Sect. 3.3,

$$U_{h0} = \mathcal{S}_0^1(\mathcal{T}) \times \mathcal{RT}^0(\mathcal{T}) \times \mathcal{P}^0(\mathcal{T}).$$

Let  $\mathcal{E}$  denote the set of edges ( $n = 2$ ) resp. faces ( $n = 3$ ). Then,  $\dim(U_{h0}) = \#\mathcal{V}_0 + \#\mathcal{E} + \#\mathcal{T} =: N$ . Numbering the nodes  $x_j$  of  $\mathcal{V}_0$ , the edges/faces  $E_j$  in  $\mathcal{E}$  and the elements  $T_j$  in  $\mathcal{T}$ , we consider the following functions:

- For  $j = 1, \dots, \#\mathcal{V}_0$  let  $v_j$  denote the nodal basis functions associated to the node  $x_j \in \mathcal{V}_0$ .
- For  $j = 1, \dots, \#\mathcal{E}$  let  $\boldsymbol{\tau}^{(j)}$  denote the Raviart–Thomas basis functions associated to the edge/face  $E_j \in \mathcal{E}$ .
- For  $j = 1, \dots, \#\mathcal{T}$  let  $\chi_j$  denote the characteristic function of the element  $T_j \in \mathcal{T}$ .

We define the basis  $(\boldsymbol{\xi}^{(j)})_{j=1}^N$  for the space  $U_{h0}$  by

$$\begin{aligned} \boldsymbol{\xi}^{(j)} &:= (v_j, 0, 0) && \text{for } j = 1, \dots, \#\mathcal{V}_0, \\ \boldsymbol{\xi}^{(\#\mathcal{V}_0+j)} &:= (0, \boldsymbol{\tau}^{(j)}, 0) && \text{for } j = 1, \dots, \#\mathcal{E}, \\ \boldsymbol{\xi}^{(\#\mathcal{V}_0+\#\mathcal{E}+j)} &:= (0, 0, \chi_j) && \text{for } j = 1, \dots, \#\mathcal{T}. \end{aligned}$$

Recall from Eq. 11 the discrete convex sets

$$\begin{aligned} K_h^s &:= \{(v_h, \boldsymbol{\tau}_h, \mu_h) \in U_{h0} : \mu_h \geq 0, v_h(x) \geq g(x) \text{ for all } x \in \mathcal{V}_0\}, \\ K_h^0 &:= \{(v_h, \boldsymbol{\tau}_h, \mu_h) \in U_{h0} : v_h(x) \geq g(x) \text{ for all } x \in \mathcal{V}_0\}, \\ K_h^1 &:= \{(v_h, \boldsymbol{\tau}_h, \mu_h) \in U_{h0} : \mu_h \geq 0\}. \end{aligned}$$

These convex subsets of  $U_{h0}$  correspond to convex subsets of  $\mathbb{R}^N$  as follows: For given obstacle function  $g \in H_0^1(\Omega) \cap C^0(\overline{\Omega})$  define the vector  $\mathbf{g} \in \mathbb{R}^N$  by

$$\mathbf{g}_j = \begin{cases} g(x_j) & \text{for } j = 1, \dots, \#\mathcal{V}_0, \\ 0 & \text{else} \end{cases}.$$

Let  $\mathcal{N} = \{1, \dots, N\}$  and define  $\mathcal{N}_\gamma^s, \mathcal{N}_\gamma^0, \mathcal{N}_\gamma^1$  by

$$\begin{aligned} \mathcal{N}_\gamma^s &:= \{1, \dots, \#\mathcal{V}_0, \#\mathcal{V}_0 + \#\mathcal{E} + 1, \dots, N\} \subset \mathcal{N}, \\ \mathcal{N}_\gamma^0 &:= \{1, \dots, \#\mathcal{V}_0\} \subset \mathcal{N}, \\ \mathcal{N}_\gamma^1 &:= \{\#\mathcal{V}_0 + \#\mathcal{E} + 1, \dots, N\} \subset \mathcal{N}. \end{aligned}$$

Then, the three sets  $K_h^s, K_h^0, K_h^1$  correspond to the sets

$$\begin{aligned} K_N^s &:= \{\mathbf{x} \in \mathbb{R}^N : \mathbf{x}_{\mathcal{N}_\gamma^s} \geq \mathbf{g}_{\mathcal{N}_\gamma^s}\}, \\ K_N^0 &:= \{\mathbf{x} \in \mathbb{R}^N : \mathbf{x}_{\mathcal{N}_\gamma^0} \geq \mathbf{g}_{\mathcal{N}_\gamma^0}\}, \end{aligned}$$

$$K_N^1 := \{ \mathbf{x} \in \mathbb{R}^N : \mathbf{x}_{\mathcal{N}_V^1} \geq 0 \}.$$

With these definitions we can now state the algebraic forms of the discrete variational inequalities:

**Discrete version of (VIa) with  $K_h^s$**

The discrete version of (VIa) with convex set  $K_h^s$  reads: find  $\mathbf{u}_h \in K_h^s$  such that

$$a_\beta(\mathbf{u}_h, \mathbf{v}_h - \mathbf{u}_h) \geq F_\beta(\mathbf{v}_h - \mathbf{u}_h) \quad \text{for all } \mathbf{v}_h \in K_h^s. \tag{19}$$

Let  $\mathbf{S}^{(s)} \in \mathbb{R}^N \times \mathbb{R}^N$  denote the Galerkin matrix of the bilinear form  $a_\beta(\cdot, \cdot)$  and let  $\mathbf{b}^{(s)} \in \mathbb{R}^N$  denote the load vector, i.e.,

$$\mathbf{S}_{jk}^{(s)} = a_\beta(\boldsymbol{\xi}^{(k)}, \boldsymbol{\xi}^{(j)}), \quad \mathbf{b}_j^{(s)} = F_\beta(\boldsymbol{\xi}^{(j)})$$

for all  $j, k = 1, \dots, N$ . Note that  $\mathbf{S}^{(s)}$  is symmetric and positive definite. Problem (19) then reads in algebraic form as: find  $\mathbf{x} \in K_N^s$  such that

$$\langle \mathbf{S}^{(s)} \mathbf{x}, \mathbf{y} - \mathbf{x} \rangle_2 \geq \langle \mathbf{b}^{(s)}, \mathbf{y} - \mathbf{x} \rangle_2 \quad \text{for all } \mathbf{y} \in K_N^s. \tag{20}$$

