

# A Penalty Scheme to Solve Constrained Non-convex Optimization Problems in $BV(\Omega)$

Carolin Natemeyer, Daniel Wachsmuth



Non-smooth and Complementarity-based Distributed Parameter Systems: Simulation and Hierarchical Optimization

Preprint Number SPP1962-178

received on October 5, 2021

Edited by SPP1962 at Weierstrass Institute for Applied Analysis and Stochastics (WIAS) Leibniz Institute in the Forschungsverbund Berlin e.V. Mohrenstraße 39, 10117 Berlin, Germany E-Mail: spp1962@wias-berlin.de

World Wide Web: http://spp1962.wias-berlin.de/

# A penalty scheme to solve constrained non-convex optimization problems in $BV(\Omega)$

Carolin Natemeyer, Daniel Wachsmuth \*

October 5, 2021

Abstract. We investigate non-convex optimization problems in  $BV(\Omega)$  with two-sided pointwise inequality constraints. We propose a regularization and penalization method to numerically solve the problem. Under certain conditions, weak limit points of iterates are stationary for the original problem. In addition, we prove optimality conditions for the original problem that contain Lagrange multipliers to the inequality constraints. Numerical experiments confirm the theoretical findings.

**Keywords.** Bounded variation, inequality constraints, optimality conditions, Lagrange multipliers, regularization scheme.

MSC 2020 classification. 49K30, 49M05, 65K10.

### **1** Introduction

Let  $\Omega \in \mathbb{R}^n$  be an open bounded set with Lipschitz boundary. We consider the possibly nonconvex optimization problem of the form

$$\min_{u \in U_{ad} \cap BV(\Omega)} f(u) + |u|_{BV(\Omega)}.$$
(1.1)

Mostly, we will work with

$$U_{ad} = \{ u \in BV(\Omega) : u_a \le u(x) \le u_b \text{ f.a.a. } x \in \Omega \},$$

$$(1.2)$$

where  $u_a, u_b \in \mathbb{R}$ . The function space setting is  $BV(\Omega)$ , i.e., the space of functions of bounded variation that consists of  $L^1(\Omega)$ -functions with weak derivative in the Banach space  $\mathcal{M}(\Omega)$  of real Borel measures on  $\Omega$ . The term  $|u|_{BV(\Omega)}$  denotes the  $BV(\Omega)$ -seminorm, which is equal to the total variation of the measure  $\nabla u$ , i.e.,  $|u|_{BV(\Omega)} = |\nabla u|(\Omega)$ . The functional  $f : L^2(\Omega) \to \mathbb{R}$ is assumed to be smooth and can be non-convex. In particular, we have in mind to choose f(u) := f(y(u)) as the reduced smooth part of an optimal control problem, incorporating the state equation. We will give more details on the assumptions on the optimal control problem in Section 2. Problem (1.1) is solvable, and existence of solutions to (1.1) can be shown by the direct method of the calculus of variations, see Theorem 2.4.

The purpose of the paper is two-fold:

(A) We prove optimality conditions for (1.1)-(1.2) that contain Lagrange multipliers to the inequality constraints  $u_a \leq u$  and  $u \leq u_b$ . Moreover, these multipliers belong to  $L^2(\Omega)$ .

<sup>\*</sup>Institut für Mathematik, Universität Würzburg, 97074 Würzburg, Germany, carolin.natemeyer@mathematik.uni-wuerzburg.de, daniel.wachsmuth@mathematik.uni-wuerzburg.de. This research was partially supported by the German Research Foundation DFG under project grant Wa 3626/3-2.

(B) We investigate an algorithmic scheme to solve (1.1)–(1.2), where weak limit points of iterates satisfy the optimality conditions from (A).

Both of these goals rely on the same approximation method. The algorithmic scheme consists of the following two parts. First, we approximate the non-differentiable total variation term by a smooth approximation and apply a continuation strategy. Second, we address the box constraints with a classical penalty method. Of course, solutions to (1.1) and appearing subproblems are not unique due to the lack of convexity, which makes the analysis challenging. In general, only stationary points of these non-convex problems can be computed. Under suitable assumptions, limit points of the generated sequences of stationary points of the subproblems and of the associated multipliers satisfy a certain necessary optimality condition for the original problem.

In addition, we apply this regularization and penalization approach to local minima of the original problem. This allows us to prove optimality conditions that contains Lagrange multipliers to the inequality constraints, see Theorem 4.3. Such a result is not available in the literature. Admittedly, we have to make the assumption that the constraints  $u_a$  and  $u_b$  are constant functions.

Regularization by total variation is nowadays a standard tool in image analysis. Following the seminal contribution [17], much research was devoted to study such kind of optimization problems. We refer to [9, 10] for a general introduction and an overview on total variation in image analysis. Optimal control problems in BV-spaces were studied for instance in [5, 8, 7]. These control problems are subject to semilinear equations, which results in non-convex control problems. Finite element discretization and convergence analysis for optimization problems in  $BV(\Omega)$  were investigated for instance in [5, 3, 11]. An extensive comparisons of algorithms to solve (1.1) with the choice  $f(u) := \frac{1}{2} ||u - g||_{L^2(\Omega)}$  can be found in [4], see also [16]. In [18], the one-sided obstacle problem in  $BV(\Omega)$  is analyzed under low regularity requirements on the obstacle. Interestingly, we could not find any results, where the existence of Lagrange multipliers to the inequality constraints in  $BV(\Omega)$  is addressed.

One natural idea to regularize (1.1) is to replace the non-differentiable BV-seminorm by a smooth approximation. This was introduced in the image processing setting in [1] with the functional

$$u\mapsto \int_\Omega \sqrt{\epsilon+|\nabla u|^2}\,\mathrm{d} x,$$

which is widely used in the literature. Our regularization method is similar, with the exception that our functional guarantees existence of solutions in  $H^1(\Omega)$ . A similar scheme was employed in the recent work [12], where a path-following inexact Newton method for a convex PDEconstrained optimal control problem in  $BV(\Omega)$  is studied.

Let us comment on the structure of this work. In Section 2 we give a brief introduction to the function space  $BV(\Omega)$  and recall some useful facts. Furthermore, we prove existence of solutions and a necessary first-order optimality condition for the optimization problem (1.1). In Section 3, we introduce the regularization scheme for (1.1) and show that limit points of the suggested smoothing and penalty method satisfy a stationary condition for the original problem, see Theorem 3.16. In Section 4, we apply the regularization scheme to derive an optimality condition for locally optimal solutions of (1.1), see Theorem 4.3. These conditions are stronger than the conditions proven in Section 2, since they contain Lagrange multipliers to the inequality constraints. Finally, we provide numerical results and details regarding the implementation of the method in Section 5.

# 2 Preliminaries and Background

In this section we want to provide some definitions and results regarding the mathematical background of the paper. For details we refer also to, e.g., [2, 11, 8, 5]. First, let us recall that  $\mathcal{M}(\Omega)$  is the dual space of  $C_0(\Omega)$ . The noram of a measure  $\mu \in \mathcal{M}(\Omega)$  is given by

$$\|\mu\|_{\mathcal{M}(\Omega)} = \sup\left\{\int_{\Omega} zd\mu \ : \ z \in C_0(\Omega), \ |z(x)| \le 1 \ \forall x \in \Omega\right\}.$$

The space of functions of bounded variation  $BV(\Omega)$  is a non-reflexive Banach space when endowed with the norm

$$||u||_{BV(\Omega)} = ||u||_{L^1(\Omega)} + |u|_{BV(\Omega)},$$

where we define the total variation of  $\nabla u$  by

$$|u|_{BV(\Omega)} := \sup\left\{\int_{\Omega} u \operatorname{div} \varphi \, \mathrm{d}x : \, \varphi \in C_0^{\infty}(\Omega)^n, \ |\varphi(x)| \le 1 \, \forall x \in \Omega\right\} = \|\nabla u\|_{\mathcal{M}(\Omega)^n}.$$
(2.1)

Here,  $|\cdot|$  denotes the Euclidean norm on  $\mathbb{R}^n$ . In the definition (2.1),  $\nabla : BV(\Omega) \to \mathcal{M}(\Omega)^n$  is a linear and continuous map. Functions in  $BV(\Omega)$  are not necessarily continuous, as an example we mention the characteristic functions of a set with sufficient regularity. If  $u \in W^{1,1}(\Omega)$ , then  $||u||_{BV(\Omega)} = ||u||_{W^{1,1}(\Omega)}$  and  $|u|_{BV(\Omega)} = ||\nabla u||_{L^1(\Omega)}$ .

The Banach space  $BV(\Omega)$  and  $|\cdot|_{BV(\Omega)}$  have some useful properties, which are recalled in the following.

**Proposition 2.1.** Let  $\Omega \subset \mathbb{R}^n$  be an open bounded set with Lipschitz boundary.

- 1. The space  $BV(\Omega)$  is continuously embedded in  $L^r(\Omega)$  for  $1 \le r \le \frac{n}{n-1}$ , while the embedding is compact for  $1 \le r < \frac{n}{n-1}$ .
- 2. Let  $(u_k) \subset BV(\Omega)$  be bounded in  $BV(\Omega)$  with  $u_k \to u$  in  $L^1(\Omega)$ . Then

$$|u|_{BV(\Omega)} \le \liminf_{k \to \infty} |u_k|_{BV(\Omega)}$$

holds.

3. For  $u \in BV(\Omega) \cap L^p(\Omega)$ :  $p \in [1, \infty)$ , there is a sequence  $(u_k) \subset C^{\infty}(\overline{\Omega})$  such that

$$u_k \to u \text{ in } L^p(\Omega) \text{ and } |u_k|_{BV(\Omega)} \to |u|_{BV(\Omega)},$$

$$(2.2)$$

that is,  $C^{\infty}(\overline{\Omega})$  is dense in  $BV(\Omega) \cap L^p(\Omega)$  with respect to the intermediate convergence (2.2).

*Proof.* (1) is [2, Thm. 10.1.3, 10.1.4]. (2) is [2, Prop. 10.1.1(1)]. (3) Can be proven analogously to [2, Theorem 10.1.2], which contains the case p = 1. A proof for  $p \in (1, \infty)$  can be obtained by replacing  $L^1$ -norms by  $L^p$ -norms in the proof by [2], which is by a standard mollification procedure.

**Notation.** Frequently, we will use the following standard notations from convex analysis. The indicator function of a convex set C is denoted by  $\delta_C$ . The normal cone of a convex set C at a point x is denoted by  $N_C(x)$ , and  $\partial h$  denotes the convex subdifferential of a convex function h. It is well known that  $\partial \delta_C(x) = N_C(x)$  holds for convex sets C. Moreover, we introduce the notation

$$J(u) := f(u) + |u|_{BV(\Omega)},$$

which will be used thoughout the paper. In addition, we will denote the positive and negative part of  $x \in \mathbb{R}$  by  $(x)_+ := \max(x, 0)$  and  $(x)_- := \min(x, 0)$ .

