

# OPTIMAL CONTROL PROBLEMS WITH CONTROL COMPLEMENTARITY CONSTRAINTS

EXISTENCE RESULTS, OPTIMALITY CONDITIONS, AND A PENALTY METHOD

Christian Clason<sup>\*</sup>    Yu Deng<sup>†</sup>    Patrick Mehlitz<sup>‡</sup>    Uwe Prüfert<sup>§</sup>

January 23, 2019

**Abstract** A special class of optimal control problems with complementarity constraints on the control functions is studied. It is shown that such problems possess optimal solutions whenever the underlying control space is a first-order Sobolev space. After deriving necessary optimality conditions of strong stationarity-type, a penalty method based on the Fischer–Burmeister function is suggested and its theoretical properties are analyzed. Finally, the numerical treatment of the problem is discussed and results of computational experiments are presented.

## 1 INTRODUCTION

Complementarity conditions appear in many mathematical optimization problems arising from real-world applications, and this phenomenon is not restricted to the finite-dimensional setting, see [23, 33, 34] and the references therein. A prominent example for a complementarity problem in function spaces is the optimal control of the obstacle problem, see [19] for an overview of existing literature. Mathematical problems with complementarity constraints (MPCCs) suffer from an inherent lack of regularity, see [36, Proposition 1.1] and [24, Lemma 3.1] for the finite- and infinite-dimensional situation, respectively, which is why the construction of suitable optimality conditions, constraint qualifications, and numerical methods is a challenging task. Using so-called *NCP-functions*, complementarity constraints can be transformed into possibly nonsmooth equality constraints that can be handled by, e.g., Newton-type methods, see [10, 14, 33] and the references therein. A satisfying overview of NCP-functions can be found in [31]. One of the most popular NCP-functions is the so-called *Fischer–Burmeister function*  $\phi: \mathbb{R}^2 \rightarrow \mathbb{R}$  given by

$$(1.1) \quad \forall a, b \in \mathbb{R}: \quad \phi(a, b) := \sqrt{a^2 + b^2} - a - b,$$

<sup>\*</sup>Universität Duisburg-Essen, Faculty of Mathematics, 45117 Essen, Germany, [christian.clason@uni-due.de](mailto:christian.clason@uni-due.de), <https://udue.de/clason>

<sup>†</sup>Technische Universität Bergakademie Freiberg, Faculty of Mathematics and Computer Science, 09596 Freiberg, Germany, [yu.deng@math.tu-freiberg.de](mailto:yu.deng@math.tu-freiberg.de), <http://www.mathe.tu-freiberg.de/nmo/mitarbeiter/yu-deng>

<sup>‡</sup>Brandenburgische Technische Universität Cottbus-Senftenberg, Institute of Mathematics, Chair of Optimal Control, 03046 Cottbus, Germany, [mehlitz@b-tu.de](mailto:mehlitz@b-tu.de), <https://www.b-tu.de/fg-optimale-steuerung/team/dr-patrick-mehlitz>

<sup>§</sup>Technische Universität Bergakademie Freiberg, Faculty of Mathematics and Computer Science, 09596 Freiberg, Germany, [uwe.pruefert@math.tu-freiberg.de](mailto:uwe.pruefert@math.tu-freiberg.de), <http://www.mathe.tu-freiberg.de/nmo/mitarbeiter/uwe-pruefert>

see [15]. Obviously, one has

$$\forall a, b \in \mathbb{R}: \quad \phi(a, b) = 0 \iff a \geq 0 \wedge b \geq 0 \wedge ab = 0,$$

which (by definition) holds for all NCP-functions. Thus, NCP-functions allow the replacement of a complementarity condition by a single equality constraint. In [33], it is shown that NCP-functions can be applied to solve complementarity problems in function space settings as well.

In this paper, an optimal control problem with complementarity constraints on the control functions is studied. Control complementarity constraints have been the subject of several recent papers including [6, 18, 25, 27]. Classically, such constraints arise from reformulating a bilevel optimal control problem with lower level control constraints as a single-level problem using lower level first-order optimality conditions, see [24, Section 5]. On the other hand, control complementarity constraints are closely related to switching conditions on the control functions, see [7–9] and the references therein. Here, it will be shown that such problems possess an optimal solution if the control space is taken as  $H^1(\Omega)$ . Recently, optimal control problems with control constraints in first-order Sobolev spaces were studied in [11, 12].

It will also be demonstrated that the Fischer–Burmeister function can be used to design penalty methods that can be exploited to find minimizers of the corresponding optimal control problem. One major advantage of this procedure is that the resulting penalized problems are unconstrained. In contrast, simply penalizing the equilibrium condition and leaving the non-negativity conditions in the constraints would lead to the appearance of Lagrange multipliers from  $H^1(\Omega)^*$  in the necessary optimality conditions of the penalized problems, which would cause some theoretical and numerical difficulties due to the presumed high regularity of the control space, see [12].

The paper is organized as follows: In the remainder of this section, the basic notation is introduced. Afterwards, the optimal control problem is formally stated and the existence of solutions is discussed in Section 2. Necessary optimality conditions of strong stationarity-type are derived in Section 3. Section 4 is dedicated to the theoretical investigation of a penalization procedure. The practical implementation of the proposed numerical method and some corresponding examples are discussed in Section 5 and Section 6, respectively. A brief summary as well as some concluding remarks are presented in Section 7.

**Basic notation** For a real Banach space  $\mathcal{X}$ ,  $\|\cdot\|_{\mathcal{X}}$  denotes its norm. The expression  $\mathcal{X}^*$  is used to represent the topological dual space of  $\mathcal{X}$ . Let  $\langle \cdot, \cdot \rangle_{\mathcal{X}} : \mathcal{X}^* \times \mathcal{X} \rightarrow \mathbb{R}$  be the associated dual pairing. For another Banach space  $\mathcal{Y}$ ,  $\mathbb{L}[\mathcal{X}, \mathcal{Y}]$  represents the Banach space of all bounded, linear operators which map from  $\mathcal{X}$  to  $\mathcal{Y}$ . For  $F \in \mathbb{L}[\mathcal{X}, \mathcal{Y}]$ ,  $F^* \in \mathbb{L}[\mathcal{Y}^*, \mathcal{X}^*]$  denotes its adjoint. If  $\mathcal{X} \subset \mathcal{Y}$  holds true while the associated identity mapping from  $\mathcal{X}$  into  $\mathcal{Y}$  is continuous, then  $\mathcal{X}$  is said to be continuously embedded into  $\mathcal{Y}$ , denoted by  $\mathcal{X} \hookrightarrow \mathcal{Y}$ .

Recall that a set  $A \subset \mathcal{X}$  is said to be weakly sequentially closed if all the limit points of weakly convergent sequences contained in  $A$  belong to  $A$  as well, and that any closed, convex set is weakly sequentially closed by Mazur’s lemma. For any  $A \subset \mathcal{X}$ , define the polar cone

$$A^\circ := \{x^* \in \mathcal{X}^* \mid \forall x \in A: \langle x^*, x \rangle_{\mathcal{X}} \leq 0\},$$

as well as the annihilator

$$A^\perp := \{x^\star \in X^\star \mid \forall x \in A: \langle x^\star, x \rangle_X = 0\}.$$

By definition,  $A^\perp = A^\circ \cap (-A)^\circ$  holds true. It is well known that  $A^\circ$  is a nonempty, closed, convex cone while  $A^\perp$  is a closed subspace of  $X^\star$ . For an arbitrary vector  $x \in X$ , set  $x^\perp := \{x\}^\perp$  for the sake of brevity.

Finally, if a function  $F: X \rightarrow Y$  is Fréchet differentiable at  $\bar{x} \in X$ , then the bounded, linear operator  $F'(\bar{x}) \in \mathbb{L}[X, Y]$  denotes its Fréchet derivative at  $\bar{x}$ .

**Function spaces** For an arbitrary bounded domain  $\Omega \subset \mathbb{R}^d$  and  $p \in [1, \infty]$ ,  $L^p(\Omega)$  denotes the usual Lebesgue space of (equivalence classes of) Lebesgue measurable functions mapping from  $\Omega$  to  $\mathbb{R}$ , which is equipped with the usual norm. It is well known that for  $p \in [1, \infty)$ , the space  $L^p(\Omega)^\star$  is isometric to  $L^{p'}(\Omega)$  for  $p' \in (1, \infty]$  such that  $1/p + 1/p' = 1$ . The associated dual pairing is given by

$$\forall u \in L^p(\Omega) \forall v \in L^{p'}(\Omega): \quad \langle v, u \rangle_{L^p(\Omega)} := \int_{\Omega} u(x)v(x)dx.$$

Recall that  $L^2(\Omega)$  is a Hilbert space whose dual  $L^2(\Omega)^\star$  will be identified with  $L^2(\Omega)$  by means of Riesz' representation theorem. For an arbitrary function  $u \in L^1(\Omega)$ ,  $\text{supp } u := \{x \in \Omega \mid u(x) \neq 0\}$  denotes the support of  $u$ . Supposing that  $A \subset \Omega$  is a Lebesgue measurable set,  $\chi_A: \Omega \rightarrow \mathbb{R}$  represents the characteristic function of  $A$  which is 1 for all  $x \in A$  and 0 else. Clearly, for a bounded domain  $\Omega$  and  $p \in [1, \infty)$ , the relation  $\|\chi_A\|_{L^p(\Omega)} = |A|^{1/p}$  is obtained where  $|A|$  denotes the Lebesgue measure of  $A$ .

The Banach space of all weakly differentiable functions from  $L^2(\Omega)$  whose weak derivatives belong to  $L^2(\Omega)$  is denoted by  $H^1(\Omega)$ . It is equipped with the usual norm

$$\forall y \in H^1(\Omega): \quad \|y\|_{H^1(\Omega)} := \left( \|y\|_{L^2(\Omega)}^2 + \sum_{i=1}^d \|\partial_{x_i} y\|_{L^2(\Omega)}^2 \right)^{1/2}.$$

Clearly,  $H^1(\Omega)$  is a Hilbert space. However, its dual  $H^1(\Omega)^\star$  will not be identified with  $H^1(\Omega)$  so that  $H^1(\Omega)$ ,  $L^2(\Omega)$ , and  $H^1(\Omega)^\star$  form a so-called Gelfand triple, i.e., they satisfy the relations  $H^1(\Omega) \hookrightarrow L^2(\Omega) \hookrightarrow H^1(\Omega)^\star$ . A detailed study of duality in Sobolev spaces can be found in [1, Section 3].

Whenever  $\Omega$  satisfies the so-called cone condition, see [1, Section 4], then the embedding  $H^1(\Omega) \hookrightarrow L^2(\Omega)$  is compact, see [1, Theorem 6.3]. In this paper,  $E \in \mathbb{L}[H^1(\Omega), L^2(\Omega)]$  is used to denote the latter.

For later use, let  $L_+^2(\Omega) \subset L^2(\Omega)$  and  $H_+^1(\Omega) \subset H^1(\Omega)$  denote the nonempty, closed, and convex cones of almost everywhere nonnegative functions in  $L^2(\Omega)$  and  $H^1(\Omega)$ , respectively.

## 2 PROBLEM SETTING AND EXISTENCE OF OPTIMAL SOLUTIONS

In this work, the model **optimal control problem with control complementarity constraints**

$$(OC^4) \quad \begin{cases} \frac{1}{2} \|D[y] - y_d\|_{\mathcal{D}}^2 + J(u, v) \rightarrow \min_{y, u, v} \\ A[y] - B[u] - C[v] = 0 \\ (u, v) \in \mathbb{C} \end{cases}$$

is studied, where for some  $\alpha_1, \alpha_2 \geq 0$  and  $\varepsilon \geq 0$ ,

$$\forall u, v \in H^1(\Omega): \quad J(u, v) := \frac{\alpha_1}{2} \|u\|_{L^2(\Omega)}^2 + \frac{\alpha_2}{2} \|v\|_{L^2(\Omega)}^2 + \frac{\varepsilon}{2} \left( \|u\|_{H^1(\Omega)}^2 + \|v\|_{H^1(\Omega)}^2 \right),$$

and  $\mathbb{C}$  denotes the complementarity set

$$\mathbb{C} := \{(w, z) \in H^1(\Omega)^2 \mid 0 \leq w(x) \perp z(x) \geq 0 \text{ a.e. on } \Omega\}.$$

Observing that  $A$  can represent a differential operator, one can interpret (OC<sup>4</sup>) as an optimal control problem with complementarity constraints on the control functions that can be used to model switching requirements on the controls. In the context of ordinary differential equations, optimal control problems with mixed control-state complementarity constraints have been studied in [6, 18, 27] recently. In [19, 20, 25], the interested reader can find some theoretical investigations of optimization problems with complementarity constraints with respect to the function spaces  $L^2(\Omega)$  and  $H_0^1(\Omega)$ . Recently, optimal control problems with switching constraints related to (OC<sup>4</sup>) have been studied in [8, 9].

For the remainder of this work, the following standing assumptions on the problem (OC<sup>4</sup>) are postulated.

**Assumption 2.1.** The domain  $\Omega \subset \mathbb{R}^d$  is nonempty, bounded, and satisfies the cone condition. Its boundary will be denoted by  $\text{bd } \Omega$ . Let the observation space  $\mathcal{D}$  as well as the state space  $\mathcal{Y}$  be Hilbert spaces. The target  $y_d \in \mathcal{D}$  will be fixed. The operator  $A \in \mathbb{L}[\mathcal{Y}, \mathcal{Y}^*]$  is an isomorphism while  $B, C \in \mathbb{L}[H^1(\Omega), \mathcal{Y}^*]$  and  $D \in \mathbb{L}[\mathcal{Y}, \mathcal{D}]$  are arbitrarily chosen. Finally,  $\varepsilon > 0$  holds.

Let  $S \in \mathbb{L}[H^1(\Omega)^2, \mathcal{D}]$  be the control-to-observation operator which maps any pair of controls  $(u, v) \in H^1(\Omega)^2$  to  $D[y]$ , where  $y \in \mathcal{Y}$  is the associated uniquely determined solution of the state equation

$$A[y] - B[u] - C[v] = 0.$$

Then,  $S$  is a well-defined continuous linear operator since  $A$  is assumed to be an isomorphism.

In the following, the existence of optimal solutions to (OC<sup>4</sup>) is discussed. First, the overall  $H^1$ -setting needed for the further theoretical treatment of (OC<sup>4</sup>) is analyzed in Section 2.1. Some comments on the setting where controls come from  $L^2(\Omega)$  are presented in Section 2.2.

## 2.1 FIRST-ORDER SOBOLEV SPACES

Since the objective function of (OC<sup>4</sup>) is continuously Fréchet differentiable, convex, and bounded from below, the only critical point for existence is the weak sequential closedness of the complementarity set  $\mathbb{C}$ .

