



Error estimates for the approximation of multibang control problems

Christian Clason¹  · Thi Bich Tram Do¹ · Frank Pörner²

Received: 12 March 2018 / Published online: 16 August 2018
© Springer Science+Business Media, LLC, part of Springer Nature 2018

Abstract

This work is concerned with optimal control problems where the objective functional consists of a tracking-type functional and an additional “multibang” regularization functional that promotes optimal control taking values from a given discrete set pointwise almost everywhere. Under a regularity condition on the set where these discrete values are attained, error estimates for the Moreau–Yosida approximation (which allows its solution by a semismooth Newton method) and the discretization of the problem are derived. Numerical results support the theoretical findings.

Keywords Multibang control · Moreau-Yosida approximation · Finite element discretization · Error estimates · Semi-smooth Newton method

1 Introduction

We consider linear-quadratic optimal control problems where the optimal control is only allowed to take values at discrete values $u_1 < \dots < u_d \in \mathbb{R}$ with $d \in \mathbb{N}$. Such problems occur, e.g., in topology optimization, nondestructive testing or medical imaging; a similar task also arises as a sub-step in segmentation or labeling problems in image processing. However, such problems are inherently nonconvex and, more importantly, not weakly lower semi-continuous and hence cannot be treated by standard techniques. A classical remedy is convex relaxation, where the nonconvex

✉ Christian Clason
christian.clason@uni-due.de

Thi Bich Tram Do
tram.do@uni-due.de

Frank Pörner
frank.poerner@mathematik.uni-wuerzburg.de

¹ Faculty of Mathematics, University Duisburg-Essen, 45117 Essen, Germany

² Department of Mathematics, University of Würzburg, Emil-Fischer-Str. 40, 97074 Würzburg, Germany

constraint $u(x) \in \{u_1, \dots, u_d\}$ is replaced by the convex constraint $u(x) \in [u_1, u_d]$, but this leads to ignoring the intermediate parameter values. In [3,5–8], it was therefore proposed to promote all desired control values using a convex *multibang* penalty

$$G(u) : L^2(\Omega) \rightarrow \mathbb{R}, \quad u \mapsto \int_{\Omega} g(u(x))dx,$$

for a suitable convex integrand $g : \mathbb{R} \rightarrow \mathbb{R}$ with a polyhedral epigraph whose vertices correspond to the desired control values u_1, \dots, u_d . We thus consider the *multibang control problem*

$$\min_{u \in L^2(\Omega)} \frac{1}{2} \|Ku - z\|_Y^2 + \alpha G(u) \tag{1.1}$$

with $\alpha > 0$, $z \in Y$ for a Hilbert space Y , and $K : L^2(\Omega) \rightarrow Y$ a linear and continuous operator (e.g., the solution operator for a linear elliptic partial differential equation). Just as in L^1 regularization for sparsity (and in linear optimization), it can be expected that minimizers are found at the vertices of G , thus yielding the desired structure. Furthermore, it was shown in [3,4,7] that this leads to a primal-dual optimality system that can be solved by a superlinearly convergent semismooth Newton method in function space [14,22] if a suitable Moreau–Yosida approximation (of the Fenchel conjugate G^* , see Proposition 2.2 below) is introduced. It turns out that this approximation can be expressed in primal form as

$$\min_{u \in L^2(\Omega)} \frac{1}{2} \|Ku - z\|_Y^2 + \alpha G(u) + \frac{\gamma}{2} \|u\|_{L^2(\Omega)}^2 \tag{1.2}$$

for a parameter $\gamma > 0$. We remark that this approach (i.e., applying the approximation to G^* instead of G) does not destroy the non-differentiability of G and hence preserves the structural properties of (1.1). Standard lower semicontinuity techniques can then be applied to show that the solutions to (1.2) converge weakly to the solution to (1.1) as $\gamma \rightarrow 0$; see [7, §4.1]. The aim of this paper is to establish strong convergence and in particular approximation error estimates for $\|\bar{u} - u_\gamma\|_{L^2(\Omega)}$.

Let us recall some literature and already known results. For the case $d = 2$ we obtain the minimization problem

$$\min_{u_1 \leq u \leq u_2} \frac{1}{2} \|Ku - z\|_Y^2. \tag{1.3}$$

and if the associated adjoint state $\bar{p}(x) \neq 0$ almost everywhere, it is well-known that \bar{u} exhibits a bang-bang structure, i.e. $\bar{u}(x) \in \{u_1, u_2\}$ almost everywhere. This problem has been studied intensively in the literature, see [20,21,23,24,26] and the references therein. Note that this list is far away from being complete. For this problem a structural assumption has been established in [24,26], which controls the behavior of the adjoint state around a singular set and guarantees that the optimal control \bar{u} exhibits a bang-bang structure. Using this assumption, error estimates for the approximation of (1.3)

can be proven; see [24]. A related question is the Moreau–Yosida approximation of state constraints; see [10,11].

If $d = 3$ and $u_1 < u_2 = 0 < u_3$, the problem (1.1) resembles the minimization problem

$$\min_{u_1 \leq u \leq u_2} \frac{1}{2} \|Ku - z\|_Y^2 + \alpha \|u\|_{L^1(\Omega)}, \tag{1.4}$$

see, e.g., [20]. The structural assumption used to prove error rates for the approximation of (1.3) can be generalized to problem (1.4). Again, approximation error estimates can be proven; see [23,24,26] and the reference therein.

We will generalize this structural assumption to the multibang control problem (1.1). We will show that this assumption is sufficient to guarantee that an optimal control \bar{u} of (1.1) satisfies $\bar{u}(x) \in \{u_1, \dots, u_d\}$ for almost all $x \in \Omega$. Furthermore, we will use this condition to prove approximation error estimates of the form

$$\|\bar{u} - u_\gamma\|_{L^2(\Omega)} = \mathcal{O}\left(\gamma^{\frac{\kappa}{2}}\right)$$

with a constant $\kappa > 0$ depending only on the structural assumption.

The paper is organized as follows. In Sect. 2 we recall some preliminary results which are needed for the convergence analysis. Our structural assumption is introduced in Sect. 3 and used to derive the approximation error estimates. This is also the main result of this paper. In Sect. 4, we establish discretization error estimates under our structural assumption. We introduce an active set method for the solution of (1.2) and show its equivalence to a semismooth Newton method in Sect. 5. Finally, numerical results to support our theoretical findings can be found in Sect. 6.

2 Preliminary results

Let $u_1 < u_2 < \dots < u_d$ be some given real numbers with $d \geq 2$, and let $\Omega \subset \mathbb{R}^n$ be a bounded domain. Following [3,5–7], we define the piecewise linear function

$$g(v) := \begin{cases} \frac{1}{2}((u_i + u_{i+1})v - u_i u_{i+1}) & \text{if } v \in [u_i, u_{i+1}], \quad 1 \leq i < d, \\ \infty & \text{else.} \end{cases}$$

As the pointwise supremum of affine functions, g is convex and continuous on the interior of its domain $\text{dom}(g) = [u_1, u_d]$. Hence, the corresponding integral functional

$$G : L^2(\Omega) \rightarrow \mathbb{R}, \quad u \mapsto \int_{\Omega} g(u(x)) dx,$$

is proper, convex and weakly lower semicontinuous as well; see, e.g., [2, Proposition 2.53].

We now consider the problem

$$\min_{u \in L^2(\Omega)} \frac{1}{2} \|Ku - z\|_Y^2 + \alpha G(u) \tag{2.1}$$

with a parameter $\alpha > 0$. Standard semi-continuity methods then yield existence of a minimizer \bar{u} , which is unique if K is injective; see [7]. We will later impose a condition which guarantees that \bar{u} exhibits a multibang structure, i.e., $\bar{u}(x) \in \{u_1, \dots, u_d\}$ for almost every $x \in \Omega$.

Let us further define the set

$$U_{\text{ad}} := \{u \in L^2(\Omega) : u_1 \leq u(x) \leq u_d\} = \text{co} \left\{ u \in L^2(\Omega) : u(x) \in \{u_1, \dots, u_d\} \right\},$$

where co denotes the convex hull. It is clear that (2.1) is equivalent to the problem

$$\min_{u \in U_{\text{ad}}} \frac{1}{2} \|Ku - z\|_Y^2 + \alpha G(u). \tag{P}$$

We will use this equivalent formulation to derive variational inequalities which will be useful in the convergence analysis. Standard convex analysis techniques then yield primal–dual optimality conditions; see, e.g., [3, 7].