**Discrete version of (VIb) with  $K_h^0$**

The discrete version of (VIb) with convex set  $K_h^0$  reads: find  $\mathbf{u}_h \in K_h^0$  such that

$$b_\beta(\mathbf{u}_h, \mathbf{v}_h - \mathbf{u}_h) \geq G_\beta(\mathbf{v}_h - \mathbf{u}_h) \quad \text{for all } \mathbf{v}_h \in K_h^0. \tag{21}$$

Let  $\mathbf{S}^{(0)} \in \mathbb{R}^N \times \mathbb{R}^N$  denote the Galerkin matrix of the bilinear form  $b_\beta(\cdot, \cdot)$  and let  $\mathbf{b}^{(0)} \in \mathbb{R}^N$  denote the load vector, i.e.,

$$\mathbf{S}_{jk}^{(0)} = b_\beta(\boldsymbol{\xi}^{(k)}, \boldsymbol{\xi}^{(j)}), \quad \mathbf{b}_j^{(0)} = G_\beta(\boldsymbol{\xi}^{(j)})$$

for all  $j, k = 1, \dots, N$ . Note that  $\mathbf{S}^{(0)}$  is non-symmetric and positive definite. Problem (21) then reads in algebraic form as: find  $\mathbf{x} \in K_N^0$  such that

$$\langle \mathbf{S}^{(0)} \mathbf{x}, \mathbf{y} - \mathbf{x} \rangle_2 \geq \langle \mathbf{b}^{(0)}, \mathbf{y} - \mathbf{x} \rangle_2 \quad \text{for all } \mathbf{y} \in K_N^0. \tag{22}$$

**Discrete version of (VIc) with  $K_h^1$**

The discrete version of (VIc) with convex set  $K_h^1$  reads: find  $\mathbf{u}_h \in K_h^1$  such that

$$c_\beta(\mathbf{u}_h, \mathbf{v}_h - \mathbf{u}_h) \geq H_\beta(\mathbf{v}_h - \mathbf{u}_h) \quad \text{for all } \mathbf{v}_h \in K_h^1. \tag{23}$$

Let  $S^{(1)} \in \mathbb{R}^N \times \mathbb{R}^N$  denote the Galerkin matrix of the bilinear form  $c_\beta(\cdot, \cdot)$  and let  $b^{(1)} \in \mathbb{R}^N$  denote the load vector, i.e.,

$$S_{jk}^{(1)} = c_\beta(\xi^{(k)}, \xi^{(j)}), \quad b_j^{(1)} = H_\beta(\xi^{(j)})$$

for all  $j, k = 1, \dots, N$ . Note that  $S^{(1)}$  is non-symmetric and positive definite. Problem (23) then reads in algebraic form as: find  $x \in K_N^1$  such that

$$\langle S^{(1)}x, y - x \rangle_2 \geq \langle b^{(1)}, y - x \rangle_2 \quad \text{for all } y \in K_N^1. \tag{24}$$

**Solver setup**

The algebraic problems (20), (22) and (24) are then solved using Algorithm 1. The initial data  $(x^0, \lambda^0)$  is chosen as the solution of

$$S^{(\star)}x^0 = b^{(\star)},$$

$$x_{\mathcal{N}_\gamma^\star}^0 = g_{\mathcal{N}_\gamma^\star}$$

and

$$\lambda^0 := \max\{0, S^{(\star)}x^0 - b^{(\star)}\},$$

where  $\star = s, \star = 0$  or  $\star = 1$ . The constant  $C$  in Algorithm 1 is chosen as  $C = 1$ . The linear systems in Line 8 of Algorithm 1 are solved using the MATLAB backslash operator.

**5.2 Error and estimator quantities**

We define the error resp. total estimator by

$$\text{err}_U := \|u - u_h\|_U, \quad \text{est}^2 := \eta^2 + \rho^2 + \text{osc}^2.$$

Note that the estimator can be decomposed into local contributions,

$$\text{est}^2 = \sum_{T \in \mathcal{T}} \text{est}(T)^2 := \sum_{T \in \mathcal{T}} \left( \|\text{div } \sigma_h + \lambda_h + \Pi_h f\|_T^2 + \|\nabla u_h - \sigma_h\|_T^2 \right. \\ \left. + (\lambda_h, (u_h - g)_+)_T + \|\nabla(g - u_h)_+\|_T^2 + \|(1 - \Pi_h)f\|_T^2 \right),$$

where  $\|\cdot\|_T$  denotes the  $L^2(T)$  norm and  $(\cdot, \cdot)_T$  the  $L^2(T)$  inner product. Moreover, we will estimate the error in the weaker norm  $\|\cdot\|_V$ . To do so we consider an upper bound given by

$$\text{err}_V^2 := \|\nabla(u - u_h)\|^2 + \|\sigma - \sigma_h\|^2 + \|\lambda - \lambda_h\|_{-1,h}^2,$$



where the evaluation of  $\|\cdot\|_{-1,h}$  is based on the discrete  $H^{-1}(\Omega)$  norm discussed in the seminal work [8]: Let  $Q_h: L^2(\Omega) \rightarrow \mathcal{S}_0^1(\mathcal{T})$  denote the  $L^2(\Omega)$  projector. Let  $\mu \in L^2(\Omega)$ . We stress that using the projection and local approximation property of  $Q_h$  yields

$$\|(1 - Q_h)\mu\|_{-1} = \sup_{0 \neq v \in H_0^1(\Omega)} \frac{\langle (1 - Q_h)\mu, (1 - Q_h)v \rangle}{\|\nabla v\|} \lesssim \|h_{\mathcal{T}}\mu\|,$$

where the involved constant depends on shape regularity of  $\mathcal{T}$ . Following [8] it holds that

$$\|\mu\|_{-1} \leq \|(1 - Q_h)\mu\|_{-1} + \|Q_h\mu\|_{-1} \lesssim \|h_{\mathcal{T}}\mu\| + \|\nabla u_h[\mu]\|$$

where  $u_h[\mu] \in \mathcal{S}_0^1(\mathcal{T})$  is the solution of

$$\langle \nabla u_h[\mu], \nabla v_h \rangle = \langle \mu, v_h \rangle \quad \text{for all } v_h \in \mathcal{S}_0^1(\mathcal{T}).$$