**Lemma 2.2.** *The set  $\mathbb{C}$  is closed.*

*Proof.* Let  $\{(u_k, v_k)\}_{k \in \mathbb{N}} \subset \mathbb{C}$  be a sequence converging to  $(\bar{u}, \bar{v}) \in H^1(\Omega)^2$ . Due to the continuity of the embedding  $H^1(\Omega) \hookrightarrow L^2(\Omega)$ , the strong convergences  $u_k \rightarrow \bar{u}$  and  $v_k \rightarrow \bar{v}$  hold in  $L^2(\Omega)$ . In particular, these convergences hold (at least along a subsequence) pointwise almost everywhere. Due to the closedness of the set  $\{(a, b) \in \mathbb{R}^2 \mid 0 \leq a \perp b \geq 0\}$ , the desired result follows.  $\square$

Although  $\mathbb{C}$  is a nonconvex set, the compactness of the embedding  $H^1(\Omega) \hookrightarrow L^2(\Omega)$  can be used in order to show that  $\mathbb{C}$  is weakly sequentially closed.

**Lemma 2.3.** *The set  $\mathbb{C}$  is weakly sequentially closed.*

*Proof.* First, a similar proof as for Lemma 2.2 shows that the complementarity set in  $L^2(\Omega)$  given by

$$(2.1) \quad \begin{aligned} \tilde{\mathbb{C}} &:= \{(w, z) \in L^2(\Omega)^2 \mid 0 \leq w(x) \perp z(x) \geq 0 \text{ a.e. on } \Omega\} \\ &= \{(w, z) \in L_+^2(\Omega)^2 \mid \langle w, z \rangle_{L^2(\Omega)} = 0\} \end{aligned}$$

is closed as well.

Next, choose a sequence  $\{(u_k, v_k)\}_{k \in \mathbb{N}} \subset \mathbb{C}$  converging weakly to  $(\bar{u}, \bar{v}) \in H^1(\Omega)^2$ . Exploiting  $u_k \rightharpoonup \bar{u}$  and  $v_k \rightharpoonup \bar{v}$  as well as the compactness of the embedding  $H^1(\Omega) \hookrightarrow L^2(\Omega)$ , there is a subsequence of  $\{(u_k, v_k)\}_{k \in \mathbb{N}}$  that converges strongly to  $(\bar{u}, \bar{v})$  in  $L^2(\Omega)^2$ . Due to the closedness of  $\tilde{\mathbb{C}}$  in  $L^2(\Omega)^2$ ,  $(\bar{u}, \bar{v}) \in \tilde{\mathbb{C}} \cap H^1(\Omega)^2$  holds, and, consequently,  $(\bar{u}, \bar{v})$  is already an element of  $\mathbb{C}$ . Thus,  $\mathbb{C}$  is weakly sequentially closed.  $\square$

As a corollary, the existence of optimal solutions to (OC<sup>4</sup>) is obtained.

**Corollary 2.4.** *The problem (OC<sup>4</sup>) possesses an optimal solution.*

*Proof.* The objective functional of (OC<sup>4</sup>) is continuously Fréchet differentiable, convex, and (due to  $\varepsilon > 0$ ) coercive. Furthermore, by Lemma 2.3, the complementarity set  $\mathbb{C}$  is weakly sequentially closed, and so is the feasible set induced by the PDE constraint. Hence, the claim follows by application of Tonelli's direct method.  $\square$

## 2.2 LEBESGUE SPACES

In the remainder of this section, the existence of optimal controls in  $L^2(\Omega)$  is investigated. In this case, the corresponding model problem is given by

$$(OC_{L^2}) \quad \begin{cases} \frac{1}{2} \|D[y] - y_d\|_{\mathcal{D}}^2 + \frac{\alpha_1}{2} \|u\|_{L^2(\Omega)}^2 + \frac{\alpha_2}{2} \|v\|_{L^2(\Omega)}^2 \rightarrow \min_{y, u, v} \\ A[y] - \tilde{B}[u] - \tilde{C}[v] = 0 \\ (u, v) \in \tilde{\mathbb{C}} \end{cases}$$

where the complementarity set  $\tilde{\mathbb{C}}$  has been defined in (2.1). Furthermore,  $\tilde{B}, \tilde{C} \in \mathbb{L}[L^2(\Omega), \mathcal{Y}^*]$  need to be chosen. As already shown in the proof of Lemma 2.3,  $\tilde{\mathbb{C}}$  is closed. However,  $\tilde{\mathbb{C}}$  is in general *not* weakly sequentially closed, as the following example shows.

**Example 2.1.** For any  $k \in \mathbb{N}$ , define the two open sets

$$\begin{aligned} P_k &:= \left\{ x \in \mathbb{R}^d \mid \prod_{j=1}^d \sin(k\pi x_j) > 0 \right\}, \\ Q_k &:= \left\{ x \in \mathbb{R}^d \mid \prod_{j=1}^d \sin(k\pi x_j) < 0 \right\}. \end{aligned}$$

Now, set  $u_k := \chi_{\Omega \cap P_k}$  and  $v_k := \chi_{\Omega \cap Q_k}$ . Obviously,  $(u_k, v_k) \in \tilde{\mathbb{C}}$  holds true for all  $k \in \mathbb{N}$ . Furthermore, the sequence  $\{(u_k, v_k)\}_{k \in \mathbb{N}} \subset L^2(\Omega)^2$  converges weakly to the point  $(\frac{1}{2}\chi_\Omega, \frac{1}{2}\chi_\Omega)$ , which does not belong to  $\tilde{\mathbb{C}}$ . Thus,  $\tilde{\mathbb{C}}$  is not weakly sequentially closed.

It may still happen that there exists an optimal solution of the complementarity-constrained problem (OC<sub>L<sup>2</sup></sub>), as illustrated by the following example. For  $\mathcal{D} := L^2(\Omega)$  and  $\mathcal{Y} := H^1(\Omega)$ , consider the elliptic optimal control problem

$$(2.2) \quad \begin{cases} \frac{1}{2} \|E[y] - y_d\|_{L^2(\Omega)}^2 + \frac{\alpha_1}{2} \|u\|_{L^2(\Omega)}^2 + \frac{\alpha_2}{2} \|v\|_{L^2(\Omega)}^2 \rightarrow \min_{y, u, v} \\ -\nabla \cdot (C \nabla y) + \mathbf{a}y = \chi_{\Omega_u} u + \chi_{\Omega_v} v & \text{a.e. on } \Omega \\ \vec{\mathbf{n}} \cdot (C \nabla y) + \mathbf{q}y = 0 & \text{a.e. on } \text{bd } \Omega \\ (u, v) \in \widetilde{\mathbb{C}} \end{cases}$$

where we recall that  $E$  represents the natural embedding from  $H^1(\Omega)$  into  $L^2(\Omega)$ ,  $\alpha_1, \alpha_2 > 0$  are constants, and  $C \in L^\infty(\Omega; S^d(\mathbb{R}))$  (where  $S^d(\mathbb{R})$  denotes the set of real symmetric  $d \times d$  matrices) satisfies the condition of uniform ellipticity, i.e.,

$$(2.3) \quad \exists c_0 > 0 \quad \forall x \in \Omega \quad \forall \xi \in \mathbb{R}^d: \quad \xi^\top C(x) \xi \geq c_0 |\xi|_2^2.$$

Moreover,  $\mathbf{a} \in L^\infty(\Omega)$  and  $\mathbf{q} \in L^\infty(\text{bd } \Omega)$  are nonnegative and satisfy  $\|\mathbf{a}\|_{L^\infty(\Omega)} + \|\mathbf{q}\|_{L^\infty(\text{bd } \Omega)} > 0$ , and  $\Omega_u, \Omega_v \subset \Omega$  are measurable sets of positive measure satisfying  $\Omega_u \cup \Omega_v = \Omega$ . Here, the PDE constraint is interpreted in the weak sense. It is well known that the associated differential operator  $A$  is elliptic, see [13, Section 6], and, thus, an isomorphism.

**Proposition 2.5.** *The problem (2.2) possesses an optimal solution.*

*Proof.* Assume without loss of generality that  $\alpha_1 \leq \alpha_2$ ; the other case can be handled analogously. Consider then the surrogate optimal control problem

$$(2.4) \quad \begin{cases} \frac{1}{2} \|E[y] - y_d\|_{L^2(\Omega)}^2 + \frac{1}{2} \|(\sqrt{\alpha_1} \chi_{\Omega_u} + \sqrt{\alpha_2} \chi_{\Omega_v \setminus \Omega_u}) z\|_{L^2(\Omega)}^2 \rightarrow \min_{y, z} \\ -\nabla \cdot (C \nabla y) + \mathbf{a}y = z & \text{a.e. on } \Omega \\ \vec{\mathbf{n}} \cdot (C \nabla y) + \mathbf{q}y = 0 & \text{a.e. on } \text{bd } \Omega \\ z \in L_+^2(\Omega). \end{cases}$$

Note that its objective is equivalent to

$$H^1(\Omega) \times L^2(\Omega) \ni (y, z) \mapsto \frac{1}{2} \|E[y] - y_d\|_{L^2(\Omega)}^2 + \frac{\alpha_1}{2} \|\chi_{\Omega_u} z\|_{L^2(\Omega)}^2 + \frac{\alpha_2}{2} \|\chi_{\Omega_v \setminus \Omega_u} z\|_{L^2(\Omega)}^2 \in \mathbb{R}.$$

The ellipticity of the underlying PDE in (2.4) implies that the associated control-to-observation operator  $\check{S}: L^2(\Omega) \rightarrow L^2(\Omega)$  is linear and continuous, see [13, Section 6.2]. Observing that  $\Omega_u \cup \Omega_v = \Omega$  holds by assumption, the reduced objective functional

$$L^2(\Omega) \ni z \mapsto \frac{1}{2} \|\check{S}[z] - y_d\|_{L^2(\Omega)}^2 + \frac{\alpha_1}{2} \|\chi_{\Omega_u} z\|_{L^2(\Omega)}^2 + \frac{\alpha_2}{2} \|\chi_{\Omega_v \setminus \Omega_u} z\|_{L^2(\Omega)}^2 \in \mathbb{R}$$

is convex, continuous, and coercive. This shows that the optimal control problem (2.4) possesses an optimal solution  $(\bar{y}, \bar{z}) \in H^1(\Omega) \times L^2(\Omega)$  with objective value  $\bar{m} \in \mathbb{R}$ .

Let  $(y, u, v) \in H^1(\Omega) \times L^2(\Omega) \times L^2(\Omega)$  be feasible to (2.2). Defining  $z := \chi_{\Omega_u} u + \chi_{\Omega_v \setminus \Omega_u} v$ ,  $(y, z)$  is feasible for (2.4). Then, the estimate

$$\begin{aligned} & \frac{1}{2} \|E[y] - y_d\|_{L^2(\Omega)}^2 + \frac{\alpha_1}{2} \|u\|_{L^2(\Omega)}^2 + \frac{\alpha_2}{2} \|v\|_{L^2(\Omega)}^2 \\ & \geq \frac{1}{2} \|E[y] - y_d\|_{L^2(\Omega)}^2 + \frac{\alpha_1}{2} \|\chi_{\Omega_u} u\|_{L^2(\Omega)}^2 + \frac{\alpha_2}{2} \|\chi_{\Omega_v \setminus \Omega_u} v\|_{L^2(\Omega)}^2 \\ & = \frac{1}{2} \|E[y] - y_d\|_{L^2(\Omega)}^2 + \frac{\alpha_1}{2} \|\chi_{\Omega_u} z\|_{L^2(\Omega)}^2 + \frac{\alpha_2}{2} \|\chi_{\Omega_v \setminus \Omega_u} z\|_{L^2(\Omega)}^2 \geq \bar{m} \end{aligned}$$

is obtained. In particular, the objective value of (2.2) is bounded from below by  $\bar{m}$ .

Define  $\bar{u} := \chi_{\Omega_u} \bar{z}$  and  $\bar{v} := \chi_{\Omega_v \setminus \Omega_u} \bar{z}$ . Then,  $(\bar{y}, \bar{u}, \bar{v})$  is feasible to (2.2) since  $\bar{y}$  is the state associated with  $\bar{z}$  and  $\chi_{\Omega_u} \bar{u} + \chi_{\Omega_v} \bar{v} = \bar{z}$  holds true. Moreover, the relation

$$\begin{aligned} & \frac{1}{2} \|E[\bar{y}] - y_d\|_{L^2(\Omega)}^2 + \frac{\alpha_1}{2} \|\bar{u}\|_{L^2(\Omega)}^2 + \frac{\alpha_2}{2} \|\bar{v}\|_{L^2(\Omega)}^2 \\ & = \frac{1}{2} \|E[\bar{y}] - y_d\|_{L^2(\Omega)}^2 + \frac{\alpha_1}{2} \|\chi_{\Omega_u} \bar{z}\|_{L^2(\Omega)}^2 + \frac{\alpha_2}{2} \|\chi_{\Omega_v \setminus \Omega_u} \bar{z}\|_{L^2(\Omega)}^2 = \bar{m} \end{aligned}$$

follows. Thus,  $(\bar{y}, \bar{u}, \bar{v})$  is an optimal solution of (2.2).  $\square$

Note that the proof of Proposition 2.5 yields a strategy for the solution of (2.2) by means of standard arguments from optimal control by solving the surrogate problem (2.4).

### 3 OPTIMALITY CONDITIONS

Consider the so-called state-reduced problem

$$(3.1) \quad \begin{cases} \frac{1}{2} \|S[u, v] - y_d\|_{\mathcal{D}}^2 + J(u, v) \rightarrow \min_{u, v} \\ (u, v) \in \mathbb{C} \end{cases}$$

which is equivalent to (OC<sup>4</sup>) by definition of the control-to-observation operator  $S$ . Using the embedding operator  $E: H^1(\Omega) \rightarrow L^2(\Omega)$ , (3.1) can be stated equivalently as

$$(3.2) \quad \begin{cases} \frac{1}{2} \|S[u, v] - y_d\|_{\mathcal{D}}^2 + J(u, v) \rightarrow \min_{u, v} \\ E[u] \in L_+^2(\Omega) \\ E[v] \in L_+^2(\Omega) \\ \langle E[u], E[v] \rangle_{L^2(\Omega)} = 0 \end{cases}$$

which is a generalized MPCC in the Banach space  $L^2(\Omega)$ . It was shown in [24, Lemma 3.1] that Robinson's constraint qualification, see [4, Section 2.3.4] for its definition, some discussion, and suitable references to the literature, does not hold at the feasible points of this problem. Moreover, since  $E$  is not surjective, the constraint qualifications needed to show that locally optimal solutions of this problem satisfy MPCC-tailored stationarity conditions (e.g., the weak or strong stationarity conditions) are not satisfied, see [24, 34] for details.

On the other hand, it is still possible to derive necessary optimality conditions for (3.1) using a standard trick from finite-dimensional MPCC theory: Define appropriate surrogate problems

which do not contain a complementarity constraint anymore and handle them with the classical KKT conditions in Banach spaces.