Proposition 2.1 *Define the sets*

$$\begin{aligned} Q_1 &:= \left\{ q : q < \frac{\alpha}{2}(u_1 + u_2) \right\}, \\ Q_i &:= \left\{ q : \frac{\alpha}{2}(u_{i-1} + u_i) < q < \frac{\alpha}{2}(u_i + u_{i+1}) \right\}, \quad 1 < i < d, \\ Q_d &:= \left\{ q : q > \frac{\alpha}{2}(u_{d-1} + u_d) \right\}, \\ Q_{i,i+1} &:= \left\{ q : q = \frac{\alpha}{2}(u_i + u_{i+1}) \right\}. \end{aligned}$$

Let $\bar{u} \in U_{\text{ad}}$ with associated adjoint state $\bar{p} := K^*(z - K\bar{u})$. Then \bar{u} is a solution to (P) if and only if

$$\bar{u}(x) \in \begin{cases} \{u_i\} & \text{if } \bar{p}(x) \in Q_i \quad 1 \leq i \leq d, \\ [u_i, u_{i+1}] & \text{if } \bar{p}(x) \in Q_{i,i+1} \quad 1 \leq i < d. \end{cases} \tag{2.2}$$

It is clear that the optimal solution \bar{u} is uniquely determined by the adjoint state on the sets $\{x \in \Omega : \bar{p}(x) \in Q_i\}$. We see furthermore that $\bar{u}(x) \in \{u_1, \dots, u_d\}$ almost everywhere on Ω if $\text{meas}\{x \in \Omega : \bar{p}(x) \in Q_{i,i+1}\} = 0$ for all $1 \leq i < d$. Hence \bar{u} has a multibang structure in this case. In the following, we will make use of this relation to construct a suitable regularity condition on these sets.

Remark 2.1 Although the dependence of the optimal controls on α is not the focus of this work – see instead the earlier works [5–8], and, in particular, [3, Section 5] – let us recall the essential features for the sake of completeness. First, note that α enters the optimality conditions (2.2) only via the case distinction for the sets Q_i and $Q_{i,i+1}$. Specifically, increasing the value of α shifts the conditions on \bar{p} so that desired control values u_i of smaller magnitude are preferred. Conversely, for $\alpha \rightarrow 0$, these conditions coincide with the well-known optimality conditions for bang-bang control problems

where only $Q_1, Q_d,$ and $Q_{1,d}$ are relevant; see, e.g., [21, Lemma 2.26]. This implies that apart from singular cases where $\text{meas}\{x \in \Omega : \bar{p}(x) = c\} \neq 0$ for some $c \in \mathbb{R}$, the value of α does not influence the “strength” of the multibang penalty in enforcing the desired control values but only the specific selection among these values.

We next introduce the Moreau–Yosida approximation of (P) with a regularization parameter $\gamma > 0,$

$$\min_{u \in U_{\text{ad}}} \frac{1}{2} \|Ku - z\|_Y^2 + \alpha G(u) + \frac{\gamma}{2} \|u\|_{L^2(\Omega)}^2. \tag{P_\gamma}$$

As for (P), arguments from convex analysis lead to the following optimality conditions; see [3,7].

Proposition 2.2 *Define the sets*

$$\begin{aligned} Q_1^\gamma &:= \left\{ q : q < \frac{\alpha}{2} \left(\left(1 + 2\frac{\gamma}{\alpha} \right) u_1 + u_2 \right) \right\}, \\ Q_i^\gamma &:= \left\{ q : \frac{\alpha}{2} \left(u_{i-1} + \left(1 + 2\frac{\gamma}{\alpha} \right) u_i \right) < q < \frac{\alpha}{2} \left(\left(1 + 2\frac{\gamma}{\alpha} \right) u_i + u_{i+1} \right) \right\}, \\ Q_{i,i+1}^\gamma &:= \left\{ q : \frac{\alpha}{2} \left(\left(1 + 2\frac{\gamma}{\alpha} \right) u_i + u_{i+1} \right) \leq q \leq \frac{\alpha}{2} \left(u_i + \left(1 + 2\frac{\gamma}{\alpha} \right) u_{i+1} \right) \right\}, \\ Q_d^\gamma &:= \left\{ q : \frac{\alpha}{2} \left(u_{d-1} + \left(1 + 2\frac{\gamma}{\alpha} \right) u_d \right) < q \right\}. \end{aligned}$$

Let $u_\gamma \in U_{\text{ad}}$ with associated adjoint state $p_\gamma := K^*(z - Ku_\gamma).$ Then u_γ is a solution to (P_\gamma) if and only if

$$u_\gamma(x) = \begin{cases} u_i & \text{if } p_\gamma(x) \in Q_i^\gamma \quad 1 \leq i \leq d, \\ \frac{1}{\gamma} (p_\gamma(x) - \frac{\alpha}{2}(u_i + u_{i+1})) & \text{if } p_\gamma(x) \in Q_{i,i+1}^\gamma \quad 1 \leq i < d. \end{cases} \tag{2.3}$$

We remark that (2.3) is the explicit pointwise characterization of $u_\gamma \in (\partial G^*)_\gamma(p_\gamma),$ where $(\partial G^*)_\gamma$ denotes the Yosida approximation of the convex subdifferential (which coincides with the Fréchet derivative of the Moreau envelope) of the Fenchel conjugate of $G,$ which justifies the term *Moreau–Yosida approximation;* see, e.g., [3, §4.1].

We can also derive purely primal first-order optimality conditions for (P) and (P_\gamma) in terms of variational inequalities using standard arguments as in, e.g., [21, Thm. 2.22].

Proposition 2.3 *Let \bar{u} and u_γ be solutions of (P) and (P_\gamma) with associated adjoint states $\bar{p} := K^*(z - K\bar{u})$ and $p_\gamma := K^*(z - Ku_\gamma),$ respectively. Then,*

$$\begin{aligned} (-\bar{p}, u - \bar{u})_{L^2(\Omega)} + \alpha G'(\bar{u}; u - \bar{u}) &\geq 0 \quad \text{for all } u \in U_{\text{ad}}, \\ (-p_\gamma + \gamma u_\gamma, u - u_\gamma)_{L^2(\Omega)} + \alpha G'(u_\gamma; u - u_\gamma) &\geq 0 \quad \text{for all } u \in U_{\text{ad}}. \end{aligned}$$

Here, $G'(\bar{u}; u - \bar{u})$ denotes the directional derivative of G at \bar{u} in direction $u - \bar{u},$ which will be characterized in the following lemma. Note that for $\bar{u}, u \in U_{\text{ad}}$ we have $u - \bar{u} \in T_{U_{\text{ad}}}(\bar{u})$ for

$$T_{U_{\text{ad}}}(u) := \left\{ v \in L^2(\Omega) : v(x) \begin{cases} \geq 0 & \text{if } u(x) = u_1 \\ \leq 0 & \text{if } u(x) = u_d \end{cases} \right\},$$

i.e., the tangential cone to U_{ad} in the point u . It thus suffices to consider directional derivatives for directions in $T_{U_{\text{ad}}}$, which helps to avoid unnecessary case distinctions in the proof. Furthermore, since $U_{\text{ad}} \subset L^\infty(\Omega)$, we only have to consider directions in $L^\infty(\Omega)$. In the following, all pointwise expressions and calculations are understood in an almost everywhere sense.

Lemma 2.1 *Let $u \in U_{\text{ad}}$ and define the sets*

$$\begin{aligned} S_i &:= \{x \in \Omega : u(x) = u_i\}, \quad i = 1, \dots, d, \\ T_i &:= \{x \in \Omega : u_i < u(x) < u_{i+1}\}, \quad i = 1, \dots, d - 1. \end{aligned}$$

The directional derivative of G in direction $v \in T_{U_{\text{ad}}}(u) \cap L^\infty(\Omega)$ is then given as

$$\begin{aligned} G'(u; v) &= \sum_{i=1}^{d-1} \int_{T_i} \frac{1}{2}(u_i + u_{i+1})v(x) \, dx \\ &+ \sum_{i=1}^d \left[\int_{S_i \cap \{v \geq 0\}} \frac{1}{2}(u_i + u_{i+1})v(x) \, dx + \int_{S_i \cap \{v < 0\}} \frac{1}{2}(u_{i-1} + u_i)v(x) \, dx \right]. \end{aligned}$$

Proof We use the definition of the directional derivative and of the sets S_i and T_i to obtain

$$\begin{aligned} G'(u; v) &:= \lim_{\rho \rightarrow 0} \frac{1}{\rho} (G(u + \rho v) - G(u)) \\ &= \lim_{\rho \rightarrow 0} \frac{1}{\rho} \left[\sum_{i=1}^{d-1} \int_{T_i} (g(u(x) + \rho v(x)) - g(u(x))) \, dx \right. \\ &\quad \left. + \sum_{i=1}^d \int_{S_i} (g(u(x) + \rho v(x)) - g(u(x))) \, dx \right]. \end{aligned}$$

We now make use of our assumption that $v \in T_{U_{\text{ad}}}(u) \cap L^\infty(\Omega)$. For such a v , we can find a $\rho > 0$ such that $u + \rho v \in U_{\text{ad}}$. Note that this is a pointwise condition, which we are going to exploit in the following. We have to differentiate between several cases.