Note that  $\|\nabla u_h[\mu]\| \leq \|\mu\|_{-1}$ . The estimate  $\|Q_h\mu\|_{-1} \lesssim \|\nabla u_h[\mu]\|$  depends on the stability of the projection  $Q_h$  in  $H^1(\Omega)$ ,  $\|\nabla Q_h v\| \lesssim \|\nabla v\|$  for  $v \in H_0^1(\Omega)$ , i.e.,

$$\begin{aligned} \|Q_h\mu\|_{-1} &= \sup_{0 \neq v \in H_0^1(\Omega)} \frac{\langle Q_h\mu, v \rangle}{\|\nabla v\|} = \sup_{0 \neq v \in H_0^1(\Omega)} \frac{\langle \mu, Q_h v \rangle}{\|\nabla v\|} \\ &= \sup_{0 \neq v \in H_0^1(\Omega)} \frac{\langle \nabla u_h[\mu], \nabla Q_h v \rangle}{\|\nabla v\|} \\ &\lesssim \sup_{0 \neq v \in H_0^1(\Omega)} \frac{\langle \nabla u_h[\mu], \nabla Q_h v \rangle}{\|\nabla Q_h v\|} = \|\nabla u_h[\mu]\|. \end{aligned}$$

Here, we use newest-vertex bisection [30] as refinement strategy where stability of the  $L^2(\Omega)$  projection is known [24].

We use an adaptive algorithm that basically consists of iterating the four steps

$$\boxed{SOLVE} \rightarrow \boxed{ESTIMATE} \rightarrow \boxed{MARK} \rightarrow \boxed{REFINE},$$

where the marking step is done with the bulk criterion, i.e., we determine a set  $\mathcal{M} \subseteq \mathcal{T}$  of (up to a constant) minimal cardinality with

$$\theta \text{ est}^2 \leq \sum_{T \in \mathcal{M}} \text{est}(T)^2.$$

For the experiments the marking parameter  $\theta$  is set to  $\frac{1}{4}$ .

Convergence rates in the figures are indicated by triangles, where the number  $\alpha$  besides the triangle denotes the experimental rate  $\mathcal{O}((\#\mathcal{T})^{-\alpha})$ . For uniform refinement we have  $h^{2\alpha} \simeq \#\mathcal{T}^{-\alpha}$ .

