# A convex analysis approach to multi-material topology optimization

C. Clason          K. Kunisch

A–8010 GRAZ,  HEINRICHSTRASSE 36,  AUSTRIA

Der Wissenschaftsfonds.

SFB sponsors:

- **Austrian Science Fund** (FWF)
- **University of Graz**
- **Graz University of Technology**
- **Medical University of Graz**
- **Government of Styria**
- **City of Graz**

# A convex analysis approach to multi-material topology optimization

Christian Clason[*]    Karl Kunisch[†]

January 14, 2016

This work is concerned with optimal control of partial differential equations where the control enters the state equation as a coefficient and should take on values only from a given discrete set of values corresponding to available materials. A "multi-bang" framework based on convex analysis is proposed where the desired piecewise constant structure is incorporated using a convex penalty term. Together with a suitable tracking term, this allows formulating the problem of optimizing the topology of the distribution of material parameters as minimizing a convex functional subject to a (nonlinear) equality constraint. The applicability of this approach is validated for two model problems where the control enters as a potential and a diffusion coefficient, respectively. This is illustrated in both cases by numerical results based on a semi-smooth Newton method.

## 1 Introduction

In this work, topology optimization consists in determining the optimal distribution of two or more given materials within a domain, where the material properties enter as the values of a spatially varying coefficient $u(x)$ into the operator of a partial differential equation. We propose to follow a direct approach and minimize a cost functional of interest subject to the constraint $u(x) \in \{u_1, \dots, u_d\}$, where $u_i$ are given parameters specific to different materials. This constraint is realized by means of the penalty functional

$$\mathcal{G}_0(u) = \int_\Omega \frac{\alpha}{2} |u(x)|^2 + \beta \prod_{i=1}^{d} |u(x) - u_i|^0 \, dx,$$

where $|0|^0 = 0$ and $|t|^0 = 1$ for $t \neq 0$, and $\alpha$ and $\beta$ are fixed parameters to be further discussed below (see Corollary 2.3). This functional was analyzed in [10] in the context of linear optimal

---

[*]Faculty of Mathematics, University Duisburg-Essen, 45117 Essen, Germany (`christian.clason@uni-due.de`)

[†]Institute of Mathematics and Scientific Computing, University of Graz, Heinrichstrasse 36, 8010 Graz, Austria, and Radon Institute, Austrian Academy of Sciences, Linz, Austria (`karl.kunisch@uni-graz.at`).

control problems. There it was shown that, under mild technical assumptions, the solutions to optimal control problems based on the convex envelope $\mathcal{G}_\Gamma$ of $\mathcal{G}_0$ have the desired property of being exactly multi-bang. This means that the solutions assume values in $\{u_1, \ldots, u_d\}$ pointwise a.e. in the control domain, provided that $\beta$ is sufficiently large. This property is related to the use of the $\ell^1$ norm in sparse optimization as the convex envelope (on the unit interval) of the $\ell^0$ "norm". Although the explicit form of $\mathcal{G}_\Gamma$ is not needed in our approach, we compute it in section 3 and remark on its relation to a direct $L^1$-type penalization of the constraint $u(x) \in \{u_1, \ldots, u_d\}$.

In this work, we focus on tracking-type functionals for multi-material optimization, i.e., we consider the optimization problem

$$(1.1) \qquad \min_{u \in U} \frac{1}{2}\|S(u) - z\|_Y^2 + \mathcal{G}_\Gamma(u),$$

where

$$U = \left\{ u \in L^2(\Omega) : u(x) \in [u_1, u_d] \quad \text{for almost all } x \in \Omega \right\}$$

is the admissible set with $u_1 < \cdots < u_d$ given, $Y$ is a Hilbert space, $z \in Y$ is the given desired state, and $S : U \to Y$ is the (nonlinear) parameter-to-state mapping.

Following [10, 9], we can derive a first-order necessary primal-dual optimality system

$$(1.2) \qquad \begin{cases} -\bar{p} = S'(\bar{u})^*(S(\bar{u}) - z), \\ \bar{u} \in \partial\mathcal{G}_0^*(\bar{p}) \end{cases}$$

(where $\partial\mathcal{G}_0^*$ is the convex subdifferential of the (convex) Fenchel conjugate of $\mathcal{G}_0$), whose Moreau–Yosida regularization is amenable to numerical solution by a superlinearly convergent semismooth Newton method. While in earlier works, we considered the case of linear $S$, the main focus here is on nonlinear, and in particular bilinear, parameter-to-state mappings. Our aim is to demonstrate that the proposed methodology provides a viable technology for solving multi-material shape and topology optimization problems without the need for computing shape or topological derivatives.

Let us very briefly point out some of the alternative approaches for topology optimization and give very selective references. Relaxation methods [1, 7, 20, 19] are amongst the earliest and most frequently used techniques. A standard approach for the two-material case consists in setting $u(x) = u_1 w(x) + u_2(1 - w(x))$ and minimizing over the set of all characteristic functions $w(x) \in \{0, 1\}$. This problem is non-convex, but its convex relaxation – minimizing over all $w(x) \in [0, 1]$ – often has a bang-bang solution, i.e., $w(x) \in \{0, 1\}$ almost everywhere. For multi-material optimization, this approach can be extended by introducing multiple characteristic functions; non-overlapping materials can be enforced by considering the third domain as an intersection of two (possibly overlapping) domains, e.g., $u(x) = u_1 w_1(x) + u_2(1 - w_1(x))w_2(x) + u_3(1 - w_1(x))(1 - w_2(x))$ for $w_1(x), w_2(x) \in [0, 1]$. For an increasing number $d$ of materials, this approach has obvious drawbacks due to the combinatorial nature and increasing non-linearity. Shape calculus techniques [20, 23] focus on the effect of smooth perturbations of the interfaces on the cost functional and have reached a high level of sophistication. From the point of view of numerical optimization, they are first-order methods and stable, with the drawback that they

mostly allow only smooth variations of the reference geometry. When combined with level-set techniques [2, 15], they are flexible enough to allow vanishing and merging of connected components, but they do not allow the creation of holes. This is allowed in the context of topological sensitivity analysis [12, 22], which investigates the effect of the creation of holes on the cost. Let us point out that in our work we do not rely in any explicit manner on knowledge of the shape or the topological derivatives. Moreover, the numerical technique that we propose is of second order rather than of gradient nature. Second-order shape or topological derivative analysis is available, but it is involved when it comes to numerical realization. Multi-material optimization for elasticity problems are further investigated in [13] by means of H-convergence methods and by phase-field methods in [8]. The work which in part is most closely related to ours is [4], see also [5, 3], where for the case of linear solution operators and two materials, the set of coefficients is expressed in terms of characteristic functions, and the resulting problem is considered in function spaces rather than in terms of subdomains and their boundaries. The first order-optimality condition is derived and formulated as a nonlinear equation for which a semi-smooth Newton method is applicable.

The general theory to be developed will be tested on two particular model problems. For the first one, the mapping $S : u \mapsto y \in H^2(\Omega)$ is the solution operator to

$$\begin{cases} -\Delta y + uy = f, \\ \qquad \partial_\nu y = 0, \end{cases}$$

for $u$ in an appropriate subset of $L^2(\Omega)$ and fixed $f \in L^2(\Omega)$. The second one is motivated by the mapping $\tilde{S} : u \mapsto y \in H_0^1(\Omega)$, where $y$ is the solution to

$$\begin{cases} -\nabla \cdot (u\nabla y) = f, \\ \qquad y = 0, \end{cases}$$

with $u$ in a subset of $L^\infty(\Omega)$. It is well known from [17] that (1.1) does not admit a solution in this case, since the differential equation is not closed under weak-$*$ convergence in $L^\infty(\Omega)$. For this reason we shall introduce a local smoothing operator $G$ and define the associated solution operator as $S = \tilde{S} \circ G$. We point out that the operator to be used in section 4 will be of local nature. It acts as smoothing of the constant values $u_i$ across interior interfaces of boundaries between different materials and will justify the use of a semi-smooth Newton method for the numerical realization.

This work is organized as follows. In section 2, existence of a solution to (1.1) is shown and the explicit form of (1.2) is derived. Section 3 is devoted to the explicit form of $G$ and its comparison to an alternative $L^1$-type penalty. The numerical solution is addressed in section 4, where the Moreau–Yosida regularization and its convergence are treated for general nonlinear mappings in section 4.1. The analysis of the semismooth Newton method for the regularized problems requires specific properties of the state equation and is therefore addressed in 4.2 separately for each model problem. Finally, numerical results are presented in section 5.

## 2 Existence and optimality conditions

We set

$$\mathcal{F} : L^2(\Omega) \to \overline{\mathbb{R}}, \qquad \mathcal{F}(u) = \frac{1}{2}\|S(u) - z\|_Y^2,$$

$$\mathcal{G}_0 : L^2(\Omega) \to \overline{\mathbb{R}}, \qquad \mathcal{G}_0(u) = \frac{\alpha}{2}\|u\|_{L^2}^2 + \beta \int_\Omega \prod_{i=1}^{d} |u(x) - u_i|^0 \, dx + \delta_U(u),$$

where $U \subset L^2(\Omega)$ is a convex and closed set $U \subset L^2(\Omega)$ and $\delta_U$ is the indicator function in the sense of convex analysis, i.e.,

$$(2.1) \qquad \delta_U(u) = \begin{cases} 0 & \text{if } u \in U, \\ \infty & \text{if } u \notin U. \end{cases}$$

For $S : U \to Y$, we assume that

(A1)   $S : U \to Y$ is weak-to-weak continuous, i.e., $\{u_n\}_{n \in \mathbb{N}} \subset U$ and $u_n \rightharpoonup u \in U$ in $L^2(\Omega)$ implies $S(u_n) \rightharpoonup S(u) \in Y$;

(A2)   $S$ is twice Fréchet differentiable.

Both assumptions are satisfied for the two model problems stated in the introduction. Now consider

$$(2.2) \qquad \min_{u \in L^2(\Omega)} \mathcal{F}(u) + \mathcal{G}(u)$$

for

$$\mathcal{G} := \mathcal{G}_0^{**},$$

where $\mathcal{G}_0^{**}$ is the biconjugate of $\mathcal{G}_0$, i.e., the Fenchel conjugate of

$$\mathcal{G}_0^* : L^2(\Omega) \to \overline{\mathbb{R}}, \qquad \mathcal{G}_0^*(q) = \sup_{u \in L^2(\Omega)} \langle q, u \rangle - \mathcal{G}_0(u).$$

Since Fenchel conjugates are always lower semicontinuous and convex, see, e.g. [6, Proposition 13.11], it follows that $\mathcal{G}$ is proper, lower semicontinuous and convex for any $\alpha > 0$ and $\beta \geq 0$. Existence of a solution to (1.1) thus follows under the stated assumptions on $S$.