(i) First, assume that $x \in T_i$ with $1 \leq i \leq d - 1$. For ρ small enough we then get $u(x) + \rho v(x) \in [u_i, u_{i+1}]$. Hence we obtain

$$\begin{aligned} g(u(x) + \rho v(x)) - g(u(x)) &= \frac{1}{2}((u_i + u_{i+1})(u(x) + \rho v(x)) - u_i u_{i+1}) \\ &\quad - \frac{1}{2}((u_i + u_{i+1})u(x) - u_i u_{i+1}) \\ &= \frac{\rho}{2}(u_i + u_{i+1})v(x). \end{aligned} \tag{2.4}$$

which yields

$$\lim_{\rho \rightarrow 0} \int_{T_i} (g(u(x) + \rho v(x)) - g(u(x))) dx = \int_{T_i} \frac{1}{2} (u_i + u_{i+1}) v(x) dx.$$

(ii) Now assume that $x \in S_i$ with $1 < i < d$. Then by definition, $u(x) = u_i$. Here we have to further differentiate between three cases.

$v(x) = 0$: Here we obtain $u(x) + \rho v(x) = u(x)$, leading to

$$g(u(x) + \rho v(x)) - g(u(x)) = 0.$$

$v(x) > 0$: Here we obtain $u(x) + \rho v(x) \in [u_i, u_{i+1}]$ for ρ small enough, leading to

$$g(u(x) + \rho v(x)) - g(u(x)) = \frac{\rho}{2} (u_i + u_{i+1}) v(x).$$

$v(x) < 0$: Here we obtain $u(x) + \rho v(x) \in [u_{i-1}, u_i]$, leading as in (2.4) to

$$g(u(x) + \rho v(x)) - g(u(x)) = \frac{\rho}{2} (u_{i-1} + u_i) v(x).$$

Combining all three cases yields

$$\begin{aligned} & \lim_{\rho \rightarrow 0} \frac{1}{\rho} \int_{S_i} (g(u(x) + \rho v(x)) - g(u(x))) dx \\ &= \int_{S_i \cap \{v \geq 0\}} \frac{1}{2} (u_i + u_{i+1}) v(x) dx + \int_{S_i \cap \{v < 0\}} \frac{1}{2} (u_{i-1} + u_i) v(x) dx. \end{aligned}$$

(iii) We are left with the special cases $x \in S_i$ for $i = 1$ and $i = d$. We only consider the case $i = 1$ as the case $i = d$ is similar. Hence we assume $x \in S_1$, which implies $u(x) = u_1$. Since $v \in T_{U_{ad}}(u)$, we have that $v(x) \geq 0$. If $v(x) > 0$, we obtain for ρ small enough that $u(x) + \rho v(x) \in T_1$ holds, leading to

$$g(u(x) + \rho v(x)) - g(u(x)) = \frac{\rho}{2} (u_1 + u_2) v(x)$$

and similar if $v(x) = 0$. This leads to

$$\lim_{\rho \rightarrow 0} \frac{1}{\rho} \int_{S_1} (g(u(x) + \rho v(x)) - g(u(x))) dx = \int_{S_1} \frac{1}{2} (u_1 + u_2) v(x) dx.$$

A similar argument for the remaining case $i = d$ finishes the proof.

□

3 Regularity assumption and error estimates

We now extend the active set condition from [24,26] to the multibang control problem. From Proposition 2.1, we see that the optimal control \bar{u} is not uniquely determined by the adjoint state \bar{p} on the *singular sets* $Q_{i,i+1}$. We therefore need to control the way in which \bar{p} “detaches” from these sets. This motivates the following assumption.

Assumption REG For the solution \bar{u} to (P) with adjoint state $\bar{p} = K^*(z - K\bar{u})$ there exists a constant $c > 0$ and $\kappa > 0$ such that

$$\text{meas} \left(\bigcup_{i=1}^{d-1} \left\{ x \in \Omega : \left| \bar{p}(x) - \frac{\alpha}{2}(u_i + u_{i+1}) \right| < \varepsilon \right\} \right) \leq c\varepsilon^\kappa$$

holds for all $\varepsilon > 0$ small enough.

Note that if \bar{u} satisfies this assumption, the sets $Q_{i,i+1}$ have Lebesgue measure zero. Hence, \bar{u} is multibang by Proposition 2.1. In addition, we have the following result, which is a direct consequence of $\text{meas}\{x \in \Omega : \bar{p}(x) \in Q_{i,i+1}\} = 0$.

Lemma 3.1 *Assume \bar{u} satisfies Assumption REG. Then $\bar{p}(x) \in Q_i$ if and only if $\bar{u}(x) = u_i$ holds almost everywhere in Ω .*

Following [9, Lemma 1.3], we can derive a sufficient condition for Assumption REG.

Theorem 3.1 *Suppose that the adjoint state $\bar{p} \in C^1(\bar{\Omega})$ and satisfies*

$$\min_{x \in K_i} |\nabla p(x)| > 0 \quad \text{for all } i = 1, \dots, d - 1,$$

where

$$K_i := \left\{ x \in \bar{\Omega} : p(x) = \frac{\alpha}{2}(u_i + u_{i+1}) \right\}.$$

Then Assumption REG holds with $\kappa = 1$.

Proof Define for $t \in \mathbb{R}$ the level sets $F_t := \{x \in \bar{\Omega} : p(x) = t\}$. Now we use a continuity argument to obtain constants $\varepsilon_0, c_0, C > 0$ such that for all $|t - \frac{\alpha}{2}(u_i + u_{i+1})| \leq \varepsilon_0$ and all $1 \leq i < d$ there holds

$$|\nabla p(x)| \geq c_0 > 0, \quad \mathcal{H}^{n-1}(F_t) \leq C,$$

where \mathcal{H}^{n-1} is the $(n-1)$ -dimensional Hausdorff measure. In the following, we denote by $\mathbb{1}_C$ the characteristic function of the set C , i.e., $\mathbb{1}_C(x) = 1$ if $x \in C$ and 0 else. We now use the co-area formula

$$\int_{\Omega} h(x)|\nabla p(x)|dx = \int_{-\infty}^{\infty} \left(\int_{p^{-1}(t)} h(x)d\mathcal{H}^{n-1}(x) \right) dt$$

with the function

$$h(x) := \mathbb{1}_{E_i}, \quad E_i := \left\{ x \in \Omega : \left| p(x) - \frac{\alpha}{2}(u_i + u_{i+1}) \right| \leq \varepsilon \right\},$$

to obtain for all $1 \leq i < d$ and $0 < \varepsilon \leq \varepsilon_0$ that

$$c_0 \operatorname{meas}(E_i) \leq \int_{E_i} |\nabla p(x)| dx = \int_{-\varepsilon}^{\varepsilon} \mathcal{H}^{n-1} \left(F_{t-\frac{\alpha}{2}(u_i+u_{i+1})} \right) dt \leq 2C\varepsilon$$

holds. Since this holds for all $1 \leq i < d$, the Assumption REG now follows with $\kappa = 1$. □

We now establish error estimates for the approximation (P_γ) of (P) . For this purpose, we first derive a stronger version of Proposition 2.3. The next result, which is similar to ones in [18,19], is the most important tool in the convergence analysis.

Lemma 3.2 *Assume that the solution \bar{u} to (P) satisfies Assumption REG. Then,*

$$(-\bar{p}, u - \bar{u})_{L^2(\Omega)} + \alpha G'(\bar{u}; u - \bar{u}) \geq c_A \|u - \bar{u}\|_{L^1(\Omega)}^{1+\frac{1}{\kappa}} \quad \forall u \in U_{\text{ad}}$$

with a constant $c_A := c_A(\kappa) > 0$.

Proof First, recall that Assumption REG implies that \bar{u} has a multibang structure. Furthermore, using Lemma 3.1 we obtain with the definition of Q_i and S_i in Proposition 2.1 and Lemma 2.1, respectively, that $\bar{u}(x) \in S_i$ if and only if $\bar{p}(x) \in Q_i$. Now we use Lemma 2.1 and the fact that $u - \bar{u} \in T_{U_{\text{ad}}}(\bar{u})$ to compute

$$\begin{aligned} & (-\bar{p}, u - \bar{u})_{L^2(\Omega)} + \alpha G'(\bar{u}; u - \bar{u}) \\ &= \int_{\{\bar{p} \in Q_1\}} \left(-\bar{p}(x) + \frac{\alpha}{2}(u_1 + u_2) \right) (u(x) - \bar{u}(x)) dx \\ &+ \int_{\{\bar{p} \in Q_d\}} \left(-\bar{p}(x) + \frac{\alpha}{2}(u_{d-1} + u_d) \right) (u(x) - \bar{u}(x)) dx \\ &+ \sum_{i=2}^{d-1} \int_{\{\bar{p} \in Q_i\} \cap \{u - \bar{u} \geq 0\}} \left(-\bar{p}(x) + \frac{\alpha}{2}(u_i + u_{i+1}) \right) (u(x) - \bar{u}(x)) dx \\ &+ \sum_{i=2}^{d-1} \int_{\{\bar{p} \in Q_i\} \cap \{u - \bar{u} < 0\}} \left(-\bar{p}(x) + \frac{\alpha}{2}(u_{i-1} + u_i) \right) (u(x) - \bar{u}(x)) dx. \end{aligned}$$

Here we have abbreviated the sets $\{\bar{p} \in Q_1\} := \{x \in \Omega : \bar{p}(x) \in Q_1\}$ and similar for the other sets. Recall that by definition, $\bar{p}(x) \in Q_1$ implies that $-\bar{p}(x) + \frac{\alpha}{2}(u_1 + u_2) > 0$. Furthermore, we know that $\bar{u}(x) = u_1$, leading to $u(x) - \bar{u}(x) = u(x) - u_1 \geq 0$.