In order to formulate an appropriate surrogate problem, let  $(\bar{u}, \bar{v}) \in H^1(\Omega)^2$  be a feasible point of (3.1) and define the measurable sets

$$(3.3) \quad I^{+0}(\bar{u}, \bar{v}) := \{x \in \Omega \mid \bar{u}(x) > 0 \wedge \bar{v}(x) = 0\},$$

$$(3.4) \quad I^{0+}(\bar{u}, \bar{v}) := \{x \in \Omega \mid \bar{u}(x) = 0 \wedge \bar{v}(x) > 0\},$$

$$(3.5) \quad I^{00}(\bar{u}, \bar{v}) := \{x \in \Omega \mid \bar{u}(x) = 0 \wedge \bar{v}(x) = 0\}.$$

Noting that  $L^2(\Omega)$  is a space of equivalence classes, it should be mentioned that these sets are well-defined up to sets of Lebesgue measure zero. This will be taken into account in the following. If  $(\bar{u}, \bar{v})$  is a locally optimal solution of (3.1), then it is also a locally optimal solution of the auxiliary problems

$$(rNLP_{\bar{u}}) \quad \left\{ \begin{array}{ll} \frac{1}{2} \|S[u, v] - y_d\|_{\mathcal{D}}^2 + J(u, v) \rightarrow \min_{u, v} & \\ u \geq 0 & \text{a.e. on } I^{+0}(\bar{u}, \bar{v}) \\ u = 0 & \text{a.e. on } I^{0+}(\bar{u}, \bar{v}) \cup I^{00}(\bar{u}, \bar{v}) \\ v \geq 0 & \text{a.e. on } I^{0+}(\bar{u}, \bar{v}) \cup I^{00}(\bar{u}, \bar{v}) \\ v = 0 & \text{a.e. on } I^{+0}(\bar{u}, \bar{v}) \end{array} \right.$$

and

$$(rNLP_{\bar{v}}) \quad \left\{ \begin{array}{ll} \frac{1}{2} \|S[u, v] - y_d\|_{\mathcal{D}}^2 + J(u, v) \rightarrow \min_{u, v} & \\ u \geq 0 & \text{a.e. on } I^{+0}(\bar{u}, \bar{v}) \cup I^{00}(\bar{u}, \bar{v}) \\ u = 0 & \text{a.e. on } I^{0+}(\bar{u}, \bar{v}) \\ v \geq 0 & \text{a.e. on } I^{0+}(\bar{u}, \bar{v}) \\ v = 0 & \text{a.e. on } I^{+0}(\bar{u}, \bar{v}) \cup I^{00}(\bar{u}, \bar{v}) \end{array} \right.$$

since their respective feasible sets are smaller than  $\mathbb{C}$  but contain  $(\bar{u}, \bar{v})$ . By standard notion, see [26, 30, 34],  $(rNLP_{\bar{u}})$  and  $(rNLP_{\bar{v}})$  are referred to as *restricted nonlinear problems*. Furthermore, the corresponding *relaxed nonlinear problem* is introduced by means of

$$(RNLP) \quad \left\{ \begin{array}{ll} \frac{1}{2} \|S[u, v] - y_d\|_{\mathcal{D}}^2 + J(u, v) \rightarrow \min_{u, v} & \\ u \geq 0 & \text{a.e. on } I^{+0}(\bar{u}, \bar{v}) \cup I^{00}(\bar{u}, \bar{v}) \\ u = 0 & \text{a.e. on } I^{0+}(\bar{u}, \bar{v}) \\ v \geq 0 & \text{a.e. on } I^{0+}(\bar{u}, \bar{v}) \cup I^{00}(\bar{u}, \bar{v}) \\ v = 0 & \text{a.e. on } I^{+0}(\bar{u}, \bar{v}). \end{array} \right.$$

Observe that the feasible points  $(u, v) \in H^1(\Omega)^2$  of (RNLP) do not necessarily satisfy the complementarity condition  $(u, v) \in \mathbb{C}$ . Combining standard techniques from finite-dimensional MPCC theory and optimization in Banach spaces, the following result is obtained, see also [34, Theorems 3.1 and 5.2]. It should be noted that due to the appearance of the two control



variables  $u$  and  $v$  in (3.1), there will be two Lagrange multipliers  $\mu$  and  $\nu$  corresponding to  $u$  and  $v$ , respectively, in the stationarity system as well. In particular, the pair  $(\mu, \nu) \in H^1(\Omega)^\star \times H^1(\Omega)^\star$  may be identified with a functional from  $(H^1(\Omega)^2)^\star$ .

**Theorem 3.1.** *Let  $(\bar{u}, \bar{v}) \in H^1(\Omega)^2$  be a locally optimal solution of (3.1). Then, there exist multipliers  $\mu, \nu \in H^1(\Omega)^\star$  satisfying*

$$(3.6a) \quad 0 = S^\star [S[\bar{u}, \bar{v}] - y_d] + J'(\bar{u}, \bar{v}) + (\mu, \nu),$$

$$(3.6b) \quad \mu \in \left\{ z \in H^1(\Omega) \left| \begin{array}{ll} z \geq 0 & \text{a.e. on } I^{+0}(\bar{u}, \bar{v}) \cup I^{00}(\bar{u}, \bar{v}) \\ z = 0 & \text{a.e. on } I^{0+}(\bar{u}, \bar{v}) \end{array} \right. \right\}^\circ,$$

$$(3.6c) \quad \langle \mu, \bar{u} \rangle_{H^1(\Omega)} = 0,$$

$$(3.6d) \quad \nu \in \left\{ z \in H^1(\Omega) \left| \begin{array}{ll} z \geq 0 & \text{a.e. on } I^{0+}(\bar{u}, \bar{v}) \cup I^{00}(\bar{u}, \bar{v}) \\ z = 0 & \text{a.e. on } I^{+0}(\bar{u}, \bar{v}) \end{array} \right. \right\}^\circ,$$

$$(3.6e) \quad \langle \nu, \bar{v} \rangle_{H^1(\Omega)} = 0.$$

*Proof.* Introducing the cones

$$\begin{aligned} \mathcal{K}_{+0} &:= \left\{ z \in H^1(\Omega) \left| \begin{array}{ll} z \geq 0 & \text{a.e. on } I^{+0}(\bar{u}, \bar{v}) \\ z = 0 & \text{a.e. on } I^{0+}(\bar{u}, \bar{v}) \cup I^{00}(\bar{u}, \bar{v}) \end{array} \right. \right\}, \\ \mathcal{K}_{0+,00} &:= \left\{ z \in H^1(\Omega) \left| \begin{array}{ll} z \geq 0 & \text{a.e. on } I^{0+}(\bar{u}, \bar{v}) \cup I^{00}(\bar{u}, \bar{v}) \\ z = 0 & \text{a.e. on } I^{+0}(\bar{u}, \bar{v}) \end{array} \right. \right\}, \end{aligned}$$

(rNLP $_{\bar{u}}$ ) is equivalent to

$$\begin{cases} \frac{1}{2} \|S[u, v] - y_d\|_{\mathcal{D}}^2 + J(u, v) \rightarrow \min_{u, v} \\ u \in \mathcal{K}_{+0} \\ v \in \mathcal{K}_{0+,00}. \end{cases}$$

Since  $(\bar{u}, \bar{v})$  is a locally optimal solution of (rNLP $_{\bar{u}}$ ), there exist multipliers  $\mu^1, \nu^1 \in H^1(\Omega)^\star$  which satisfy the corresponding KKT conditions

$$(3.7) \quad \begin{cases} 0 = S^\star [S[\bar{u}, \bar{v}] - y_d] + J'(\bar{u}, \bar{v}) + (\mu^1, \nu^1), \\ \mu^1 \in \mathcal{K}_{+0}^\circ \cap \bar{u}^\perp, \\ \nu^1 \in \mathcal{K}_{0+,00}^\circ \cap \bar{v}^\perp, \end{cases}$$

see [4, Theorem 3.9]. Considering (rNLP $_{\bar{v}}$ ) in a similar way, there exist  $\mu^2, \nu^2 \in H^1(\Omega)^\star$  which satisfy

$$(3.8) \quad \begin{cases} 0 = S^\star [S[\bar{u}, \bar{v}] - y_d] + J'(\bar{u}, \bar{v}) + (\mu^2, \nu^2), \\ \mu^2 \in \mathcal{K}_{+0,00}^\circ \cap \bar{u}^\perp, \\ \nu^2 \in \mathcal{K}_{0+}^\circ \cap \bar{v}^\perp, \end{cases}$$

where

$$\begin{aligned}\mathcal{K}_{+0,00} &:= \left\{ z \in H^1(\Omega) \left| \begin{array}{ll} z \geq 0 & \text{a.e. on } I^{+0}(\bar{u}, \bar{v}) \cup I^{00}(\bar{u}, \bar{v}) \\ z = 0 & \text{a.e. on } I^{0+}(\bar{u}, \bar{v}) \end{array} \right. \right\}, \\ \mathcal{K}_{0+} &:= \left\{ z \in H^1(\Omega) \left| \begin{array}{ll} z \geq 0 & \text{a.e. on } I^{0+}(\bar{u}, \bar{v}) \\ z = 0 & \text{a.e. on } I^{+0}(\bar{u}, \bar{v}) \cup I^{00}(\bar{u}, \bar{v}) \end{array} \right. \right\}.\end{aligned}$$

Combining the respective first condition in (3.7) and (3.8) yields  $\mu^1 = \mu^2$  and  $\nu^1 = \nu^2$ . Since  $\mathcal{K}_{+0,00}^\circ \cap \bar{u}^\perp$  is a subset of  $\mathcal{K}_{+0}^\circ \cap \bar{u}^\perp$  while  $\mathcal{K}_{0+,00}^\circ \cap \bar{v}^\perp$  is a subset of  $\mathcal{K}_{0+}^\circ \cap \bar{v}^\perp$ , the desired result is obtained by setting  $\mu := \mu^2$  and  $\nu := \nu^1$ .  $\square$

Note that the system (3.6) coincides with the KKT conditions of (RNLP). In this regard, it is reasonable to call the conditions (3.6) a *strong stationarity-type system*.

**Remark 3.2.** It is difficult to give an explicit characterization of the multipliers  $\mu, \nu \in H^1(\Omega)^\star$ . Assume that  $\Omega$  has a Lipschitz boundary. Introducing  $\mathcal{H}_{\mathcal{A}} := \{z \in H^1(\Omega) \mid z = 0 \text{ a.e. on } \mathcal{A}\}$  for a fixed measurable set  $\mathcal{A} \subset \Omega$  and using the relation  $H_+^1(\Omega)^\circ = H^1(\Omega)^\star \cap \mathcal{M}_-(\bar{\Omega})$ , see [12, Lemma 3.1], it holds that

$$\mu \in (H_+^1(\Omega) \cap \mathcal{H}_{I^{0+}(\bar{u}, \bar{v})})^\circ = \text{cl} \left( H^1(\Omega)^\star \cap \mathcal{M}_-(\bar{\Omega}) + \mathcal{H}_{I^{0+}(\bar{u}, \bar{v})}^\perp \right)$$

where  $\mathcal{M}_-(\bar{\Omega})$  denotes the set of all finite, nonpositive Borel measures on  $\bar{\Omega}$ . A similar result can be obtained to characterize  $\nu$ . However, due to the appearance of the closure as well as the annihilated subspace associated with  $\mathcal{H}_{I^{0+}(\bar{u}, \bar{v})}$ , this characterization is of limited practical use; in particular, it cannot be deduced that  $\mu$  and  $\nu$  are measures. Applying the machinery of capacity theory, see [3, 4], a more advanced approach to the characterization of  $\mu$  and  $\nu$  can be attempted. For this purpose, one could strengthen the constraints in (rNLP $_{\bar{u}}$ ), (rNLP $_{\bar{v}}$ ), and (RNLP) to hold *quasi-everywhere* on the respective subdomains, i.e., the respective conditions hold up to sets of  $H^1$ -capacity zero. Then, one needs to find explicit expressions for the polar cone associated with sets of type

$$\left\{ z \in H^1(\Omega) \left| \begin{array}{ll} z \geq 0 & \text{quasi-everywhere on } \mathcal{A} \\ z = 0 & \text{quasi-everywhere on } \Omega \setminus \mathcal{A} \end{array} \right. \right\}$$

where  $\mathcal{A} \subset \Omega$  is measurable. The price one has to pay when using this approach is a less restrictive stationarity system than (3.6). In particular, the polar cones from (3.6b) and (3.6d) would be replaced by larger ones.

In order to state necessary optimality conditions of strong stationarity-type that avoid the appearance of multipliers and allow a numerical implementation, one can exploit the definition of the polar cone in the system (3.6).

**Corollary 3.3.** *Let  $(\bar{u}, \bar{v}) \in H^1(\Omega)^2$  be a locally optimal solution of (3.1). Then, the condition*

$$0 = \langle S[\bar{u}, \bar{v}] - y_d, S[\bar{u}, \bar{v}] \rangle_{\mathcal{D}} + J'(\bar{u}, \bar{v})[\bar{u}, \bar{v}]$$

holds, and for any pair  $(z_u, z_v) \in H_+^1(\Omega) \times H_+^1(\Omega)$ ,

$$\left. \begin{aligned} \text{supp } z_u &\subset I^{+0}(\bar{u}, \bar{v}) \cup I^{00}(\bar{u}, \bar{v}) \\ \text{supp } z_v &\subset I^{0+}(\bar{u}, \bar{v}) \cup I^{00}(\bar{u}, \bar{v}) \end{aligned} \right\} \implies \langle S^*[S[\bar{u}, \bar{v}] - y_d] + J'(\bar{u}, \bar{v}), (z_u, z_v) \rangle_{H^1(\Omega)^2} \geq 0.$$

*Proof.* Due to [Theorem 3.1](#), there exist  $\mu, \nu \in H^1(\Omega)^*$  satisfying (3.6). Testing (3.6a) with  $(\bar{u}, \bar{v})$  while exploiting (3.6c), (3.6e), and the definition of the adjoint operator, the first statement of the corollary follows.

The second statement is a consequence of (3.6a), (3.6b), and (3.6d).  $\square$

**Remark 3.4.** According to standard terminology for MPCCs, the necessary optimality conditions (3.6) are of strong stationarity-type, see, e.g., [34, Definition 5.1] and [35, Definition 2.7]. Recall that a feasible point  $(\bar{u}, \bar{v}) \in H^1(\Omega)^2$  of (3.1) and thus of (3.2) is a strongly stationary point of (3.2) in the sense of [34, Definition 5.1] if and only if there are multipliers  $(\mu, \nu) \in L^2(\Omega)^2$  satisfying

$$(3.9a) \quad 0 = S^*[S[\bar{u}, \bar{v}] - y_d] + J'(\bar{u}, \bar{v}) + (E, E)^*[\mu, \nu],$$

$$(3.9b) \quad \mu = 0 \quad \text{a.e. on } I^{+0}(\bar{u}, \bar{v}),$$

$$(3.9c) \quad \nu = 0 \quad \text{a.e. on } I^{0+}(\bar{u}, \bar{v}),$$

$$(3.9d) \quad \mu \leq 0 \wedge \nu \leq 0 \quad \text{a.e. on } I^{00}(\bar{u}, \bar{v}),$$

see also [25, Definition 4.1]. If  $\mathbb{C}$  is replaced by  $\tilde{\mathbb{C}}$  and  $\varepsilon = 0$  is taken in the definition of  $J$  (in which case  $E$  is the identity mapping), the systems (3.6) and (3.9) are equivalent. However, for  $\mathbb{C}$  and  $\varepsilon > 0$ , the necessary optimality conditions (3.6) are weaker than (3.9), which can be seen as follows: It is clear that whenever  $(\tilde{\mu}, \tilde{\nu}) \in L^2(\Omega)^2$  satisfy the classical strong stationarity conditions (3.9), then the multipliers  $\mu := E^*[\tilde{\mu}]$  and  $\nu := E^*[\tilde{\nu}]$  satisfy (3.6). On the other hand, by means of [Theorem 3.1](#), the multipliers appearing in the system (3.6) may come from  $H^1(\Omega)^* \setminus L^2(\Omega)$  in general.