**Proposition 2.1.** *There exists a solution $\bar{u} \in U$ to (1.1) for any $\alpha > 0$ and $\beta \geq 0$.*

*Proof.* Due to Assumption (A1), the tracking term $\mathcal{F}$ is weakly lower semicontinuous and bounded from below. Similarly, $\mathcal{G}_0$ is bounded from below by 0, which implies that $\mathcal{G}_0^{**} \geq 0$ as well, see, e.g. [6, Proposition 13.14]. Since $U$ is a compact subset of $L^2(\Omega)$, we have

$$U = \operatorname{dom}\mathcal{G}_0 \subset \operatorname{dom}\mathcal{G}_0^{**} \subset \overline{\operatorname{dom}\mathcal{G}_0} = \overline{U} = U,$$

see, e.g., [6, Proposition 13.40], and hence that $\mathcal{G} = \mathcal{G}_0^{**}$ is coercive. This implies that $\mathcal{F} + \mathcal{G}$ is proper, weakly lower semicontinous and coercive, and application of Tonelli's direct method yields existence of a minimizer. $\qquad \square$

We next derive first-order necessary optimality conditions of primal-dual type.

**Proposition 2.2.** *Let $\bar{u} \in U$ be a local minimizer of (2.2). Then there exists a $\bar{p} \in L^2(\Omega)$ satisfying*

(2.3)
$$\begin{cases} -\bar{p} = S'(\bar{u})^*(S(\bar{u}) - z), \\ \bar{u} \in \partial \mathcal{G}^*(\bar{p}). \end{cases}$$

*Proof.* Let $\bar{u} \in U$ be a local minimizer, i.e., for $t > 0$ small enough and any $u \in U$ there holds

(2.4)
$$\mathcal{F}(\bar{u}) + \mathcal{G}(\bar{u}) \leq \mathcal{F}(\bar{u} + t(u - \bar{u})) + \mathcal{G}(\bar{u} + t(u - \bar{u})).$$

Since $\mathcal{G}$ is convex, we have

$$\mathcal{G}(\bar{u} + t(u - \bar{u})) = \mathcal{G}(tu + (1 - t)\bar{u}) \leq t\mathcal{G}(u) + (1 - t)\mathcal{G}(\bar{u}),$$

which implies

$$\mathcal{G}(tu + (1 - t)\bar{u}) - \mathcal{G}(\bar{u}) \leq t(\mathcal{G}(u) - \mathcal{G}(\bar{u})).$$

Inserting this in (2.4) and rearranging yields

$$\mathcal{F}(\bar{u} + t(u - \bar{u})) - \mathcal{F}(\bar{u}) + t(\mathcal{G}(u) - \mathcal{G}(\bar{u})) \geq 0.$$

Since $\mathcal{F}$ is Fréchet-differentiable due to Assumption (A2), we can divide by $t > 0$ and let $t \to 0$ to obtain

$$\langle \mathcal{F}'(\bar{u}), u - \bar{u} \rangle + \mathcal{G}(u) - \mathcal{G}(\bar{u}) \geq 0$$

for every $u \in U$, i.e.,

$$\bar{p} := -\mathcal{F}'(\bar{u}) \in \partial \mathcal{G}(\bar{u}).$$

Since $\mathcal{G}$ is convex, this is equivalent to $\bar{u} \in \partial \mathcal{G}^*(\bar{p})$. Applying the chain rule for Fréchet derivatives to $\mathcal{F}$ then yields the desired optimality conditions. $\square$

The question of optimality of solutions to Problem (2.2) with respect to the non-convex functional $\mathcal{F} + \mathcal{G}_0$ has been addressed (for linear $S$) in [10]; here we only remark that since $\mathcal{G} = \mathcal{G}_0^{**} \leq \mathcal{G}_0$ and $\mathcal{G}(u) = \mathcal{G}_0(u)$ for $u(x) \in \{u_1, \ldots, u_d\}$ almost everywhere (see section 3 below), it follows that if a (local) minimizer $\bar{u}$ of (2.2) satisfies $\bar{u}(x) \in \{u_1, \ldots, u_d\}$ almost everywhere, we have for all $u \in U$ (sufficiently close to $\bar{u}$) that

$$\mathcal{F}(u) + \mathcal{G}_0(u) \geq \mathcal{F}(u) + \mathcal{G}(u) \geq \mathcal{F}(\bar{u}) + \mathcal{G}(\bar{u}) = \mathcal{F}(\bar{u}) + \mathcal{G}_0(\bar{u}),$$

i.e., $\bar{u}$ is a (local) minimizer of $\mathcal{F} + \mathcal{G}_0$ as well.

Since $\mathcal{G}^* = (\mathcal{G}_0^{**})^* = \mathcal{G}_0^{***} = \mathcal{G}_0^*$, see, e.g., [6, Proposition 13.14 (iii)], we can make use of the following characterization from [10, § 2.1].

**Corollary 2.3.** *If $\alpha$ and $\beta$ satisfy the relation*

$$(2.5) \qquad \frac{\alpha}{2}(u_{i+1} - u_i) \leq \sqrt{2\alpha\beta} \quad \text{for all } 1 \leq i < d,$$

*then $u \in \partial\mathcal{G}^*(p)$ if and only if for almost all $x \in \Omega$,*

$$(2.6) \qquad u(x) \in \begin{cases} \{u_1\} & p(x) < \frac{\alpha}{2}(u_1 + u_2), \\ \{u_i\} & \frac{\alpha}{2}(u_{i-1} + u_i) < p(x) < \frac{\alpha}{2}(u_i + u_{i+1}), \qquad 1 < i < d, \\ \{u_d\} & p(x) > \frac{\alpha}{2}(u_{d-1} + u_d), \\ [u_i, u_{i+1}] & p(x) = \frac{\alpha}{2}(u_i + u_{i+1}), \qquad 1 \leq i < d. \end{cases}$$

Thus, with (2.5) holding, $u(x)$ coincides with one of the preassigned control values $u_i$, except in the singular cases when $p(x) = \frac{\alpha}{2}(u_i + u_{i+1})$ for some $i$. If, on the other hand, (2.5) is not satisfied, then $u = \frac{1}{\alpha}p$ may hold on subsets $\hat{\Omega}$ of nontrivial measure. In this case we call $u|_{\hat{\Omega}}$ a *free arc*, and refer to [10] for details.

# 3 Relation to $L^1$ penalization

We now compare the penalty $\mathcal{G}$ to a direct $L^1$ penalization of $u(x) - u_i$, $i \in \{1, \ldots, d\}$. First, we give an explicit characterization of $\mathcal{G} = \mathcal{G}_0^{**}$. Since $\mathcal{G}_0$ is defined via the integral of a pointwise function of $u(x)$, we can compute the Fenchel conjugate and its subdifferential pointwise as well; see, e.g., [11, Props. IV.1.2, IX.2.1], [6, Prop. 16.50]. It therefore suffices to consider

$$g_0 : \mathbb{R} \to \overline{\mathbb{R}}, \qquad g_0(v) = \frac{\alpha}{2}|v|^2 + \beta \prod_{i=1}^{d} |v - u_i|^0 + \delta_{[u_1, u_d]}(v),$$

where $\delta_{[u_1, u_d]}$ is again the indicator function in the sense of convex analysis, cf. (2.1).

To compute $g_0^{**}$ we make use of the fact that the biconjugate coincides with the lower convex envelope (or Gamma-regularization)

$$g_\Gamma(v) = \sup\{a(v) : a : \mathbb{R} \to \mathbb{R} \text{ is affine and } a \leq g_0\},$$

see, e.g., [21, Theorem 2.2.4 (a)]. We assume again that (2.5) holds.

First, note that $g_0(u_i) = \frac{\alpha}{2}u_i^2$ for all $1 \leq i \leq d$, which implies that $g_\Gamma(u_i) \leq \frac{\alpha}{2}u_i^2$. Now consider a single interval $[u_i, u_{i+1}]$ for $1 \leq i < d$. Obviously, a candidate for $g_\Gamma(v)$ in $v \in \{u_i, u_{i+1}\}$ is given by the linear interpolant $g_i$ of $g_0(u_i)$ and $g_0(u_{i+1})$, i.e.,

$$g_i(v) = \frac{\alpha}{2}\left((u_i + u_{i+1})v - u_i u_{i+1}\right).$$

This function in fact satisfies the conditions for $g_\Gamma$ also for $v \in (u_i, u_{i+1})$, which follows from the fact that on this open interval, the quadratic function

$$(g_0 - g_i)(v) = \frac{\alpha}{2}\left(v^2 - (u_i + u_{i+1})v + u_i u_{i+1}\right) + \beta$$

6

has a unique minimizer (since $\alpha > 0$) in its critical point $\bar{v} = \frac{1}{2}(u_i + u_{i+1})$, where

$$(g_0 - g_i)(v) = \frac{\alpha}{2}\left(-\frac{1}{4}(u_i + u_{i+1})^2 + u_i u_{i+1}\right) + \beta$$
$$= -\frac{\alpha}{8}(u_{i+1} - u_i)^2 + \beta \geq 0$$

by (2.5). Hence, $g_i(v) \leq g_0(v)$ for all $v \in [u_i, u_{i+1}]$ with equality in $v \in \{u_i, u_{i+1}\}$.