We similarly obtain on Q_d that $-\bar{p}(x) + \frac{\alpha}{2}(u_{d-1} + u_d) < 0$ and $u(x) - \bar{u}(x) = u(x) - u_d \leq 0$. Finally, if $\bar{p}(x) \in Q_i$ for $1 < i < d$, we obtain that

$$\frac{\alpha}{2}(u_{i-1} + u_i) < \bar{p}(x) < \frac{\alpha}{2}(u_i + u_{i+1}),$$

which leads to

$$-\bar{p}(x) + \frac{\alpha}{2}(u_i + u_{i+1}) > 0 \quad \text{and} \quad -\bar{p}(x) + \frac{\alpha}{2}(u_{i-1} + u_i) < 0.$$

This allows us to write

$$\begin{aligned} & (-\bar{p}, u - \bar{u})_{L^2(\Omega)} + \alpha G'(\bar{u}; u - \bar{u}) \\ &= \int_{\{\bar{p} \in Q_1\}} \left| -\bar{p}(x) + \frac{\alpha}{2}(u_1 + u_2) \right| |u(x) - \bar{u}(x)| dx \\ &+ \int_{\{\bar{p} \in Q_d\}} \left| -\bar{p}(x) + \frac{\alpha}{2}(u_{d-1} + u_d) \right| |u(x) - \bar{u}(x)| dx \\ &+ \sum_{i=2}^{d-1} \int_{\{\bar{p} \in Q_i\} \cap \{u - \bar{u} \geq 0\}} \left| -\bar{p}(x) + \frac{\alpha}{2}(u_i + u_{i+1}) \right| |u(x) - \bar{u}(x)| dx \\ &+ \sum_{i=2}^{d-1} \int_{\{\bar{p} \in Q_i\} \cap \{u - \bar{u} < 0\}} \left| -\bar{p}(x) + \frac{\alpha}{2}(u_{i-1} + u_i) \right| |u(x) - \bar{u}(x)| dx. \end{aligned}$$

Now let $\varepsilon > 0$ and consider the set

$$Q_1^\varepsilon := \left\{ q : q \leq \frac{\alpha}{2}(u_1 + u_2) - \varepsilon \right\} \subset Q_1.$$

Let $\bar{p}(x) \in Q_1^\varepsilon$. Together with $-\bar{p}(x) + \frac{\alpha}{2}(u_1 + u_2) > 0$, this implies that

$$\left| -\bar{p}(x) + \frac{\alpha}{2}(u_1 + u_2) \right| = -\bar{p}(x) + \frac{\alpha}{2}(u_1 + u_2) \geq \varepsilon,$$

leading to

$$\begin{aligned} \int_{\{\bar{p} \in Q_1\}} \left| -\bar{p} + \frac{\alpha}{2}(u_1 + u_2) \right| |u - \bar{u}| dx &\geq \int_{\{\bar{p} \in Q_1^\varepsilon\}} \left| -\bar{p} + \frac{\alpha}{2}(u_1 + u_2) \right| |u - \bar{u}| dx \\ &\geq \varepsilon \int_{\{\bar{p} \in Q_1^\varepsilon\}} |u - \bar{u}| dx. \end{aligned}$$

We similarly define

$$Q_d^\varepsilon := \left\{ q \geq \frac{\alpha}{2}(u_{d-1} + u_d) + \varepsilon \right\},$$

leading to

$$\int_{\{\bar{p} \in Q_d\}} \left| -\bar{p}(x) + \frac{\alpha}{2}(u_{d-1} + u_d) \right| |u(x) - \bar{u}(x)| dx \geq \varepsilon \int_{\{\bar{p} \in Q_d^\varepsilon\}} |u(x) - \bar{u}(x)| dx,$$

as well as for $1 < i < d$

$$Q_i^\varepsilon := \left\{ \varepsilon + \frac{\alpha}{2}(u_{i-1} + u_i) \leq q \leq \frac{\alpha}{2}(u_i + u_{i+1}) - \varepsilon \right\} \subset Q_i.$$

The latter leads to

$$\begin{aligned} \left| -\bar{p}(x) + \frac{\alpha}{2}(u_i + u_{i+1}) \right| &= -\bar{p}(x) + \frac{\alpha}{2}(u_i + u_{i+1}) \geq \varepsilon, \\ \left| -\bar{p}(x) + \frac{\alpha}{2}(u_{i-1} + u_i) \right| &= \bar{p}(x) - \frac{\alpha}{2}(u_{i-1} + u_i) \geq \varepsilon \end{aligned}$$

and therefore

$$\begin{aligned} &\int_{\{\bar{p} \in Q_i\} \cap \{u - \bar{u} \geq 0\}} \left| -\bar{p}(x) + \frac{\alpha}{2}(u_i + u_{i+1}) \right| |u(x) - \bar{u}(x)| dx \\ &+ \int_{\{\bar{p} \in Q_i\} \cap \{u - \bar{u} < 0\}} \left| -\bar{p}(x) + \frac{\alpha}{2}(u_{i-1} + u_i) \right| |u(x) - \bar{u}(x)| dx \\ &\geq \varepsilon \int_{\{\bar{p} \in Q_i^\varepsilon\} \cap \{u - \bar{u} \geq 0\}} |u(x) - \bar{u}(x)| dx + \varepsilon \int_{\{\bar{p} \in Q_i^\varepsilon\} \cap \{u - \bar{u} < 0\}} |u(x) - \bar{u}(x)| dx \\ &= \varepsilon \int_{\{\bar{p} \in Q_i^\varepsilon\}} |u(x) - \bar{u}(x)| dx. \end{aligned}$$

We now combine all these estimates to obtain

$$\begin{aligned} &(-\bar{p}, u - \bar{u})_{L^2(\Omega)} + \alpha G'(\bar{u}; u - \bar{u}) \\ &\geq \varepsilon \sum_{i=1}^d \int_{\{\bar{p} \in Q_i^\varepsilon\}} |u(x) - \bar{u}(x)| dx \\ &= \varepsilon \sum_{i=1}^d \left(\int_{\{\bar{p} \in Q_i\}} |u(x) - \bar{u}(x)| dx - \int_{\{\bar{p} \in Q_i\} \setminus \{\bar{p} \in Q_i^\varepsilon\}} |u(x) - \bar{u}(x)| dx \right) \\ &= \varepsilon \|u - \bar{u}\|_{L^1(\Omega)} - \varepsilon \sum_{i=1}^d \int_{\{\bar{p} \in Q_i\} \setminus \{\bar{p} \in Q_i^\varepsilon\}} |u(x) - \bar{u}(x)| dx \end{aligned}$$

$$\geq \varepsilon \|u - \bar{u}\|_{L^1(\Omega)} - \varepsilon \|u - \bar{u}\|_{L^\infty(\Omega)} \sum_{i=1}^d \int_{\{\bar{p} \in Q_i\} \setminus \{\bar{p} \in Q_i^e\}} 1 \, dx,$$

where we have used the L^∞ -boundedness of $u - \bar{u}$ in the last step. We now use Assumption REG to estimate the remaining sum, yielding

$$\sum_{i=1}^d \int_{\{\bar{p} \in Q_i\} \setminus \{\bar{p} \in Q_i^e\}} 1 \, dx = \text{meas} \left(\bigcup_{i=1}^{d-1} \left\{ x \in \Omega : \left| \bar{p}(x) - \frac{\alpha}{2}(u_i + u_{i+1}) \right| < \varepsilon \right\} \right) \leq c\varepsilon^\kappa.$$

Summarizing, we have for a constant $c > 1$ that

$$(-\bar{p}, u - \bar{u})_{L^2(\Omega)} + \alpha G'(\bar{u}; u - \bar{u}) \geq \varepsilon \|u - \bar{u}\|_{L^1(\Omega)} - c\varepsilon^{\kappa+1},$$

and hence setting

$$\varepsilon := c^{-\frac{2}{\kappa}} \|u - \bar{u}\|_{L^1(\Omega)}^{\frac{1}{\kappa}}$$

finishes the proof. □

We now have everything at hand to prove approximation error estimates.