**Remark 3.5.** In this section, only the property of  $S$  to be a bounded, linear operator has been exploited. Thus, the optimality conditions obtained in [Theorem 3.1](#) and [Corollary 3.3](#) are applicable in many different situations, e.g., in case where  $S$  is the control-to-observation operator associated with a linear elliptic equation where  $u$  and  $v$  only operate on some subdomain, or for a linear parabolic equation where the controls  $u$  and  $v$  only depend on time. The latter problems are closely related to the switching-constrained problems examined in [7–9].

It should be noted that similar necessary optimality conditions can be derived if  $S: H^1(\Omega)^2 \rightarrow \mathcal{D}$  is Fréchet differentiable but not necessarily linear.

## 4 PENALIZATION OF COMPLEMENTARITY CONSTRAINTS

In order to find optimal solutions of (OC<sup>4</sup>), an obvious idea would be to penalize the violation of the equilibrium condition

$$(4.1) \quad u(x)v(x) = 0 \quad \text{a.e. on } \Omega$$

in (OC<sup>4</sup>). This is related to the approaches used in [7–9] for the treatment of switching-constrained optimal control problems. However, the resulting penalized problem would still involve inequality constraints for the controls in  $H^1(\Omega)$ , and thus the associated KKT conditions would involve Lagrange multipliers from  $H^1(\Omega)^* \cap \mathcal{M}_-(\overline{\Omega})$ , see [12, Section 5] for details. This, however, may provoke theoretical and numerical difficulties that should be avoided here.

To get around these issues, the penalization of the overall complementarity constraint using the Fischer–Burmeister function is proposed here, which leads to penalized problems in which the only constraint is the state equation.

#### 4.1 PENALTY TERM

Let  $\phi : \mathbb{R}^2 \rightarrow \mathbb{R}$  denote the Fischer–Burmeister function introduced in (1.1) and let the mapping  $\Phi : L^2(\Omega)^2 \rightarrow L^2(\Omega)$  be the associated Nemytskii operator defined by

$$\forall (w, z) \in L^2(\Omega)^2 \quad \forall x \in \Omega: \quad \Phi(w, z)(x) := \phi(w(x), z(x)).$$

This operator is well-defined since for all  $w, z \in L^2(\Omega)$ , one has

$$\begin{aligned} \|\Phi(w, z)\|_{L^2(\Omega)} &\leq \left( \int_{\Omega} (w^2(x) + z^2(x)) dx \right)^{1/2} + \|w\|_{L^2(\Omega)} + \|z\|_{L^2(\Omega)} \\ &\leq \left( \int_{\Omega} (|w(x)| + |z(x)|)^2 dx \right)^{1/2} + \|w\|_{L^2(\Omega)} + \|z\|_{L^2(\Omega)} \\ &\leq 2 \left( \|w\|_{L^2(\Omega)} + \|z\|_{L^2(\Omega)} \right) < +\infty, \end{aligned}$$

i.e.,  $\Phi$  maps from  $L^2(\Omega)^2$  to  $L^2(\Omega)$ , see also [33, Section 3.3]. For a detailed introduction to the theory of superposition operators in Lebesgue spaces, the interested reader is referred to [2, 17].

The violation of the complementarity constraint  $(u, v) \in \mathbb{C}$  can then be penalized using the functional  $F : H^1(\Omega)^2 \rightarrow \mathbb{R}_0^+$  defined by

$$(4.2) \quad \forall (u, v) \in H^1(\Omega)^2: \quad F(u, v) := \frac{1}{2} \int_{\Omega} \phi^2(u(x), v(x)) dx = \frac{1}{2} \|\Phi(E[u], E[v])\|_{L^2(\Omega)}^2.$$

Recall that  $E \in \mathbb{L} [H^1(\Omega), L^2(\Omega)]$  represents the natural embedding  $H^1(\Omega) \hookrightarrow L^2(\Omega)$ .

It is obvious that  $\Phi$  cannot be Fréchet differentiable since  $\phi$  is not smooth. In contrast,  $F$  is a continuously Fréchet differentiable mapping.

**Lemma 4.1.** *Let  $(\bar{u}, \bar{v}) \in H^1(\Omega)^2$  be arbitrarily chosen. Then,  $F$  is continuously Fréchet differentiable at  $(\bar{u}, \bar{v})$ . The associated Fréchet derivative is given by*

$$\forall (\delta^u, \delta^v) \in H^1(\Omega)^2: \quad F'(\bar{u}, \bar{v})[\delta^u, \delta^v] = \int_{\Omega} \phi(\bar{u}(x), \bar{v}(x)) (\eta_{\bar{u}}(x) \delta^u(x) + \eta_{\bar{v}}(x) \delta^v(x)) dx,$$

where  $\eta_{\bar{u}}, \eta_{\bar{v}} \in L^\infty(\Omega)$  are defined by

$$(4.3a) \quad \forall x \in \Omega: \quad \eta_{\bar{u}}(x) = \begin{cases} \frac{\bar{u}(x)}{\sqrt{\bar{u}(x)^2 + \bar{v}(x)^2}} - 1 & \text{if } x \notin I^{00}(\bar{u}, \bar{v}), \\ 0 & \text{if } x \in I^{00}(\bar{u}, \bar{v}), \end{cases}$$

$$(4.3b) \quad \forall x \in \Omega: \quad \eta_{\bar{v}}(x) = \begin{cases} \frac{\bar{v}(x)}{\sqrt{\bar{u}(x)^2 + \bar{v}(x)^2}} - 1 & \text{if } x \notin I^{00}(\bar{u}, \bar{v}), \\ 0 & \text{if } x \in I^{00}(\bar{u}, \bar{v}), \end{cases}$$

and  $I^{00}(\bar{u}, \bar{v})$  is defined by (3.5).

*Proof.* Let  $f: \mathbb{R}^2 \rightarrow \mathbb{R}$  be given by

$$\forall (a, b) \in \mathbb{R}^2: \quad f(a, b) := \frac{1}{2}\phi(a, b)^2.$$

One can check that  $f$  is continuously differentiable with gradient

$$\forall (a, b) \in \mathbb{R}^2: \quad \nabla f(a, b) = \begin{cases} \phi(a, b) \begin{pmatrix} \frac{a}{\sqrt{a^2 + b^2}} - 1 \\ \frac{b}{\sqrt{a^2 + b^2}} - 1 \end{pmatrix} & \text{if } (a, b) \neq (0, 0), \\ \begin{pmatrix} 0 \\ 0 \end{pmatrix} & \text{if } (a, b) = (0, 0). \end{cases}$$

Clearly, the Nemytskii-operator  $\mathcal{F}$  associated with  $f$  maps from  $L^2(\Omega)^2$  to  $L^1(\Omega)$ , since  $\Phi$  maps  $L^2(\Omega)^2$  to  $L^2(\Omega)$ . Noting that  $a/\sqrt{a^2 + b^2} \in [-1, 1]$  and  $b/\sqrt{a^2 + b^2} \in [-1, 1]$  hold for all  $(a, b) \in \mathbb{R}^2 \setminus \{(0, 0)\}$ , the Nemytskii operator associated with  $\nabla f$  maps from  $L^2(\Omega)^2$  to  $L^2(\Omega)^2$ . Applying [17, Theorems 4 and 7],  $\mathcal{F}: L^2(\Omega)^2 \rightarrow L^1(\Omega)$  is continuously Fréchet differentiable. Furthermore,

$$\forall x \in \Omega: \quad \mathcal{F}'(w, z)[\delta_w, \delta_z](x) = \nabla_a f(w(x), z(x))\delta_w(x) + \nabla_b f(w(x), z(x))\delta_z(x)$$

for any  $(w, z), (\delta_w, \delta_z) \in L^2(\Omega)^2$ .

Define  $L \in \mathbb{L}[L^1(\Omega), \mathbb{R}]$  by  $L[w] := \int_{\Omega} w(x)dx$ . Then,  $F = L \circ \mathcal{F} \circ (E, E)$ . Since all involved mappings are continuously Fréchet differentiable, the assertion of the lemma follows by exploiting the chain rule for Fréchet differentiable functions, see [32, Theorem 2.20].  $\square$

**Remark 4.2.** As the penalty functional  $F$  is smooth, it cannot lead to exact penalization of the complementarity constraints, see, e.g., [16, Theorem 5.9]. Although Section 6 demonstrates that a penalty method using  $F$  behaves well in numerical practice, in principle any other NCP-function, see [31] for an overview, can be used to construct similar penalty methods.

One possible alternative would be to use  $F_1: H^1(\Omega)^2 \rightarrow \mathbb{R}_0^+$  given by

$$\forall (u, v) \in H^1(\Omega)^2: \quad F_1(u, v) := \int_{\Omega} |\phi(u(x), v(x))|dx = \|\tilde{\Phi}(E[u], E[v])\|_{L^1(\Omega)},$$

where  $\tilde{\Phi}: L^2(\Omega)^2 \rightarrow L^1(\Omega)$  is the mapping  $E_{L^2 \rightarrow L^1} \circ \Phi$  where  $E_{L^2 \rightarrow L^1}$  represents the continuous embedding  $L^2(\Omega) \hookrightarrow L^1(\Omega)$ . This leads to a nonsmooth but Lipschitz continuous mapping.

Another approach would be to exploit the so-called *smoothed* Fischer–Burmeister function  $\phi_\theta: \mathbb{R}^2 \rightarrow \mathbb{R}$  given by

$$\forall (a, b) \in \mathbb{R}^2: \quad \phi_\theta(a, b) := \sqrt{a^2 + b^2 + 2\theta} - a - b,$$

which is continuously differentiable for any  $\theta > 0$ , see [22]. Using [17, Theorems 4 and 7], one can check that the associated Nemytskii operator  $\tilde{\Phi}_\theta: L^2(\Omega)^2 \rightarrow L^1(\Omega)$  is continuously Fréchet differentiable. Define  $F_{1,\theta}: H^1(\Omega)^2 \rightarrow \mathbb{R}_0^+$  by means of

$$\forall (u, v) \in H^1(\Omega)^2: \quad F_{1,\theta}(u, v) := \int_{\Omega} |\phi_\theta(u(x), v(x))| dx = \|\tilde{\Phi}_\theta(E[u], E[v])\|_{L^1(\Omega)}.$$

Clearly,  $F_{1,0}$  corresponds to  $F_1$ . For  $\theta > 0$  this approach can be seen as a mixture of a penalty and a smoothing method. However, it needs to be noted that  $F_{1,\theta}$  is nonsmooth even for positive values of  $\theta$ .

#### 4.2 EXISTENCE, CONVERGENCE RESULTS, AND OPTIMALITY CONDITIONS

Using the penalty functional  $F$  defined in (4.2) to penalize the complementarity constraints in (3.1) leads to the family of penalized problems

$$(P_k) \quad \frac{1}{2} \|S[u, v] - y_d\|_{\mathcal{D}}^2 + J(u, v) + \sigma_k F(u, v) \rightarrow \min_{u, v},$$

where  $\{\sigma_k\}_{k \in \mathbb{N}} \subset \mathbb{R}^+$  is a sequence of positive real numbers tending to infinity as  $k \rightarrow \infty$ . The first question is about the existence of solutions of  $(P_k)$ .

**Proposition 4.3.** *For any  $\sigma_k > 0$ , the penalized problem  $(P_k)$  possesses an optimal solution.*

*Proof.* Let  $\{(u_l, v_l)\}_{l \in \mathbb{N}} \subset H^1(\Omega)^2$  be a minimizing sequence for  $(P_k)$  and let  $\bar{m} \in \overline{\mathbb{R}}$  be the corresponding infimal value. Since  $J$  is, due to  $\varepsilon > 0$ , coercive and bounded from below, this sequence is bounded in  $H^1(\Omega)^2$  and, thus, possesses a weakly convergent subsequence (without relabeling) with weak limit  $(\bar{u}, \bar{v}) \in H^1(\Omega)^2$ . Due to the compactness of  $H^1(\Omega) \hookrightarrow L^2(\Omega)$ , the strong convergences  $u_l \rightarrow \bar{u}$  and  $v_l \rightarrow \bar{v}$  hold in  $L^2(\Omega)$ . Noting that the operator  $\Phi$  is continuous on  $L^2(\Omega)^2$ , see [17, Theorem 4], it follows that

$$\lim_{l \rightarrow \infty} F(u_l, v_l) = F(\bar{u}, \bar{v}).$$

Thus, the continuity of  $S$  and the weak lower semicontinuity of norms imply that

$$\begin{aligned} \frac{1}{2} \|S[\bar{u}, \bar{v}] - y_d\|_{\mathcal{D}}^2 + J(\bar{u}, \bar{v}) + \sigma_k F(\bar{u}, \bar{v}) \\ \leq \liminf_{l \rightarrow \infty} \left( \frac{1}{2} \|S[u_l, v_l] - y_d\|_{\mathcal{D}}^2 + J(u_l, v_l) \right) + \sigma_k \lim_{l \rightarrow \infty} F(u_l, v_l) \\ = \liminf_{l \rightarrow \infty} \left( \frac{1}{2} \|S[u_l, v_l] - y_d\|_{\mathcal{D}}^2 + J(u_l, v_l) + \sigma_k F(u_l, v_l) \right) = \bar{m}, \end{aligned}$$

i.e.,  $(\bar{u}, \bar{v})$  is a global minimizer of  $(P_k)$ .  $\square$

Next, the convergence of solutions of  $(P_k)$  as  $\sigma_k \rightarrow \infty$  is addressed.

**Proposition 4.4.** *Fix a sequence  $\{\sigma_k\}_{k \in \mathbb{N}} \subset \mathbb{R}^+$  tending to infinity as  $k \rightarrow \infty$ . For any  $k \in \mathbb{N}$ , let  $(u_k, v_k) \in H^1(\Omega)^2$  be a global minimizer of  $(P_k)$ . Then,  $\{(u_k, v_k)\}_{k \in \mathbb{N}}$  contains a subsequence converging strongly in  $H^1(\Omega)^2$  to a point  $(\bar{u}, \bar{v}) \in \mathbb{C}$  such that  $(\bar{y}, \bar{u}, \bar{v})$ , where  $\bar{y} \in \mathcal{Y}$  is the state associated with  $(\bar{u}, \bar{v})$ , is an optimal solution of  $(OC^4)$ .*

*Moreover, any subsequence of  $\{(u_k, v_k)\}_{k \in \mathbb{N}}$  converging weakly to some  $(\bar{u}, \bar{v})$  in  $H^1(\Omega)^2$  produces a global minimizer of  $(OC^4)$  in the above sense.*

*Proof.* For any  $k \in \mathbb{N}$ , the estimate

$$\frac{1}{2} \|S[u_k, v_k] - y_d\|_{\mathcal{D}}^2 + J(u_k, v_k) + \sigma_k F(u_k, v_k) \leq \frac{1}{2} \|y_d\|_{\mathcal{D}}^2$$

follows from the feasibility of  $(0, 0) \in H^1(\Omega)^2$  for  $(P_k)$ . Thus, since  $J$  is coercive and bounded from below while  $F$  only takes nonnegative values,  $\{(u_k, v_k)\}_{k \in \mathbb{N}}$  is bounded and therefore contains a weakly convergent subsequence (which, as all further subsequences, will not be relabeled). Recalling the compactness of  $H^1(\Omega) \hookrightarrow L^2(\Omega)$ , the sequence  $\{(u_k, v_k)\}_{k \in \mathbb{N}}$  converges strongly to  $(\bar{u}, \bar{v})$  in  $L^2(\Omega)^2$  and thus pointwise almost everywhere at least along a subsequence. Furthermore, the relation

$$0 \leq \|\Phi(E[u_k], E[v_k])\|_{L^2(\Omega)} \leq \sqrt{\frac{1}{\sigma_k}} \|y_d\|_{\mathcal{D}} \rightarrow 0$$

is obtained as  $k \rightarrow \infty$ . Consequently, at least along a subsequence,  $\{\Phi(E[u_k], E[v_k])\}_{k \in \mathbb{N}}$  converges pointwise a.e. to 0. By definition of  $\Phi$ ,  $(\bar{u}, \bar{v}) \in \mathbb{C}$  follows.