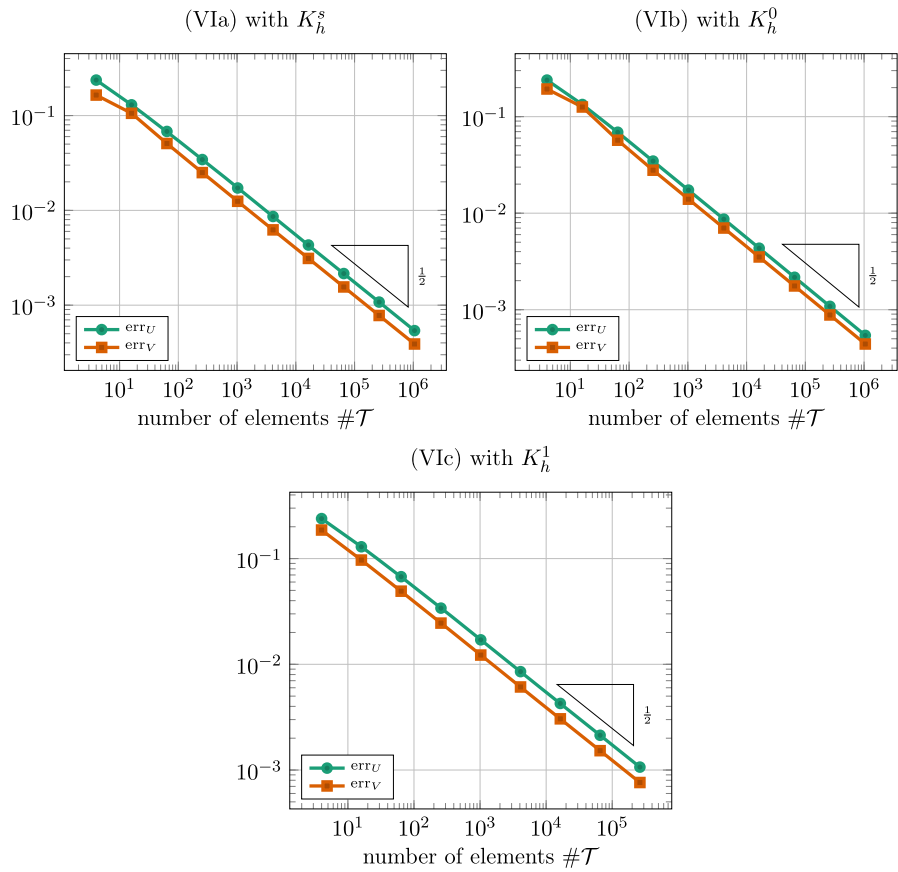


Fig. 1 Convergence rates for the problem from Sect. 5.3

### 5.3 Smooth solution

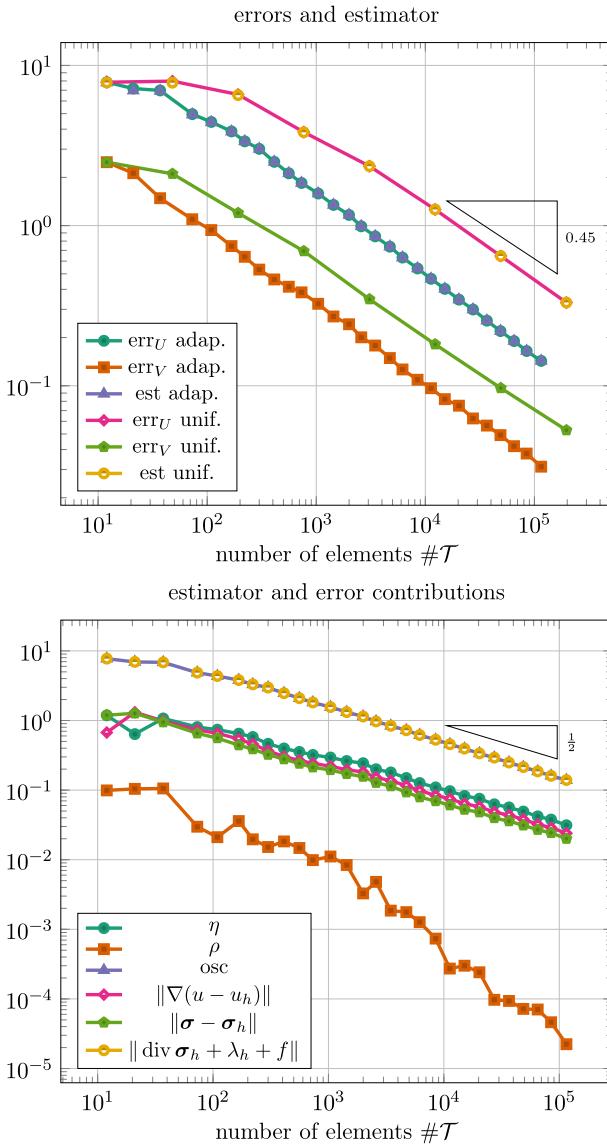
Let  $\Omega = (0, 1)^2$ ,  $u(x, y) = (1 - x)x(1 - y)y$ ,

$$f(x, y) := \begin{cases} 0 & x < \frac{1}{2} \\ -\Delta u(x, y) & x \geq \frac{1}{2} \end{cases}.$$

Then,  $u$  solves the obstacle problem (1) with data  $f$  and obstacle

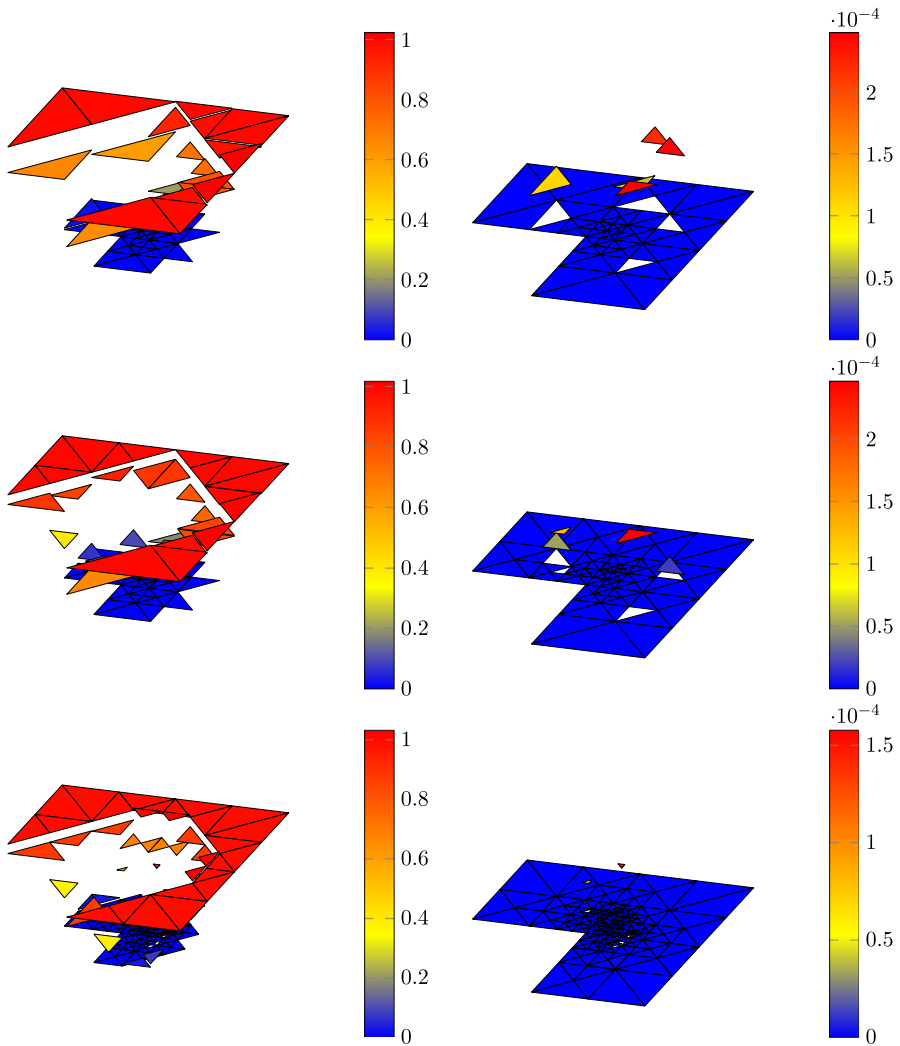
$$g(x, y) = \begin{cases} (1 - x)x(1 - y)y & x \leq \frac{1}{2} \\ \tilde{g}(x)(1 - y)y & x \in (\frac{1}{2}, \frac{3}{4}), \\ 0 & x \geq \frac{3}{4} \end{cases},$$

where  $\tilde{g}$  is the unique polynomial of degree 3 such that  $g$  and  $\nabla g$  are continuous at the lines  $x = \frac{1}{2}, \frac{3}{4}$ . In particular,  $g \in H^2(\Omega)$ . Note that  $\lambda = -\Delta u - f \in H^1(\mathcal{T})$ . Figure 1



**Fig. 2** Convergence rates for the problem from Sect. 5.4. The upper plot shows the total errors and estimators for uniform and adaptive refinement. The lower plot compares the error and estimator contributions in the case of adaptive refinements

shows that the convergence rates for the solutions of the discrete variational inequalities (VIa)–(VIc) based on the convex sets  $K_h^s, K_h^0, K_h^1$  are optimal. This perfectly fits to our theoretic considerations in Theorems 13–15. Additionally, we plot  $err_V$  which is in all cases slightly smaller than  $err_U$  but of the same order. Note that since  $\lambda$  is a



**Fig. 3** Approximation  $\lambda_h$  (left) and distribution of the estimator contribution  $\rho^2$  (right) for the example from Sect. 5.4

$\mathcal{T}$ -elementwise polynomial, an inverse inequality shows that  $h\|\lambda - \lambda_h\| \lesssim \|\lambda - \lambda_h\|_{-1}$  and thus  $\text{err}_V$  is equivalent to  $\|\mathbf{u} - \mathbf{u}_h\|_V$ .