Theorem 3.2 *Let \bar{u} be a solution of (P) with corresponding state $\bar{y} := K\bar{u}$ and assume that Assumption REG is satisfied. Furthermore, let u_γ be the solution of (P $_\gamma$) for $\gamma > 0$ with corresponding state $y_\gamma := Ku_\gamma$. Then there exists a constant $c > 0$ such that*

$$\frac{1}{\gamma} \|y_\gamma - \bar{y}\|_Y^2 + \frac{1}{\gamma} \|u_\gamma - \bar{u}\|_{L^1(\Omega)}^{1+\frac{1}{\kappa}} + \|u_\gamma - \bar{u}\|_{L^2(\Omega)}^2 \leq c\gamma^\kappa.$$

Proof First note that G is a convex function and hence that

$$G'(\bar{u}; u_\gamma - \bar{u}) + G'(u_\gamma; \bar{u} - u_\gamma) \leq 0.$$

We thus obtain from Proposition 2.3 and Lemma 3.2 that

$$\begin{aligned} (-\bar{p}, u - \bar{u})_{L^2(\Omega)} + \alpha G'(\bar{u}; u - \bar{u}) &\geq c_A \|u_\gamma - \bar{u}\|_{L^1(\Omega)}^{1+\frac{1}{\kappa}} \quad \forall u \in U_{\text{ad}}, \\ (-p_\gamma, u - u_\gamma)_{L^2(\Omega)} + \alpha G'(u_\gamma; u - u_\gamma) + \gamma (u_\gamma, u - u_\gamma)_{L^2(\Omega)} &\geq 0 \quad \forall u \in U_{\text{ad}}. \end{aligned}$$

Inserting $u = u_\gamma$ and $u = \bar{u}$ into two above inequalities, respectively, and then adding both yields

$$\begin{aligned} (-\bar{p} + p_\gamma, u_\gamma - \bar{u})_{L^2(\Omega)} + \alpha (G'(\bar{u}; u_\gamma - \bar{u}) + G'(u_\gamma; \bar{u} - u_\gamma)) + \gamma (u_\gamma, \bar{u} - u_\gamma)_{L^2(\Omega)} \\ \geq c_A \|u_\gamma - \bar{u}\|_{L^1(\Omega)}^{1+\frac{1}{\kappa}}. \end{aligned}$$

We now use the definition of $\bar{p} = K^*(z - K\bar{u})$ and $p_\gamma = K^*(z - Ku_\gamma)$ to deduce that

$$(-\bar{p} + p_\gamma, u_\gamma - \bar{u})_{L^2(\Omega)} = -\|y_\gamma - \bar{y}\|_Y^2.$$

Hence, by adding $\gamma\|\bar{u} - u_\gamma\|_{L^2(\Omega)}^2$ to the inequality above and rearranging terms, we obtain that

$$\begin{aligned} \|y_\gamma - \bar{y}\|_Y^2 + c_A\|u_\gamma - \bar{u}\|_{L^1(\Omega)}^{1+\frac{1}{\kappa}} + \gamma\|u_\gamma - \bar{u}\|_{L^2(\Omega)}^2 &\leq \alpha(G'(\bar{u}; u_\gamma - \bar{u}) + G'(u_\gamma; \bar{u} - u_\gamma)) \\ &\quad + \gamma(\bar{u}, \bar{u} - u_\gamma)_{L^2(\Omega)} \\ &\leq \gamma(\bar{u}, \bar{u} - u_\gamma)_{L^2(\Omega)} \\ &\leq c\gamma\|u_\gamma - \bar{u}\|_{L^1(\Omega)} \\ &\leq \frac{c_A}{2}\|u_\gamma - \bar{u}\|_{L^1(\Omega)}^{1+\frac{1}{\kappa}} + c\gamma^{\kappa+1}, \end{aligned}$$

where we have used Young’s inequality in the last step. The stated inequality now follows immediately. \square

4 Discretization error estimates

In practice, the exact operator K is not realizable, and a discretization $K_h : L^2(\Omega) \rightarrow Y_h$ with finite dimensional range Y_h must be employed. Denote by $u_{\gamma,h}$ the solution of the discrete problem

$$\min_{u \in U_{ad}} \frac{1}{2}\|K_h u - z\|_Y^2 + \alpha G(u) + \frac{\gamma}{2}\|u\|_{L^2(\Omega)}^2 \tag{P_{\gamma,h}}$$

with corresponding state $y_{\gamma,h} := K_h u_{\gamma,h}$ and adjoint state $p_{\gamma,h} := K_h^*(z - y_{\gamma,h})$. If K is the solution operator of an elliptic partial differential equation and K_h its finite element discretization as in the next section, $(P_{\gamma,h})$ can be interpreted as a variational discretization [12,13].

We assume that for all $h > 0$, the estimate

$$\|(K - K_h)u_{\gamma,h}\|_Y + \|(K^* - K_h^*)(y_{\gamma,h} - z)\|_{L^2(\Omega)} \leq \delta(h), \tag{4.1}$$

holds uniformly for all $\gamma > 0$ with a monotonically increasing function $\delta : \mathbb{R}_0^+ \rightarrow \mathbb{R}$ such that $\delta(0) = 0$. Note that this approximation condition only needs to be satisfied for the solutions to the discretized problem $(P_{\gamma,h})$. However, as in [23] the condition can also be replaced by a corresponding uniform condition for the solution to the continuous problem (P_γ) .

Now, we follow [23, Proposition 1.8] and estimate the discretization error for the solution to (P_γ) .

Theorem 4.1 *For all $\gamma > 0$ and $h \geq 0$ there holds*

$$\|y_\gamma - y_{\gamma,h}\|_Y^2 + \gamma\|u_\gamma - u_{\gamma,h}\|_{L^2(\Omega)}^2 \leq (1 + \gamma^{-1})\delta(h)^2.$$

Proof With $u_{\gamma,h}$ and u_γ solutions to $(P_{\gamma,h})$ and (P_γ) , respectively, we have from Proposition 2.3 that

$$\begin{aligned} (-p_{\gamma,h} + \gamma u_{\gamma,h}, u_\gamma - u_{\gamma,h})_{L^2(\Omega)} + \alpha G'(u_{\gamma,h}; u_\gamma - u_{\gamma,h}) &\geq 0, \\ (-p_\gamma + \gamma u_\gamma, u_{\gamma,h} - u_\gamma)_{L^2(\Omega)} + \alpha G'(u_\gamma; u_{\gamma,h} - u_\gamma) &\geq 0. \end{aligned}$$

Adding these two inequalities, substituting

$$p_{\gamma,h} = -K_h^*(K_h u_{\gamma,h} - z), \quad p_\gamma = -K^*(K u_\gamma - z),$$

and using the convexity of G then yields

$$\begin{aligned} (K_h^*(K_h u_{\gamma,h} - z) + \gamma u_{\gamma,h}, u_\gamma - u_{\gamma,h}) + (K^*(K u_\gamma - z) + \gamma u_\gamma, u_{\gamma,h} - u_\gamma) \\ \geq -\alpha (G'(u_{\gamma,h}; u_\gamma - u_{\gamma,h}) + G'(u_\gamma; u_{\gamma,h} - u_\gamma)) \geq 0. \end{aligned}$$

We thus obtain that

$$\begin{aligned} \gamma \|u_{\gamma,h} - u_\gamma\|_{L^2(\Omega)}^2 &\leq (K_h^*(y_{\gamma,h} - z) - K^*(y_\gamma - z), u_\gamma - u_{\gamma,h}) \\ &\leq ((K_h^* - K^*)(y_{\gamma,h} - z), u_\gamma - u_{\gamma,h}) \\ &\quad + (K^*(y_{\gamma,h} - y_\gamma), u_\gamma - u_{\gamma,h}). \end{aligned}$$

The rest of the proof follows similarly to the proof of [23, Proposition 1.6]. The first term on the right-hand side is estimated by the Cauchy–Schwarz inequality and the inequality (4.1) as

$$((K_h^* - K^*)(y_{\gamma,h} - z), u_\gamma - u_{\gamma,h}) \leq \frac{\gamma}{2} \|u_{\gamma,h} - u_\gamma\|_{L^2(\Omega)}^2 + \frac{1}{2\gamma} \delta(h)^2.$$

Rewriting the second term and using again the Cauchy–Schwarz inequality combined with the inequality (4.1), we obtain

$$\begin{aligned} (K^*(y_{\gamma,h} - y_\gamma), u_\gamma - u_{\gamma,h}) &= -\|y_\gamma - y_{\gamma,h}\|_Y^2 + (y_\gamma - y_{\gamma,h}, (K_h - K)u_{\gamma,h}) \\ &\leq -\frac{1}{2} \|y_\gamma - y_{\gamma,h}\|_Y^2 + \frac{1}{2} \delta(h)^2. \end{aligned}$$

Adding these two estimates, we finally arrive at

$$\frac{1}{2} \|y_\gamma - y_{\gamma,h}\|_Y^2 + \frac{\gamma}{2} \|u_\gamma - u_{\gamma,h}\|_{L^2(\Omega)}^2 \leq \left(\frac{1}{2} + \frac{1}{2\gamma}\right) \delta(h)^2.$$

□

Combining the approximation error estimate from Theorem 3.2 and the discretization error estimate from Theorem 4.1, we immediately obtain the following result.