Now choose  $(u, v) \in \mathbb{C}$  arbitrarily. Since this point is feasible to  $(P_k)$ , it follows for any  $k \in \mathbb{N}$  that

$$\begin{aligned} \frac{1}{2} \|S[u, v] - y_d\|_{\mathcal{D}}^2 + J(u, v) &\geq \frac{1}{2} \|S[u_k, v_k] - y_d\|_{\mathcal{D}}^2 + J(u_k, v_k) + \sigma_k F(u_k, v_k) \\ &\geq \frac{1}{2} \|S[u_k, v_k] - y_d\|_{\mathcal{D}}^2 + J(u_k, v_k). \end{aligned}$$

Thus, using the weak lower semicontinuity of the functionals, one obtains

$$\begin{aligned} \frac{1}{2} \|S[\bar{u}, \bar{v}] - y_d\|_{\mathcal{D}}^2 + J(\bar{u}, \bar{v}) &\leq \liminf_{k \rightarrow \infty} \left( \frac{1}{2} \|S[u_k, v_k] - y_d\|_{\mathcal{D}}^2 + J(u_k, v_k) \right) \\ &\leq \limsup_{k \rightarrow \infty} \left( \frac{1}{2} \|S[u_k, v_k] - y_d\|_{\mathcal{D}}^2 + J(u_k, v_k) \right) \\ &\leq \limsup_{k \rightarrow \infty} \left( \frac{1}{2} \|S[u_k, v_k] - y_d\|_{\mathcal{D}}^2 + J(u_k, v_k) + \sigma_k F(u_k, v_k) \right) \\ &\leq \frac{1}{2} \|S[u, v] - y_d\|_{\mathcal{D}}^2 + J(u, v) \end{aligned}$$

for all  $(u, v) \in \mathbb{C}$ . Consequently,  $(\bar{u}, \bar{v})$  is a global minimizer of the state-reduced problem (3.1). Choosing  $u := \bar{u}$  and  $v := \bar{v}$  in the above estimate, one obtains

$$\frac{1}{2} \|S[u_k, v_k] - y_d\|_{\mathcal{D}}^2 + J(u_k, v_k) \rightarrow \frac{1}{2} \|S[\bar{u}, \bar{v}] - y_d\|_{\mathcal{D}}^2 + J(\bar{u}, \bar{v}),$$

and  $J(u_k, v_k) \rightarrow J(\bar{u}, \bar{v})$  follows by Lemma A.1. Since  $u_k \rightarrow \bar{u}$  and  $v_k \rightarrow \bar{v}$  in  $L^2(\Omega)$ , the definition of  $J$  and  $\varepsilon > 0$  imply that

$$\|u_k\|_{H^1(\Omega)}^2 + \|v_k\|_{H^1(\Omega)}^2 \rightarrow \|\bar{u}\|_{H^1(\Omega)}^2 + \|\bar{v}\|_{H^1(\Omega)}^2.$$

Now, applying [Lemma A.1](#) once more yields

$$\|u_k\|_{H^1(\Omega)}^2 \rightarrow \|\bar{u}\|_{H^1(\Omega)}^2, \quad \|v_k\|_{H^1(\Omega)}^2 \rightarrow \|\bar{v}\|_{H^1(\Omega)}^2.$$

Combining this with the weak convergences  $u_k \rightharpoonup \bar{u}$  and  $v_k \rightharpoonup \bar{v}$  in  $H^1(\Omega)$ , the convergences  $u_k \rightarrow \bar{u}$  and  $v_k \rightarrow \bar{v}$  in  $H^1(\Omega)$  follow since the latter is a Hilbert space. This yields the first assertion.

If  $\{(u_k, v_k)\}_{k \in \mathbb{N}}$  contains a subsequence converging weakly to some  $(\bar{u}, \bar{v}) \in H^1(\Omega)^2$  in  $H^1(\Omega)^2$ , then the above arguments can be partially repeated to show that  $(\bar{u}, \bar{v})$  is a global minimizer of (3.1). This completes the proof.  $\square$

An obvious advantage of  $(P_k)$  is that it is a smooth and unconstrained problem, allowing the straightforward derivation of necessary optimality conditions. Hence, the following result is a direct consequence of Fermat's rule and [Lemma 4.1](#).

**Proposition 4.5.** *For fixed  $\sigma_k > 0$ , let  $(u_k, v_k) \in H^1(\Omega)^2$  be a locally optimal solution of  $(P_k)$ . Then, the corresponding functions  $\eta_{u_k}, \eta_{v_k} \in L^\infty(\Omega)$  defined as in (4.3) satisfy*

$$0 = S^* [S[u_k, v_k] - y_d] + J'(u_k, v_k) \sigma_k (E, E)^* [\Phi(E[u_k], E[v_k]) \eta_{u_k}, \Phi(E[u_k], E[v_k]) \eta_{v_k}].$$

**Remark 4.6.** Similar results as in this section can be shown for the penalty terms induced by the nonsmooth functionals  $F_1$  and  $F_{1, \theta_k}$  given in [Remark 4.2](#) using the continuity of the associated Nemytskii operators  $\tilde{\Phi}$  and  $\tilde{\Phi}_{\theta_k}$  as well as calculus rules for Clarke's generalized derivative, see [5]. Obtaining a convergence result as in [Proposition 4.4](#) for  $F_{1, \theta_k}$  additionally requires to choose  $\sigma_k$  and  $\theta_k$  such that  $\sigma_k \sqrt{\theta_k} \rightarrow 0$  as  $k \rightarrow \infty$ .

**Remark 4.7.** Using the boundedness of the solutions and passing to subsequences, it is possible by pointwise inspection to take the limit  $k \rightarrow \infty$  in the optimality system from [Proposition 4.5](#) and derive the existence of multipliers  $\mu, \nu \in H^1(\Omega)^*$  which satisfy the polarity relations from [Theorem 3.1](#) with respect to the index sets  $I^{+0}(\bar{u}, \bar{v})$  and  $I^{0+}(\bar{u}, \bar{v})$ . This can be seen as a natural extension of the so-called weak stationarity concept, see [25, Definition 4.1], to (3.1). However, it does not seem to be possible to infer the polarity relations for  $\mu$  and  $\nu$  on  $I^{00}(\bar{u}, \bar{v})$  found in the strong stationarity system from [Theorem 3.1](#). Noting that our penalty approach is related to Scholtes' relaxation technique for the numerical solution of finite-dimensional MPCCs which yields so-called Clarke-stationary points in general, see [21, Section 3.1] for details, this observation does not seem to be too surprising since Clarke-stationarity is much weaker than strong stationarity.

## 5 NUMERICAL TREATMENT

This section deals with the numerical implementation of the penalization technique described in [Section 4](#) following a “first-discretize-then-optimize approach” based on a finite element discretization. In order to concentrate on the complementarity constraint, the state equation is chosen as the elliptic model problem

$$(PDE) \quad \begin{cases} -\nabla \cdot (C \nabla y) + ay = bu + cv & \text{a.e. on } \Omega \\ \vec{n} \cdot (C \nabla y) = 0 & \text{a.e. on } \text{bd } \Omega. \end{cases}$$



Here,  $\Omega \subset \mathbb{R}^d$  is a domain with Lipschitz boundary  $\text{bd}(\Omega)$ ,  $C \in L^\infty(\Omega; S^d(\mathbb{R}))$  satisfies the condition of uniform ellipticity (2.3), and the functions  $\mathbf{a}, \mathbf{b}, \mathbf{c} \in L^\infty(\Omega)$  do not vanish while  $\mathbf{a}$  is additionally nonnegative, see also Section 2.2. Set  $\mathcal{D} := L^2(\Omega)$  and  $\mathcal{Y} := H^1(\Omega)$ . The operator  $D := E$  represents the natural embedding  $H^1(\Omega) \hookrightarrow L^2(\Omega)$ . Note that the weak formulation of the associated state equation can be written in the abstract form  $A[y] - B[u] - C[v] = 0$ , where the bounded, linear operators  $A, B, C \in \mathbb{L}[H^1(\Omega), H^1(\Omega)^*]$  are given for all  $y, u, v, w \in H^1(\Omega)$  as

$$\begin{aligned}\langle A[y], w \rangle_{H^1(\Omega)} &:= \int_{\Omega} (C(x) \nabla y(x)) \cdot \nabla w(x) dx + \int_{\Omega} \mathbf{a}(x) y(x) w(x) dx, \\ \langle B[u], w \rangle_{H^1(\Omega)} &:= \int_{\Omega} \mathbf{b}(x) u(x) w(x) dx, \\ \langle C[v], w \rangle_{H^1(\Omega)} &:= \int_{\Omega} \mathbf{c}(x) v(x) w(x) dx.\end{aligned}$$

It can be checked that the operator  $A$  is elliptic and self-adjoint under the postulated assumptions, see, e.g., [13, Section 6]. The operators  $B$  and  $C$  are self-adjoint as well.

### 5.1 FINITE ELEMENT DISCRETIZATION

While the discretization of  $(P_k)$  is rather standard, some notation needs to be introduced for the sake of the following subsection. Let the domain  $\Omega$  be discretized by a suitable tessellation  $\Omega_\Delta$ , where  $n_p$  denotes the number of vertices and  $n_e$  the number of elements in  $\Omega_\Delta$ . All functions from  $H^1(\Omega)$  ( $y, u, v$ , and  $p$ ) are represented by finite elements from  $\mathcal{P}^1(\Omega_\Delta)$ . The corresponding coefficient vectors are denoted by  $\vec{y}, \vec{u}, \vec{v}$ , and  $\vec{p}$ , respectively. The set of test functions  $H^1(\Omega)$  is represented by the same basis functions.

The coefficient functions  $C, \mathbf{a}, \mathbf{b}$ , and  $\mathbf{c}$  as well as the desired state  $y_d$  are assumed to be chosen from  $L^\infty(\Omega)$  and discretized by functions from  $\mathcal{P}^0(\Omega_\Delta)$ ; their discrete approximations are denoted by  $C, \vec{a}, \vec{b}, \vec{c}$ , and  $\vec{y}_d$ , respectively. The matrix  $E_{10} \in \mathbb{R}^{n_e \times n_p}$  realizes the discrete projection of  $\mathcal{P}^1$  approximations into  $\mathcal{P}^0$  and corresponds to the natural embedding operator  $E: H^1(\Omega) \rightarrow L^2(\Omega)$ . The mass matrices  $M_0(1)$  and  $M_1(1)$  correspond to the finite element spaces  $\mathcal{P}^0(\Omega_\Delta)$  and  $\mathcal{P}^1(\Omega_\Delta)$ , respectively. The stiffness matrix associated with the constant coefficient 1 (i.e.,  $C$  is the identity in  $\mathbb{R}^{d \times d}$ ) is denoted by  $K(1)$ . A detailed description of this discretization and the specific forms of these matrices can be found in [11].

The main difficulty when discretizing  $(P_k)$  lies in the handling of the penalty term  $F(u, v)$ . Since the Fischer–Burmeister function is penalized with respect to the space  $L_2(\Omega)$ , the mass matrix  $M_0(1)$  can be used to evaluate integrals over all elements. Interpreting powers and square roots of a vector in a componentwise fashion, a reasonable discretization of  $F(u, v)$  is given by

$$\tilde{F}(\vec{u}, \vec{v}) = \frac{1}{2} \left( \sqrt{(E_{10}\vec{u})^2 + (E_{10}\vec{v})^2} - E_{10}\vec{u} - E_{10}\vec{v} \right)^\top M_0(1) \left( \sqrt{(E_{10}\vec{u})^2 + (E_{10}\vec{v})^2} - E_{10}\vec{u} - E_{10}\vec{v} \right)$$

for all  $\vec{u}, \vec{v} \in \mathbb{R}^{n_p}$ . The appearance of  $E_{10}$  is motivated by the proof of Lemma 4.1, where the penalty functional  $F$  has been represented as the composition of three differentiable mappings: The natural embedding  $E: H^1(\Omega) \rightarrow L^2(\Omega)$ , the Nemytskii-operator associated with the squared Fischer–Burmeister function (as a mapping from  $L^2(\Omega)^2$  to  $L^1(\Omega)$ ), and a linear integral operator.

This discretization strategy leads to the finite-dimensional problem associated with  $(P_k)$  given by

$$(5.1) \quad \begin{cases} \frac{1}{2}(E_{10}\vec{y} - \vec{y}_d)^\top M_0(1)(E_{10}\vec{y} - \vec{y}_d) + \frac{\alpha_1}{2}\vec{u}^\top M_1(1)\vec{u} + \frac{\alpha_2}{2}\vec{v}^\top M_1(1)\vec{v} \\ + \frac{\varepsilon}{2}\vec{u}^\top (M_1(1) + K(1))\vec{u} + \frac{\varepsilon}{2}\vec{v}^\top (M_1(1) + K(1))\vec{v} + \sigma_k \tilde{F}(\vec{u}, \vec{v}) \rightarrow \min_{\vec{y}, \vec{u}, \vec{v}} \\ (M_1(\vec{a}) + K(C))\vec{y} - M_1(\vec{b})\vec{u} - M_1(\vec{c})\vec{v} = 0. \end{cases}$$

For the optimality conditions, one first observes that the quadratic function  $\tilde{F}$  is differentiable everywhere and that its derivative at  $(\vec{u}, \vec{v})$  is given by

$$(5.2) \quad \tilde{F}'(\vec{u}, \vec{v}) = \begin{pmatrix} E_{10}^\top \text{diag}(T_u(\vec{u}, \vec{v})) M_0(1) \left( \sqrt{(E_{10}\vec{u})^2 + (E_{10}\vec{v})^2} - E_{10}\vec{u} - E_{10}\vec{v} \right) \\ E_{10}^\top \text{diag}(T_v(\vec{u}, \vec{v})) M_0(1) \left( \sqrt{(E_{10}\vec{u})^2 + (E_{10}\vec{v})^2} - E_{10}\vec{u} - E_{10}\vec{v} \right) \end{pmatrix},$$

where the vectors  $T_u(\vec{u}, \vec{v}), T_v(\vec{u}, \vec{v}) \in \mathbb{R}^{n_e}$  are defined for all  $i \in \{1, \dots, n_e\}$  as

$$T_u(\vec{u}, \vec{v})_i := \begin{cases} \frac{(E_{10}\vec{u})_i}{\sqrt{(E_{10}\vec{u})_i^2 + (E_{10}\vec{v})_i^2}} - 1 & \text{if } (E_{10}\vec{u})_i \neq 0 \text{ or } (E_{10}\vec{v})_i \neq 0, \\ 0 & \text{if } (E_{10}\vec{u})_i = (E_{10}\vec{v})_i = 0, \end{cases}$$

$$T_v(\vec{u}, \vec{v})_i := \begin{cases} \frac{(E_{10}\vec{v})_i}{\sqrt{(E_{10}\vec{u})_i^2 + (E_{10}\vec{v})_i^2}} - 1 & \text{if } (E_{10}\vec{u})_i \neq 0 \text{ or } (E_{10}\vec{v})_i \neq 0, \\ 0 & \text{if } (E_{10}\vec{u})_i = (E_{10}\vec{v})_i = 0. \end{cases}$$

Note that the case  $(E_{10}\vec{u})_i = (E_{10}\vec{v})_i = 0$  corresponds to the *biactive* case, i.e., where the discretized controls  $\vec{u}$  and  $\vec{v}$  (interpreted in the discretized counterpart of  $L^2(\Omega)$ , i.e., elementwise) are zero at the same time.