### 5.4 Manufactured solution on L-shaped domain

We consider the same problem as given in [4, Section 5.2], where  $g = 0$ ,  $\Omega = (-2, 2)^2 \setminus [0, 2]^2$  and

$$f(r, \varphi) := -r^{2/3} \sin(2/3\varphi)(\gamma'(r)/r + \gamma''(r)) - 4/3r^{-1/3}\gamma'(r) \sin(2/3\varphi) - \delta(r),$$

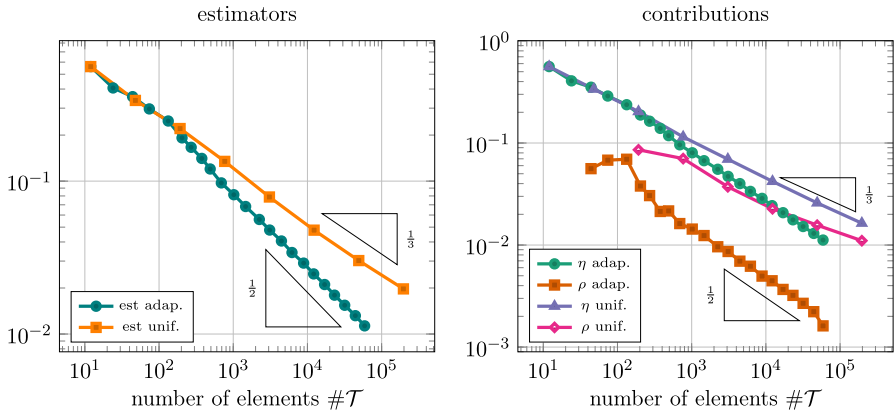


Fig. 4 Experimental convergence rates for the problem from Sect. 5.5

where  $(r, \varphi)$  denote polar coordinates. With  $r_* = 2(r - 1/4)$ ,  $\gamma, \delta$  are given by

$$\gamma(r) := \begin{cases} 1 & r_* < 0, \\ -6r_*^5 + 15r_*^4 - 10r_*^3 + 1 & 0 \leq r_* < 1, \\ 0 & 1 \leq r_*, \end{cases} \quad \delta(r) := \begin{cases} 0 & r \leq 5/4, \\ 1 & r > 5/4. \end{cases}$$

The exact solution then reads  $u(r, \varphi) = r^{2/3} \sin(2/3\varphi)\gamma(r)$ . Note that  $u$  has a generic singularity at the reentrant corner. We consider the discrete version of (VIa), where solutions are sought in the convex set  $K_h^s$ . We conducted various tests with  $\beta$  between 1 and 100 and the results were in all cases comparable. For the results displayed here we have used  $\beta = 3$ . Figure 2 displays convergence rates in the case of uniform and adaptive mesh-refinement. We note that in the first plot the lines for  $err_U$  and  $est$  are almost identical. In the second plot we compare the contributions of the overall error and estimator in the adaptive case. The lines for  $osc$  and  $\|div \sigma_h + \lambda_h + f\|$  are almost identical. This means that the estimator contribution  $\|div \sigma_h + \lambda_h + \Pi_h f\|$  in  $\eta$  is negligible and  $osc$  is dominating the overall estimator. We observe from the first plot that  $err_V$  is much smaller than  $err_U$  but has the same rate of convergence. In the uniform case we see that the errors and estimators approximately converge at rate 0.45. One would expect a smaller rate due to the singularity. However, in this example the solution has a large gradient so that the algorithm first refines the regions where the gradient resp.  $f$  is large. This preasymptotic behavior was also observed in [4, Section 5.2]. Nevertheless, adaptivity yields a significant error reduction.

Figure 3 shows the approximation  $\lambda_h$  (left column) and the distribution of the estimator contribution  $\rho^2$  (right column) on some adaptively refined meshes.

### 5.5 Unknown solution

For our final experiment, we choose  $\Omega = (-1, 1)^2 \setminus [-1, 0]^2$ ,  $f = 1$ , and the pyramid-like obstacle  $g(x) = \max\{0, \text{dist}(x, \partial\Omega_u) - \frac{1}{4}\}$ , where  $\Omega_u = (0, 1)^2$ . The solution