Theorem 4.2 *If \bar{u} satisfies Assumption REG, then*

$$\frac{1}{\gamma} \|y_\gamma - \bar{y}\|_Y^2 + \|u_{\gamma,h} - \bar{u}\|_{L^2(\Omega)}^2 \leq c \left(\gamma^{-1} (1 + \gamma^{-1}) \delta(h)^2 + \gamma^\kappa \right)$$

holds for all $\gamma > 0$ and $h \geq 0$.

5 Active set method for the regularized problem

Let us now consider the special case where $y = Ku$ is given as the unique solution of the partial differential equation

$$\begin{cases} Ay = u & \text{in } \Omega, \\ y = 0 & \text{on } \partial\Omega. \end{cases} \tag{5.1}$$

with A being a linear second-order elliptic differential operator, e.g., $A = -\Delta$. In this case, the optimality conditions from Proposition 2.2 can be solved using a superlinearly convergent semi-smooth Newton method in function space; see [3,6,7].

We recall that (2.3) can be written as $u_\gamma = H_\gamma(p_\gamma)$ for $H_\gamma : L^r(\Omega) \rightarrow L^2(\Omega)$ with $r \geq 2$,

$$[H_\gamma(p)](x) = \begin{cases} u_i & \text{if } p(x) \in Q_i^\gamma \ 1 \leq i \leq d, \\ \frac{1}{\gamma} (p(x) - \frac{\alpha}{2}(u_i + u_{i+1})) & \text{if } p(x) \in Q_{i,i+1}^\gamma \ 1 \leq i < d, \end{cases}$$

where $p_\gamma \in H_0^1(\Omega)$ is the solution to the adjoint equation

$$\begin{cases} A^*p = z - y_\gamma & \text{in } \Omega, \\ p = 0 & \text{on } \partial\Omega, \end{cases} \tag{5.2}$$

and y_γ is the solution to (5.1) with $u = u_\gamma$. From the natural $H_0^1(\Omega)$ regularity of solutions to (5.2), the Sobolev embedding $H_0^1(\Omega) \hookrightarrow L^r(\Omega)$ for some $r > 2$, and the general theory of semi-smooth Newton methods in function space [22], we deduce that the superposition operator H_γ is Newton differentiable from $L^r(\Omega)$ to $L^2(\Omega)$ with

$$[D_N H_\gamma(p)h](x) = \begin{cases} \frac{1}{\gamma}h(x) & \text{if } p(x) \in Q_{i,i+1}^\gamma, \\ 0 & \text{else.} \end{cases}$$

A Newton step for the solution of (P_γ) can therefore be formulated as

$$\begin{pmatrix} -\text{Id} & A & 0 \\ 0 & \text{Id} & A^* \\ 0 & A & -D_N H_\gamma(p^k) \end{pmatrix} \begin{pmatrix} u^{k+1} - u^k \\ y^{k+1} - y^k \\ p^{k+1} - p^k \end{pmatrix} = - \begin{pmatrix} Ay^k - u^k \\ A^*p^k + y^k - z \\ Ay^k - H_\gamma(p^k) \end{pmatrix} \tag{5.3}$$

In [3], this was reduced to a symmetric system in (y, p) . Here, we instead consider an equivalent primal active set formulation that has proven to be more robust for small values of γ and h . In a slight abuse of notation, we introduce

$$Q_i^k := \{x \in \Omega : p^k(x) \in Q_i^\gamma\}, \quad 1 \leq i \leq d,$$

and similarly for $Q_{i,i+1}^k$. The following algorithm is an extension of the one proposed in [20] for $G(u) = \|u\|_{L^1(\Omega)}$.

Algorithm 1 Choose initial data u^0, p^0 and parameters α, γ , set $k = 0$ and compute the sets Q_i^0 for $1 \leq i \leq d$ and $Q_{i,i+1}^0$ for $1 \leq i < d$.

1. Solve for $(u^{k+1}, y^{k+1}, p^{k+1}, \lambda^{k+1})$ satisfying

$$\begin{cases} Ay^{k+1} - u^{k+1} = 0, \\ A^* p^{k+1} + y^{k+1} - z = 0, \\ -p^{k+1} + \gamma u^{k+1} + \alpha \lambda^{k+1} = 0, \end{cases} \tag{5.4a}$$

$$\begin{aligned} \left(1 - \sum_{i=1}^d \mathbb{1}_{Q_i^k}\right) \lambda^{k+1} + \left(1 - \sum_{i=1}^{d-1} \mathbb{1}_{Q_{i,i+1}^k}\right) u^{k+1} \\ = \sum_{i=1}^d \mathbb{1}_{Q_i^k} u_i + \frac{1}{2} \sum_{i=1}^{d-1} \mathbb{1}_{Q_{i,i+1}^k} (u_i + u_{i+1}). \end{aligned} \tag{5.4b}$$

2. Compute the sets Q_i^{k+1} for $1 \leq i \leq d$ and $Q_{i,i+1}^{k+1}$ for $1 \leq i < d$.
3. If $Q_i^k = Q_i^{k+1}$ for $1 \leq i \leq d$ and $Q_{i,i+1}^k = Q_{i,i+1}^{k+1}$ for $1 \leq i < d$, then go to step 4. Otherwise set $k = k + 1$ and go to step 2.
4. STOP: u^{k+1} is a solution of (P_γ) .

The stopping criterion yields solutions of (P_γ) .

Lemma 5.1 *If*

$$\begin{aligned} Q_i^k &= Q_i^{k+1} \quad 1 \leq i \leq d, \\ Q_{i,i+1}^k &= Q_{i,i+1}^{k+1} \quad 1 \leq i < d, \end{aligned}$$

then the solution (u^{k+1}, p^{k+1}) computed from (5.4) satisfy (2.3). In particular, u^{k+1} is a solution to (P_γ) .

Proof Since for fixed Q_i^k and $Q_{i,i+1}^k$ the solution of (5.4) is unique, we have $(u^k, y^k, p^k) = (u^{k+1}, y^{k+1}, p^{k+1})$. Inserting this into (5.4b) and comparing with (2.3) yields the claim. \square

We now show that Algorithm 1 coincides with a semi-smooth Newton method, which implies locally superlinear convergence.

Theorem 5.1 *The active set step (5.4) is equivalent to the semi-smooth Newton step (5.3).*

Proof Clearly, the first two equations of (5.3) are equivalent to the first two equation of (5.4a). It therefore remains to consider the last equation, which is given by

$$A(y^{k+1} - y^k) - D_N H_\gamma(p^k)(p^{k+1} - p^k) = -Ay^k + H_\gamma(p^k). \tag{5.5}$$

Let us define the function

$$\lambda^{k+1}(x) := \begin{cases} -\frac{1}{\alpha}(-p^{k+1}(x) + \gamma u^{k+1}) & \text{if } x \in Q_i^k, \\ \frac{1}{2}(u_i + u_{i+1}) & \text{if } x \in Q_{i,i+1}^k. \end{cases}$$

We now make a case distinction pointwise almost everywhere.

- (i) If $x \in Q_i^k$, (5.5) reduces to $[Ay^{k+1}](x) = u_i$, and from the first line of (5.3) we obtain $u^{k+1}(x) = u_i$.
- (ii) If $x \in Q_{i,i+1}^k$, (5.5) shows that

$$\gamma u^{k+1}(x) - p^{k+1}(x) + \frac{\alpha}{2}(u_i + u_{i+1}) = \gamma u^{k+1}(x) - p^{k+1}(x) + \alpha \lambda^{k+1}(x) = 0.$$

Hence the third row of (5.3) is equivalent to (5.4b). In both cases, we obtain from the definition of λ^{k+1} that

$$-p^{k+1} + \gamma u^{k+1} + \alpha \lambda^{k+1} = 0,$$

which finally gives (5.4a) and therefore the claimed equivalence. □

6 Numerical results

In this section we present some numerical results and convergence rates. Let $\Omega \subset \mathbb{R}^d$ be a bounded Lipschitz domain and K be the operator mapping u to the weak solution y of

$$\begin{cases} -\Delta y = u & \text{in } \Omega, \\ y = 0 & \text{on } \partial\Omega. \end{cases} \tag{6.1}$$

The operator K_h is correspondingly defined via the Galerkin approximation of (6.1) using linear finite elements on a triangulation of Ω , which is chosen in such a way that the approximation condition (4.1) is satisfied; see [23]. For the multibang penalty, we take $(u_1, \dots, u_5) = (-2, -1, 0, 1, 2)$ and $\alpha = 2$. We implemented Algorithm 1 in Python using DOLFIN [16,17], which is part of the open-source computing platform FEniCS [1,15]. The linear system (5.4) arising from the active set step is solved using the sparse direct solver `spsolve` from SciPy. The code used to obtain the following results can be downloaded from <https://github.com/clason/multibangestimates>.