#### 2.1 Standing assumption

In order to prove existence of solutions of (1.1) and to analyze the regularization scheme later on, we need some assumptions on the ingredients of the optimal control problem (1.1). Let us start with collecting those in the following paragraph.

#### Assumption A.

- (A1)  $f: L^2(\Omega) \to \mathbb{R}$  is bounded from below and weakly lower semicontinuous.
- (A2)  $f + \delta_{U_{ad}}$  is weakly coercive in  $L^2(\Omega)$ , i.e., for all sequences  $(u_k)$  with  $u_k \in U_{ad}$  and  $\|u_k\|_{L^2(\Omega)} \to \infty$  it follows  $f(u_k) \to +\infty$ .
- (A3)  $U_{ad}$  is a convex and closed subset of  $L^2(\Omega)$  with  $U_{ad} \cap BV(\Omega) \neq \emptyset$ .
- (A4)  $f: L^2(\Omega) \to \mathbb{R}$  is continuously Fréchet differentiable.
- (A5)  $U_{ad} := \{ u \in L^2(\Omega) : u_a \le u(x) \le u_b \text{ f.a.a. } x \in \Omega \}$  with  $u_a, u_b \in \mathbb{R}$  and  $u_a < u_b$ .

Here, assumptions (A1)-(A3) will be used to prove existence of solutions of (1.1). Condition (A4) is necessary to derive necessary optimality conditions. The assumption (A5) will be used in Section 3 to prove boundedness of Lagrange multipliers associated to the inequality constraints in  $U_{ad}$ .

Example 2.2. We consider

$$f(u) := \int_{\Omega} L(x, y_u(x)) \, \mathrm{d}x$$

where  $y_u \in H_0^1(\Omega)$  is defined as the unique weak solution to the elliptic partial differential equation

$$(Ay)(x) + d(x, y(x)) = u(x)$$
 a.e. in  $\Omega$ .

Let us assume that A is an uniformly elliptic operator with bounded coefficients and L, d are Carathéodory functions, continuously differentiable with respect to y such that derivatives are bounded on bounded sets and that d is monotonically increasing with respect to y. Then it is well known that f is covered by Assumption A, see for instance [6].

Example 2.3. Another example is given by the functional

$$f(u) := \int_{\Omega} \int_{I} L(x, t, y_u(x, t)) \,\mathrm{d}x \,\mathrm{d}t.$$

with  $y_u \in L^2(I, H^1_0(\Omega)), I := (0, T), T > 0, \Omega \subset \mathbb{R}^n$ , as the solution of the parabolic equation

$$\partial_t y(x,t) + (Ay)(x,t) + d(x,t,y(x,t)) = u(x,t)$$
 a.e. in  $\Omega \times I$ .

Assuming again an uniformly elliptic operator A and measurable functions L, d of class  $C^2$  w.r.t. y with bounded derivatives, such that d is monotonically increasing, the functional f satisfies Assumption A. We refer to [20, Chapter 5], [7].

#### **2.2** Existence of solutions of (1.1)

Next, we show that under suitable assumptions on the function f, Problem (1.1) has at least one solution in  $L^2(\Omega) \cap BV(\Omega)$ .

**Theorem 2.4.** Let Assumptions (A1)–(A3) be satisfied. Then (1.1) has a solution  $u \in L^2(\Omega) \cap BV(\Omega)$ .

Proof. The proof is standard. We recall it by following the lines of the proof of [5, Theorem 2.1]. Consider a minimizing sequence  $(u_k) \subset L^2(\Omega) \cap BV(\Omega)$ . Since f is bounded from below by (A1),  $(|u_k|_{BV(\Omega)})$  is bounded. By (A2),  $(u_k)$  is bounded in  $L^2(\Omega)$ , and hence  $(u_k)$  is bounded in  $BV(\Omega)$ . By Proposition 2.1, there is a subsequence  $(u_{k_n})$  and  $\bar{u} \in L^2(\Omega) \cap BV(\Omega)$  with  $u_{k_n} \rightharpoonup \bar{u}$  in  $L^2(\Omega)$  and  $u_{k_n} \rightarrow \bar{u}$  in  $L^1(\Omega)$ . Due to (A3),  $U_{ad}$  is weakly closed in  $L^2(\Omega)$ , and  $\bar{u} \in U_{ad}$  follows. By weak lower semicontinuity (A1) and Proposition 2.1, we obtain

$$J(\bar{u}) \leq \liminf_{k_n \to \infty} J(u_{k_n}) = \inf_{u \in L^2(\Omega) \cap BV(\Omega)} J(u),$$

therefore,  $\bar{u}$  is a solution.

#### 2.3 Necessary optimality conditions

Next, we provide a first-order necessary optimality condition for (1.1). A similar result with proof can be found in [5, Theorem 2.3].

**Theorem 2.5.** Let Assumptions (A3)-(A4) be satisfied. Let  $\bar{u} \in BV(\Omega) \cap U_{ad}$  be locally optimal for (1.1) with respect to  $BV(\Omega) \cap L^2(\Omega)$ , i.e., there is r > 0 such that  $J(\bar{u}) \leq J(u)$  for all  $u \in U_{ad}$ with  $\|u - \bar{u}\|_{BV(\Omega)} + \|u - \bar{u}\|_{L^2(\Omega)} < r$ . Then there is  $\lambda \in \partial \| \cdot \|_{\mathcal{M}(\Omega)} (\nabla \bar{u}) \subset (\mathcal{M}(\Omega)^n)^*$  such that

$$-\nabla f(\bar{u}) \in -\operatorname{div} \lambda + N_{U_{ad}}(\bar{u}) \ in \ (BV(\Omega) \cap L^2(\Omega))^*,$$
(2.3)

where  $-\operatorname{div}: (\mathcal{M}(\Omega)^n)^* \to (BV(\Omega) \cap L^2(\Omega))^*$  is the adjoint operator of  $\nabla: BV(\Omega) \cap L^2(\Omega) \to \mathcal{M}(\Omega)^n$ , and the normal cone  $N_{U_{ad}}$  of  $U_{ad}$  is determined with respect to  $BV(\Omega) \cap L^2(\Omega)$ :

$$N_{U_{ad}} = \{ \mu \in (BV(\Omega) \cap L^2(\Omega))^* : \langle \mu, u - \bar{u} \rangle \le 0 \quad \forall u \in U_{ad} \cap BV(\Omega) \}.$$

*Proof.* Let us define  $X := BV(\Omega) \cap L^2(\Omega)$ . By standard arguments, we find

$$-\nabla f(\bar{u}) \in \partial(|\cdot|_{BV(\Omega)} + \delta_{U_{ad}})(\bar{u}) \subset X^*.$$

Recall,  $|u|_{BV(\Omega)} = ||\nabla u||_{\mathcal{M}(\Omega)} = (||\cdot||_{\mathcal{M}(\Omega)} \circ \nabla)(u)$ . Following [5, Theorem 2.3], let  $-\operatorname{div}$ :  $(\mathcal{M}(\Omega)^n)^* \to X^*$  denote the adjoint operator of  $\nabla : X \to \mathcal{M}(\Omega)^n$ , i.e.,

$$\langle \operatorname{div} \lambda, u \rangle = - \langle \lambda, \nabla u \rangle \quad \forall u \in X, \ \lambda \in (\mathcal{M}(\Omega)^n)^*.$$

By the sum and chain rules for the convex subdifferential, (2.3) can be rewritten as

$$-\nabla f(\bar{u}) \in -\operatorname{div}\left(\partial \|\cdot\|_{\mathcal{M}(\Omega)}(\nabla \bar{u})\right) + N_{U_{ad}}(\bar{u}),$$

which is the claim.

#### 3 Regularization scheme

In this section, we introduce the regularization scheme for (1.1). We will use a smoothing of the *BV*-norm as well as a penalization of the constraints  $u \in U_{ad}$ . In order to approximate the *BV*-norm by smooth functions, we introduce the following family  $(\psi_{\epsilon})_{\epsilon>0}$  of smooth functions with  $\psi_{\epsilon}(t) \approx |t|$ . For  $\epsilon > 0$  we define  $\psi_{\epsilon} : \mathbb{R}^n \to \mathbb{R}^n$  by

$$\psi_{\epsilon}(t) := \sqrt{\epsilon + |t|^2} + \epsilon |t|^2, \qquad (3.1)$$

where  $|\cdot|$  denotes the Euclidian norm in  $\mathbb{R}^n$ . As a first direct consequence of the above definition we have:

**Lemma 3.1.** For  $\epsilon > 0$ ,  $(\psi_{\epsilon})_{\epsilon > 0}$  is a family of twice continuously differentiable functions from  $\mathbb{R}^n$  to  $\mathbb{R}^n$  with the following properties:

- (1)  $t \mapsto \psi_{\epsilon}(t)$  is convex.
- (2)  $\psi_{\epsilon}(t) \ge |t| + \epsilon |t|^2$  and  $\psi'_{\epsilon}(t)t \ge 0$  for all  $t \in \mathbb{R}^n$ .
- (3) For all  $t \in \mathbb{R}^n$ ,  $\epsilon \mapsto \psi_{\epsilon}(t)$  is monotonically increasing and  $\psi_{\epsilon}(t) \to |t|$  as  $\epsilon \searrow 0$ .
- (4)  $t \mapsto \psi'_{\epsilon}(t)t$  is coercive, i.e.,  $\psi'_{\epsilon}(t)t \to \infty$  as  $|t| \to \infty$ .
- (5)  $\psi'_{\epsilon}(t)t \ge |t| \sqrt{\epsilon} \text{ for all } t \in \mathbb{R}^n.$

*Proof.* Properties (1)-(4) are immediate consequences of the definition. (5) can be proven as follows:

$$\psi_{\epsilon}'(t)t - |t| \ge \frac{|t|^2 - |t|\sqrt{\epsilon + |t|^2}}{\sqrt{\epsilon + |t|^2}} = \frac{-\epsilon|t|^2}{\sqrt{\epsilon + |t|^2}(|t|^2 + |t|\sqrt{\epsilon + |t|^2})} \ge \frac{-\epsilon|t|^2}{\sqrt{\epsilon + |t|^2}|t|^2} \ge -\sqrt{\epsilon}.$$

The BV-minimization problem (1.1) is then approximated by

$$\min_{u \in H^1(\Omega)} f(u) + \int_{\Omega} \psi_{\epsilon}(\nabla u) \, \mathrm{d}x \quad \text{s.t. } u \in U_{ad}.$$
 (P<sub>\epsilon</sub>)

The choice (3.1) of  $\psi_{\epsilon}$  guarantees the existence of solutions of this problem in  $H^1(\Omega)$ . Note that the standard approximation  $|t| \approx \sqrt{\epsilon + |t|^2}$  does Lemma 3.1(2), which ensures that  $u \mapsto \int_{\Omega} \psi_{\epsilon}(\nabla u) \, dx$  is weakly coercive in  $H^1(\Omega)$ . The existence of solutions  $u \in H^1(\Omega)$  (and the fact that  $\nabla u$  is a measurable function) is important for the subsequent analysis.