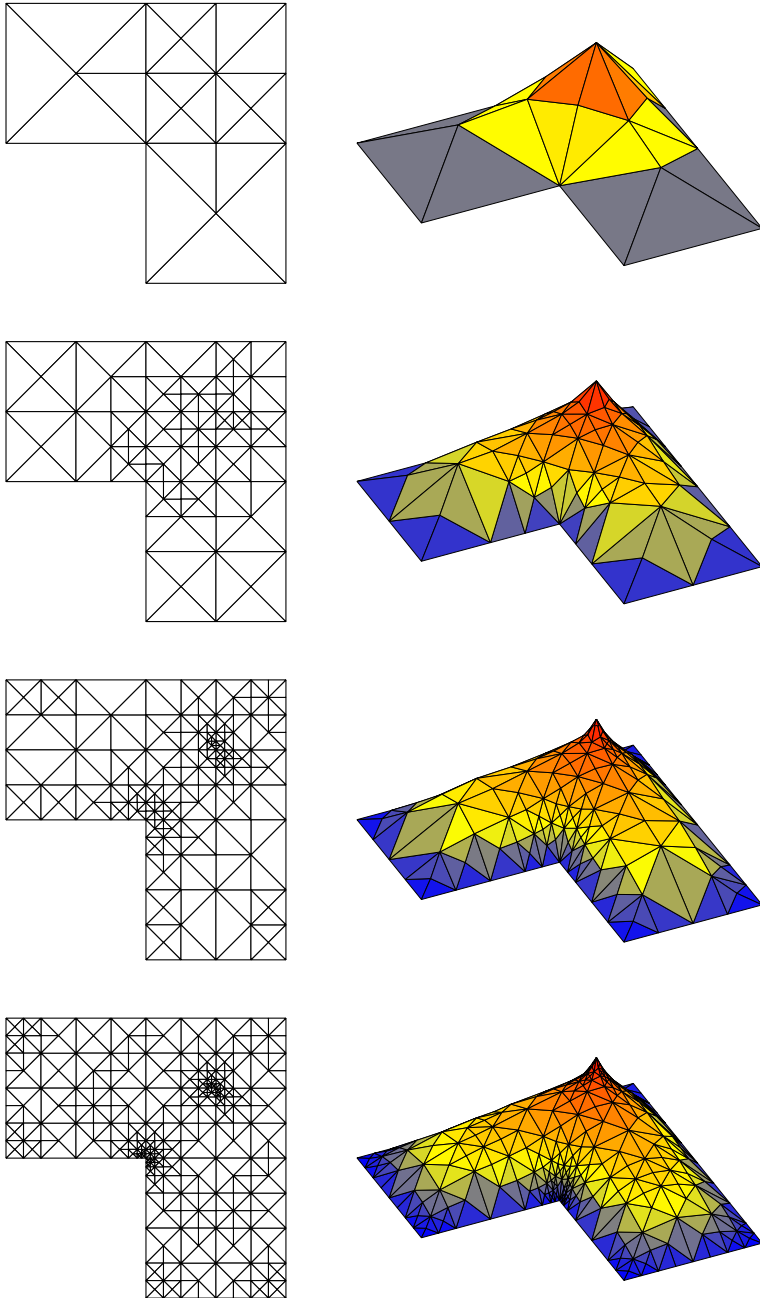


Fig. 5 Adaptively refined meshes and corresponding solution component  $u_h$  for the problem from Sect. 5.5

in this case is unknown. We solve the discrete version of (VIa) with convex set  $K_h^S$ . Since  $f$  is constant we have  $\text{osc} = 0$ . Figure 4 shows the overall estimator (left) and its contributions (right). We observe that uniform refinement leads to the reduced rate  $\frac{1}{3}$ , whereas for adaptive refinement we recover the optimal rate. Heuristically, we expect the solution to have a singularity at the reentrant corner as well as in the contact regions. This would explain the reduced rates. Figure 5 visualizes meshes produced by the adaptive algorithm and corresponding solution components  $u_h$ . We observe strong refinements towards the corner  $(0, 0)$  and around the point  $(\frac{1}{2}, \frac{1}{2})$ , which coincides with the tip of the pyramid obstacle.

**Acknowledgements** This work was supported by CONICYT through FONDECYT project “Least-squares methods for obstacle problems” under Grant 11170050.

### Appendix A: Non-convexity of functional $J$

Recall that the functional  $J(\cdot; f, g)$  is convex if and only if

$$a_1(\mathbf{u} - \mathbf{v}, \mathbf{u} - \mathbf{v}) \geq 0 \quad \text{for all } \mathbf{u}, \mathbf{v} \in K^S := \{(w, \boldsymbol{\chi}, v) \in U : w \geq g, v \geq 0\}.$$

In the following we construct a simple example that shows that the above inequality does not hold in general, thus  $J$  is not convex resp.  $a(\cdot, \cdot) = a_1(\cdot, \cdot)$  is not coercive.

To that end, let  $u \in H_0^1(\Omega)$  denote the solution of  $\Delta u = 1$  in the square domain  $\Omega = (0, d)^2$ . Then,  $u \leq 0$  in  $\Omega$ . Choose the obstacle as  $g = u$  (or  $g \leq u$ ). Note that  $\mathbf{v} := (0, 0, 0) \in K^S$  and that  $\mathbf{u} := (u, \boldsymbol{\sigma}, \Delta u) := (u, \nabla u, \Delta u) \in K^S$ . We have that

$$\|1\|_{-1} = \sup_{0 \neq v \in H_0^1(\Omega)} \frac{\langle 1, v \rangle}{\|\nabla v\|} = \sup_{0 \neq v \in H_0^1(\Omega)} \frac{-(\nabla u, \nabla v)}{\|\nabla v\|} = \|\nabla u\|.$$

Therefore  $\|1\|_{-1}^2 = \|\nabla u\|^2 = -\langle 1, u \rangle$ . Using this we infer that

$$a_1(\mathbf{u} - \mathbf{v}, \mathbf{u} - \mathbf{v}) = \|\text{div } \boldsymbol{\sigma} + 1\|^2 + \|\nabla u - \boldsymbol{\sigma}\|^2 + \langle 1, u \rangle = \|2\|^2 - \|1\|_{-1}^2.$$

Hence,  $a_1(\mathbf{u} - \mathbf{v}, \mathbf{u} - \mathbf{v}) < 0$  if and only if  $\|2\| < \|1\|_{-1}$ . Clearly,  $\|2\| = 2|\Omega|^{1/2} = 2d$ . We investigate the scaling of the negative order norm. Let  $\widehat{\Omega}$  denote the unit square  $(0, 1)^2$ . Then,

$$\|1\|_{-1} = \sup_{0 \neq v \in H_0^1(\Omega)} \frac{\langle 1, v \rangle}{\|\nabla v\|} = |\Omega| \sup_{0 \neq \widehat{v} \in H_0^1(\widehat{\Omega})} \frac{\langle 1, \widehat{v} \rangle_{\widehat{\Omega}}}{\|\nabla \widehat{v}\|_{\widehat{\Omega}}} = |\Omega| \|1\|_{-1, \widehat{\Omega}} =: |\Omega| C.$$

Finally,  $\|2\| < \|1\|_{-1}$  if

$$2d < Cd^2$$

which holds for sufficiently large  $d$ .