Example 1: $\kappa = 1$ We first consider $\Omega = (0, 1)$ and define

$$\begin{aligned} \bar{p}(x) := & \left(\frac{27}{2}x\right) \mathbb{1}_{[0, \frac{2}{9}]}(x) \\ & + \left(-72 + \frac{3123x}{2} - 13122x^2 + 54675x^3 - 111537x^4\right. \\ & \left. + \frac{177147}{2}x^5\right) \mathbb{1}_{[\frac{2}{9}, \frac{3}{9}]}(x) \\ & + (9 - 18x) \mathbb{1}_{[\frac{3}{9}, \frac{6}{9}]}(x) \end{aligned}$$

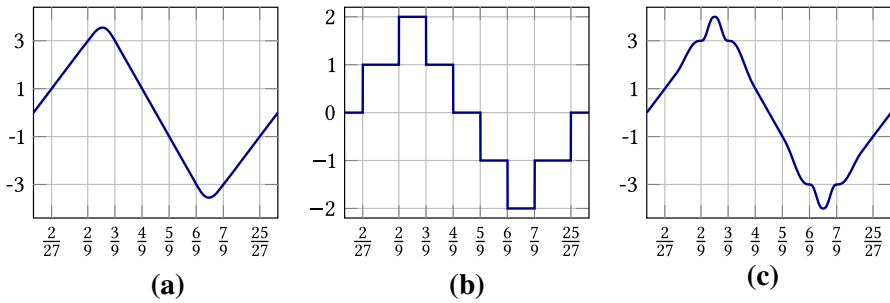


Fig. 1 Constructed optimal adjoint states \bar{p} and optimal control \bar{u} . **a** adjoint state, Example 1. **b** control, Examples 1 and 2. **c** adjoint state, Example 2

$$\begin{aligned}
 & + \left(-20079 + 136062x - 367416x^2 + 494262x^3 \right. \\
 & \quad \left. - \frac{662661}{2}x^4 + \frac{177147}{2}x^5 \right) \mathbb{1}_{\left[\frac{6}{9}, \frac{7}{9}\right)}(x) \\
 & + \left(-\frac{27}{2} + \frac{27}{2}x \right) \mathbb{1}_{\left[\frac{7}{9}, 1\right]}(x), \\
 \bar{u}(x) & := \mathbb{1}_{\left[\frac{2}{27}, \frac{2}{9}\right)}(x) + 2\mathbb{1}_{\left[\frac{2}{9}, \frac{3}{9}\right)}(x) + \mathbb{1}_{\left[\frac{3}{9}, \frac{4}{9}\right)}(x) \\
 & \quad - \mathbb{1}_{\left[\frac{5}{9}, \frac{6}{9}\right)}(x) - 2\mathbb{1}_{\left[\frac{6}{9}, \frac{7}{9}\right)}(x) - \mathbb{1}_{\left[\frac{7}{9}, \frac{25}{27}\right)}(x), \\
 \bar{y}(x) & := \sin(2\pi x) \\
 e_\Omega & := -\Delta \bar{y} - \bar{u}, \\
 z & := -Ke_\Omega - \Delta \bar{p} + \bar{y},
 \end{aligned}$$

see Fig. 1a and b. Note that $\bar{p}, \bar{y} \in C^2(\bar{\Omega})$, and that \bar{u} and \bar{p} satisfy the optimality conditions in Proposition 2.1. Hence, (\bar{u}, \bar{p}) are a solution to (P) . From Theorem 3.1 we further deduce that Assumption REG is satisfied with $\kappa = 1$.

We now compute the solution of $(P_{\gamma,h})$ for different values of h , where Ω is divided into equidistant elements with mesh size h . From Theorem 3.2 we expect that the numerical convergence rate

$$\kappa_{\gamma,h} := \frac{1}{\log(2)} \log \left(\frac{\|u_{\gamma/2,h} - \bar{u}\|_{L^2(\Omega)}^2}{\|u_{\gamma,h} - \bar{u}\|_{L^2(\Omega)}^2} \right)$$

satisfies $\kappa_{\gamma,h} \geq \kappa = 1$. We compute $\kappa_{\gamma,h}$ for different but fixed mesh sizes h . Due to the discretization error, we expect a certain saturation effect for small γ ; see Theorem 4.2. Note that for $d = 2$, it is known that Assumption REG is not only sufficient for convergence rates similar to Theorem 3.2 but also necessary for high convergence rates; see [25]. Hence, we expect that $\kappa_{\gamma,h} \approx 1$, which can be observed from Table 1a and Fig. 2a. In addition, the discretization error dominates for small γ as expected.

Example 2: $\kappa < 1$ We also consider an example where Assumption REG is only satisfied with $\kappa < 1$. The idea is to violate the assumption of the sufficient condition presented in Theorem 3.1. We modify the adjoint state \bar{p} from Example 1 to

Table 1 Computed numerical order of convergence for different h

$\gamma \setminus h$	10^{-4}	10^{-5}	10^{-6}
(a) Example 1			
2^{-4}	1.0143	1.0142	1.0141
2^{-6}	1.0028	1.0008	1.0007
2^{-8}	1.0211	1.0004	0.9998
2^{-10}	0.9295	1.0038	0.9989
2^{-12}	0.6828	1.0049	0.9954
2^{-14}	0.0	0.9592	0.9917
2^{-16}	0.0	-0.0096	0.9701
2^{-18}	0.0	0.0	0.1308
(b) Example 2			
2^{-4}	0.4679	0.4679	0.4679
2^{-6}	0.3993	0.3992	0.3992
2^{-8}	0.3668	0.3665	0.3664
2^{-10}	0.3509	0.3518	0.3513
2^{-12}	0.3379	0.3470	0.3453
2^{-14}	0.3293	0.3496	0.3424
2^{-16}	0.2986	0.3649	0.3413
2^{-18}	0.1774	0.4122	0.3274

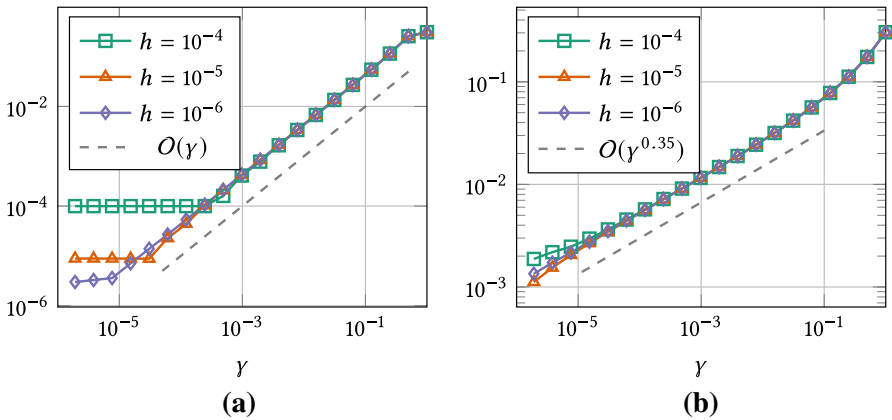


Fig. 2 Discretization and approximation error $\|u_{\gamma,h} - \bar{u}\|_{L^2(\Omega)}^2$ for different γ and h . **a** Example 1. **b** Example 2

$$\begin{aligned}
\bar{p}(x) := & \left(\frac{27}{2}x\right) \mathbb{1}_{\left[0, \frac{3}{27}\right)}(x) \\
& + \left(266085x^5 - \frac{433593}{2}x^4 + \frac{135765}{2}x^3 - \frac{20437}{2}x^2 + \frac{6812}{9}x - \frac{1703}{81}\right) \mathbb{1}_{\left[\frac{3}{27}, \frac{2}{9}\right)}(x) \\
& + \left(11334492x^5 - 14168034x^4 + 7054821x^3 - \frac{3498235}{2}x^2\right. \\
& \left. + \frac{1943450}{9}x - \frac{860051}{81}\right) \mathbb{1}_{\left[\frac{2}{9}, \frac{5}{18}\right)}(x) \\
& + \left(-11334492x^5 + 17316666x^4 - 10553301x^3 + \frac{6413635}{2}x^2\right. \\
& \left. - \frac{1457650}{3}x + \frac{528697}{18}\right) \mathbb{1}_{\left[\frac{5}{18}, \frac{3}{9}\right)}(x) \\
& + \left(-\frac{709317}{2}x^5 + 696195x^4 - \frac{1085913}{2}x^3 + 210182x^2 - \frac{121150}{3}x + \frac{27761}{9}\right) \mathbb{1}_{\left[\frac{3}{9}, \frac{4}{9}\right)}(x) \\
& + (-18x + 9) \mathbb{1}_{\left[\frac{4}{9}, \frac{5}{9}\right)}(x) \\
& + \left(-\frac{707859}{2}x^5 + \frac{2149821}{2}x^4 - \frac{2604285}{2}x^3 + \frac{1573075}{2}x^2 - \frac{710804}{3}x + \frac{256331}{9}\right) \mathbb{1}_{\left[\frac{5}{9}, \frac{6}{9}\right)}(x) \\
& + \left(-11340324x^5 + 39376206x^4 - 54660123x^3 + \frac{75835981}{2}x^2\right. \\
& \left. - \frac{39434798}{3}x + \frac{16396175}{9}\right) \mathbb{1}_{\left[\frac{6}{9}, \frac{13}{18}\right)}(x) \\
& + \left(11340324x^5 - 42526134x^4 + 63759915x^3 - \frac{95552197}{2}x^2\right. \\
& \left. + \frac{161022862}{9}x - \frac{433967467}{162}\right) \mathbb{1}_{\left[\frac{13}{18}, \frac{7}{9}\right)}(x) \\
& + \left(265356x^5 - \frac{2221101}{2}x^4 + \frac{3712707}{2}x^3 - 1549124x^2\right. \\
& \left. + \frac{11616563}{18}x - \frac{17395339}{162}\right) \mathbb{1}_{\left[\frac{7}{9}, \frac{8}{9}\right)}(x) \\
& + \left(\frac{27}{2}x - \frac{27}{2}\right) \mathbb{1}_{\left[\frac{8}{9}, 1\right)}(x),
\end{aligned}$$

see Fig. 1c, while the remaining functions remain unchanged. Note that for, e.g., $\hat{x} := \frac{2}{9}$, we obtain $p'(\hat{x}) = 0$ and $p(\hat{x}) = 3$, which violates the assumption of Theorem 3.1. Hence we expect that $\kappa < 1$ holds, resulting in a much slower convergence speed; see Theorem 3.2. This is corroborated by our numerical results: We obtain $\kappa_{\gamma, h} \approx 0.35 < 1$, which can be seen in Table 1b and Fig. 2b. Due to the slower convergence speed, we do not observe a saturation effect for the chosen range of γ and h .