Due to the presence of the inequality constraints, Problem  $(P_{\epsilon})$  is difficult to solve. Following existing approaches in the literature, see, e.g., [14, 19], we will use a smooth penalization of these constraints. We define the smooth function  $\max_{\rho}$  by

$$\max_{\rho}(x) := \begin{cases} \max(0, x) & \text{if } |x| \ge \frac{1}{2\rho}, \\ \frac{\rho}{2}(x + \frac{1}{2\rho})^2 & \text{if } |x| \le \frac{1}{2\rho}, \end{cases}$$
(3.2)

where  $\rho > 0$ . Due to the inequalities

$$0 \le \max(0, x) \le \max_{\rho}(x) \le \max(x, 0) + \frac{1}{2\rho} \quad \forall x \in \mathbb{R},$$
(3.3)

it can be considered as an approximation of  $x \mapsto \max(x, 0) = (x)_+$ . In addition, one verifies  $\max(0, x) \leq t^{-1} \max_{\rho}(tx)$  for all  $x \in \mathbb{R}$  and t > 0.

Let us introduce

$$M_{\rho}(x) := \int_{-\infty}^{x} \max_{\rho}(t) \,\mathrm{d}t. \tag{3.4}$$

Using this function, we define the penalized problem by

$$\min_{u \in H^1(\Omega)} f(u) + \int_{\Omega} \psi_{\epsilon}(\nabla u) \, \mathrm{d}x + \int_{\Omega} \frac{1}{\rho} \left( M_{\rho}(\rho(u_a - u)) + M_{\rho}(\rho(u - u_b)) \right) \, \mathrm{d}x. \tag{$P_{\epsilon,\rho}$}$$

If  $u \in H^1(\Omega)$  is a local solution of  $(P_{\epsilon,\rho})$  then it satisfies

$$\int_{\Omega} \nabla f(u)v + \psi'_{\epsilon}(\nabla u)\nabla v - \lambda^{a}_{\rho}(u)v + \lambda^{b}_{\rho}(u)v \,\mathrm{d}x = 0 \quad \forall v \in H^{1}(\Omega),$$
(3.5)

where we used the abbreviations

$$\lambda_{\rho}^{a}(u) := \max_{\rho}(\rho(u_{a} - u)), \ \lambda_{\rho}^{b}(u) := \max_{\rho}(\rho(u - u_{b})).$$
(3.6)

Existence of solutions of  $(P_{\epsilon,\rho})$  and necessity of (3.5) for local optimality can be proven by standard arguments.

**Corollary 3.2.** Let assumptions (A1)–(A5) be satisfied. Then the equation (3.5) admits a solution  $u \in H^1(\Omega)$ .

We will investigate the behavior of a penalty and smoothing method to solve (1.1). Since  $(P_{\epsilon,\rho})$  is a non-convex problem, it is unrealistic to assume that one can compute global solutions. Instead, the iterates will be chosen as stationary points of  $(P_{\epsilon,\rho})$ . Hence, we are interested in the behavior of stationary points  $u_{\epsilon,\rho}$  and corresponding multipliers  $\lambda_{\rho}^{a}(u_{\epsilon,\rho})$ ,  $\lambda_{\rho}^{b}(u_{\epsilon,\rho})$  as  $\rho \to \infty$  and  $\epsilon \searrow 0$ .

The resulting method then reads as follows.

Algorithm 3.3. Choose  $\epsilon_0 \in (0, 1), \rho_0 > 0$  and  $u_0 \in H^1(\Omega)$ .

1. Compute  $u_k$  as solution to

$$\int_{\Omega} \nabla f(u)v + \psi'_{\epsilon_k}(\nabla u)\nabla v - \lambda^a_{\rho_k}(u)v + \lambda^b_{\rho_k}(u)v \,\mathrm{d}x = 0 \quad \forall v \in H^1(\Omega),$$
(3.7)

where  $\lambda_{o}^{a}(u), \lambda_{o}^{b}(u)$  are defined in (3.6).

- 2. Choose  $\rho_{k+1} > \rho_k$ ,  $\epsilon_{k+1} < \epsilon_k$ .
- 3. If a suitable stopping criterion is satisfied: Stop. Else set k := k + 1 and go to Step 1.

In view of Corollary 3.2, the algorithm is well-defined. In the following, we assume that the algorithm generates an infinite sequence of iterates  $(u_k, \lambda_{\rho_k}^a(u_k), \lambda_{\rho_k}^b(u_k))$ . Here, we are interested to prove that weak limit points are stationary, i.e., they satisfy the optimality condition (2.3) for (1.1). Throughout the subsequent analysis, we assume that assumptions (A1)–(A5) are satisfied

#### 3.1 A-priori bounds

In order to investigate the sequences of iterates  $(u_k)$  and its (weak) limit points, it is reasonable to derive bounds of iterates  $u_k$ , i.e., solutions of (3.5), first. To this end, we will study solutions  $u \in H^1(\Omega)$  of the nonlinear variational equation

$$\int_{\Omega} \psi_{\epsilon}'(\nabla u) \nabla v - \lambda_{\rho}^{a}(u)v + \lambda_{\rho}^{b}(u)v \,\mathrm{d}x = \int_{\Omega} gv \,\mathrm{d}x \quad \forall v \in H^{1}(\Omega)$$
(3.8)

for  $\epsilon > 0$ ,  $\rho > 0$ , and  $g \in L^2(\Omega)$ . The functions  $\lambda_{\rho}^a(u)$ ,  $\lambda_{\rho}^b(u)$  are defined in (3.6) as

$$\lambda_{\rho}^{a}(u) := \max_{\rho}(\rho(u_{a}-u)), \ \lambda_{\rho}^{b}(u) := \max_{\rho}(\rho(u-u_{b})).$$

Let us start with the following lemma. It shows that the supports of the multipliers  $\lambda_{\rho}^{a}(u)$ ,  $\lambda_{\rho}^{b}(u)$  do not overlap if the penalty parameter is large enough.

**Lemma 3.4.** Let  $u \in H^1(\Omega)$  be a solution of (3.8) to  $g \in L^2(\Omega)$ . Suppose  $\rho^2 \geq \frac{1}{u_b - u_a}$ . Then it holds  $\lambda_{\rho}^a(u) \cdot \lambda_{\rho}^b(u) = 0$  a.e. on  $\Omega$ . Proof. Let  $x \in \Omega$  such that  $\lambda_{\rho}^{a}(u)(x) \cdot \lambda_{\rho}^{b}(u)(x) \neq 0$ . Then it holds  $u(x) < u_{a}(x) + \frac{1}{2\rho^{2}}$  and  $u(x) > u_{b}(x) - \frac{1}{2\rho^{2}}$ . Consequently,  $\rho^{2} < \frac{1}{u_{b} - u_{a}}$  follows.

Under the assumptions that the bounds  $u_a$  and  $u_b$  are constant, we can prove the following series of helpful results regarding the boundedness of iterates. Let us start with the boundedness of the multiplier sequences. This result is inspired by related results for the  $H^1$ -obstacle problem, see, e.g., [13, Lemma 5.1], see also [19, Lemma 2.3]. The results for  $H^1$ -obstacle problems require the assumption  $\Delta u_a, \Delta u_b \in L^2(\Omega)$ . It is not clear to us, how the following proof can be generalized to non-constant obstacles  $u_a, u_b$ .

**Lemma 3.5.** Let  $u \in H^1(\Omega)$  be a solution of (3.8) to  $g \in L^2(\Omega)$ . Suppose  $\rho^2 \ge \frac{1}{u_b - u_a}$ . Then it holds

$$\|\lambda_{\rho}^{a}(u)\|_{L^{2}(\Omega)} + \|\lambda_{\rho}^{b}(u)\|_{L^{2}(\Omega)} \leq 2\|g\|_{L^{2}(\Omega)}.$$

*Proof.* To show boundedness of the multipliers  $\lambda_{\rho}^{a}(u), \lambda_{\rho}^{b}(u)$  in  $L^{2}(\Omega)$ , we test the optimality condition (3.8) with  $\lambda_{\rho}^{a}(u)$  and  $\lambda_{\rho}^{b}(u)$ , respectively. We get for  $\lambda_{\rho}^{a}(u) = \max_{\rho}(\rho(u_{a} - u))$ 

$$\|\lambda_{\rho}^{a}(u)\|_{L^{2}(\Omega)}^{2} = \int_{\Omega} \psi_{\epsilon}'(\nabla u) \nabla(\max_{\rho}(\rho(u_{a}-u))) \,\mathrm{d}x - \int_{\Omega} g\lambda_{\rho}^{a}(u) \,\mathrm{d}x + \int_{\Omega} \lambda_{\rho}^{b}(u)\lambda_{\rho}^{a}(u) \,\mathrm{d}x.$$

Due to Lemma 3.4, the last term is zero. It remains to analyze the first term. Here, we find

$$\int_{\Omega} \psi_{\epsilon}'(\nabla u) \nabla \max_{\rho}(\rho(u_a - u)) \, \mathrm{d}x = \int_{\Omega} \psi_{\epsilon}'(\nabla u) \rho \max_{\rho}'(\rho(u_a - u)) \nabla(-u) \, \mathrm{d}x \le 0,$$

where we used Lemma 3.1(2) and  $\max_{\rho} \geq 0$ . This proves  $\|\lambda_{\rho}^{a}(u)\|_{L^{2}(\Omega)} \leq \|g\|_{L^{2}(\Omega)}$ . Similarly,  $\|\lambda_{\rho}^{b}(u)\|_{L^{2}(\Omega)} \leq \|g\|_{L^{2}(\Omega)}$  can be proven.

**Corollary 3.6.** Let  $u \in H^1(\Omega)$  be a solution of (3.8) to  $g \in L^2(\Omega)$ . Suppose  $\rho^2 \geq \frac{1}{u_b - u_a}$ . Then it holds

$$||(u - u_b)_+||_{L^2(\Omega)} + ||(u_a - u)_+||_{L^2(\Omega)} \le 2\rho^{-1} ||g||_{L^2(\Omega)}.$$

*Proof.* Due to the definition of  $\max_{\rho}$ , we have  $\max(x, 0) \leq \rho^{-1} \max_{\rho}(\rho x)$  for all  $x \in \mathbb{R}$ . This implies

$$\|(u-u_b)_+\|_{L^2(\Omega)} \le \rho^{-1} \|\max_{\rho}(\rho(u-u_b))\|_{L^2(\Omega)} = \rho^{-1} \|\lambda_{\rho}^b(u)\|_{L^2(\Omega)},$$

and the claim follows by Lemma 3.5 above.