## Appendix B: Proof of Lemma 12

We use  $T_C := T_C(v)$  and  $T_{NC} := T_{NC}(v)$ . If  $|T_{NC}| = 0$ , then  $v = 0$  on  $T$  and the first inequality is trivial. Note that also  $\nabla v = 0$ . Therefore, the second inequality is trivial as well. From now on we thus assume  $|T_{NC}| > 0$ . Using that  $v = 0$  on  $T_C$  and the identity

$$v(\mathbf{x}) - v(\mathbf{y}) = (\mathbf{x} - \mathbf{y}) \cdot \int_0^1 \nabla v(s\mathbf{x} + (1-s)\mathbf{y}) ds \quad \text{for all } \mathbf{x}, \mathbf{y} \in T$$

we obtain that

$$\begin{aligned} \|v\|_T^2 &= \|v\|_{T_{NC}}^2 = \int_{T_{NC}} v(\mathbf{x})^2 d\mathbf{x} = \frac{1}{|T_C|} \int_{T_{NC}} \int_{T_C} (v(\mathbf{x}) - v(\mathbf{y}))^2 d\mathbf{y} d\mathbf{x} \\ &= \frac{1}{|T_C|} \int_{T_{NC}} \int_{T_C} \left( (\mathbf{x} - \mathbf{y}) \cdot \int_0^1 \nabla v(s\mathbf{x} + (1-s)\mathbf{y}) ds \right)^2 \\ &\leq \frac{h_T^2}{|T_C|} \int_T \int_T \int_0^1 |\nabla v(s\mathbf{x} + (1-s)\mathbf{y})|^2 ds d\mathbf{y} d\mathbf{x} \\ &= 2 \frac{h_T^2}{|T_C|} \int_T \int_{1/2}^1 \int_T |\nabla v(s\mathbf{x} + (1-s)\mathbf{y})|^2 d\mathbf{x} ds d\mathbf{y}, \end{aligned}$$

where in the ultimate step we have used symmetry in  $\mathbf{x}$  and  $\mathbf{y}$ . With the substitution  $\mathbf{z} = \phi(\mathbf{x}) = s\mathbf{x} + (1-s)\mathbf{y} \in T$  we further get that

$$\begin{aligned} \int_T \int_{1/2}^1 \int_T |\nabla v(s\mathbf{x} + (1-s)\mathbf{y})|^2 d\mathbf{x} ds d\mathbf{y} &= \int_T \int_{1/2}^1 s^{-n} \int_{\phi(T)} |\nabla v(\mathbf{z})|^2 d\mathbf{z} ds d\mathbf{y} \\ &\leq \int_T \int_{1/2}^1 s^{-n} \int_T |\nabla v(\mathbf{z})|^2 d\mathbf{z} ds d\mathbf{y} \\ &= |T| \frac{1 - (1/2)^{1-n}}{1-n} \|\nabla v\|_T^2. \end{aligned}$$

Putting altogether this proves the first inequality.

For the second inequality we use the fact that the gradient on level sets vanishes, see [13, Theorem 3.3]. This means that  $\nabla v = 0$  a.e. in  $T_C$ . Then, the same lines of proof as above (with  $v$  replaced by the components of  $\nabla v$ ) show the second inequality, which finishes the proof.

## Appendix C: Proof of Lemma 19

We consider the decompositions

$$\mathcal{J}^{k-1} = \underbrace{\{j \in \mathcal{J}^{k-1} : \lambda_j^k > 0\}}_{=: \mathcal{J}_1^{k-1}} \cup \underbrace{\{j \in \mathcal{J}^{k-1} : \lambda_j^k = 0\}}_{=: \mathcal{J}_2^{k-1}},$$



$$\mathcal{I}^{k-1} = \underbrace{\{i \in \mathcal{I}^{k-1} : i \in \mathcal{N}_\omega \text{ or } i \in \mathcal{N}_\gamma \text{ with } \mathbf{x}_i^k \geq \mathbf{g}_i\}}_{=: \mathcal{I}_1^{k-1}} \cup \underbrace{\{i \in \mathcal{I}^{k-1} \cap \mathcal{N}_\gamma : \mathbf{x}_i^k < \mathbf{g}_i\}}_{=: \mathcal{I}_2^{k-1}}.$$

For the decomposition of  $\mathcal{J}^k$  note that from Lines 8–9 of Algorithm 1 we have that  $\lambda_{\mathcal{N}_\gamma \cap \mathcal{I}^{k-1}} = 0$  and  $\mathbf{x}_{\mathcal{J}^k} = \mathbf{g}_{\mathcal{J}^k}$ . This yields

$$\mathcal{J}^k = \{j \in \mathcal{J}^{k-1} : \lambda_j^k > 0\} \cup \{i \in \mathcal{I}^{k-1} \cap \mathcal{N}_\gamma : \mathbf{x}_i^k < \mathbf{g}_i\} = \mathcal{J}_1^{k-1} \cup \mathcal{I}_2^{k-1}$$

If  $\mathcal{J}^k = \mathcal{J}^{k-1}$ , then, since all decompositions are disjoint,

$$\mathcal{I}_2^{k-1} = \emptyset = \mathcal{J}_2^{k-1}.$$

This also means that  $\widehat{\lambda}_j^k > 0$  if and only if  $\lambda_j^k > 0$  with  $\mathbf{x}_j^k = \mathbf{g}_j$ . Thus,

$$\begin{aligned} \lambda_{\mathcal{N}_\gamma}^k &= \max\{0, \lambda_{\mathcal{N}_\gamma}^k - C(\mathbf{x}_{\mathcal{N}_\gamma} - \mathbf{g}_{\mathcal{N}_\gamma})\}, \\ \lambda_{\mathcal{N}_\omega}^k &= 0, \end{aligned}$$

which implies that (18) is satisfied for  $\mathbf{x} = \mathbf{x}^k$ ,  $\lambda = \lambda^k$  or equivalently  $\mathbf{x} = \mathbf{x}^k$  solves (17).