To obtain a global function, we define $\bar{g} : [u_1, u_d] \to \mathbb{R}$ via

$$\bar{g}(v) := g_i(v) \qquad \text{for } v \in [u_i, u_{i+1}], \quad 1 \leq i < d.$$

It remains to verify that for each fixed $i$, we have $g_j(v) \leq g_i(v)$ for all $j \neq i$ and $v \in [u_i, u_{i+1}]$. A short computation shows that $g_j(u_i) \leq g_i(u_i)$. Moreover, due to the ordering of the $u_i$ we have

$$g_j'(v) = \frac{\alpha}{2}(u_j + u_{j+1}) > \frac{\alpha}{2}(u_{i+1} + u_{i+2}) = g_i'(v)$$

for all $j > i$ and similarly $g_i'(v) < g_j'(v)$ for all $j < i$. This implies that $g_j(v) \leq g_i(v)$ for all $j \neq i$ and $v \in [u_i, u_{i+1}]$. Using again that $\operatorname{dom} g_\Gamma = \operatorname{dom} g_0 = [u_1, u_d]$ since the interval is closed, we obtain

$$g_0^{**}(v) = g_\Gamma(v) = \bar{g}(v) + \delta_{[u_1, u_d]}(v)$$
$$= \begin{cases} \frac{\alpha}{2}\left((u_i + u_{i+1})v - u_i u_{i+1}\right) & v \in [u_i, u_{i+1}], \quad 1 \leq i < d, \\ \infty & v \in \mathbb{R} \setminus [u_1, u_d]. \end{cases}$$

and hence

$$G(u) = \int_\Omega g_\Gamma(u(x)) \, dx.$$

From the above, we have that $g_\Gamma$ is the unique continuous and piecewise (on $[u_i, u_{i+1}]$) affine function with $g_\Gamma(u_i) = \frac{\alpha}{2}u_i^2$. It is not surprising that using such a function in optimization promotes solutions lying in the "kinks" (cf. sparse optimization using $\ell_1$-type norms, where the only "kink" is at $v = 0$). Other penalties $h$ with a similar piecewise affine structure can be constructed by prescribing different values for $h(u_i)$, although the obvious choice $h(u_i) = \alpha|u_i|$ results in a shifted $\ell_1$ norm which has only one "kink" at $v = \min_i |u_i|$ and hence does not have the desired structure.

An alternative to this piecewise affine construction is the direct $\ell^1$-penalization of the deviation, i.e., choosing

$$h(v) = \alpha \sum_{i=1}^{d} |v - u_i| + \delta_{[u_1, u_d]}(v).$$

(Note that the product $\prod_{i=1}^d |v - u_i|$ is a polynomial of order $d$ and hence in general is not convex.) We first point out that the value $h(u_i)$ depends on all $u_j$, $1 \leq j \leq d$, (and in particular, on $d$) rather than on $u_i$ only, which may be undesirable; see Figure 1. To further illustrate the
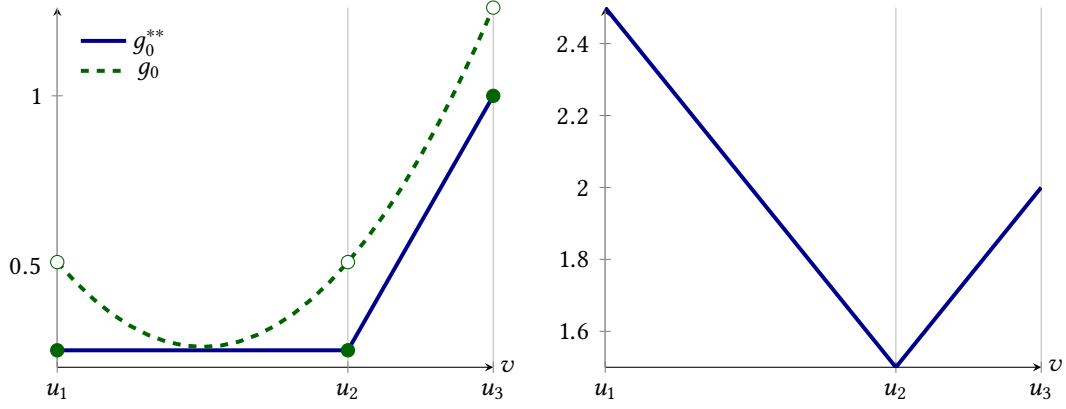
Figure 1: Plot of $g_0^{**}$ and $g_0$ (left), $h$ (right) for $d = 3$, $(u_1, u_2, u_3)$, $\alpha = 0.5$, $\beta = 0.26$ (satisfying (2.5))

practical difference between using $g_\Gamma$ and $h$, we compute the corresponding subdifferential $\partial h^*$ which would appear in (2.3). First, we determine the Fenchel conjugate

$$h^*(q) = \sup_{v \in [u_1, u_d]} vq - \alpha \sum_{i=1}^{d} |v - u_i|.$$

Since the function to be maximized is continuous and piecewise affine on $\mathbb{R}$, the supremum must be attained at $\bar{v} = u_i$ for some $1 \leq i \leq d$. Making use of the fact that the $u_i$ are ordered, we obtain that $h^*(q)$ must be equal to one of the functions

$$h_i^*(q) = qu_i - \alpha\left(\sum_{j=1}^{i-1}(u_i - u_j) + \sum_{j=i+1}^{d}(u_j - u_i)\right)$$

$$= u_i(q + \alpha(d + 1 - 2i)) + \alpha\sum_{j=1}^{i-1} u_j - \alpha\sum_{j=i+1}^{d} u_j$$

(with the convention that empty sums evaluate to 0). It remains to determine the supremum over $1 \leq i \leq d$ based on the value of $q$. For this, we first compare $h_i^*(q)$ with $h_{i+1}^*(q)$. Simple rearrangement of terms shows that $h_i^*(q) \leq h_{i+1}^*(q)$ if and only if

$$\alpha(2i - d)(u_{i+1} - u_i) \leq q(u_{i+1} - u_i).$$

Since $u_{i+1} > u_i$, we deduce that this is the case if and only if $q \geq \alpha(2i - d)$. Hence, the supremum is attained for the largest $i$ for which $q \geq \alpha(2i - d)$. This yields

$$h^*(q) = \begin{cases} u_1(q + \alpha(d - 1)) - \alpha\sum_{j=2}^{d} u_j & \frac{1}{\alpha}q < 2 - d, \\ u_i(q + \alpha(d + 1 - 2i)) - \alpha\sum_{j=1}^{i-1} u_j + \alpha\sum_{j=i+1}^{d} u_j & 2(i-1) - d \leq \frac{1}{\alpha}q < 2i - d, \ 1 < i < d, \\ u_d(q - \alpha(d + 1)) + \alpha\sum_{j=1}^{d-1} u_j & \frac{1}{\alpha}q \geq d - 2. \end{cases}$$
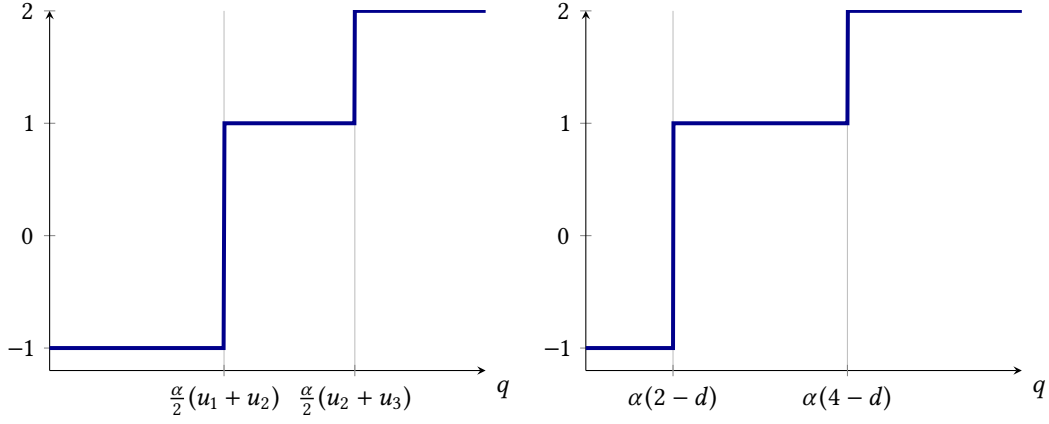
8

Figure 2: Plot of $\partial g^*$ (left), $\partial h^*$ (right) for $d = 3$, $(u_1, u_2, u_3)$, $\alpha = 0.5$, $\beta = 0.26$

Since $h^*$ is continuous and piecewise differentiable, we have that the convex subdifferential is given by

$$\partial h^*(q) = \begin{cases} \{u_1\} & \frac{1}{\alpha}q < 2 - d, \\ \{u_i\} & 2(i-1) - d < \frac{1}{\alpha}q < 2i - d, \quad 1 < i < d, \\ \{u_d\} & \frac{1}{\alpha}q > d - 2, \\ [u_i, u_{i+1}] & \frac{1}{\alpha}q = 2i - d, \quad 1 \le i < d. \end{cases}$$

Comparing this with Corollary 2.3, we see that the case distinction is independent of $u_i$, but rather depends on $d$ only, with the individual cases always being intervals of length $2\alpha$. In particular, for fixed $q$, the value $\partial h^*(q)$ changes if the number of parameters $d$ is increased, independent of the magnitude of the additional parameters. Furthermore, since the distribution of intervals is symmetric around the origin, $h$ tends to favor for increasing $\alpha$ those $u_i$ closer to the "middle parameter" $u_{d/2}$, rather than those of smaller magnitude as is the case for $g_0^{**}$; see Figure 2.

## 4 Numerical solution

For the numerical solution, we follow the approach described in [9] for linear parameter-to-state mappings, where we replace $\partial \mathcal{G}^*$ by its Moreau–Yosida regularization and apply a semi-smooth Newton method with backtracking line search and continuation. In this section, we describe the necessary modifications for nonlinear mappings, arguing in terms of the functional instead of the optimality system. We first introduce the regularization and discuss its convergence to the original problem for general nonlinear mappings in section 4.1. The explicit form and well-posedness of the Newton step (from which superlinear convergence follows) requires exploiting the structure of the mapping, hence we discuss it separately for each model problem in section 4.2.

## 4.1 Regularization

Since $\mathcal{F}$ is not convex, we cannot proceed directly to the regularized system. Instead, we start by considering for $\gamma > 0$ the regularized problem

$$(4.1) \qquad \min_{u \in L^2(\Omega)} \mathcal{F}(u) + \mathcal{G}(u) + \frac{\gamma}{2}\|u\|_{L^2(\Omega)}^2.$$

By the same arguments as in the proof of Proposition 2.1, we obtain the existence of a minimizer $u_\gamma \in U$. We now address convergence of $u_\gamma$ as $\gamma \to 0$.