**Corollary 3.7.** Let  $u \in H^1(\Omega)$  be a solution of (3.8) to  $g \in L^2(\Omega)$ . Suppose  $\rho^2 \ge \frac{1}{u_b - u_a}$ . Then it holds

$$|u||_{L^{2}(\Omega)} \leq 2\rho^{-1} ||g||_{L^{2}(\Omega)} + ||\max(|u_{a}|, |u_{b}|)||_{L^{2}(\Omega)}$$

*Proof.* The claim is a consequence of Corollary 3.6 and the identity

$$u = (u - u_b)_+ - (u_a - u)_+ + \operatorname{proj}_{[u_a, u_b]}(u).$$
(3.9)

**Lemma 3.8.** Let  $u \in H^1(\Omega)$  be a solution of (3.8) to  $g \in L^2(\Omega)$ . Suppose  $\rho^2 \ge \frac{1}{u_b - u_a}$ . Then it holds

$$\|\nabla u\|_{L^{1}(\Omega)} \leq 3\|g\|_{L^{2}(\Omega)}\|u\|_{L^{2}(\Omega)} + \sqrt{\epsilon}|\Omega|.$$

*Proof.* We test the optimality condition (3.5) with u and use Lemma 3.1(5) to get the estimate

$$\begin{split} \int_{\Omega} |\nabla u| - \sqrt{\epsilon} \, \mathrm{d}x &\leq \int_{\Omega} \psi_{\epsilon}'(\nabla u) \nabla u \, \mathrm{d}x = \left| \int_{\Omega} gu - \lambda_{\rho}^{a}(u)u + \lambda_{\rho}^{b}(u)u \, \mathrm{d}x \right| \\ &\leq (\|g\|_{L^{2}(\Omega)} + \|\lambda_{\rho}^{a}(u)\|_{L^{2}(\Omega)} + \|\lambda_{\rho}^{b}(u)\|_{L^{2}(\Omega)}) \|u\|_{L^{2}(\Omega)}. \end{split}$$

The claim follows with the estimate of Lemma 3.5.

#### **3.2** Preliminary convergence results

As next step, we derive convergence properties of solutions  $u_k \in H^1(\Omega)$  of

$$\int_{\Omega} \psi_{\epsilon_k}'(\nabla u_k) \nabla v - \lambda_{\rho_k}^a(u_k) v + \lambda_{\rho_k}^b(u_k) v \, \mathrm{d}x = \int_{\Omega} g_k v \, \mathrm{d}x \quad \forall \, v \in H^1(\Omega)$$
(3.10)

where

$$\epsilon_k \searrow 0, \ \rho_k \to +\infty, \ g_k \rightharpoonup g \text{ in } L^2(\Omega).$$
 (3.11)

From the results of the previous section, we immediately obtain that  $(u_k)$  is bounded in  $BV(\Omega)$ , and  $(\lambda_{\rho_k}^a(u_k))$  and  $(\lambda_{\rho_k}^b(u_k))$  are bounded in  $L^2(\Omega)$ . Moreover, strong limit points of  $(u_k)$  satisfy the inequality constraints in (1.2) due to Corollary 3.6. In a first result, we need to lift the strong convergence of (a subsequence of)  $(u_k)$  in  $L^1(\Omega)$ , which is a consequence of the compact embedding  $BV(\Omega) \hookrightarrow L^r(\Omega), r < \frac{n}{n-1}$ , to strong convergence in  $L^2(\Omega)$ .

**Lemma 3.9.** Assume (3.11). Let  $u_k$ ,  $k \in \mathbb{N}$ , be a solution of (3.10). Suppose  $u_k \to u$  in  $L^1(\Omega)$ . Then  $u_k \to u$  in  $L^2(\Omega)$ .

*Proof.* We use again the identity (3.9):

$$u_k = (u_k - u_b)_+ - (u_a - u_k)_+ + \operatorname{proj}_{[u_a, u_b]}(u_k).$$

By Corollary 3.6, the first two terms converge to zero in  $L^2(\Omega)$ , which implies  $u_a \leq u \leq u_b$ almost everywhere. The sequence  $(\operatorname{proj}_{[u_a,u_b]}(u_k))$  converges to  $\operatorname{proj}_{[u_a,u_b]}(u) = u$  in  $L^1(\Omega)$  and is bounded in  $L^{\infty}(\Omega)$ . By Hölder inequality it converges in  $L^2(\Omega)$ .

We are now in the position to prove existence of suitably converging subsequence under assumption (3.11).

**Theorem 3.10.** Assume (3.11). Let  $(u_k)$ ,  $k \in \mathbb{N}$ , be a family of solutions of (3.10). Then there is a subsequence such that  $u_{k_n} \to u^*$ ,  $\lambda^a_{\rho_{k_n}}(u_{k_n}) \rightharpoonup \lambda^a$ , and  $\lambda^b_{\rho_{k_n}}(u_{k_n}) \rightharpoonup \lambda^b$  in  $L^2(\Omega)$ .

Proof. By Lemmas 3.5, Corollary 3.7, and Lemma 3.8,  $(u_k)$  is bounded in  $BV(\Omega)\cap L^2(\Omega)$ , and  $(\lambda_{\rho_k}^a(u_k))$  and  $(\lambda_{\rho_k}^b(u_k))$  are bounded in  $L^2(\Omega)$ . Then we can choose  $(u_k)$  as a subsequence that converges strongly in  $L^1(\Omega)$  by Proposition 2.1. Due to Lemma 3.9 this convergence is strong in  $L^2(\Omega)$ . Now, extracting additional weakly converging subsequences from  $(\lambda_{\rho_k}^a(u_k))$  and  $(\lambda_{\rho_k}^b(u_k))$  finishes the proof.

The next result shows that limit points of  $(u_k, \lambda^a_{\rho_k}(u_k), \lambda^b_{\rho_k}(u_k))$  satisfy the usual complementarity conditions.

**Lemma 3.11.** Assume (3.11). Let  $u_k$ ,  $k \in \mathbb{N}$ , be a solution of (3.10). Let  $u_k \to u^*$ ,  $\lambda^a_{\rho_k}(u_k) \rightharpoonup \lambda^a$ , and  $\lambda^b_{\rho_k}(u_k) \rightharpoonup \lambda^b$  in  $L^2(\Omega)$ . Then it follows that

$$(\lambda^{b}, u^{*} - u_{b}) = 0 \text{ and } (\lambda^{a}, u_{a} - u^{*}) = 0.$$

*Proof.* Due to (3.11),  $(g_k)$  is bounded in  $L^2(\Omega)$ . We deduce using (3.3)

$$\begin{aligned} |(\lambda_{\rho_k}^b(u_k), u_k - u_b)| &\leq \int_{\Omega} \max_{\rho_k} (\rho_k(u_k - u_b)) |u_k - u_b| \, \mathrm{d}x \\ &\leq \int_{\Omega} \left( \rho_k(u_k - u_b)_+ + \frac{1}{2\rho_k} \right) |u_k - u_b| \, \mathrm{d}x \\ &= \rho_k \int_{\Omega} (u_k - u_b)_+ (u_k - u_b) \, \mathrm{d}x + \int_{\Omega} \frac{1}{2\rho_k} |u_k - u_b| \, \mathrm{d}x \\ &= \rho_k ||(u_k - u_b)_+||_{L^2(\Omega)}^2 + \frac{1}{2\rho_k} ||u_k - u_b||_{L^1(\Omega)}. \end{aligned}$$

Due to Corollary 3.6 and the boundedness of  $(u_k)$  in  $L^2(\Omega)$ , both expressions tend to zero for  $k \to \infty$ . This implies

$$(\lambda^{b}, u^{*} - u_{b}) = \lim_{k \to \infty} (\lambda^{b}_{\rho_{k}}(u_{k}), u_{k} - u_{b}) = 0.$$

The same argumentation yields the claim for  $(\lambda^a, u_a - u^*)$ .

We will now show that weak limit points of solutions to (3.10) satisfy a stationary condition similar to the one for the original problem (1.1). We will utilize this result twice: first we apply it to iterates of Algorithm 3.3, second we will use it to prove a optimality condition for (1.1)that has a different structure than that of Theorem 2.5.

**Theorem 3.12.** Assume (A1)-(A5) and (3.11). Let  $(u_k)$ ,  $k \in \mathbb{N}$ , be a family of solutions of (3.10). Let  $u_k \to u^*$ ,  $\lambda^a_{\rho_k}(u_k) \rightharpoonup \lambda^a$ , and  $\lambda^b_{\rho_k}(u_k) \rightharpoonup \lambda^b$  in  $L^2(\Omega)$ . Then it holds

$$u^{*} \in U_{ad}$$
  
 $\lambda^{a} \ge 0, \qquad \lambda^{b} \ge 0,$ 
 $\lambda^{a}, u_{a} - u^{*}) = 0, \qquad (\lambda^{b}, u^{*} - u_{b}) = 0.$ 
(3.12)

In addition, there is  $\mu^* \in L^{\infty}(\Omega)^n$  with div  $\mu^* \in L^2(\Omega)$  such that

(

$$-\operatorname{div}\mu^* - \lambda^a + \lambda^b = g$$

and

$$-\operatorname{div}\mu^* \in \partial(|\cdot|_{BV(\Omega)})(u^*)$$

Moreover, there is  $\lambda^* \in \partial \| \cdot \|_{\mathcal{M}(\Omega)}(\nabla u^*) \subset (\mathcal{M}(\Omega)^n)^*$  with div  $\lambda^* \in L^2(\Omega)$  such that

$$-\operatorname{div}\lambda^* - \lambda^a + \lambda^b = g.$$

*Proof.* The system (3.12) is a consequence of Corollary 3.6, Lemma 3.11, and the non-negativity of max<sub> $\rho$ </sub>. In order to pass to the limit in (3.10), we need to analyze the term involving  $\psi'_{\epsilon_k}(\nabla u_k)$ . Here, we argue similar as in the proof of [8, Theorem 10]. Let  $v \in C_c^{\infty}(\Omega)$ . Then, we find

$$\int_{\Omega} \psi_{\epsilon_k}'(\nabla u_k) \nabla v \, \mathrm{d}x = \int_{\Omega} \frac{\nabla u_k}{\sqrt{\epsilon_k + |\nabla u_k|^2}} \nabla v + 2\epsilon \nabla u_k \nabla v \, \mathrm{d}x = \int_{\Omega} \frac{\nabla u_k}{\sqrt{\epsilon_k + |\nabla u_k|^2}} \nabla v - 2\epsilon u_k \Delta v \, \mathrm{d}x.$$