Combining (5.1) and (5.2), it is now possible to obtain the following KKT system for the problem (5.1):

$$(5.3a) \quad E_{10}^\top M_0(1)E_{10}\vec{y} - E_{10}^\top M_0(1)\vec{y}_d - (M_1(\vec{a}) + K(C))\vec{p} = 0$$

$$(5.3b) \quad [\alpha_1 M_1(1) + \varepsilon (M_1(1) + K(1))]\vec{u} + \sigma_k \tilde{F}'_u(\vec{u}, \vec{v}) + M_1(\vec{b})\vec{p} = 0$$

$$(5.3c) \quad [\alpha_2 M_1(1) + \varepsilon (M_1(1) + K(1))]\vec{v} + \sigma_k \tilde{F}'_v(\vec{u}, \vec{v}) + M_1(\vec{c})\vec{p} = 0$$

$$(5.3d) \quad -(M_1(\vec{a}) + K(C))\vec{y} + M_1(\vec{b})\vec{u} + M_1(\vec{c})\vec{v} = 0.$$

Recall that  $\vec{p}$  represents the discretized adjoint state and can also be considered as a multiplier related to the discretized state equation. Since the function  $\tilde{F}'$  is nonsmooth but Lipschitz continuous, the nonlinear system (5.3) can be solved using a damped semismooth Newton-type method, see [29]. Note that the domain of nonsmoothness associated with the mapping  $\tilde{F}': \mathbb{R}^{n_p} \times \mathbb{R}^{n_p} \rightarrow \mathbb{R}^{n_p} \times \mathbb{R}^{n_p}$  is given by

$$\{(\vec{u}, \vec{v}) \in \mathbb{R}^{n_p} \times \mathbb{R}^{n_p} \mid \exists i \in \{1, \dots, n_e\}: (E_{10}\vec{u})_i = (E_{10}\vec{v})_i = 0\}.$$

A particular Newton derivative can then be chosen as an element of Clarke's generalized Jacobian, see [5], associated with  $\tilde{F}'$  at  $(\vec{u}, \vec{v})$  that is zero at indices corresponding to biactive components of  $(E_{10}\vec{u}, E_{10}\vec{v})$ . This choice will be used in the proposed method.

Next, due to the well-known local convergence behavior of Newton's method, the initialization of  $\vec{u}$  and  $\vec{v}$  for the numerical solution of (5.3) has to be taken into consideration. For that purpose, consider the (infinite-dimensional) problem

$$(OCNC) \quad \begin{cases} \frac{1}{2} \|E[y] - y_d\|_{L^2(\Omega)}^2 + J(u, v) \rightarrow \min_{y, u, v} \\ -\nabla \cdot (C \nabla y) + ay = bu + cv & \text{a.e. on } \Omega \\ \vec{n} \cdot (C \nabla y) = 0 & \text{a.e. on } \text{bd } \Omega \\ u, v \geq 0 & \text{a.e. on } \Omega \end{cases}$$

which results from (OC<sup>4</sup>) by omitting the equilibrium condition (4.1) and merely imposing nonnegativity constraints. Note that (OCNC) is convex and can be solved globally by combining a penalty algorithm and a semismooth Newton method, see [12]. The associated global minimizer is uniquely determined. If its solution already satisfies the equilibrium condition (4.1), then a global minimizer of (OC<sup>4</sup>) has already been detected. The discretized counterpart of (OCNC) can be derived similarly as stated above. The associated (discrete) optimal solution  $(\vec{y}_0, \vec{u}_0, \vec{v}_0)$  will be used as the starting vector of the semismooth Newton-type method. An abstract description of the proposed numerical method for the computational solution of (OC<sup>4</sup>) is presented in Algorithm 1. In step S2 of this algorithm,  $\|\cdot\|_M$  denotes a weighted Euclidean norm which represents the discretized  $H^1$ -norm, see [12] for details.

---

**Algorithm 1** Abstract algorithm

---

- S0** Let  $\{\sigma_k\}_{k \in \mathbb{N}}$  be a sequence of positive penalty parameters with  $\sigma_k \rightarrow \infty$  as  $k \rightarrow \infty$ . Let a tolerance  $\text{eps} > 0$  be given. Let  $(\vec{y}_0, \vec{u}_0, \vec{v}_0)$  be the (discrete) optimal solution associated with (OCNC). Compute  $\vec{p}_0$  as a solution of the discretized adjoint equation with source  $E_{10}\vec{y}_0 - \vec{y}_d$ . Set  $k := 1$ .
- S1** Solve the discretized KKT system (5.3) for fixed  $\sigma_k$  by a damped, semismooth Newton-type method with starting point  $(\vec{y}_{k-1}, \vec{u}_{k-1}, \vec{v}_{k-1}, \vec{p}_{k-1})$ . Let  $(\vec{y}_k, \vec{u}_k, \vec{v}_k, \vec{p}_k)$  be the associated solution.
- S2** If  $\|(\vec{u}_k, \vec{v}_k) - (\vec{u}_{k-1}, \vec{v}_{k-1})\|_M < \text{eps}$  holds true, then return  $(\vec{u}_k, \vec{v}_k)$ . Otherwise, set  $k := k + 1$  and go to S1.
- 

## 5.2 CHECKING STRONG STATIONARITY

It has to be noted that in step S1 of Algorithm 1, one generally only computes critical points to (5.1). Since the penalty functional  $F$  defined in (4.2) is not convex, these cannot be guaranteed to be global minimizers of (5.1) and therefore the convergence result of Proposition 4.4 does not apply. It is therefore sensible to verify whether the output is at least a strongly stationary point of (3.1) in the sense of Corollary 3.3, since the local minimizers of (3.1) can be found among its strongly stationary points. Note that available first-order methods for the numerical solution of complementarity problems mainly compute so-called Clarke- or Mordukhovich-stationary points and that these stationarity notions are weaker than strong stationarity, see,

e.g., [21] for a discussion of the finite-dimensional situation. Thus, checking strong stationarity is recommendable even if a directly discretized version of (3.1) is solved using the available techniques from finite-dimensional MPCC-theory. A possible approach for verifying strong stationarity is described in the following.

Let  $(y, u, v) \in H^1(\Omega)^3$  be feasible to (OC<sup>4</sup>). If this point is a local minimizer, then Corollary 3.3 implies that

$$(5.4) \quad \langle y - y_d, y \rangle_{L^2(\Omega)} + \alpha_1 \langle u, u \rangle_{L^2(\Omega)} + \alpha_2 \langle v, v \rangle_{L^2(\Omega)} + \varepsilon \langle u, u \rangle_{H^1(\Omega)} + \varepsilon \langle v, v \rangle_{H^1(\Omega)} = 0$$

and that

$$(5.5) \quad \langle y - y_d, z_y \rangle_{L^2(\Omega)} + \alpha_1 \langle u, z_u \rangle_{L^2(\Omega)} + \alpha_2 \langle v, z_v \rangle_{L^2(\Omega)} + \varepsilon \langle u, z_u \rangle_{H^1(\Omega)} + \varepsilon \langle v, z_v \rangle_{H^1(\Omega)} \geq 0$$

for any pair  $(z_u, z_v) \in H_+^1(\Omega)^2$  with

$$\text{supp } z_u \subset I^{+0}(u, v) \cup I^{00}(u, v), \quad \text{supp } z_v \subset I^{0+}(u, v) \cup I^{00}(u, v),$$

where  $z_y \in H^1(\Omega)$  is the solution of the state equation  $A[z_y] - B[z_u] - C[z_v] = 0$ .

Using the same discretization technique as described in Section 5.1, a discrete counterpart to (5.4) is

$$(5.6) \quad \Theta := \vec{y}^\top E_{10}^\top M_0(1) E_{10} \vec{y} - \vec{y}^\top E_{10}^\top M_0(1) \vec{y}_d + \alpha_1 \vec{u}^\top M_1(1) \vec{u} + \alpha_2 \vec{v}^\top M_1(1) \vec{v} \\ + \varepsilon \vec{u}^\top (K(1) + M_1(1)) \vec{u} + \varepsilon \vec{v}^\top (K(1) + M_1(1)) \vec{v} = 0.$$

Clearly, a certain tolerance for the violation of (5.6) needs to be imposed in practice.

The numerical verification of condition (5.5) requires an appropriate choice of discrete test functions  $\vec{z}_u, \vec{z}_v$  for given discretized controls  $(\vec{u}, \vec{v})$  in the finite element space  $\mathcal{P}^1(\Omega_\Delta)$ . Considering the employed finite element discretization of (3.1), one particular choice is from the set of basis functions associated with  $\mathcal{P}^1(\Omega_\Delta)$ . Since the support of each of these “hat functions” covers all elements adjoining a single vertex, a corresponding elementwise approximation of the set  $I^{+0}(u, v)$ ,  $I^{0+}(u, v)$ , and  $I^{00}(u, v)$ , see (3.3), (3.4), and (3.5), respectively, is required as well. This can be defined using the projection of  $\vec{u}, \vec{v}$  from  $\mathcal{P}^1(\Omega_\Delta)$  to  $\mathcal{P}^0(\Omega_\Delta)$  using the matrix  $E_{10}$ , which will be denoted by  $\vec{u}^0 := E_{10} \vec{u}$  and  $\vec{v}^0 := E_{10} \vec{v}$ , respectively. This leads to the corresponding discrete sets

$$I^{+0}(\vec{u}, \vec{v}) := \{i \in \{1, \dots, n_e\} \mid \vec{u}_i^0 > 0 \text{ and } \vec{v}_i^0 = 0\}, \\ I^{00}(\vec{u}, \vec{v}) := \{i \in \{1, \dots, n_e\} \mid \vec{u}_i^0 = 0 \text{ and } \vec{v}_i^0 = 0\}, \\ I^{0+}(\vec{u}, \vec{v}) := \{i \in \{1, \dots, n_e\} \mid \vec{u}_i^0 = 0 \text{ and } \vec{v}_i^0 > 0\}.$$

For any pair of basis vectors  $(\vec{z}_u, \vec{z}_v)$  whose support is contained in  $I^{+0}(\vec{u}, \vec{v}) \cup I^{00}(\vec{u}, \vec{v})$  and  $I^{0+}(\vec{u}, \vec{v}) \cup I^{00}(\vec{u}, \vec{v})$ , respectively, one can then check whether

$$(5.7) \quad \Sigma(\vec{z}_u, \vec{z}_v) := \vec{z}_y^\top E_{10}^\top M_0(1) E_{10} \vec{y} - \vec{z}_y^\top E_{10}^\top M_0(1) \vec{y}_d + \alpha_1 \vec{u}^\top M_1(1) \vec{z}_u + \alpha_2 \vec{v}^\top M_1(1) \vec{z}_v \\ + \varepsilon \vec{u}^\top (K(1) + M_1(1)) \vec{z}_u + \varepsilon \vec{v}^\top (K(1) + M_1(1)) \vec{z}_v \geq 0,$$

where the state  $\vec{z}_y$  associated with  $(\vec{z}_u, \vec{z}_v)$  is obtained via

$$(M_1(\vec{a}) + K(C)) \vec{z}_y = M_1(\vec{b}) \vec{z}_u + M_1(\vec{c}) \vec{z}_v.$$

In numerical practice, a certain tolerance with respect to negative values of  $\Sigma(\vec{z}_u, \vec{z}_v)$  is necessary since [Algorithm 1](#) involves a penalty procedure and hence yields, in general, only *almost* feasible points for (OC<sup>4</sup>). Rather than testing for nonnegativity, it is thus checked whether  $\Sigma(\vec{z}_u, \vec{z}_v)$  is larger than a given negative tolerance.

## 6 NUMERICAL EXAMPLES

The proposed numerical method from [Section 5](#) is illustrated by means of three experiments. These examples are of academical nature and constructed in such a way that the different features of the stationarity test are visualized. In the first example, [Algorithm 1](#) computes globally optimal controls, and thus the results of the corresponding stationarity test provide a first benchmark for a *numerically passed* stationarity test. Examples 2 and 3 provide nontrivial situations where the stationarity test is passed and failed, respectively. Recall that whenever the stationarity test fails, the considered point cannot be a local minimizer of the underlying complementarity-constrained program, see [Corollary 3.3](#).

Let  $\Omega := (0, 1)^2 \subset \mathbb{R}^2$ . For all examples in this section, let  $\mathbf{C}$  be the identity matrix in  $\mathbb{R}^{2 \times 2}$  and let  $\mathbf{a} \equiv 1$ ,  $\mathbf{b} = \chi_{\Omega_u}$ , as well as  $\mathbf{c} = \chi_{\Omega_v}$  hold where  $\Omega_u := \{(x_1, x_2) \in \Omega \mid x_2 < 0.25\}$  and  $\Omega_v := \{(x_1, x_2) \in \Omega \mid x_2 > 0.75\}$  are fixed subdomains of  $\Omega$ . The values  $\alpha_1 = \alpha_2 = 0$  are fixed for this section. Furthermore,  $\varepsilon := 10^{-8}$  is used for all experiments. The implementation is carried out using the object oriented finite element MATLAB class library OOPDE, see [\[28\]](#).

In order to construct examples where the controls are independent of  $x_2$ , cf. [\[7, Section 6\]](#) where parabolic problems were considered and the controls only depend on time, the problem (OC<sup>4</sup>) will be equipped with the additional restrictions

$$(6.1) \quad \partial_{x_2} u = \partial_{x_2} v = 0 \quad \text{a.e. on } \Omega.$$

These constraints realize controls depending only on  $x_1$  and being constant with respect to  $x_2$  while allowing to use the same finite element space for the discretization of  $u$ ,  $v$ , and  $y$ . Note that the additional constraints do not influence the complementarity constraints (which are now imposed on  $\Omega$  rather than  $(0, 1)$ ). Due to these additional gradient constraints, structured grids on the discretized domain  $\Omega_\Delta$  are preferentially used for the following examples. On unstructured grids, which can be created by local refinement of an arbitrary set of triangles of a structured mesh, the use of basis functions from  $\mathcal{P}^1(\Omega_\Delta)$  forces the resulting controls to be globally affine, see [\[11, Section 7.2\]](#) for details. This issue can be solved by choosing basis functions from  $\mathcal{P}^2(\Omega_\Delta)$ . A detailed discussion of optimal control problems with gradient constraints can be found in [\[11\]](#).