**Proposition 4.1.** *The family $\{u_\gamma\}_{\gamma>0}$ of global minimizers to (4.1) contains at least one subsequence $\{u_{\gamma_n}\}_{n \in \mathbb{N}}$ converging to a global minimizer of (2.2) as $n \to \infty$. Furthermore, for any such subsequence the convergence is strong.*

*Proof.* Since $U$ is bounded, the set $\{u_\gamma\}_{\gamma>0}$ contains a subsequence $\{u_{\gamma_n}\}_{n \in \mathbb{N}}$ with $\gamma_n \to 0$ converging weakly to some $\bar{u}$. Furthermore, it follows that $\lim_{n \to \infty} \frac{\gamma_n}{2}\|u_{\gamma_n}\|_{L^2(\Omega)}^2 = 0$. By the weak lower semicontinuity of $\mathcal{J} := \mathcal{F} + \mathcal{G}$ and the optimality of $u_{\gamma_n}$, we thus have for any $u \in U$ that

$$\mathcal{J}(\bar{u}) \le \liminf_{n \to \infty} \mathcal{J}(u_{\gamma_n}) = \liminf_{n \to \infty} \mathcal{J}(u_{\gamma_n}) + \frac{\gamma_n}{2}\|u_{\gamma_n}\|_{L^2(\Omega)}^2$$
$$\le \mathcal{J}(u) + \lim_{n \to \infty} \frac{\gamma_n}{2}\|u\|_{L^2(\Omega)}^2 = \mathcal{J}(u),$$

i.e., $\bar{u}$ is a global minimizer of (2.2).

To show strong convergence, it suffices to show $\limsup_{n \to \infty} \|u_{\gamma_n}\| \le \|\bar{u}\|$. This follows from

$$\mathcal{J}(u_{\gamma_n}) + \frac{\gamma_n}{2}\|u_{\gamma_n}\|_{L^2(\Omega)}^2 \le \mathcal{J}(\bar{u}) + \frac{\gamma_n}{2}\|\bar{u}\|_{L^2(\Omega)}^2 \le \mathcal{J}(u_{\gamma_n}) + \frac{\gamma_n}{2}\|\bar{u}\|_{L^2(\Omega)}^2$$

for every $n \in \mathbb{N}$ due to the optimality of $u_\gamma$ and $\bar{u}$. Hence, $\|u_{\gamma_n}\|_{L^2(\Omega)} \to \|\bar{u}\|_{L^2(\Omega)}$, which together with weak convergence implies strong convergence in the Hilbert space $L^2(\Omega)$ of the subsequence. $\qquad\qquad\square$

Arguing as in the proof of Proposition 2.2, we obtain the abstract first-order necessary optimality conditions

$$\begin{cases} -p_\gamma = \mathcal{F}'(u_\gamma), \\ \quad u_\gamma \in \partial(\mathcal{G}_\gamma)^*(p_\gamma), \end{cases}$$

where

$$\mathcal{G}_\gamma(u) := \mathcal{G}(u) + \frac{\gamma}{2}\|u\|_{L^2(\Omega)}^2.$$

We now use that $(\mathcal{G} + \frac{\gamma}{2}\|\cdot\|_{L^2(\Omega)}^2)^*$ is equal to the infimal convolution of $\mathcal{G}^*$ and $\frac{1}{2\gamma}\|\cdot\|_{L^2(\Omega)}^2$, which in turn coincides with the Moreau envelope of $\mathcal{G}^*$; see, e.g., [6, Proposition 13.21]. Furthermore, the Moreau envelope is Fréchet-differentiable with Lipschitz-continuous gradient which coincides with the Moreau–Yosida regularization $(\partial \mathcal{G}^*)_\gamma$ of $\partial \mathcal{G}^*$; see, e.g., [6, Proposition 12.29]. We can

therefore make use of the pointwise characterization of $H_\gamma := (\partial \mathcal{G}^*)_\gamma = \partial (\mathcal{G}_\gamma)^*$ from [9, Appendix A.2], assuming again that (2.5) holds, to obtain

$$(4.2) \qquad [H_\gamma(p)](x) = \begin{cases} u_i & p(x) \in Q_i^\gamma, & 1 \le i \le d, \\ \frac{1}{\gamma}\left(p(x) - \frac{\alpha}{2}(u_i + u_{i+1})\right) & p(x) \in Q_{i,i+1}^\gamma, & 1 \le i < d. \end{cases}$$

where

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

to obtain the explicit primal-dual first-order necessary conditions

$$(4.3) \qquad \begin{cases} -p_\gamma = S'(u_\gamma)^*(S(u_\gamma) - z), \\ \phantom{-}u_\gamma = H_\gamma(p_\gamma). \end{cases}$$

Comparing (4.2) to (2.6), we observe that the Moreau–Yosida regularization is of local nature, acting along interfaces between regions with different material parameters.

Since $H_\gamma$ is a superposition operator defined by a Lipschitz continuous and piecewise differentiable scalar function, $H_\gamma$ is Newton-differentiable from $L^r(\Omega) \to L^2(\Omega)$ for any $r > 2$; see, e.g., [14, Example 8.12] or [24, Theorem 3.49]. Its Newton derivative at $p$ in direction $h$ is given pointwise almost everywhere by

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

## 4.2 Semismooth Newton method

We now wish to apply a semismooth Newton method to (4.3). For this purpose, we need to argue that $p_\gamma \in V$ for some $V \hookrightarrow L^r(\Omega)$ with $r > 2$ and show uniform invertibility of the Newton step. Since the control-to-state mapping is nonlinear, this requires exploiting its concrete structure. We thus directly consider the specific model problems.

### 4.2.1 Potential problem

We first express (4.3) in equivalent form by introducing the state $y_\gamma = S(u_\gamma) \in H^1(\Omega)$, i.e., satisfying for $u = u_\gamma$

$$(4.4) \qquad \begin{cases} -\Delta y + uy = f & \text{in } \Omega, \\ \partial_\nu y = 0 & \text{on } \partial\Omega. \end{cases}$$

In the following, we assume that $\Omega \subset \mathbb{R}^N$, $N \leq 3$, is sufficiently regular such that for any $f \in L^2(\Omega)$ and any $u \in U = U_M := \{u \in L^2(\Omega) : u_1 \leq u \leq M \text{ a.e.}\}$, the solution to (4.4) satisfies $y \in H^2(\Omega)$ together with the uniform a priori estimate

$$(4.5) \qquad \qquad \|y\|_{H^2(\Omega)} \leq C_M \|f\|_{L^2(\Omega)}.$$

We also consider for given $u \in U_M$ and $y \in H^2(\Omega)$ the adjoint equation

$$(4.6) \qquad \qquad \begin{cases} -\Delta w + uw = -(y - z) & \text{in } \Omega, \\ \partial_\nu w = 0 & \text{on } \partial\Omega, \end{cases}$$

whose solution $w \in H^2(\Omega)$ also satisfies the uniform a priori estimate (4.5). Due to the Sobolev embedding theorem, we have that the solutions $y$ and $w$ are also bounded in $L^\infty(\Omega)$ uniformly with respect to $u \in U_M$.

By standard Lagrangian calculus, we can now write $p_\gamma = y_\gamma w_\gamma$, where $w_\gamma \in H^1(\Omega)$ is the solution to (4.6) with $u = u_\gamma$ and $y = y_\gamma$. We further eliminate $u_\gamma$ using the second equation of (4.3) to obtain the reduced system

$$(4.7) \qquad \qquad \begin{cases} -\Delta w_\gamma + H_\gamma(-y_\gamma w_\gamma)w_\gamma + y_\gamma = z, \\ -\Delta y_\gamma + H_\gamma(-y_\gamma w_\gamma)y_\gamma = f. \end{cases}$$

Due the regularity of $y_\gamma$ and $p_\gamma$, we can consider this as an equation in $L^2(\Omega) \times L^2(\Omega)$ for $(y_\gamma, p_\gamma) \in H^2(\Omega) \times H^2(\Omega)$. By the Sobolev embedding theorem, we have $y_\gamma w_\gamma \in L^\infty(\Omega)$, and hence that the system (4.7) is semismooth. By the chain rule, the Newton derivative of $H_\gamma(-yw)$ with respect to $y$ in direction $\delta y$ is given by

$$D_{N,y} H_\gamma(-yw)\delta y = -\frac{1}{\gamma}\chi(-yw)\, w\, \delta y,$$

where $\chi(-yw)$ is the characteristic function of the inactive set

$$\mathcal{S}_\gamma(-yw) := \bigcup_{i=1}^{d-1}\left\{x \in \Omega : -y(x)w(x) \in Q_{i,i+1}^\gamma\right\}.$$

Similarly,

$$D_{N,w} H_\gamma(-yw)\delta w = -\frac{1}{\gamma}\chi(-yw)\, y\, \delta w.$$

For convenience, we set $\chi^k := \chi(-y^k w^k)$. A Newton step consists in solving

$$(4.8) \quad \begin{pmatrix} 1 - \frac{1}{\gamma}\chi^k(w^k)^2 & -\Delta + H_\gamma(-y^k w^k) - \frac{1}{\gamma}\chi^k y^k w^k \\ -\Delta + H_\gamma(-y^k w^k) - \frac{1}{\gamma}\chi^k y^k w^k & -\frac{1}{\gamma}\chi^k(y^k)^2 \end{pmatrix}\begin{pmatrix} \delta y \\ \delta w \end{pmatrix}$$
$$= -\begin{pmatrix} -\Delta w^k + H_\gamma(-y^k w^k)w^k + y^k - z \\ -\Delta y^k + H_\gamma(-y^k w^k)y^k - f \end{pmatrix}$$

and setting $y^{k+1} = y^k + \delta y$ and $w^{k+1} = w^k + \delta w$.

To show local superlinear convergence, it remains to prove uniformly bounded invertibility of (4.8). We proceed in several steps. First, we consider the off-diagonal terms in (4.8).