Let us define  $\mu_k \in L^{\infty}(\Omega)$  by

$$\mu_k := \frac{\nabla u_k}{\sqrt{\epsilon_k + |\nabla u_k|^2}}$$

Clearly, the sequence  $(\mu_k)$  is bounded in  $L^{\infty}(\Omega)^n$ , and there exists a subsequence converging weak-star in  $L^{\infty}(\Omega)^n$ . W.l.o.g. we can assume  $\mu_k \rightharpoonup^* \mu^*$  in  $L^{\infty}(\Omega)^n$ . Since  $(u_k)$  is bounded in  $L^2(\Omega)$ , we obtain

$$\lim_{k \to \infty} \int_{\Omega} \psi'_{\epsilon_k}(\nabla u_k) \nabla v \, \mathrm{d}x = \lim_{k \to \infty} \int_{\Omega} \mu_k \nabla v - 2\epsilon u_k \Delta v \, \mathrm{d}x = \int_{\Omega} \mu^* \nabla v \, \mathrm{d}x$$

Then we can pass to the limit in (3.10) to find

$$\int_{\Omega} \mu^* \nabla v - \lambda^a v + \lambda^b v \, \mathrm{d}x = \int_{\Omega} g v \, \mathrm{d}x$$

which is satisfied for all  $v \in C_c^{\infty}(\Omega)$ . This implies

$$-\operatorname{div} \mu^* = g + \lambda^a - \lambda^b \in L^2(\Omega)$$

Let now  $v \in C^{\infty}(\overline{\Omega})$ . By convexity of  $\psi_{\epsilon}$ , we have

$$\int_{\Omega} \psi_{\epsilon_k}(\nabla u_k) + \psi'_{\epsilon_k}(\nabla u_k)(\nabla v - \nabla u_k) \, \mathrm{d}x \le \int_{\Omega} \psi_{\epsilon_k}(\nabla v) \, \mathrm{d}x.$$
(3.13)

Here, we find  $\int_{\Omega} \psi_{\epsilon_k}(\nabla v) \, \mathrm{d}x \to \|\nabla v\|_{L^1(\Omega)}$  and

$$\liminf_{k \to \infty} \int_{\Omega} \psi_{\epsilon_k}(\nabla u_k) \, \mathrm{d}x \ge \liminf_{k \to \infty} \|\nabla u_k\|_{L^1(\Omega)} \ge |u^*|_{BV(\Omega)},$$

cf., Proposition 2.1. Here, we used that  $(u_k)$  is bounded in  $BV(\Omega)$  due to Lemma 3.8 and (3.11). Using the equation (3.10), we find

$$\int_{\Omega} \psi_{\epsilon_k}'(\nabla u_k) (\nabla v - \nabla u_k) \, \mathrm{d}x = \int_{\Omega} (g_k + \lambda_{\rho_k}^a(u_k) - \lambda_{\rho_k}^b(u_k)) (v - u_k) \, \mathrm{d}x$$
$$\to \int_{\Omega} (g + \lambda^a - \lambda^b) (v - u^*) \, \mathrm{d}x = \int_{\Omega} -\operatorname{div} \mu^* (v - u^*) \, \mathrm{d}x.$$

Then we can pass to the limit in (3.13) to obtain

$$|u^*|_{BV(\Omega)} + \int_{\Omega} -\operatorname{div} \mu^*(v - u^*) \,\mathrm{d}x \le |v|_{BV(\Omega)}$$

for all  $v \in C^{\infty}(\overline{\Omega})$ . Due to the density result of Proposition 2.1 with respect to intermediate convergence (2.2), the inequality holds for all  $v \in BV(\Omega) \cap L^2(\Omega)$ . Consequently,  $-\operatorname{div} \mu^* \in \partial(|\cdot|_{BV(\Omega)})(u^*) \subset (BV(\Omega) \cap L^2(\Omega))^*$ . Using the chain rule as in Theorem 2.5, we find

$$-\operatorname{div} \mu^* \in -\operatorname{div} \left(\partial \|\cdot\|_{\mathcal{M}(\Omega)}(\nabla u^*)\right),$$

which proves the existence of  $\lambda^*$  with the claimed properties.

#### **3.3** Convergence of iterates

We are now going to apply the results of the previous two sections to the iterates of Algorithm 3.3. In terms of (3.10), we have to set  $g_k := \nabla f(u_k)$ . As can be seen from, e.g., Theorem 3.12, the boundedness of  $(\nabla f(u_k))$  in  $L^2(\Omega)$  will be crucial for any convergence analysis. Unfortunately, this boundedness can only be guaranteed in exceptional cases. Here, we prove it under the assumption that  $\nabla f$  is globally Lipschitz continuous. In Section 3.4 we show that convexity of f or global optimality of  $u_k$  is sufficient.

**Lemma 3.13.** Let  $\nabla f : L^2(\Omega) \to L^2(\Omega)$  be globally Lipschitz continuous with modulus  $L_f$ . Assume  $\rho_k \to \infty$ . Then  $(u_k)$  and  $(\nabla f(u_k))$  are bounded in  $L^2(\Omega)$ .

*Proof.* Due to the Lipschitz continuity of  $\nabla f$ , we have

$$\|\nabla f(u_k)\|_{L^2(\Omega)} \le L_f \|u_k\|_{L^2(\Omega)} + \|\nabla f(0)\|_{L^2(\Omega)}.$$

By Corollary 3.7, we find for k sufficiently large

$$\begin{aligned} \|u_k\|_{L^2(\Omega)} &\leq 2\rho_k^{-1} \|\nabla f(u_k)\|_{L^2(\Omega)} + \|\max(|u_a|, |u_b|)\|_{L^2(\Omega)} \\ &\leq 2\rho_k^{-1} L_f \|u_k\|_{L^2(\Omega)} + 2\rho_k^{-1} \|\nabla f(0)\|_{L^2(\Omega)} + \|\max(|u_a|, |u_b|)\|_{L^2(\Omega)}. \end{aligned}$$

If k is such that  $2\rho_k^{-1}L_f < \frac{1}{2}$ , then  $||u_k||_{L^2(\Omega)} \le 4\rho_k^{-1}||\nabla f(0)||_{L^2(\Omega)} + 2||\max(|u_a|, |u_b|)||_{L^2(\Omega)}$ , which proves the claim.

The next observation is a simple consequence of previous results and shows the close relation between boundedness of  $(u_k)$  in  $L^2(\Omega)$  and  $BV(\Omega)$  and the boundedness of  $(\nabla f(u_k))$  in  $L^2(\Omega)$ .

**Lemma 3.14.** Assume  $\rho_k \to \infty$ . Then the following statements are equivalent:

- (1)  $(\nabla f(u_k))$  is bounded in  $L^2(\Omega)$ ,
- (2)  $(u_k)$  is bounded in  $L^2(\Omega)$  and  $BV(\Omega)$ ,
- (3)  $\{u_k: k \in \mathbb{N}\}$  is pre-compact in  $L^2(\Omega)$ .

Proof. (1)  $\Rightarrow$  (2): The boundedness of  $(u_k)$  is a direct consequence of Corollary 3.7 and Lemma 3.8. (2)  $\Rightarrow$  (3) follows from Proposition 2.1 and Lemma 3.9. (3)  $\Rightarrow$  (1): Since  $u \mapsto \nabla f(u)$  is continuous from  $L^2(\Omega)$  to  $L^2(\Omega)$  by (A4), the set { $\nabla f(u_k) : k \in \mathbb{N}$ } is pre-compact in  $L^2(\Omega)$  and thus bounded.

Similarly to Theorem 3.10, we have the following result on the existence of converging subsequences.

**Theorem 3.15.** Suppose  $\epsilon_k \searrow 0$  and  $\rho_k \to \infty$ . Let  $(u_k)$  solve (3.5). Assume that  $(\nabla f(u_k))$  is bounded in  $L^2(\Omega)$ . Then there is a subsequence such that  $u_{k_n} \to u^*$ ,  $\lambda^a_{\rho_{k_n}}(u_{k_n}) \rightharpoonup \lambda^a$ , and  $\lambda^b_{\rho_{k_n}}(u_{k_n}) \rightharpoonup \lambda^b$  in  $L^2(\Omega)$ .

*Proof.* This result can be proven with similar arguments as Theorem 3.10.

We finally arrive at the following convergence result for iterates of Algorithm 3.3 which is a consequence of Theorem 3.12.

**Theorem 3.16.** Assume (A1)-(A5). Suppose  $\epsilon_k \searrow 0$  and  $\rho_k \to \infty$ . Let  $(u_k)$  solve (3.5). Assume that there is a subsequence with  $u_{k_n} \to u^*$ ,  $\lambda^a_{k_n} \to \lambda^a$ , and  $\lambda^b_{k_n} \to \lambda^b$  in  $L^2(\Omega)$ . Then it holds

$$\lambda^a \ge 0, \qquad \lambda^b \ge 0,$$
  
 $(\lambda^a, u_a - u^*) = 0, \quad (\lambda^b, u^* - u_b) = 0.$ 

In addition, there is  $\mu^* \in L^{\infty}(\Omega)^n$  with div  $\mu^* \in L^2(\Omega)$  such that

$$-\operatorname{div} \mu^* - \lambda^a + \lambda^b = -\nabla f(u^*)$$

and

$$-\operatorname{div} \mu^* \in \partial(|\cdot|_{BV(\Omega)})(u^*).$$

Moreover, there is  $\lambda^* \in \partial \| \cdot \|_{\mathcal{M}(\Omega)}(\nabla u^*) \subset (\mathcal{M}(\Omega)^n)^*$  with div  $\lambda^* \in L^2(\Omega)$  such that

$$-\operatorname{div}\lambda^* - \lambda^a + \lambda^b = g.$$

*Proof.* By assumption, we have  $\nabla f(u_{k_n}) \to \nabla f(u^*)$ . The proof is now a direct consequence of Theorem 3.12.

#### 3.4 Global solutions

The next theorem shows that global optimality is sufficient to obtain boundedness of iterates. We note that if f is convex, solutions of (3.5) are global solutions to the penalized problem  $(P_{\epsilon,\rho})$ .

**Theorem 3.17.** Assume (A1)-(A5). Suppose  $\epsilon_k \searrow 0$  and  $\rho_k \to \infty$ . Suppose  $(u_k)$  is the corresponding sequence of global solutions to the penalized problems  $(P_{\epsilon,\rho})$ . Then  $(u_k)$  is bounded in  $BV(\Omega) \cap L^2(\Omega)$ .