To compare results, the solutions of the control problem (OCNC) without complementarity constraints (equipped with the additional constraints (6.1)) will be considered. Recall that optimal controls  $(u, v) \in H^1(\Omega)^2$  of (OCNC) additionally fulfilling the equilibrium condition (4.1) solve (OC<sup>4</sup>) as well, and that these controls are used as starting points for solving (OC<sup>4</sup>). Since the computed controls are nearly constant with respect to  $x_2$ , only  $u(x_1, 0)$  and  $v(x_1, 0)$  are plotted for the sake of easier comparison. To evaluate the satisfaction of the complementarity conditions, the maximal absolute value of the Fischer–Burmeister function applied componentwise to  $(\vec{u}^0, \vec{v}^0)$  is reported. Furthermore,  $\Sigma(\vec{z}_u, \vec{z}_v)$  from (5.7) is checked with a tolerance

$$(6.2) \quad \text{tol} := 0.01 \left| \min_{(\vec{z}_u, \vec{z}_v) \text{ feasible test pair}} \Sigma(\vec{z}_u, \vec{z}_v) \right|,$$

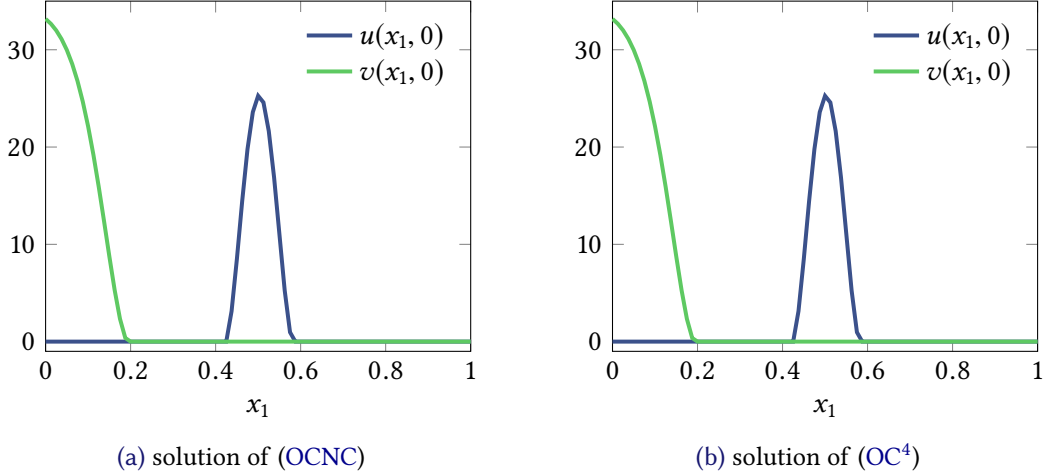


Figure 1: Example 1: computed controls

and the number as well as distribution of pairs  $(\vec{z}_u, \vec{z}_v)$  for which  $\Sigma(\vec{z}_u, \vec{z}_v) > \text{tol}$  (“numerically positive”),  $|\Sigma(\vec{z}_u, \vec{z}_v)| \leq \text{tol}$  (“numerically zero”), or  $\Sigma(\vec{z}_u, \vec{z}_v) < -\text{tol}$  (“numerically negative”) holds is given.

**Example 1** In this example, the desired state is given by the discontinuous function

$$y_d(x) := \begin{cases} 3 & \text{for } x \in [(0.25, 0.75) \times (0, 0.25)] \cup [(0, 0.5) \times (0.75, 1)] \\ 1 & \text{otherwise.} \end{cases}$$

The optimal controls of problem (OCNC) are already (numerically) complementary, see Figure 1a, and thus provide a globally optimal solution of (OC<sup>4</sup>). Correspondingly, they coincide with the controls computed for (OC<sup>4</sup>), see Figure 1b, for which the maximal absolute value of the Fischer–Burmeister function is  $2.08 \cdot 10^{-5}$ . With the tolerance chosen as  $\text{tol} = 2.27 \cdot 10^{-10}$ , 5789 pairs are labeled as numerically positive, 228 as numerically zero, and 544 as numerically negative, see Figure 2b. Thus, only 8.3% of all tested pairs belong to the latter category. Note that  $\Theta = -1.65 \cdot 10^{-7}$  holds for the constant defined in (5.6).

Observing that Algorithm 1 computes the globally optimal solution of (OC<sup>4</sup>) in this example, the above data represent an *approximately* passed stationarity test.

**Example 2** Here, the desired state is chosen to be the (weak) solution of the elliptic boundary value problem

$$\begin{cases} -\Delta y(x) = 0 & \text{a.e. on } \Omega \\ y(x) = 2 \max\{0; x_1 \cos(0.75\pi x_1)\} & \text{a.e. on } \Gamma_1 \\ y(x) = 0.25 & \text{a.e. on } \Gamma_2 \\ \vec{n}(x) \cdot \nabla y(x) = 0 & \text{a.e. on } \Gamma_3 \end{cases}$$

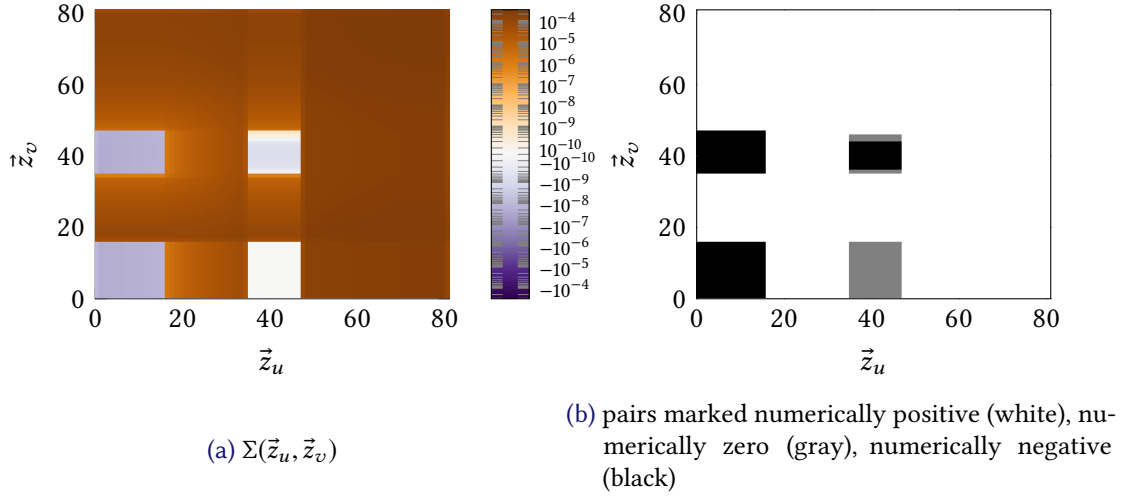


Figure 2: Example 1: values of stationarity test and distribution of failed pairs

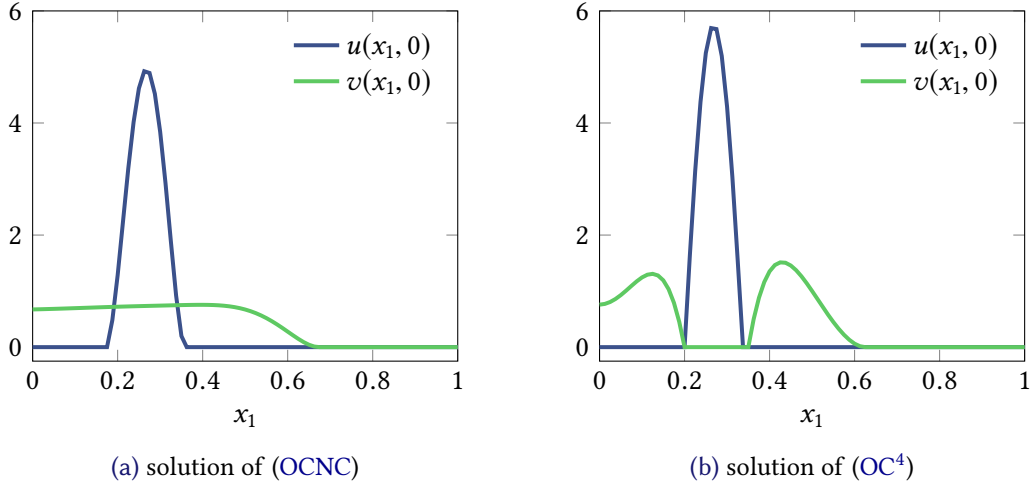


Figure 3: Example 2: computed controls

where  $\Gamma_1 := [0, 1] \times \{0\}$ ,  $\Gamma_2 := [0, 1] \times \{1\}$ , and  $\Gamma_3 := \{0, 1\} \times [0, 1]$  are fixed. The optimal controls of the associated problem (OCNC) do not fulfill the complementarity condition but already provide a biactive set, see Figure 3a. On the other hand, the computed solution for (OC<sup>4</sup>) approximately satisfies the complementarity condition, see Figure 3b, with a maximal absolute value of the Fischer–Burmeister function of approximately  $3.58 \cdot 10^{-6}$ . The minimal value of  $\Sigma(\vec{z}_u, \vec{z}_v)$  was approximately  $-1.62 \cdot 10^{-6}$ , cf. Figure 4a. Accordingly, the tolerance for the stationarity test was chosen as  $\text{tol} = 1.617 \cdot 10^{-8}$ . This leads to 4000 pairs  $(\vec{z}_u, \vec{z}_v)$  marked as “numerically positive”, 2256 as “numerically zero”, and 305 as “numerically negative” and thus failing the strong stationarity test (5.7), see Figure 4b. These amount to approximately 4.7% of the total number 6561 of pairs. Note that pairs where the stationarity test fails correlate with those basis

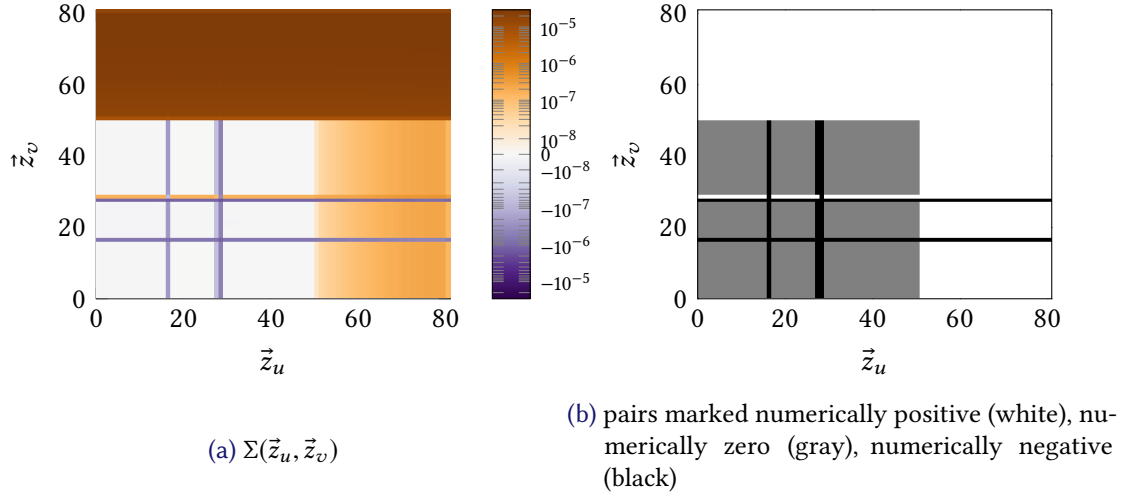


Figure 4: Example 2: values of stationarity test and distribution of failed pairs

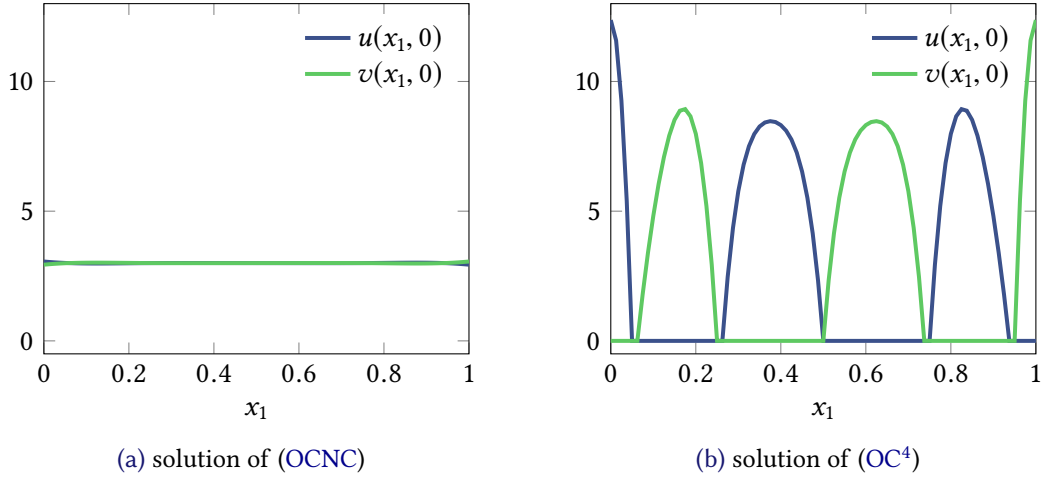


Figure 5: Example 3: computed controls

functions associated with nodes where the subdomains  $I^{+0}(\vec{u}, \vec{v})$  and  $I^{0+}(\vec{u}, \vec{v})$  meet. Finally,  $\Theta = -2.01 \cdot 10^{-9}$  holds.

**Example 3** In the last experiment, the desired state is given by  $y_d \equiv 1.5$ . The optimal controls for the problem (OCNC) are nearly constant functions, see Figure 5a. The controls for the problem (OC<sup>4</sup>) computed via Algorithm 1 are complementary, see Figure 5b. The maximal absolute value of the Fischer–Burmeister function is  $2.02 \cdot 10^{-6}$ . Using the tolerance  $\text{tol} = 1.11 \cdot 10^{-7}$  leads to 0 numerically positive, 5328 numerically zero, and 1233 numerically negative pairs, see Figure 6b. These are more than 18.5% of all tested pairs. In this example,  $\Theta = -3.35 \cdot 10^{-10}$  holds true.