**Lemma 4.2.** *For any $\gamma > 0$ and $y, w \in H^2(\Omega)$, the linear operator $B : H^2(\Omega) \to L^2(\Omega)$,*

$$B = -\Delta + H_\gamma(-yw) - \frac{1}{\gamma}\chi(-yw)yw,$$

*is uniformly invertible, and there exists a constant $C > 0$ independent of $y, w$ such that*

$$\|B^{-1}\|_{\mathcal{L}(L^2(\Omega), H^2(\Omega))} \le C.$$

*Proof.* We first note that by definition, $[H_\gamma(p)](x) \in [u_1, u_d]$ for any $p \in L^2(\Omega)$. Furthermore, on the inactive set $S_\gamma(-yw)$ we have, again by definition,

$$u_1 \le \frac{\alpha}{2\gamma}(u_1 + u_2) + u_1 \le \frac{1}{\gamma}(-yw)(x) \le \frac{\alpha}{2\gamma}(u_{d-1} + u_d) + u_d \le (1 + \frac{\alpha}{\gamma})u_d.$$

Thus, $H_\gamma(-yw) - \frac{1}{\gamma}\chi(-yw)yw \in U_M$ for $M = (2 + \frac{\alpha}{\gamma})u_d$, and the claim follows from the a priori estimate (4.5). $\qquad\square$

**Proposition 4.3.** *For $\gamma > 0$, let $(y_\gamma, w_\gamma) \in H^2(\Omega) \times H^2(\Omega)$ be a solution to (4.7) with $w_\gamma$ satisfying $\|w_\gamma\|_{L^\infty(\Omega)} < \sqrt{\gamma}$. Furthermore, let $U(y_\gamma)$ be a bounded neighborhood of $y_\gamma$ in $H^2(\Omega)$, and let $U(w_\gamma)$ be a bounded neighborhood of $w_\gamma$ in $H^2(\Omega)$ such that $\|w\|_{L^\infty(\Omega)} \le \sqrt{\gamma}$ for any $w \in U(w_\gamma)$. Then there exists a constant $C > 0$ such that for any $(y, w) \in U(y_\gamma) \times U(w_\gamma)$ and any $r_1, r_2 \in L^2(\Omega)$, there exists a unique solution $(\delta y, \delta w) \in H^2(\Omega) \times H^2(\Omega)$ to*

$$(4.9) \qquad \begin{pmatrix} 1 - \frac{1}{\gamma}\chi(-yw)w^2 & B \\ B & -\frac{1}{\gamma}\chi(-yw)y^2 \end{pmatrix} \begin{pmatrix} \delta y \\ \delta w \end{pmatrix} = \begin{pmatrix} r_1 \\ r_2 \end{pmatrix}$$

*satisfying*

$$\|\delta y\|_{H^2(\Omega)} + \|\delta w\|_{H^2(\Omega)} \le C\left(\|r_1\|_{L^2(\Omega)} + \|r_2\|_{L^2(\Omega)}\right).$$

*Proof.* We exploit the invertibility of $B$ to obtain the required bounds on $\delta y$ and $\delta w$. For the sake of convenience, we set $\omega := S_\gamma(-yw)$ and $h := 1 - \frac{1}{\gamma}\chi(-yw)w^2$. As a first step, we introduce the following bilinear form on $L^2(\omega) \times L^2(\omega)$:

$$a_\omega(w_1, w_2) := (w_1, w_2)_{L^2(\omega)} + \left(hB^{-1}(\frac{1}{\sqrt{\gamma}}yE_\omega w_1), B^{-1}(\frac{1}{\sqrt{\gamma}}yE_\omega w_2)\right)_{L^2(\Omega)},$$

where $E_\omega$ denotes the extension by zero operator from $\omega$ to $\Omega$. Due to the assumption on $w$, we have that $h$ ia nonnegative. Thus the second term on the right hand side of the above equation is non-negative as well. Hence $a_\omega$ is symmetric, continuous and elliptic on $L^2(\omega)$ (uniformly on the set of admissible $(y, w)$). This implies the existence of a unique solution $\delta\tilde{w} \in L^2(\omega)$ to

$$(4.10) \qquad a_\omega(\delta\tilde{w}, \tilde{w}) = \left(\frac{1}{\sqrt{\gamma}}yB^{-1}\left(r_1 - hB^{-1}r_2\right), \tilde{w}\right)_{L^2(\omega)} \qquad \text{for all } \tilde{w} \in L^2(\omega)$$

satisfying

$$\|\delta\tilde{w}\|_{L^2(\omega)} \le C\left(\|r_1\|_{L^2(\Omega)} + \|r_2\|_{L^2(\Omega)}\right).$$

(Here and below, $C$ is a generic constant that may change its value between occurences but does not depend on $y$ and $w$.)

Next we consider the auxiliary equation

(4.11) $$B\delta y = r_2 + \frac{1}{\sqrt{\gamma}} y E_\omega \delta \tilde{w}.$$

From Lemma 4.2 we obtain a unique solution $\delta y \in H^2(\Omega)$ to (4.11) satisfying

$$\|\delta y\|_{H^2(\Omega)} \le C \left( \|r_2\|_{L^2(\Omega)} + \frac{1}{\sqrt{\gamma}} \|\delta \tilde{w}\|_{L^2(\omega)} \right) \le C \left( \|r_1\|_{L^2(\Omega)} + \|r_2\|_{L^2(\Omega)} \right),$$

using that $y \in U(y_\gamma)$ is uniformly bounded in $L^\infty(\Omega)$. Given $\delta y \in H^2(\Omega)$, the first equation of (4.9) now admits a unique solution $\delta w \in H^2(\Omega)$ satisfying

$$\|\delta w\|_{H^2(\Omega)} \le C \left( \|r_1\|_{L^2(\Omega)} + \|\delta y\|_{L^2(\Omega)} \right) \le C \left( \|r_1\|_{L^2(\Omega)} + \|r_2\|_{L^2(\Omega)} \right),$$

using the uniform boundedness of $w \in U(w_\gamma)$ in $L^\infty(\Omega)$.

To complete the proof, it remains to verify that $\delta w = \frac{1}{\sqrt{\gamma}} y \delta \tilde{w}$ on $\omega$. For this purpose we note that by the first of equation of (4.9) and (4.11),

$$\delta w + B^{-1} \left( h B^{-1} \left( \frac{1}{\sqrt{\gamma}} y E_\omega \delta \tilde{w} \right) \right) = B^{-1} \left( r_1 - h B^{-1} r_2 \right).$$

Taking the inner product of this equation in $L^2(\omega)$ with $\frac{1}{\gamma} y E_\omega w_2$ for arbitrary $w_2 \in L^2(\omega)$ and subtracting (4.10), we arrive at

$$\left( \frac{1}{\gamma} y \delta w - \delta \tilde{w}, w_2 \right)_{L^2(\omega)} = 0 \qquad \text{for all } w_2 \in L^2(\omega).$$

Inserting into (4.11) now verifies the second equation of (4.9). $\qquad\qquad\square$

We remark that according to the a priori estimate (4.5), the required smallness of $w_\gamma$ corresponds to smallness of the tracking error $\|y_\gamma - z\|_{L^2(\Omega)}$. In the following we give an alternative sufficient condition for the uniform continuous invertibility of the Newton iteration matrix (4.9) that does not rely on the smallness of $w_\gamma$. For this purpose, we set $\omega_\gamma := S_\gamma(-y_\gamma w_\gamma)$ and define

$$\partial \omega_\gamma := \bigcup_{i=1}^{d-1} \left\{ x \in \Omega : -y_\gamma(x) w_\gamma(x) \in \partial Q_{i,i+1}^\gamma \right\}.$$

We also introduce the compact self-adjoint operator

$$C : L^2(\omega_\gamma) \to L^2(\omega_\gamma), \qquad C = \left( B^{-1}(\frac{1}{\sqrt{\gamma}} y E_{\omega_\gamma}) \right)^* (h_\gamma \, \mathrm{Id}) \left( B^{-1}(\frac{1}{\sqrt{\gamma}} y E_{\omega_\gamma}) \right),$$

where $h_\gamma = 1 - \frac{1}{\gamma} \chi(-y_\gamma w_\gamma) w_\gamma^2$ and $B = B(y_\gamma, w_\gamma)$. We require the following two assumptions.

(H1)  $-1 \notin \sigma(C)$,

(H2)  $|\partial \omega_\gamma| = 0$.

14

**Proposition 4.4.** *For $\gamma > 0$, let $(y_\gamma, w_\gamma) \in H^2(\Omega) \times H^2(\Omega)$ be a solution to (4.7) satisfying (H1) and (H2). Then there exists a neighborhood $U(y_\gamma) \times U(w_\gamma)$ of $(y_\gamma, w_\gamma)$ in $H^2(\Omega) \times H^2(\Omega)$ such that the conclusion of Proposition 4.3 holds.*

*Proof.* By (H1) and as a consequence of the proof of Proposition 4.3, the system matrix in (4.9) is continuously invertible in $(y_\gamma, w_\gamma)$. Since the set of continuously invertible operators between Hilbert spaces is open with respect to the topology of the operator norm (see, e.g., [25, Theorem 6.2.3]), the claim will be established once we have argued that the system matrix, considered as an operator from $H^2(\Omega) \times H^2(\Omega)$ to $L^2(\Omega) \times L^2(\Omega)$, depends continuously in the operator norm on $(y, w) \in H^2(\Omega) \times H^2(\Omega)$ in a neighborhood of $(y_\gamma, w_\gamma)$. For this purpose, we first argue that $p := -yw \mapsto \chi(p)$ is continuous from $C(\overline{\Omega})$ to $L^2(\Omega)$ in a neighborhood of $p_\gamma := -y_\gamma w_\gamma$. For $\varepsilon > 0$ sufficiently small, we set

$$\partial \mathcal{S}_\gamma^\varepsilon := \bigcup_{i=1}^{d-1} \left\{ x \in \Omega : \text{dist}\left(p_\gamma(x), \partial Q_{i,i+1}^\gamma\right) < \varepsilon \right\}.$$

The family $\{\partial \mathcal{S}_\gamma^\varepsilon\}_{\varepsilon > 0}$ is monotone with respect to set inclusion and satisfies

$$\lim_{\varepsilon \to 0} \left|\partial \mathcal{S}_\gamma^\varepsilon\right| = \left|\lim_{\varepsilon \to 0} \partial \mathcal{S}_\gamma^\varepsilon\right| = |\partial \mathcal{S}_\gamma| = 0.$$

For any $\varepsilon > 0$ and any $p \in C(\overline{\Omega})$ such that $\|p - p_\gamma\|_{C(\overline{\Omega})} < \frac{\varepsilon}{2}$, we thus have

$$\|\chi(p) - \chi(p_\gamma)\|_{L^2(\Omega)}^2 = \int_{\Omega \setminus \partial \mathcal{S}_\gamma^\varepsilon} |\chi(p)(x) - \chi(p_\gamma)(x)|^2 \, dx + \int_{\partial \mathcal{S}_\gamma^\varepsilon} |\chi(p)(x) - \chi(p_\gamma)(x)|^2 \, dx$$

$$= 0 + \left|\partial \mathcal{S}_\gamma^\varepsilon\right| \to 0 \qquad \text{for } \varepsilon \to 0,$$

since $\text{dist}\left(p(x), \partial Q_{i,i+1}^\gamma\right) < \frac{\varepsilon}{2}$ on $\Omega \setminus \partial \mathcal{S}_\gamma^\varepsilon$ due to the choice of $p$. Due to the continuous embedding $H^2(\Omega) \hookrightarrow C(\overline{\Omega})$, there exists $\eta = \eta(\varepsilon)$ such that $\|y - y_\gamma\|_{H^2(\Omega)} < \eta$ and $\|w - w_\gamma\|_{H^2(\Omega)} < \eta$ implies $\|yw - y_\gamma w_\gamma\|_{C(\overline{\Omega})} < \frac{\varepsilon}{2}$. Hence $yw \to \chi(-yw)$ is continuous from $H^2(\Omega) \times H^2(\Omega)$ to $L^2(\Omega)$.