*Proof.* We introduce the notation

$$j_{\epsilon,\rho}(u) := f(u) + \int_{\Omega} \psi_{\epsilon}(\nabla u) \,\mathrm{d}x + \int_{\Omega} \frac{1}{\rho} \left( M_{\rho}(\rho(u_a - u)) + M_{\rho}(\rho(u - u_b)) \right) \,\mathrm{d}x.$$

Set  $\tilde{u} := \frac{1}{2}(u_a + u_b) \in H^1(\Omega)$ . Then  $j_{\epsilon_k,\rho_k}(\tilde{u}) = f(\tilde{u})$  for  $\rho_k$  large enough. Let  $u_k$  be a global minimizer of  $j_{\epsilon_k,\rho_k}$ . This implies

$$f(u_k) + \int_{\Omega} \psi_{\epsilon_k}(\nabla u_k) \,\mathrm{d}x + \int_{\Omega} \frac{1}{\rho_k} \left( M_{\rho_k}(\rho_k(u_a - u)) + M_{\rho_k}(\rho_k(u - u_b)) \right) \,\mathrm{d}x \le f(\tilde{u}).$$

Since f is bounded from below, there is K > 0 such that

$$\int_{\Omega} \psi_{\epsilon_k}(\nabla u_k) \,\mathrm{d}x + \int_{\Omega} \frac{1}{\rho_k} \left( M_{\rho_k}(\rho_k(u_a - u)) + M_{\rho_k}(\rho_k(u - u_b)) \right) \,\mathrm{d}x \le K.$$

This proves that  $(\nabla u_k)$  is bounded in  $L^1(\Omega)$  by Lemma 3.1(2). By construction, we have

$$M_{\rho}(x) = \int_{-\infty}^{x} \max_{\rho}(t) \, \mathrm{d}t \ge \int_{-\infty}^{x} \max(t, 0) \, \mathrm{d}t = \frac{1}{2} \max(0, x)^{2}$$

This implies

$$\frac{\rho_k}{2} \left( \|(u_k - u_b)_+\|_{L^2(\Omega)}^2 + \|(u_a - u)_+\|_{L^2(\Omega)}^2 \right) \le K,$$

and the boundedness of  $(u_k)$  in  $L^2(\Omega)$  is now a consequence of identity (3.9).

# 4 Optimality condition by regularization

Let us assume  $\bar{u} \in BV(\Omega) \cap L^2(\Omega)$  is locally optimal to (1.1). In this section, we want to show that there is a sequence of solutions  $(u_{\rho,\epsilon})$  of certain regularized problems converging to  $\bar{u}$ . This will allow us to prove optimality conditions for  $\bar{u}$  that are similar to the systems obtained in Theorems 3.12 and 3.16. Again, we work under the assumptions (A1)–(A5).

The solution  $\bar{u}$  satisfies the necessary optimality condition

$$-\nabla f(\bar{u}) \in \partial \left( |\cdot|_{BV(\Omega)} \right) (\bar{u}) + N_{U_{ad}}(\bar{u}) \text{ in } (BV(\Omega) \cap L^2(\Omega))^*, \tag{4.1}$$

see also Theorem 2.5. It is easy to see that (4.1) implies that  $\bar{u}$  is the unique solution to the linearized, strictly convex problem

$$\min_{u \in BV(\Omega) \cap L^2(\Omega)} \nabla f(\bar{u}) \cdot u + |u|_{BV(\Omega)} + \frac{1}{2} ||u - \bar{u}||_{L^2(\Omega)}^2 + \delta_{U_{ad}}(u).$$
(4.2)

In fact, let  $u^* \in BV(\Omega) \cap L^2(\Omega)$  be the solution of (4.2). Then we have the following optimality condition

$$-\nabla f(\bar{u}) \in \partial\left(|\cdot|_{BV(\Omega)}\right)(u^*) + (u^* - \bar{u}) + N_{U_{ad}}(u^*) \text{ in } (BV(\Omega) \cap L^2(\Omega))^*,$$

which is satisfied by  $u^* := \bar{u}$ . Let us approximate (4.2) by the family of unconstrained convex problems

$$\min_{u \in H^1(\Omega)} \nabla f(\bar{u}) \cdot u + \int_{\Omega} \psi_{\epsilon}(\nabla u) \, \mathrm{d}x + \frac{1}{2} \|u - \bar{u}\|_{L^2(\Omega)}^2 + \frac{1}{\rho} \int_{\Omega} M_{\rho}(\rho(u_a - u)) + M_{\rho}(\rho(u - u_b)) \, \mathrm{d}x.$$
(4.3)

The optimality condition for the unique solution  $u_{\epsilon,\rho}$  to (4.3) is given by

$$\int_{\Omega} \nabla f(\bar{u})v + \psi'_{\epsilon}(\nabla u_{\epsilon,\rho})\nabla v + (u_{\epsilon,\rho} - \bar{u})v - \lambda^{a}_{\rho}(u_{\epsilon,\rho})v + \lambda^{b}_{\rho}(u_{\epsilon,\rho})v \,\mathrm{d}x = 0.$$
(4.4)

for all  $v \in H^1(\Omega)$ .

**Corollary 4.1.** Suppose  $\epsilon_k \searrow 0$  and  $\rho_k \to \infty$ . Suppose  $(u_k)$  is the corresponding sequence of global solutions to the penalized problems (4.3). Then  $(u_k)$  is bounded in  $BV(\Omega) \cap L^2(\Omega)$ .

*Proof.* The claim follows by a similar argumentation as in the proof of Theorem 3.17.  $\Box$ 

**Lemma 4.2.** Suppose  $\epsilon_k \searrow 0$  and  $\rho_k \to \infty$ . Let  $(u_k)$  be the corresponding sequence of global solutions to the penalized problems (4.3). Then  $u_k \to \bar{u}$  in  $L^2(\Omega)$ , and the sequences  $(\lambda^a_{\rho_k}(u_k))$  and  $(\lambda^b_{\rho_k}(u_k))$  are bounded in  $L^2(\Omega)$ .

Proof. Due to Corollary 4.1,  $(u_k)$  is bounded in  $BV(\Omega) \cap L^2(\Omega)$ . Suppose for the moment  $u_k \to u^*$  in  $L^1(\Omega)$  and  $u_k \rightharpoonup u^*$  in  $L^2(\Omega)$ . By Lemma 3.9 applied to  $g_k := -\nabla f(\bar{u}) - (u_k - \bar{u})$  and  $g := -\nabla f(\bar{u}) - (u^* - \bar{u})$ , we obtain  $u_k \to u^*$  in  $L^2(\Omega)$ . By Theorem 3.10, the corresponding sequences  $(\lambda^a_{\rho_k}(u_k))$  and  $(\lambda^b_{\rho_k}(u_k))$  are bounded in  $L^2(\Omega)$ . Suppose  $\lambda^a_{\rho_k}(u_k) \rightharpoonup \lambda^a$  and  $\lambda^b_{\rho_k}(u_k) \rightharpoonup \lambda^b$  in  $L^2(\Omega)$ . By Theorem 3.12, we have

$$-\nabla f(\bar{u}) - (u^* - \bar{u}) + \lambda^a - \lambda^b \in \partial(|\cdot|_{BV(\Omega)})(u^*).$$

Due to the complementarity conditions (3.12) of Theorem 3.12, we get

$$-\nabla f(\bar{u}) - (u^* - \bar{u}) \in \partial \left( |\cdot|_{BV(\Omega)} + \delta_{U_{ad}} \right) (u^*).$$

Since  $-\nabla f(\bar{u}) \in \partial \left( |\cdot|_{BV(\Omega)} + \delta_{U_{ad}} \right) (\bar{u})$ , we have by the monotonicity of the subdifferential

$$(-\nabla f(\bar{u}) - (u^* - \bar{u}) - (-\nabla f(\bar{u})), \ u^* - \bar{u}) \ge 0,$$

which implies  $u^* = \bar{u}$ .

With similar arguments, we can show that every subsequence of  $(u_k)$  contains another subsequence that converges in  $L^2(\Omega)$  to  $\bar{u}$ . Hence, the convergence of the whole sequence follows.  $\Box$ 

This convergence result enables us to prove that  $\bar{u}$  satisfies an optimality condition similar to those of Theorems 3.12 and 3.16.

**Theorem 4.3.** Assume (A1)-(A5). Let  $\bar{u}$  be locally optimal for (1.1). Then there is

$$\lambda^* \in \partial \| \cdot \|_{\mathcal{M}(\Omega)}(\nabla \bar{u}) \subset (\mathcal{M}(\Omega)^n)^*$$

with div  $\lambda^* \in L^2(\Omega)$  such that

$$-\operatorname{div}\lambda^* - \lambda^a + \lambda^b = \nabla f(\bar{u})$$

and

$$\lambda^a \ge 0, \qquad \lambda^b \ge 0,$$
$$(\lambda^a, u_a - u^*) = 0, \quad (\lambda^b, u^* - u_b) = 0.$$

*Proof.* We define  $(u_k)$  as global solutions to the penalized problems (4.3) to parameter sequences  $\epsilon_k \searrow 0$  and  $\rho_k \to \infty$ . Due to Lemma 4.2, we have  $u_k \to \bar{u}$  in  $L^2(\Omega)$ . Define  $g_k := -\nabla f(\bar{u}) - (u_k - \bar{u})$  and  $g := -\nabla f(\bar{u})$ . Now, the claim follows by Theorem 3.12.

Clearly, the optimality conditions of Theorem 4.3 are stronger than those of Theorem 2.5. However, the proofs above only work on the strong assumptions that the bounds  $u_a$  and  $u_b$  are constant functions. Here, it is not clear to us, under which assumptions the above techniques carry over to non-constant  $u_a$  and  $u_b$ .

# 5 Numerical tests

In this section, the suggested algorithm is tested with selected examples. To this end, we implemented Algorithm 3.3 in python using FEnicCS, [15]. Our examples are carried out in the optimal control setting. In particular, f is given by the reduced tracking type functional

$$f(u) := \frac{1}{2} \|S(u) - y_d\|_{L^2(\Omega)}^2,$$

where S is the weak solution operator of some elliptic partial differential equation (PDE) specified below. To solve the partial differential equation, the domain is divided into a regular triangular mesh, and the PDE as well as the control are discretized with piecewise linear finite elements. If not mentioned otherwise, the computations are done on a 128 by 128 grid, which results in a mesh size of h = 0.022.