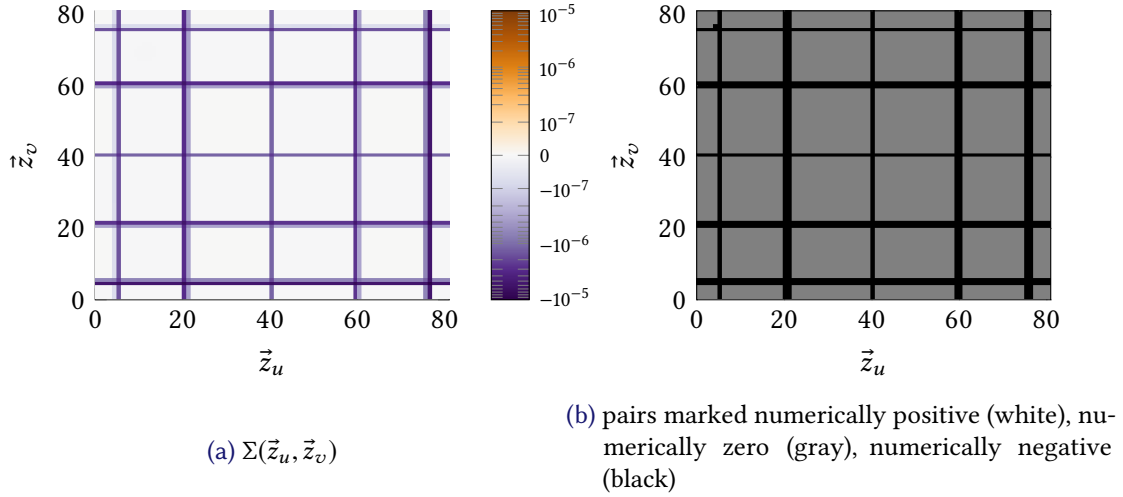


Figure 6: Example 3: values of stationarity test and distribution of failed pairs

	Example 1	Example 2	Example 3
$y_d$	in $L^2(\Omega)$	in $H^1(\Omega)$	constant
$\varepsilon$	$10^{-8}$	$10^{-8}$	$10^{-8}$
complementarity	$2.08 \cdot 10^{-5}$	$3.58 \cdot 10^{-6}$	$2.02 \cdot 10^{-6}$
tol	$2.27 \cdot 10^{-10}$	$1.62 \cdot 10^{-8}$	$1.11 \cdot 10^{-7}$
$\Theta$	$-1.65 \cdot 10^{-7}$	$-2.01 \cdot 10^{-9}$	$-3.35 \cdot 10^{-10}$
num. neg. pairs	8.3%	4.7%	18.5%
stationarity test	passed	passed	failed

Table 1: summary of experiments

**Summary** The results of the numerical experiments are summarized in Table 1, where “complementarity” refers to the maximal absolute value of the elementwise Fischer–Burmeister function. Noting that Experiment 1 provides a benchmark for a passed stationarity test, a computed solution of (OC<sup>4</sup>) is considered as *approximately passing* the strong stationarity test if  $|\Theta| \leq \sqrt{\text{tol}}$  holds for  $\Theta$  defined in (5.6) and the tolerance defined in (6.2), while the number of numerically negative tested pairs is at most 10% of the total number of tested pairs.

It has to be mentioned that more experiments with the same parameter settings of the above three examples were implemented for unstructured grids. In Algorithm 1, the inner iteration implements a damped Newton method to compute the optimal solution of the KKT system (5.3) with the fixed penalty parameter  $\sigma_k$ , which increases in every outer loop. All experiments show that there is no significant correlation between the number of (inner) Newton iterations and the mesh size. However, the solutions calculated on unstructured grids differ significantly from those ones obtained on structured grids, and this phenomenon is not restricted to the use of basis functions from  $\mathcal{P}^1(\Omega_\Delta)$ . The reason behind this fact may be the inherent nonconvexity

of the optimal control problem (OC<sup>4</sup>), which causes the existence of several local minimizers (and thus strongly stationary points). This also explains the observed fact that the output of [Algorithm 1](#) heavily relies on the initial guess for the controls.

## 7 CONCLUSIONS

Optimal control problems with complementarity constraints on the controls admit solutions if the controls are chosen from a first-order Sobolev space. Although necessary optimality conditions of strong stationarity-type can be derived in this case, the explicit characterization of the associated Lagrange multipliers is difficult and remains the topic of further research. However, a penalty method based on the Fischer–Burmeister function can be formulated that ensures convergence to a global minimizers of the original complementarity-constrained problem. In theory, this requires computing global minimizers of the penalized problems, and it has to be investigated whether an adapted method based on KKT points is theoretically possible. Nevertheless, numerical examples illustrate that combined with a computable check for a discrete strong stationarity-type condition, this approach leads to a numerical procedure that in many cases results in nearly strongly stationary points. In light of prominent literature which deals with the numerical treatment of finite-dimensional complementarity problems, see [\[21\]](#) and the references therein, this seems to be the best to be hoped for.

## APPENDIX A A HELPFUL LEMMA

In the proof of [Proposition 4.4](#), the following lemma is used twice.

**Lemma A.1.** *Let  $\{\alpha_k\}_{k \in \mathbb{N}}, \{\beta_k\}_{k \in \mathbb{N}} \subset \mathbb{R}$  be sequences such that  $\alpha_k + \beta_k \rightarrow \alpha + \beta$  holds where  $\alpha, \beta \in \mathbb{R}$  satisfy*

$$\alpha \leq \liminf_{k \rightarrow \infty} \alpha_k, \quad \beta \leq \liminf_{k \rightarrow \infty} \beta_k.$$

*Then, the convergences  $\alpha_k \rightarrow \alpha$  and  $\beta_k \rightarrow \beta$  are valid.*

*Proof.* The assumptions imply that

$$\begin{aligned} \alpha &\leq \liminf_{k \rightarrow \infty} \alpha_k \leq \limsup_{k \rightarrow \infty} \alpha_k = \limsup_{k \rightarrow \infty} (\alpha_k + \beta_k - \beta_k) \\ &= \lim_{k \rightarrow \infty} (\alpha_k + \beta_k) + \limsup_{k \rightarrow \infty} (-\beta_k) \leq \alpha + \beta - \beta = \alpha, \end{aligned}$$

which implies that  $\alpha_k \rightarrow \alpha$ . Now,  $\beta_k \rightarrow \beta$  follows from  $\alpha_k + \beta_k \rightarrow \alpha + \beta$ . □

## ACKNOWLEDGMENTS

The authors sincerely thank Frank Heyde for fruitful discussions about the explicit form of the generalized second-order derivative of the discretized squared Fischer–Burmeister function. Furthermore, the authors appreciate the comments of two anonymous reviewers which helped to improve the presentation of the obtained results. This work is partially supported by the DFG

grants *Parameter Identification in Models With Sharp Phase Transitions* and *Analysis and Solution Methods for Bilevel Optimal Control Problems* under the respective grant numbers CL 487/2-1 and DE 650/10-1 within the Priority Program SPP 1962 (Non-smooth and Complementarity-based Distributed Parameter Systems: Simulation and Hierarchical Optimization).

## REFERENCES

- [1] ADAMS & FOURNIER, *Sobolev Spaces*, Elsevier Science, 2003.
- [2] APPELL & ZABREJKO, *Nonlinear Superposition Operators*, Cambridge University Press, 1990.
- [3] ATTOUCH, BUTTAZZO & MICHAILLE, *Variational Analysis in Sobolev and BV Spaces*, Society for Industrial & Applied Mathematics (SIAM), 2006, DOI: [10.1137/1.9781611973488](https://doi.org/10.1137/1.9781611973488).
- [4] BONNANS & SHAPIRO, *Perturbation Analysis of Optimization Problems*, Springer, 2000, DOI: [10.1007/978-1-4612-1394-9](https://doi.org/10.1007/978-1-4612-1394-9).
- [5] CLARKE, *Optimization and Nonsmooth Analysis*, SIAM, 1990, DOI: [10.1137/1.9781611971309](https://doi.org/10.1137/1.9781611971309).
- [6] CLARKE & de PINHO, Optimal control problems with mixed constraints, *SIAM Journal on Control and Optimization* 48 (2010), 4500–4524, DOI: [10.1137/090757642](https://doi.org/10.1137/090757642).
- [7] CLASON, ITO & KUNISCH, A convex analysis approach to optimal controls with switching structure for partial differential equations, *ESAIM: Control, Optimisation and Calculus of Variations* 22 (2016), 581–609, DOI: [10.1051/cocv/2015017](https://doi.org/10.1051/cocv/2015017).
- [8] CLASON, RUND & KUNISCH, Nonconvex penalization of switching control of partial differential equations, *Systems & Control Letters* 106 (2017), 1–8, DOI: [10.1016/j.sysconle.2017.05.006](https://doi.org/10.1016/j.sysconle.2017.05.006).
- [9] CLASON, RUND, KUNISCH & BARNARD, A convex penalty for switching control of partial differential equations, *Systems & Control Letters* 89 (2016), 66–73, DOI: [10.1016/j.sysconle.2015.12.013](https://doi.org/10.1016/j.sysconle.2015.12.013).
- [10] DE LUCA, FACCHINEI & KANZOW, A theoretical and numerical comparison of some semismooth algorithms for complementarity problems, *Computational Optimization and Applications* 16 (2000), 173–205, DOI: [10.1023/A:1008705425484](https://doi.org/10.1023/A:1008705425484).
- [11] DENG, MEHLITZ & PRÜFERT, On an optimal control problem with gradient constraints, tech. rep. SPP-1962-050, DFG Priority Programme 1962, 2018, URL: <https://spp1962.wias-berlin.de/preprints/050.pdf>.
- [12] DENG, MEHLITZ & PRÜFERT, Optimal control in first-order Sobolev spaces with inequality constraints, *Computational Optimization and Applications* (2019), 1–30, DOI: [10.1007/s10589-018-0053-8](https://doi.org/10.1007/s10589-018-0053-8).
- [13] EVANS, *Partial Differential Equations*, American Mathematical Society, 2010, DOI: [10.1112/blms/20.4.375](https://doi.org/10.1112/blms/20.4.375).
- [14] FACCHINEI, FISCHER & KANZOW, Regularity properties of a semismooth reformulation of variational inequalities, *SIAM Journal on Optimization* 8 (1998), 850–869, DOI: [10.1137/S1052623496298194](https://doi.org/10.1137/S1052623496298194).

- [15] FISCHER, A special Newton-type optimization method, *Optimization* 24 (1992), 269–284, DOI: [10.1080/02331939208843795](https://doi.org/10.1080/02331939208843795).
- [16] GEIGER & KANZOW, *Theorie und Numerik restringierter Optimierungsaufgaben*, Springer, 2002.
- [17] GOLDBERG, KAMPOWSKY & TRÖLTZSCH, On Nemytskij operators in  $L_p$ -spaces of abstract functions, *Mathematische Nachrichten* 155 (1992), 127–140, DOI: [10.1002/mana.19921550110](https://doi.org/10.1002/mana.19921550110).
- [18] GUO & YE, Necessary optimality conditions for optimal control problems with equilibrium constraints, *SIAM Journal on Control and Optimization* 54 (2016), 2710–2733, DOI: [10.1137/15M1013493](https://doi.org/10.1137/15M1013493).
- [19] HARDER & WACHSMUTH, Comparison of optimality systems for the optimal control of the obstacle problem, *GAMM-Mitteilungen* 40 (2018), 312–338, DOI: [10.1002/gamm.201740004](https://doi.org/10.1002/gamm.201740004).
- [20] HARDER & WACHSMUTH, The limiting normal cone of a complementarity set in Sobolev spaces, *Optimization* 67 (2018), 1579–1603, DOI: [10.1080/02331934.2018.1484467](https://doi.org/10.1080/02331934.2018.1484467).
- [21] HOHEISEL, KANZOW & SCHWARTZ, Theoretical and numerical comparison of relaxation methods for mathematical programs with complementarity constraints, *Mathematical Programming* 137 (2013), 257–288, DOI: [10.1007/s10107-011-0488-5](https://doi.org/10.1007/s10107-011-0488-5).
- [22] KANZOW, Some noninterior continuation methods for linear complementarity problems, *SIAM Journal on Matrix Analysis and Applications* 17 (1996), 851–868, DOI: [10.1137/S0895479894273134](https://doi.org/10.1137/S0895479894273134).
- [23] LUO, PANG & RALPH, *Mathematical Programs with Equilibrium Constraints*, Cambridge University Press, 1996, DOI: [10.1017/CBO9780511983658](https://doi.org/10.1017/CBO9780511983658).
- [24] MEHLITZ & WACHSMUTH, Weak and strong stationarity in generalized bilevel programming and bilevel optimal control, *Optimization* 65 (2016), 907–935, DOI: [10.1080/02331934.2015.1122007](https://doi.org/10.1080/02331934.2015.1122007).
- [25] MEHLITZ & WACHSMUTH, The limiting normal cone to pointwise defined sets in Lebesgue spaces, *Set-Valued and Variational Analysis* 26 (2018), 449–467, DOI: [10.1007/s11228-016-0393-4](https://doi.org/10.1007/s11228-016-0393-4).
- [26] PANG & FUKUSHIMA, Complementarity constraint qualifications and simplified B-stationarity conditions for mathematical programs with equilibrium constraints, *Computational Optimization and Applications* 13 (1999), 111–136, DOI: [10.1023/A:1008656806889](https://doi.org/10.1023/A:1008656806889).
- [27] PANG & STEWART, Differential variational inequalities, *Mathematical Programming A* 113 (2008), 345–424, DOI: [10.1007/s10107-006-0052-x](https://doi.org/10.1007/s10107-006-0052-x).
- [28] PRÜFERT, *OOPDE: An object oriented toolbox for finite elements in Matlab*, TU Bergakademie Freiberg, 2015, URL: <http://www.mathe.tu-freiberg.de/files/personal/255/oopde-quickstart-guide-2015.pdf>.
- [29] QI & SUN, A survey of some nonsmooth equations and smoothing Newton methods, in: *Progress in Optimization: Contributions from Australasia*, Springer US, 1999, 121–146, DOI: [10.1007/978-1-4613-3285-5\\_7](https://doi.org/10.1007/978-1-4613-3285-5_7).

- [30] SCHEEL & SCHOLTES, Mathematical programs with complementarity constraints: stationarity, optimality, and sensitivity, *Mathematics of Operations Research* 25 (2000), 1–22, DOI: [10.1287/moor.25.1.1.15213](https://doi.org/10.1287/moor.25.1.1.15213).
- [31] SUN & QI, On NCP-functions, *Computational Optimization and Applications* 13 (1999), 201–220, DOI: [10.1023/A:1008669226453](https://doi.org/10.1023/A:1008669226453).
- [32] TRÖLTZSCH, *Optimal Control of Partial Differential Equations: Theory, Methods and Applications*, American Mathematical Society, 2010, DOI: [10.1090/gsm/112](https://doi.org/10.1090/gsm/112).
- [33] ULBRICH, *Semismooth Newton Methods for Variational Inequalities and Constrained Optimization Problems in Function Spaces*, MOS-SIAM, 2011, DOI: [10.1137/1.9781611970692](https://doi.org/10.1137/1.9781611970692).
- [34] WACHSMUTH, Mathematical programs with complementarity constraints in Banach spaces, *Journal on Optimization Theory and Applications* 166 (2015), 480–507, DOI: [10.1007/s10957-014-0695-3](https://doi.org/10.1007/s10957-014-0695-3).
- [35] YE, Necessary and sufficient optimality conditions for mathematical programs with equilibrium constraints, *J. Math. Anal. Appl.* 307 (2005), 350–369, DOI: [10.1016/j.jmaa.2004.10.032](https://doi.org/10.1016/j.jmaa.2004.10.032).
- [36] YE, ZHU & ZHU, Exact penalization and necessary optimality conditions for generalized bilevel programming problems, *SIAM Journal on Optimization* 7 (1997), 481–507, DOI: [10.1137/S1052623493257344](https://doi.org/10.1137/S1052623493257344).