7 Conclusions

For optimal control problems with a convex penalty promoting minimizers that pointwise almost everywhere take on values from a given discrete set, Moreau–Yosida approximation allows the solution by a superlinearly convergent semi-smooth Newton method. On a structural assumption on the behavior of the adjoint state near singular sets, convergence rates as the approximation parameter $\gamma \rightarrow 0$ can be derived. The same assumption also yields discretization error estimates for fixed $\gamma > 0$. Numerical experiments corroborate the predicted rate.

This work can be extended in a number of directions. First, an active set condition similar to Assumption REG was derived in [19] for the approximation of bang-bang control of a semilinear equation and could be adapted to the multibang control setting. Of particular interest would be the extension to problems where the control enters into the principal part of an elliptic equation as in the case of topology optimization problems [5,7].

On the other hand, the applicability of the multibang penalty G to the regularization of inverse problems was demonstrated in [3]. There, a condition related to Assumption REG was used to derive strong convergence as $\alpha \rightarrow 0$, albeit without rates; and a natural question is whether the more quantitative Assumption REG would allow obtaining such rates at least in $L^2(\Omega)$. Finally, combined regularization, approximation, and discretization estimates for the convergence $(\alpha, \gamma, h) \rightarrow 0$ would be highly useful.

Acknowledgements This work was funded by the German Research Foundation (DFG) under Grants Cl 487/1-1 and Wa 3626/1-1.

References

1. Alnæs, M.S., Blechta, J., Hake, J., Johansson, A., Kehlet, B., Logg, A., Richardson, C., Ring, J., Rognes, M.E., Wells, G.N.: The FEniCS project version 1.5. Arch. Numer. Softw. **3**(100), 9–23 (2015). <https://doi.org/10.11588/ans.2015.100.20553>
2. Barbu, V., Precupanu, T.: Convexity and Optimization in Banach Spaces, fourth edn. Springer Monographs in Mathematics. Springer, Dordrecht (2012). <https://doi.org/10.1007/978-94-007-2247-7>
3. Clason, C., Do, T.B.T.: Convex regularization of discrete-valued inverse problems. In: Hofmann, B., Leitão, A., Zubelli, J. (eds.) New Trends in Parameter Identification for Mathematical Models, Trends in Mathematics. Springer, Berlin (2018). https://doi.org/10.1007/978-3-319-70824-9_2
4. Clason, C., Ito, K., Kunisch, K.: A convex analysis approach to optimal controls with switching structure for partial differential equations. ESAIM Control Optim. Calc. Var. **22**(2), 581–609 (2016). <https://doi.org/10.1051/cocv/2015017>
5. Clason, C., Kruse, F., Kunisch, K.: Total variation regularization of multi-material topology optimization. ESAIM Math. Model. Numer. Anal. **52**(1), 275–303 (2018). <https://doi.org/10.1051/m2an/2017061>
6. Clason, C., Kunisch, K.: Multi-bang control of elliptic systems. Annales de l'Institut Henri Poincaré (C) Analyse Non Linéaire **31**(6), 1109–1130 (2014). <https://doi.org/10.1016/j.anihpc.2013.08.005>
7. Clason, C., Kunisch, K.: A convex analysis approach to multi-material topology optimization. ESAIM Math. Modell. Numer. Anal. **50**(6), 1917–1936 (2016). <https://doi.org/10.1051/m2an/2016012>
8. Clason, C., Tameling, C., Wirth, B.: Vector-valued multibang control of differential equations. SIAM J. Control Optim. **56**(3), 2295–2326 (2018). <https://doi.org/10.1137/16M1104998>
9. Deckelnick, K., Hinze, M.: A note on the approximation of elliptic control problems with bang-bang controls. Comput. Optim. Appl. **51**(2), 931–939 (2012). <https://doi.org/10.1007/s10589-010-9365-z>
10. Hintermüller, M., Hinze, M.: Moreau-Yosida regularization in state constrained elliptic control problems: error estimates and parameter adjustment. SIAM J. Numer. Anal. **47**(3), 1666–1683 (2009). <https://doi.org/10.1137/080718735>
11. Hintermüller, M., Schiela, A., Wollner, W.: The length of the primal-dual path in Moreau–Yosida-based path-following methods for state constrained optimal control. SIAM J. Optim. **24**(1), 108–126 (2014). <https://doi.org/10.1137/120866762>
12. Hinze, M.: A variational discretization concept in control constrained optimization: The linear-quadratic case. Comput. Optim. Appl. **30**(1), 45–61 (2005). <https://doi.org/10.1007/s10589-005-4559-5>
13. Hinze, M., Matthes, U.: A note on variational discretization of elliptic Neumann boundary control. Control Cybern. **38**(3), 577–591 (2009)

14. Ito, K., Kunisch, K.: Lagrange multiplier approach to variational problems and applications. In: *Advances in Design and Control*, vol. 15. SIAM, Philadelphia, PA (2008). <https://doi.org/10.1137/1.9780898718614>
15. Logg, A., Mardal, K.A., Wells, G.N.: *Automated Solution of Differential Equations by the Finite Element Method*. Springer, Berlin (2012). <https://doi.org/10.1007/978-3-642-23099-8>
16. Logg, A., Wells, G.N.: Dofin: Automated finite element computing. *ACM Trans. Math. Softw.* (2010). <https://doi.org/10.1145/1731022.1731030>
17. Logg, A., Wells, G.N., Hake, J.: DOLFIN: a C++/Python Finite Element Library, chap. 10. Springer (2012). https://doi.org/10.1007/978-3-642-23099-8_10
18. Pörner, F., Wachsmuth, D.: An iterative Bregman regularization method for optimal control problems with inequality constraints. *Optimization* **65**(12), 2195–2215 (2016). <https://doi.org/10.1080/02331934.2016.1238082>
19. Pörner, F., Wachsmuth, D.: Tikhonov regularization of optimal control problems governed by semi-linear partial differential equations. *Math. Control Relat. Fields* **8**(1), 315–335 (2018). <https://doi.org/10.3934/mcrf.2018013>
20. Stadler, G.: Elliptic optimal control problems with L^1 -control cost and applications for the placement of control devices. *Comput. Optim. Appl.* **44**(2), 159–181 (2009). <https://doi.org/10.1007/s10589-007-9150-9>
21. Tröltzsch, F.: *Optimal control of partial differential equations: theory, methods and applications*. American Mathematical Society (2010). <https://doi.org/10.1090/gsm/112>. Translated from the German by Jürgen Sprekels
22. Ulbrich, M.: *Semismooth newton methods for variational inequalities and constrained optimization problems in function spaces*, MOS-SIAM Series on Optimization, vol. 11. SIAM, Philadelphia, PA (2011). <https://doi.org/10.1137/1.9781611970692>
23. Wachsmuth, D.: Adaptive regularization and discretization of bang-bang optimal control problems. *Electron. Trans. Numer. Anal.* **40**, 249–267 (2013)
24. Wachsmuth, D., Wachsmuth, G.: Regularization error estimates and discrepancy principle for optimal control problems with inequality constraints. *Control Cybern.* **40**(4), 1125–1158 (2011)
25. Wachsmuth, D., Wachsmuth, G.: Necessary conditions for convergence rates of regularizations of optimal control problems. In: *System modeling and optimization, IFIP Adv. Inf. Commun. Technol.*, vol. 391, pp. 145–154. Springer, Heidelberg (2013). https://doi.org/10.1007/978-3-642-36062-6_15
26. Wachsmuth, G., Wachsmuth, D.: Convergence and regularization results for optimal control problems with sparsity functional. *ESAIM Control Optim. Calc. Var.* **17**(3), 858–886 (2011). <https://doi.org/10.1051/cocv/2010027>