In a similar manner, one argues continuity of $H_\gamma$ from $H^2(\Omega) \times H^2(\Omega)$ to $L^2(\Omega)$, since the pointwise case distinction in the definition (4.2) can equivalently be expressed via the sum of characteristic functions. It follows from these considerations that the system matrix in (4.9) as an operator from $H^2(\Omega) \times H^2(\Omega)$ to $L^2(\Omega) \times L^2(\Omega)$ depends continuous on $(y, w) \in H^2(\Omega) \times H^2(\Omega)$. $\qquad\square$

Semismoothness of (4.7) together with Proposition 4.3 or Proposition 4.4 now implies local convergence of the Newton iteration; see, e.g., [14, Theorem 8.6].

**Theorem 4.5.** *Under the assumptions of either Proposition 4.3 or Proposition 4.4, if $(y^0, w^0)$ is sufficiently close in $H^2(\Omega) \times H^2(\Omega)$ to a solution $(y_\gamma, w_\gamma)$ to (4.7), the semismooth Newton iteration (4.9) converges superlinearly in $H^2(\Omega) \times H^2(\Omega)$ to $(y_\gamma, w_\gamma)$.*

### 4.2.2 Diffusion problem

We now consider the optimization of the leading coefficient. Here we are immediately faced with the difficulty that the state equation is not closed with respect to weak convergence of $u$ in $L^2(\Omega)$ or even weak-$*$ convergence in $L^\infty(\Omega)$; in particular, we cannot expect (A1) to hold. This is a classical difficulty concerning the identification of diffusion coefficients when only pointwise bounds are available. In this respect we recall results from [17] where, for given data $z$, and inhomogeneities $f$ and $g$, examples for non-existence of solutions to the problem

$$\min_{0 < u_1 \leq u \leq u_2} \int_\Omega |y(u) - z|^2 \, dx \qquad \text{s.t } -\nabla \cdot (u \nabla y) = f, \quad y|_{\partial \Omega} = g,$$

are given, as well as the notion of H- and G-convergence [18]. To address this difficulty and thus to ensure (A2), we propose to introduce a *local* bounded smoothing operator $G : L^2(\Omega) \to L^2(\Omega)$ with the property that its restrictions satisfy $G \in \mathcal{L}(L^s(\Omega), W^{1,s}(\Omega))$ and $G^* \in \mathcal{L}(W^{1,s}(\Omega), W^{1,s}(\Omega))$ for $s \in (n, \infty)$ and $G(U_M) \subset U_M$. This choice of $s$ guarantees that $W^{1,s}(\Omega)$ embeds compactly into $C(\overline{\Omega})$ and that $W^{1,s}(\Omega)$ is a Banach algebra. For example, we can choose $G$ as local averaging, i.e.,

$$(4.12) \qquad [Gu](x) = \frac{1}{|B_\rho|} \int_{B_\rho} u(x + \xi) \, d\xi,$$

where $B_\rho$ is a ball with radius $\rho > 0$ and center at the origin, and $u$ is extended by $u_1$ outside of $\Omega$.

The corresponding state equation is

$$(4.13) \qquad \begin{cases} -\nabla \cdot (Gu \, \nabla y) = f & \text{in } \Omega, \\ \qquad\qquad\qquad y = 0 & \text{on } \partial \Omega. \end{cases}$$

We assume that $\Omega \subset \mathbb{R}^N$, $N \leq 3$, is sufficiently regular such that for any $f \in L^s(\Omega)$ and any $u \in U = U_M$ defined as above, the solution to (4.13) satisfies $y \in W^{2,s}(\Omega) \cap H_0^1(\Omega)$ together with the uniform a priori estimate

$$(4.14) \qquad \|y\|_{W^{2,s}(\Omega)} \leq C_M \|f\|_{L^s(\Omega)}.$$

This is the natural $W^{2,s}(\Omega)$ regularity estimate for strongly elliptic equations, see [16, page 191]. Here we use that the set $G(U_M)$ is bounded in $W^{1,s}(\Omega)$ and hence that elements in $G(U_M)$ have a uniform modulus of continuity (which affects the constant $C_M$). Setting $S : u \mapsto y$ in (4.13) and $Y = L^2(\Omega)$, the assumptions (A1) and (A1) are satisfied. Digressing for a moment, we recall that our solutions to (2.2) and (4.1) still depend on $G$, and in particular in the case of (4.12), they depend on $\rho$. Let us denote this dependence by $u_\rho$. Then as $\rho \to 0$, these solution converge weakly in $L^s(\Omega)$ and G-converge to a – possibly different – limit which both satisfies the constraints involved in $U$ and appears as diffusion coefficient in the state equation; see, e.g., [1, Chapter 1.3].

We next turn for given $z \in L^s(\Omega)$ and any $u \in U_M$ and $y \in W^{2,s}(\Omega)$ to the adjoint equation

$$(4.15) \qquad \begin{cases} -\nabla \cdot (Gu \, \nabla w) = -(y - z) & \text{in } \Omega, \\ \qquad\qquad\qquad w = 0 & \text{on } \partial \Omega, \end{cases}$$

whose solution $w \in W^{2,s}(\Omega) \cap H_0^1(\Omega)$ also satisfies the uniform a priori estimate (4.14). We note that the solutions $y$ and $w$ satisfy $\nabla y \cdot \nabla w \in W^{1,s}(\Omega)$.

Using the solution $y_\gamma$ to (4.13) for $u = u_\gamma$ and the solution $w_\gamma$ to (4.15) for $u = u_\gamma$ and $y = y_\gamma$, we can write $p_\gamma = -G^*(\nabla y_\gamma \cdot \nabla w_\gamma) \in W^{1,s}(\Omega)$ and thus express (4.3) equivalently as

$$\begin{cases} -\nabla \cdot (Gu_\gamma \nabla w_\gamma) + y_\gamma = z, \\ u_\gamma - H_\gamma(-G^*(\nabla y_\gamma \cdot \nabla w_\gamma)) = 0, \\ -\nabla \cdot (Gu_\gamma \nabla y_\gamma) = f. \end{cases}$$

After eliminating $u_\gamma$ using the second equation, the reduced system has the form

$$(4.16) \quad \begin{cases} -\nabla \cdot \left( \left( GH_\gamma(-G^*(\nabla y_\gamma \cdot \nabla w_\gamma)) \right) \nabla w_\gamma \right) + y_\gamma = z, \\ -\nabla \cdot \left( \left( GH_\gamma(-G^*(\nabla y_\gamma \cdot \nabla w_\gamma)) \right) \nabla y_\gamma \right) = f. \end{cases}$$

We consider this again as an equation in $L^s(\Omega) \times L^s(\Omega)$ for $(y_\gamma, p_\gamma) \in (W^{2,s}(\Omega) \cap H_0^1(\Omega)) \times (W^{2,s}(\Omega) \cap H_0^1(\Omega))$, and interpret $H_\gamma$ as bounded linear operator from $W^{1,s}(\Omega)$ to $L^s(\Omega)$. This renders system (4.16) semismooth. Appealing again to the chain rule for Newton derivatives and introducing $\chi = \chi(-G^*(\nabla y \cdot \nabla w))$, we obtain the Newton system

$$(4.17) \quad \begin{pmatrix} \mathrm{Id} + A^k(w^k, \cdot, w^k) & -\nabla \cdot \left( Gu^k \, \nabla \cdot \right) + A^k(y^k, \cdot, w^k) \\ -\nabla \cdot \left( Gu^k \, \nabla \cdot \right) + A^k(w^k, \cdot, y^k) & A^k(y^k, \cdot, y^k) \end{pmatrix} \begin{pmatrix} \delta y \\ \delta w \end{pmatrix}$$
$$= -\begin{pmatrix} -\nabla \cdot \left( Gu^k \, \nabla w^k \right) + y^k - z \\ -\nabla \cdot \left( Gu^k \, \nabla y^k \right) - f \end{pmatrix},$$

where we have set $u^k := H_\gamma(-G^*(\nabla y^k \cdot \nabla w^k))$ and

$$A^k(v_1, v_2, v_3) := \nabla \cdot \left( G\left( \tfrac{1}{\gamma} \chi^k G^*(\nabla v_1 \cdot \nabla v_2) \right) \nabla v_3 \right).$$

Note that for all $y, w, \delta y, \delta w \in H^2(\Omega)$,

$$\left( A^k(y, \delta y, w), \delta w \right)_{L^2(\Omega)} = \left( A^k(w, \delta w, y), \delta y \right)_{L^2(\Omega)}.$$

It remains to provide sufficient conditions for the uniform bounded invertibility of the system matrix in (4.17). For this purpose we specify the critical set $\partial \omega_\gamma$ for the present case:

$$\partial \omega_\gamma := \bigcup_{i=1}^{d-1} \left\{ x \in \Omega : -G^*(\nabla y_\gamma(x) \cdot \nabla w_\gamma(x)) \in \partial Q_{i,i+1}^\gamma \right\}.$$

**Theorem 4.6.** *Let $(y_\gamma, w_\gamma)$ denote a solution to (4.16), assume that $|\partial \omega_\gamma| = 0$, and that the system matrix (4.17) evaluated at $(y_\gamma, w_\gamma)$ is continuous invertible as an operator from $(W^{2,s} \cap H_0^1(\Omega))^2$ to $(L^s(\Omega))^2$. Then, if $(y^0, w^0)$ is sufficiently close in $(W^{2,s} \cap H_0^1(\Omega))^2$ to $(y_\gamma, w_\gamma)$, the semismooth Newton iteration (4.9) converges superlinearly to $(y_\gamma, w_\gamma)$.*