## References

1. Attia, F.S., Cai, Z., Starke, G.: First-order system least squares for the Signorini contact problem in linear elasticity. *SIAM J. Numer. Anal.* **47**(4), 3027–3043 (2009)
2. Banz, L., Schröder, A.: Biorthogonal basis functions in *hp*-adaptive FEM for elliptic obstacle problems. *Comput. Math. Appl.* **70**(8), 1721–1742 (2015)
3. Banz, L., Stephan, E.P.: A posteriori error estimates of *hp*-adaptive IPDG-FEM for elliptic obstacle problems. *Appl. Numer. Math.* **76**, 76–92 (2014)
4. Bartels, S., Carstensen, C.: Averaging techniques yield reliable a posteriori finite element error control for obstacle problems. *Numer. Math.* **99**(2), 225–249 (2004)
5. Bochev, P., Gunzburger, M.: Least-squares finite element methods. In: *International Congress of Mathematicians*, vol. III, pp. 1137–1162. Eur. Math. Soc., Zürich (2006)
6. Bochev, P.B., Gunzburger, M.D.: *Least-Squares Finite Element Methods*. Applied Mathematical Sciences, vol. 166. Springer, New York (2009)
7. Braess, D.: A posteriori error estimators for obstacle problems—another look. *Numer. Math.* **101**(3), 415–421 (2005)
8. Bramble, J.H., Lazarov, R.D., Pasciak, J.E.: A least-squares approach based on a discrete minus one inner product for first order systems. *Math. Comput.* **66**(219), 935–955 (1997)
9. Burman, E., Hansbo, P., Larson, M.G., Stenberg, R.: Galerkin least squares finite element method for the obstacle problem. *Comput. Methods Appl. Mech. Eng.* **313**, 362–374 (2017)
10. Chen, Z., Nocketto, R.H.: Residual type a posteriori error estimates for elliptic obstacle problems. *Numer. Math.* **84**(4), 527–548 (2000)
11. Chouly, F., Hild, P.: A Nitsche-based method for unilateral contact problems: numerical analysis. *SIAM J. Numer. Anal.* **51**(2), 1295–1307 (2013)
12. Drouot, G., Hild, P.: Optimal convergence for discrete variational inequalities modelling Signorini contact in 2D and 3D without additional assumptions on the unknown contact set. *SIAM J. Numer. Anal.* **53**(3), 1488–1507 (2015)

13. Evans, L.C., Gariepy, R.F.: Measure Theory and Fine Properties of Functions. Textbooks in Mathematics, revised edn. CRC Press, Boca Raton (2015)
14. Falk, R.S.: Error estimates for the approximation of a class of variational inequalities. *Math. Comput.* **28**, 963–971 (1974)
15. Führer, T., Heuer, N., Stephan, E.P.: On the DPG method for Signorini problems. *IMA J. Numer. Anal.* **38**(4), 1893–1926 (2018)
16. Glowinski, R.: Numerical methods for nonlinear variational problems. In: Scientific Computation. Springer, Berlin (2008). Reprint of the 1984 original
17. Glowinski, R., Lions, J.-L., Trémolières, R.: Numerical Analysis of Variational Inequalities, Volume 8 of Studies in Mathematics and Its Applications. North-Holland Publishing Co., Amsterdam (1981). Translated from the French
18. Gustafsson, T., Stenberg, R., Videman, J.: Mixed and stabilized finite element methods for the obstacle problem. *SIAM J. Numer. Anal.* **55**(6), 2718–2744 (2017)
19. Gustafsson, T., Stenberg, R., Videman, J.: On finite element formulations for the obstacle problem—mixed and stabilised methods. *Comput. Methods Appl. Math.* **17**(3), 413–429 (2017)
20. Gustafsson, T., Stenberg, R., Videman, J.: A stabilised finite element method for the plate obstacle problem. *BIT* **59**(1), 97–124 (2019)
21. Hintermüller, M., Ito, K., Kunisch, K.: The primal–dual active set strategy as a semismooth Newton method. *SIAM J. Optim.* **13**(3), 865–888 (2003), 2002
22. Hoppe, R.H.W., Kornhuber, R.: Adaptive multilevel methods for obstacle problems. *SIAM J. Numer. Anal.* **31**(2), 301–323 (1994)
23. Kärkkäinen, T., Kunisch, K., Tarvainen, P.: Augmented Lagrangian active set methods for obstacle problems. *J. Optim. Theory Appl.* **119**(3), 499–533 (2003)
24. Karkulik, M., Pavlicek, D., Praetorius, D.: On 2D newest vertex bisection: optimality of mesh-closure and  $H^1$ -stability of  $L_2$ -projection. *Constr. Approx.* **38**(2), 213–234 (2013)
25. Kinderlehrer, D., Stampacchia, G.: An Introduction to Variational Inequalities and Their Applications, Volume 31 of Classics in Applied Mathematics. Society for Industrial and Applied Mathematics (SIAM), Philadelphia (2000). Reprint of the 1980 original
26. Krause, R., Müller, B., Starke, G.: An adaptive least-squares mixed finite element method for the Signorini problem. *Numer. Methods Partial Differ. Equ.* **33**(1), 276–289 (2017)
27. Nochetto, R.H., Siebert, K.G., Veerer, A.: Pointwise a posteriori error control for elliptic obstacle problems. *Numer. Math.* **95**(1), 163–195 (2003)
28. Nochetto, R.H., Siebert, K.G., Veerer, A.: Fully localized a posteriori error estimators and barrier sets for contact problems. *SIAM J. Numer. Anal.* **42**(5), 2118–2135 (2005)
29. Rodrigues, J.-F.: Obstacle Problems in Mathematical Physics, Volume 134 of North-Holland Mathematics Studies. North-Holland Publishing Co., Amsterdam (1987). *Notas de Matemática [Mathematical Notes]*, 114
30. Stevenson, R.: The completion of locally refined simplicial partitions created by bisection. *Math. Comput.* **77**(261), 227–241 (2008)
31. Veerer, A.: Efficient and reliable a posteriori error estimators for elliptic obstacle problems. *SIAM J. Numer. Anal.* **39**(1), 146–167 (2001)
32. Weiss, A., Wohlmuth, B.I.: A posteriori error estimator for obstacle problems. *SIAM J. Sci. Comput.* **32**(5), 2627–2658 (2010)

**Publisher's Note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.