Let us define  $j_{\epsilon,\rho}: H^1(\Omega) \to \mathbb{R}$  by

$$j_{\epsilon,\rho}(u) := f(u) + \int_{\Omega} \psi_{\epsilon}(\nabla u) \, \mathrm{d}x + \int_{\Omega} \frac{1}{\rho} \left( M_{\rho}(\rho(u_a - u)) + M_{\rho}(\rho(u - u_b)) \right) \, \mathrm{d}x$$

with  $M_{\rho}$  as defined in (3.4). It is given in our tests by the specific choice

$$M_{\rho}(x) := \begin{cases} \frac{1}{2}x^2 + \frac{1}{24\rho^2} & \text{if } x > \frac{1}{2\rho}, \\ \frac{\rho}{6}(x + \frac{1}{2\rho})^3 & \text{if } |x| < \frac{1}{2\rho}, \\ 0 & \text{otherwise.} \end{cases}$$

let us recall that we use the following function to approximate the BV-seminorm:

$$\psi_{\epsilon} = \sqrt{\epsilon + t^2} + \epsilon t^2.$$

Concerning the continuation strategy for the parameters  $\epsilon$  and  $\rho$  in Algorithm 3.3, we set  $\epsilon_0 := 0.5$  in the initialization and decrease  $\epsilon$  by factor 0.5 after each iteration. The penalty parameter is increased by factor 2 after every iteration and is initialized with  $\rho_0 := 2$ . Algorithm 3.3 is stopped if the following termination criterion is satisfied:

$$R_k^{\rho} \le 10^{-4} \text{ and } R_k^{\epsilon} \le 10^{-3},$$
(5.1)

where the residuals  $R_k^{\rho}$ ,  $R_k^{\epsilon}$  are given by

$$R_k^{\rho} := \|(u_a - u_k)_+\|_{L^2(\Omega)} + \|(u_k - u_b)_+\|_{L^2(\Omega)} + (\lambda_k^a, u_a - u_k) + (\lambda_k^b, u_k - u_b).$$

and

$$R_k^{\epsilon} := \|\nabla u_k\|_{L^1(\Omega)} - \langle \mu_k, \nabla u_k \rangle$$

with  $\mu_k := \frac{\nabla u_k}{\sqrt{\epsilon_k + |\nabla u_k|^2}}$  as in the proof of Theorem 3.12. Here, the residuum  $R_k^{\rho}$  measures the violation of the box-constraints and of the complementarity condition in Theorem 3.16. Let us discuss the choice of the residuum  $R^{\epsilon}$ . It can be interpreted as a residual in the subgradient inequality. Since  $\|\mu_k\|_{L^{\infty}(\Omega)^n} \leq 1$ , we have  $\int_{\Omega} \mu_k \cdot \nabla v \, dx \leq \|\nabla v\|_{L^1(\Omega)}$  for  $v \in W^{1,1}(\Omega)$ . Hence,  $R_k^{\epsilon} \geq 0$ . This implies

$$\langle \mu_k, \nabla v - \nabla u_k \rangle \leq \|\nabla v\|_{L^1(\Omega)} - \|\nabla u_k\|_{L^1(\Omega)} + R_k^{\epsilon} \quad \forall v \in W^{1,1}(\Omega).$$

Hence,  $\mu_k$  can be interpreted is an element of the  $\varepsilon$ -subdifferential to the error level  $\varepsilon := R_k^{\epsilon}$ .

#### 5.1 Globalized Newton Method for the subproblems

To solve the variational subproblems of form  $(P_{\epsilon,\rho})$ , i.e.,

$$\min_{u\in H^1(\Omega)} j_{\epsilon,\rho},$$

we use a globalized Newton method. Let us recall the notation

$$\lambda^{a}(u) := \max_{\rho}(\rho(u_{a} - u)),$$
$$\lambda^{b}(u) := \max_{\rho}(\rho(u - u_{b})).$$

and introduce

$$\Lambda^{a}(u) := -\max_{\rho}'(\rho(u_{a} - u)) = \begin{cases} -1 & \text{if } \rho(u_{a} - u) > \frac{1}{2\rho}, \\ -\rho\left(\rho(u_{a} - u) + \frac{1}{2\rho}\right) & \text{if } |\rho(u_{a} - u)| < \frac{1}{2\rho}, \\ 0 & \text{otherwise}, \end{cases}$$

and

$$\Lambda^{b}(u) := \max_{\rho}'(\rho(u - u_{b})) = \begin{cases} 1 & \text{if } \rho(u - u_{b}) > \frac{1}{2\rho}, \\ \rho\left(\rho(u - u_{b}) + \frac{1}{2\rho}\right) & \text{if } |\rho(u - u_{b})| < \frac{1}{2\rho}, \\ 0 & \text{otherwise.} \end{cases}$$

The Newton method with a line search strategy is given as follows:

Algorithm 5.1 (Global Newton method). Set  $k = 0, \rho > 0, \epsilon \in (0, 1), u_a, u_b \in \mathbb{R}, \eta > 0, p > 2, \phi \in (0, 1), \tau \in (0, \frac{1}{2})$ . Choose  $u_0 \in H^1(\Omega)$ .

1. Compute the search direction  $w_k$  by solving

$$j''(u_k)w = -\nabla j(u_k)(u), \qquad (5.2)$$

where

$$\nabla j(u) := -\operatorname{div}(\psi_{\epsilon}(\nabla u)) + \nabla f(u) - \max_{\rho}(\rho(u_a - u)) + \max_{\rho}(\rho(u - u_b))$$

and

$$j''(u)d = -\operatorname{div}\left(\psi_{\epsilon}''(\nabla u)\nabla d\right) + f''(u)d - \rho\Lambda^{a}(u)d + \rho\Lambda^{b}(u)d.$$

If  $\nabla j(u_k) \cdot w_k \leq -\eta ||w_k||^p$ : set  $w_k := -\nabla j(u_k)$ .

2. (line search) Find  $\sigma_k := \max\{\phi^l : l = 0, 1, 2, ...\}$  such that

$$j(u_k + \sigma_k w_k) - j_k(u_k) \le \tau \sigma_k \nabla J(u_k) \cdot w_k$$

- 3. Set  $u_{k+1} := u_k + \sigma w_k$ .
- 4. If a suitable stopping criteria is satisfied: Stop.
- 5. Set k := k + 1 and go to step 1.

Let us provide details regarding the implementation of Algorithm (5.1). In the initialization of Algorithm 5.1, we set  $\phi = 0.5$ ,  $\tau = 10^{-4}$ ,  $\eta = 10^{-8}$  and p = 2.1. In addition, we employed the following termination criterion:

If 
$$||u_{k+1} - u_k||_{L^2(\Omega)} + ||y_{k+1} - y_k||_{L^2(\Omega)} + ||p_{k+1} - p_k||_{L^2(\Omega)} < 10^{-10}$$
: Stop

That is, if there is no sufficient change between consecutive iterates, we assume that the method resulted in a stationary (minimal) point of the subproblem.

#### 5.2 Example 1: linear elliptic PDE

First, we consider the optimal control problem

$$\min_{u \in BV(\Omega)} \frac{1}{2} \|y - y_d\|_{L^2(\Omega)}^2 + \beta |u|_{BV(\Omega)}$$
(5.3)

subject to

 $-\Delta y = u \text{ on } \Omega, \ y = 0 \text{ on } \partial \Omega$ 

and the box constraints

 $u_a \leq u(x) \leq u_b$  f.a.a.  $x \in \Omega$ .

Note that (5.3) as well as the subproblem

$$\min_{u \in BV(\Omega)} J_{\epsilon,\rho}(y,u) := \frac{1}{2} \|y - y_d\|_{L^2(\Omega)}^2 + \beta \int_{\Omega} \psi_{\epsilon}(\nabla u) \, \mathrm{d}x \\
+ \int_{\Omega} \frac{1}{\rho} \left( M_{\rho}(\rho(u_a - u)) + M_{\rho}(\rho(u - u_b)) \right) \, \mathrm{d}x \qquad (5.4)$$
s.t.  $-\Delta y = u \text{ on } \Omega, y = 0 \text{ on } \partial\Omega,$ 

are convex and uniquely solvable. Let us introduce the adjoint state  $p \in H_0^1(\Omega)$  as the solution of the partial differential equation

$$-\Delta p = y - y_d$$
 on  $\Omega$ ,  $y = 0$  on  $\partial \Omega$ .

Applying Algorithm 5.1 to the reduced functional of problem (5.4) results in the following system of equations that has to be solved in each Newton step:

$$G(y, p, u)(\delta y, \delta p, \delta u) = F(y, p, u),$$

where F is given by

$$F(y, p, u) := \begin{pmatrix} -\Delta y - u \\ -\Delta p - (y - y_d) \\ p - \beta \operatorname{div}(\psi'_{\epsilon}(\nabla u)) - \max_{\rho}(\rho(u_a - u)) + \max_{\rho}(\rho(u - u_b)) \end{pmatrix}$$

The equation F(y, p, u) = 0 is the optimality system to problem (5.4). The derivative of F in direction  $(\delta y, \delta p, \delta u) \in H_0^1(\Omega) \times H_0^1(\Omega) \times H^1(\Omega)$  is given by

$$G(y, p, u)(\delta y, \delta p, \delta u) = \begin{pmatrix} -\Delta \delta y - \delta u \\ -\Delta \delta p - \delta y \\ \delta p - \beta \operatorname{div}(\psi_{\epsilon}''(\nabla u))\nabla \delta u - \rho \Lambda^a \delta u + \rho \Lambda^b \delta u. \end{pmatrix}$$



Figure 1: Optimal control u.

The solution of the Newton step (5.2) is then given by  $w := \delta u$ . We adapt the example problem data from [11]. Here,  $\Omega = [-1, 1]^2$  and

$$y_d := \begin{cases} 1 & \text{on } (-0.5, 0.5)^2 \\ 0 & \text{otherwise} \end{cases}$$

In the computations we set  $-u_a = u_b = 10$  and  $\beta = 0.0001$ . This example (without additional box constraints) was also used in [12].

In Table 1, we see the convergence behavior of iterates. Here, the errors

$$E_u := ||u_k - u_{ref}||_{L^2(\Omega)}, \quad E_J := |J_k - J_{ref}|$$

are presented, where  $u_{ref}$  is the final iterate after the algorithm terminated at step k = 19 and  $J_{ref} := J(u_{ref})$ . Furthermore, we observe  $R^{\epsilon} = O(\sqrt{\epsilon})$  and  $R^{\rho} = O(\frac{1}{\rho})$  as  $\epsilon \sim 2^{-k}$  and  $\rho_k \sim 2^k$ .

| k  | $E_u$ | $E_J$               | $R_k^\epsilon$      | $R_k^{ ho}$          |
|----|-------|---------------------|---------------------|----------------------|
| 12 | 1.11  | $6.0 \cdot 10^{-4}$ | $8.0 \cdot 10^{-3}$ | $1.3 \cdot 10^{-9}$  |
| 13 | 0.80  | $3.5\cdot10^{-4}$   | $5.9\cdot10^{-3}$   | $6.7 \cdot 10^{-10}$ |
| 14 | 0.56  | $1.9 \cdot 10^{-4}$ | $4.2 \cdot 10^{-3}$ | $3.4 \cdot 10^{-10}$ |
| 15 | 0.34  | $1.0\cdot10^{-4}$   | $3.0\cdot10^{-3}$   | $1.7\cdot10^{-10}$   |
| 16 | 0.17  | $5.5\cdot10^{-5}$   | $2.1\cdot 10^{-3}$  | $8.3\cdot10^{-11}$   |
| 17 | 0.07  | $2.4\cdot 10^{-5}$  | $1.5\cdot 10^{-3}$  | $4.2\cdot10^{-11}$   |
| 18 | 0.02  | $8.2 \cdot 10^{-6}$ | $1.1 \cdot 10^{-3}$ | $2.1 \cdot 10^{-11}$ |
| 19 |       |                     | $7.6\cdot10^{-4}$   | $1.1 \cdot 10^{-11}$ |

Table 1: Computed errors during the final iterations.