*Proof.* It suffices to argue that the system matrix depends continuously on $(y, w) \in (W^{2,s}(\Omega) \cap H_0^1(\Omega))^2$ in a neighborhood of $(y_\gamma, w_\gamma)$ considered as operators in $\mathcal{L}((W^{2,s}(\Omega) \cap H_0^1(\Omega))^2, L^s(\Omega)^2)$. For this purpose we consider the operator

$$(W^{2,s}(\Omega) \cap H_0^1(\Omega))^2 \ni (y, w) \mapsto A(w, \cdot, w) \in \mathcal{L}(W^{2,s}(\Omega) \cap H_0^1(\Omega), L^s(\Omega)),$$

where $A$ still depends on $\chi = \chi(-G^*(\nabla y \cdot \nabla w))$. First we argue exactly as in the proof of Proposition 4.4 that

$$(W^{2,s}(\Omega) \cap H_0^1(\Omega))^2 \ni (y, w) \mapsto \chi = \chi(-G^*(\nabla y \cdot \nabla w)) \in L^s(\Omega)$$

is continuous. Next we observe that

$$W^{2,s}(\Omega) \cap H_0^1(\Omega) \ni w \mapsto G^*(\nabla w \cdot \nabla \cdot) \in \mathcal{L}(W^{2,s}(\Omega) \cap H_0^1(\Omega), W^{1,s}(\Omega))$$

is continuous, and consequently

$$(W^{2,s}(\Omega) \cap H_0^1(\Omega))^2 \ni (y, w) \mapsto G(\tfrac{1}{\gamma} \chi G^*(\nabla w \cdot \nabla \cdot)) \in \mathcal{L}(W^{2,s}(\Omega) \cap H_0^1(\Omega), L^s(\Omega))$$

is continuous as well. From here we can conclude that $(y, w) \mapsto A(w, \cdot, w)$ is continuous from $(W^{2,s}(\Omega) \cap H_0^1(\Omega))^2$ to $\mathcal{L}((W^{2,s}(\Omega) \cap H_0^1(\Omega)), L^s(\Omega))$. We argue similarly for $A(w, \cdot, y), A(y, \cdot, w)$ and $A(y, \cdot, y)$, which establishes the claim. $\qquad\square$

Returning to the assumption on the well-posedness of the system matrix at $(y_\gamma, w_\gamma)$, we now argue that this is indeed the case if $w_\gamma$ is sufficiently small in the $W^{2,s}(\Omega)$ norm, i.e., for small residual problems. For $w = 0$, the system matrix in (4.17) has the form

$$\begin{pmatrix} \mathrm{Id} & -\nabla \cdot (u_1 \nabla \cdot) \\ -\nabla \cdot (u_1 \nabla \cdot) & 0 \end{pmatrix}$$

since $u_\gamma = GH_\gamma(0) = Gu_1 = u_1$ because $Gu = u$ for $u$ constant. This operator is clearly continuously invertible. A perturbation argument as in the proof of Theorem 4.6 implies continuous invertibility also for $(y_\gamma, w_\gamma)$ if $\|w_\gamma\|_{W^{2,s}(\Omega)}$ is sufficiently small.

# 5 Numerical examples

We illustrate the behavior of the proposed approach with numerical examples modeling a simple material design problem for the potential and the diffusion equation, in which a reference binary material distribution $u_r$ (i.e., using only two values: matrix or void, and material) has already been obtained. The goal is now to obtain a comparable behavior using additionally available materials of intermediate density (and hence presumably lower cost) by solving the multi-material optimization problem (2.2) with target $z = y_r$ (the solution to the state equation corresponding to the reference coefficient $u_r$) and an extended list $u_b$ of feasible material parameters containing the two original values. Here, the tracking term $\mathcal{F}$ penalizes the deviation from the reference state, while the "multi-bang" term $\mathcal{G}$ both promotes the desired discrete structure and favors materials with lower density; the trade-off between the two goals

is controlled by the parameter $\alpha$. We point out that not strictly enforcing attainment of the target allows parameter distributions that are different from the original binary distribution (which is only recovered in the limit $\alpha \to 0$). For each example, we report on the deviation from the reference state as well as on the achieved total material cost reduction (as measured by the difference of the $L^2$ norms of the reference and computed coefficients).

The multi-material optimization problem (2.2) is solved using the described regularized semismooth Newton method. To address the local convergence of Newton methods and to avoid having to choose the Moreau–Yosida regularization parameter $\gamma$ a priori, a continuation strategy is applied where the problem is solved starting with a large $\gamma^0 = 1$ and the initial guess $(y_0, p_0) = (0, 0)$. The regularization parameter is then successively reduced via $\gamma^{k+1} = \gamma^k/2$, taking the previous solution as a starting point. The iteration is terminated if $\gamma = 10^{-12}$ is reached or more than 300 Newton iterations are performed. This is combined with a non-monotone backtracking line seach based on the residual of the optimality system (4.3), starting with a step length of 1 and using a reduction factor of $1/2$, where a minimal step length of $10^{-6}$ is accepted even if it leads to a (small) increase in the residual norm.

The partial differential equations are discretized using finite differences on a uniform grid of $128 \times 128$ grid points. Our Matlab implementation of the described algorithm can be downloaded from https://github.com/clason/multimaterialcontrol.

## 5.1 Potential problem

We first consider the design problem associated with equation (4.4), where we fix $\Omega = [-1, 1]^2$ and

$$f(x_1, x_2) = \sin(\pi x_1)\cos(\pi x_2).$$

The reference material parameter is

$$(5.1) \qquad u_r(x_1, x_2) = \begin{cases} 2.5 & \text{if } 1/4 < |x|^2 < \frac{3}{4} \text{ and } x_1 > \frac{1}{10}, \\ 2.5 & \text{if } 1/4 < |x|^2 < \frac{3}{4} \text{ and } x_1 < -\frac{1}{10}, \\ 1.5 & \text{else}, \end{cases}$$

see Figure 3a. We then solve the multi-material design problem for the target $z = y_r$ with the extended feasible parameter set $\{1, 1.5, 2, 2.5\}$ for different values of $\alpha$ using the described algorithm. In all cases, after some initial reduced steps were taken for $\gamma < 5 \cdot 10^{-5}$, the Newton iteration entered a superlinear phase and converged after at most three iterations. Depending on $\gamma$, the total number of Newton iterations was between 5 and 28. The algorithm always terminated at $\gamma \approx 10^{-12}$ because the minimal value of $\gamma$ was reached. The final material distributions $u_\gamma$ for $\alpha \in \{10^{-5}, 10^{-6}, 10^{-7}\}$ are shown in Figure 3b–d. As can be seen, at almost all points, only the feasible parameter values are attained, where lower values of $\alpha$ lead to increased use of higher density materials. The relative tracking error $e_T := \|y_\gamma - y_r\|_{L^2}/\|y_r\|_{L^2}$ as well as the relative total material cost reduction $e_M := (\|u_r\|_{L^2} - \|u_\gamma\|_{L^2})/\|u_r\|_{L^2}$ for each value of $\alpha$ are given in Table 1a.

(a) reference coefficient $u_r$

(b) optimal coefficient $u_\gamma$ for $\alpha = 10^{-5}$

(c) optimal coefficient $u_\gamma$ for $\alpha = 10^{-6}$

(d) optimal coefficient $u_\gamma$ for $\alpha = 10^{-7}$
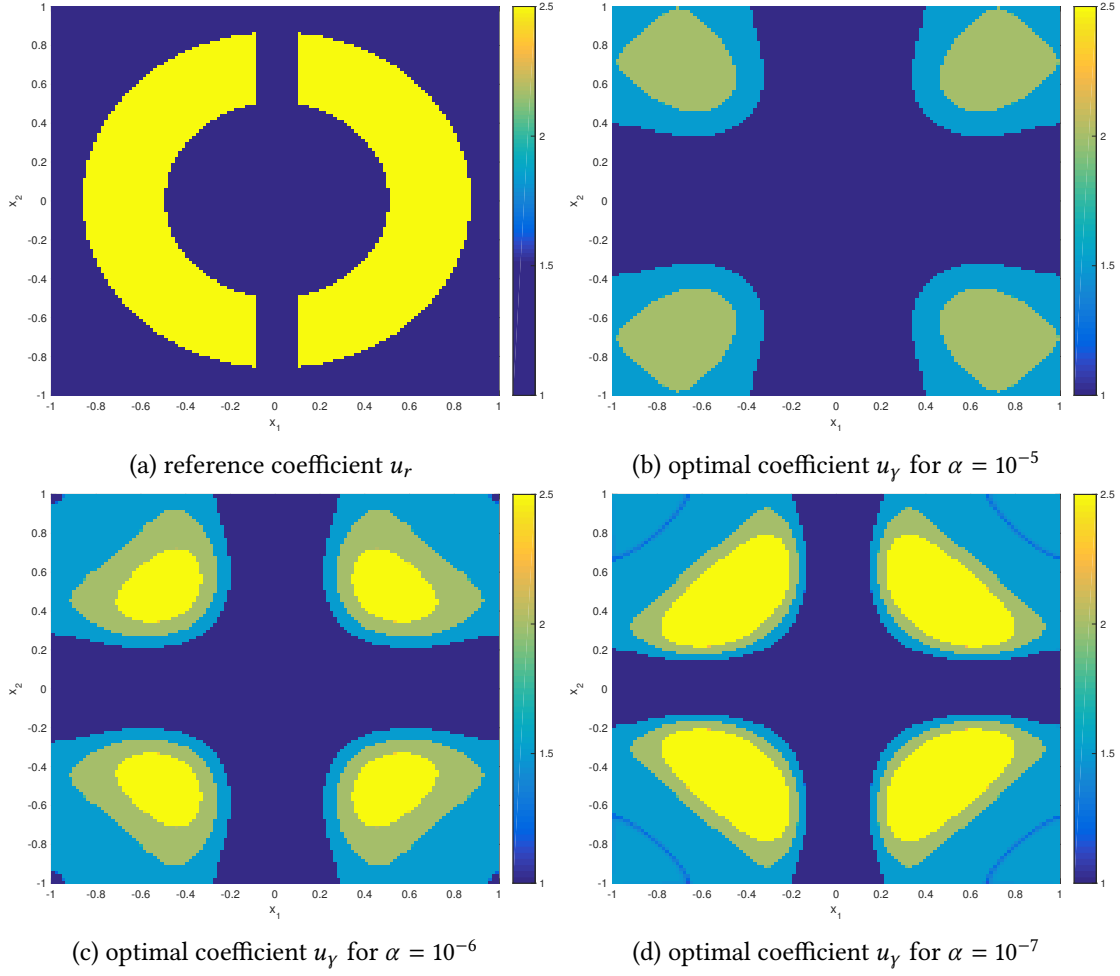
Figure 3: Results for potential problem