Figure 1 shows the optimal control. The result is in agreement with the results obtained in [12]. In Figure 2 the computed optimal controls are depicted for the unconstrained case (left), i.e., constraints are inactive during the computation process. The right plot shows the optimal control u, when lower and upper bound are set to  $u_a = -5$  and  $u_b = 18$ .



Figure 2: Optimal control u for different choices of  $u_a, u_b$ .

#### 5.3 Example 2: Semilinear elliptic optimal control problem

Let us now consider the following problem with semilinear state equation. That is, we study the minimization problem

$$\min_{u \in U_{-d}} f_{sl}(u) + \beta |u|_{BV(\Omega)},$$

where  $f_{sl}$  is given by the standard tracking type functional  $u \mapsto ||y_u - y_d||^2_{L^2(\Omega)}$ , and  $y_u$  is the weak solution of the semilinear elliptic state equation

$$-\Delta y + y^3 = u$$
 in  $\Omega$ ,  $y = 0$  on  $\partial \Omega$ .

The adjoint state  $p \in H_0^1$  is given now as solution to the equation

$$-\Delta p + 3y^2 p = y - y_d$$
 in  $\Omega$ ,  $p = 0$  on  $\partial \Omega$ .

For this example, system of equations (5.2) for the state  $y \in H_0^1(\Omega)$ , the adjoint state  $p \in H_0^1(\Omega)$ and the control variable  $u \in H^1(\Omega)$  is given by

$$G(y, p, u)(\delta y, \delta p, \delta u) = F(y, p, u)$$

with

$$F(y, p, u) := \begin{pmatrix} -\Delta y + y^3 - u \\ -\Delta p + 3y^2 p - (y - y_d) \\ p - \beta \operatorname{div}(\psi'_{\epsilon}(\nabla u)) - \max_{\rho}(\rho(u_a - u)) + \max_{\rho}(\rho(u - u_b)) \end{pmatrix}.$$

The derivative in direction  $(\delta y, \delta p, \delta u) \in H^1_0(\Omega) \times H^1_0(\Omega) \times H^1(\Omega)$  is given with

$$G(y, p, u)(\delta y, \delta p, \delta u) = \begin{pmatrix} -\Delta \delta y + 3y^2 \delta y - \delta u \\ -\Delta \delta p + 3y^2 \delta p - \delta y + 6yp \delta y \\ \delta p - \beta \operatorname{div}(\psi_{\epsilon}''(\nabla u)) \nabla \delta u - \rho \Lambda^a \delta u + \rho \Lambda^b \delta u. \end{pmatrix}$$

The data is given as in Example 1. The optimal control is depicted in Figure 3. It is close to the solution of Example 1. Let us consider the performance of the algorithm on different levels of discretization for this example. Table 2 shows the number of outer iterations ( $\sharp$ it), as well as the total number of newton iterations ( $\sharp$ newt) needed until the stopping criterion (5.1) holds for increasing meshsizes. The last column shows the final objective value  $J_{\epsilon,\rho}$ . The residuals  $R^{\epsilon}$ and  $R^{\rho}$  behaved as in Example 1.

| h     | ‡it | #newt | $\epsilon_{final}$ | $ ho_{final}$ | $J_{\epsilon,\rho}$ |
|-------|-----|-------|--------------------|---------------|---------------------|
| 0.088 | 16  | 182   | $2^{-16}$          | $2^{16}$      | 0.0596              |
| 0.044 | 19  | 201   | $2^{-19}$          | $2^{19}$      | 0.0685              |
| 0.022 | 19  | 314   | $2^{-19}$          | $2^{19}$      | 0.0737              |
| 0.011 | 19  | 486   | $2^{-19}$          | $2^{19}$      | 0.0767              |

Table 2: Number of iterations and newton steps for different mesh-sizes.



Figure 3: Optimal control u for the semilinear problem.

#### 5.4 Experiments with non-constant constraints

So far our analysis and numerical experiments are restricted to the case where  $u_a, u_b$  are constant functions. This assumption was needed to show the boundedness of multipliers  $\lambda_k^a(u), \lambda_k^b(u)$  in  $L^2(\Omega)$  in Lemma 3.5, which is crucial for the final result Theorem 3.16. For this section we tested Algorithm 3.3 also for non-constant functions  $u_a, u_b \in L^{\infty}(\Omega)$ .

Here, we consider again the linear optimal control problem and data from Example 1 with different choices for  $u_a, u_b$ :

(i) 
$$u_a := -100, \ u_b(x_1, x_2) := 8\sin(\pi x_1)\sin(\pi x_2),$$
 (5.5)

(*ii*) 
$$u_a := -100, \ u_b(x_1, x_2) := -4(x_1 - 0.5)^2 - 4x_2^2 + 10,$$
 (5.6)

In Figure 4 the behavior of the quantity  $\|\lambda_k^a(u)\|_{L^2(\Omega)}^2 + \|\lambda_k^b(u)\|_{L^2(\Omega)}^2$  is plotted along the iterations, i.e., for increasing  $\rho_k$ , for different discretization levels. In Figure 5, the respective solution plots are shown. While the multipliers seem to be bounded for one example, their norm grows with  $\rho$  (and thus with  $\epsilon$ ) for the other example. Clearly, more research has to be done to develop necessary and sufficient conditions for the boundedness of the multipliers.

# Acknowledgement

The authors are grateful to Gerd Wachsmuth for an inspiring discussion that led to an improvement of Theorem 3.12 and subsequent results.



Figure 4: The  $L^2$ -norm of multipliers  $\lambda_k^a(u)$ ,  $\lambda_k^b(u)$  for Example (5.5) (left) and (5.6) (right).



Figure 5: Optimal control u for Example (5.5) (left) and (5.6) (right).

# References

- R. Acar and C. R. Vogel. Analysis of bounded variation penalty methods for ill-posed problems. *Inverse Problems*, 10(6):1217–1229, 1994.
- [2] H. Attouch, G. Buttazzo, and G. Michaille. Variational analysis in Sobolev and BV spaces, volume 6 of MPS/SIAM Series on Optimization. Society for Industrial and Applied Mathematics (SIAM), Philadelphia, PA; Mathematical Programming Society (MPS), Philadelphia, PA, 2006. Applications to PDEs and optimization.
- [3] S. Bartels. Total variation minimization with finite elements: convergence and iterative solution. *SIAM J. Numer. Anal.*, 50(3):1162–1180, 2012.
- [4] S. Bartels and M. Milicevic. Iterative finite element solution of a constrained total variation regularized model problem. *Discrete Contin. Dyn. Syst. Ser. S*, 10(6):1207–1232, 2017.
- [5] E. Casas, K. Kunisch, and C. Pola. Regularization by functions of bounded variation and applications to image enhancement. *Appl. Math. Optim.*, 40(2):229–257, 1999.
- [6] E. Casas, R. Herzog, and G. Wachsmuth. Optimality conditions and error analysis of semilinear elliptic control problems with L<sup>1</sup> cost functional. SIAM J. Optim., 22(3):795– 820, 2012.
- [7] E. Casas, F. Kruse, and K. Kunisch. Optimal control of semilinear parabolic equations by BV-functions. SIAM J. Control Optim., 55(3):1752–1788, 2017.

- [8] E. Casas and K. Kunisch. Analysis of optimal control problems of semilinear elliptic equations by BV-functions. Set-Valued Var. Anal., 27(2):355–379, 2019.
- [9] A. Chambolle, V. Caselles, D. Cremers, M. Novaga, and T. Pock. An introduction to total variation for image analysis. In *Theoretical foundations and numerical methods for sparse* recovery, volume 9 of Radon Ser. Comput. Appl. Math., pages 263–340. Walter de Gruyter, Berlin, 2010.
- [10] T. F. Chan, S. Esedoglu, F. E. Park, and A. M. Yip. Total variation image restoration: Overview and recent developments. In N. Paragios, Y. Chen, and O. D. Faugeras, editors, *Handbook of Mathematical Models in Computer Vision*, pages 17–31. Springer, 2006.
- [11] C. Clason and K. Kunisch. A duality-based approach to elliptic control problems in nonreflexive Banach spaces. ESAIM Control Optim. Calc. Var., 17(1):243–266, 2011.
- [12] D. Hafemeyer and F. Mannel. A path-following inexact Newton method for optimal control in BV, 2020.
- [13] D. Kinderlehrer and G. Stampacchia. An introduction to variational inequalities and their applications, volume 88 of Pure and Applied Mathematics. Academic Press, Inc. [Harcourt Brace Jovanovich, Publishers], New York-London, 1980.
- [14] K. Kunisch and D. Wachsmuth. Sufficient optimality conditions and semi-smooth Newton methods for optimal control of stationary variational inequalities. ESAIM Control Optim. Calc. Var., 18(2):520–547, 2012.
- [15] H. P. Langtangen and A. Logg. Solving PDEs in Python, volume 3 of Simula SpringerBriefs on Computing. Springer, Cham, 2016. The FEniCS tutorial I.
- [16] M. Milicevic. Finite Element Discretization and Iterative Solution of Total Variation Regularized Minimization Problems and Application to the Simulation of Rate-Independent Damage Evolutions. PhD thesis, Universität Freiburg, 2018.
- [17] L. I. Rudin, S. Osher, and E. Fatemi. Nonlinear total variation based noise removal algorithms. *Phys. D*, 60(1-4):259–268, 1992. Experimental mathematics: computational issues in nonlinear science (Los Alamos, NM, 1991).
- [18] C. Scheven and T. Schmidt. On the dual formulation of obstacle problems for the total variation and the area functional. Ann. Inst. H. Poincaré Anal. Non Linéaire, 35(5):1175– 1207, 2018.
- [19] A. Schiela and D. Wachsmuth. Convergence analysis of smoothing methods for optimal control of stationary variational inequalities with control constraints. *ESAIM Math. Model. Numer. Anal.*, 47(3):771–787, 2013.
- [20] F. Tröltzsch. Optimal control of partial differential equations, volume 112 of Graduate Studies in Mathematics. American Mathematical Society, Providence, RI, 2010. Theory, methods and applications, Translated from the 2005 German original by Jürgen Sprekels.