## 5.2  Diffusion problem

For the design problem associated with equation (4.13), we set $f \equiv 10$ and $u_r$ as given in (5.1). The smoothing operator $G$ is taken as averaging over the local five-point stencil; the smoothed reference coefficient $Gu_r$ is shown in Figure 4a to facilitate comparison. For the multimaterial design problem, we choose the extended feasible parameter set $\{1.5, 1.75, 2, 2.25, 2.5\}$ and $\alpha \in \{10^{-2}, 10^{-3}, 10^{-6}\}$ (the last value to illustrate the behavior for $\alpha \to 0$). In these cases, the algorithm terminated prematurely due to reaching the maximal number of Newton iterations at $\gamma^* \approx 4.8 \cdot 10^{-7}$, $\gamma^* \approx 6.0 \cdot 10^{-8}$, and $\gamma^* \approx 9.3 \cdot 10^{-10}$, respectively. The behavior of the Newton method is similar as in the potential problem, although the required number of Newton iterations now increases significantly as $\gamma$ is decreased due to the line search leading to smaller step lengths (including, e.g., for $\alpha = 10^{-3}$ in total six non-monotone steps due to the minimal step length being reached). The corresponding material coefficients $Gu_\gamma$ from the last successful iteration at $\gamma = 2\gamma^*$ are shown in Figure 4b–d. Although the multi-bang structure is no longer perfect, it

(a) reference coefficient $Gu_r$

(b) optimal coefficient $Gu_\gamma$ for $\alpha = 10^{-2}$

(c) optimal coefficient $Gu_\gamma$ for $\alpha = 10^{-3}$

(d) optimal coefficient $Gu_\gamma$ for $\alpha = 10^{-6}$
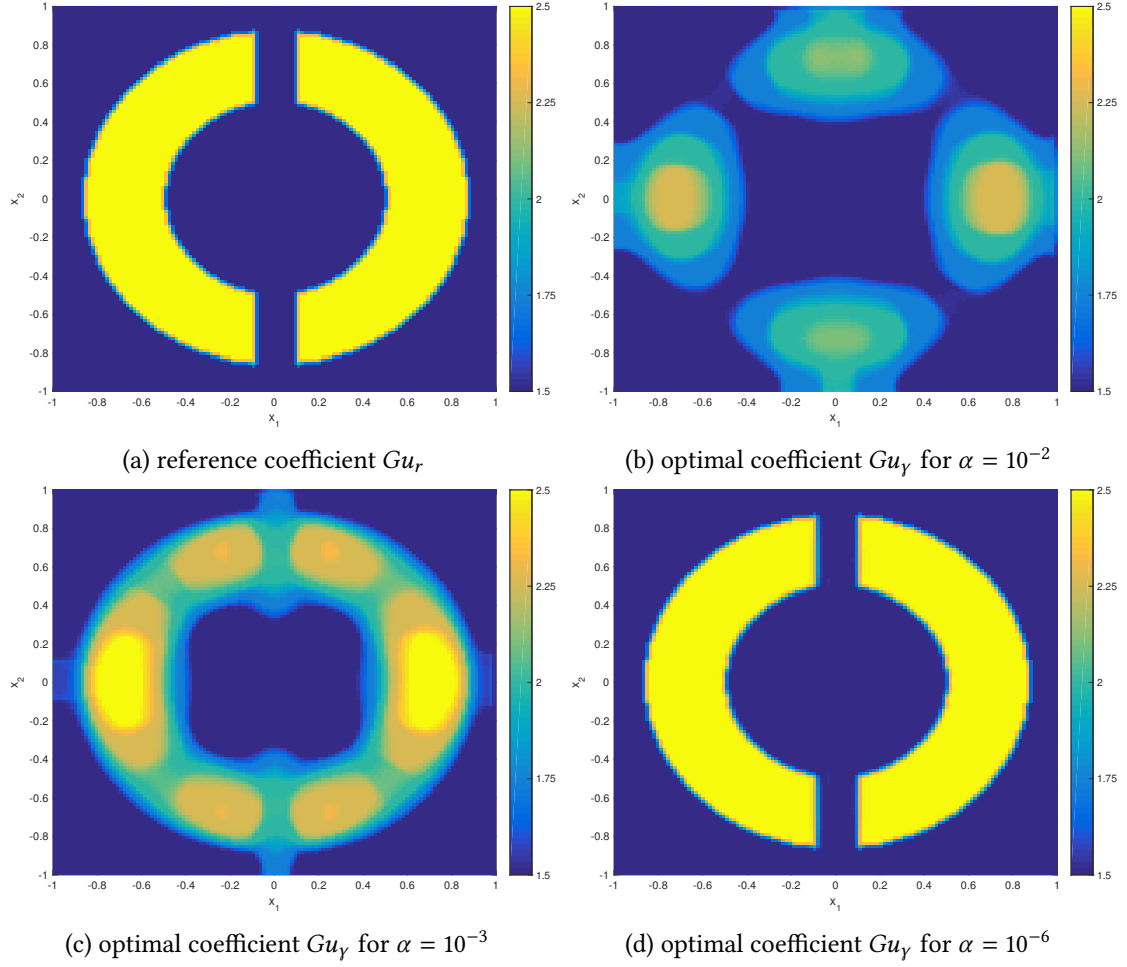
Figure 4: Results for diffusion problem

can be observed that the penalty is successful in promoting the desired parameter values even in the presence of the smoothing operator $G$. Figure 4d also indicates that the original binary reference distribution $u_r$ is recovered for $\alpha \to 0$. Finally, the relative tracking errors and relative material cost reductions for these values of $\alpha$ are given in Table 1b.

Table 1: Relative tracking error $e_T$ and material cost reduction $e_M$ for different values of $\alpha$

(a) Potential problem

| $\alpha$ | $10^{-5}$ | $10^{-6}$ | $10^{-7}$ |
|---|---|---|---|
| $e_T$ | $2.95 \cdot 10^{-2}$ | $8.28 \cdot 10^{-3}$ | $2.01 \cdot 10^{-3}$ |
| $e_M$ | $2.89 \cdot 10^{-1}$ | $1.82 \cdot 10^{-1}$ | $1.10 \cdot 10^{-1}$ |

(b) Diffusion problem

| $\alpha$ | $10^{-1}$ | $10^{-2}$ | $10^{-6}$ |
|---|---|---|---|
| $e_T$ | $4.96 \cdot 10^{-2}$ | $1.15 \cdot 10^{-2}$ | $5.29 \cdot 10^{-5}$ |
| $e_M$ | $1.16 \cdot 10^{-2}$ | $4.61 \cdot 10^{-1}$ | $7.29 \cdot 10^{-4}$ |

# 6 Conclusion

A convex analysis approach is presented for the determination of piecewise constant coefficients in a partial differential equation where the constants range over a predetermined discrete set. Since the subdomains where the coefficient is constant are not specified a priori, this constitutes a topology optimization problem. Two model applications are analyzed in detail. For the case where the unknown coefficient enters into the potential term, the numerical results are very encouraging. If the unknown parameter enters into the diffusion term, regularization is required that has a smoothing effect on the solutions, and thus the numerical results are less "crisp". In practice, this could be addressed by a post-processing step, either by standard thresholding or by evaluating the unregularized subdifferential at the computed optimal dual variable, i.e., taking an appropriate selection $\tilde{u} \in \partial \mathcal{G}^*(p_\gamma)$. Since the considered problems resemble inverse coefficient problems, it comes as no surprise that the diffusion problem is more ill-posed than the potential problem.

In future work, we plan to return to the diffusion problem and to formulate the multi-topology optimization problem based on a bounded variation framework using a functional including the total variation seminorm. It may also be of interest to search for other types of functionals which serve the purpose of multi-material topology optimization. In particular, we note that the currently used formulation in (1.1) favors values $u(x) = u_i$ with small magnitude over other ones. Depending on the practical relevance of the $u_i$, this may not be a desired effect. In this case, functionals should be constructed that favor different criteria (e.g., the weight or the price of different materials) while still keeping the "multi-bang" property feature of promoting controls with values only from the given set.

## Acknowledgment

## References

[1] G. ALLAIRE, *Shape Optimization by the Homogenization Method*, Springer, 2002.

[2] G. ALLAIRE, F. JOUVE, AND A.-M. TOADER, *Structural optimization using sensitivity analysis and a level-set method*, J. Comput. Phys., 194 (2004), pp. 363–393.

[3] S. AMSTUTZ, *A semismooth Newton method for topology optimization*, Nonlinear Anal., 73 (2010), pp. 1585–1595.

[4] ——, *Analysis of a level set method for topology optimization*, Optim. Methods Softw., 26 (2011), pp. 555–573.

[5] S. AMSTUTZ AND H. ANDRÄ, *A new algorithm for topology optimization using a level-set method*, Journal of Computational Physics, 216 (2006), pp. 573 – 588.

[6] H. H. BAUSCHKE AND P. L. COMBETTES, *Convex Analysis and Monotone Operator Theory in Hilbert Spaces*, CMS Books in Mathematics/Ouvrages de Mathématiques de la SMC, Springer, New York, 2011.

[7] M. P. BENDSØE AND O. SIGMUND, *Topology Optimization*, Springer-Verlag, Berlin, 2003.

[8] L. BLANK, M. H. FARSHBAF-SHAKER, H. GARCKE, C. RUPPRECHT, AND V. STYLES, *Multi-material phase field approach to structural topology optimization*, in Trends in PDE Constrained Optimization, G. Leugering et al., eds., vol. 165 of International Series of Numerical Mathematics, Springer International Publishing, 2014, pp. 231–246.

[9] C. CLASON, K. ITO, AND K. KUNISCH, *A convex analysis approach to optimal controls with switching structure for partial differential equations*, ESAIM: Control, Optimisation and Calculus of Variations, forthcoming (2015).

[10] C. CLASON AND K. KUNISCH, *Multi-bang control of elliptic systems*, Annales de l'Institut Henri Poincaré (C) Analyse Non Linéaire, 31 (2014), pp. 1109–1130.

[11] I. EKELAND AND R. TÉMAM, *Convex Analysis and Variational Problems*, vol. 28 of Classics Appl. Math., SIAM, Philadelphia, 1999.

[12] S. GARREAU, P. GUILLAUME, AND M. MASMOUDI, *The topological asymptotic for PDE systems: the elasticity case*, SIAM J. Control Optim., 39 (2001), pp. 1756–1778 (electronic).

[13] J. HASLINGER, M. KOČVARA, G. LEUGERING, AND M. STINGL, *Multidisciplinary free material optimization*, SIAM Journal on Applied Mathematics, 70 (2010), pp. 2709–2728.

[14] K. ITO AND K. KUNISCH, *Lagrange Multiplier Approach to Variational Problems and Applications*, vol. 15 of Advances in Design and Control, SIAM, Philadelphia, PA, 2008.

[15] K. ITO, K. KUNISCH, AND Z. LI, *Level-set function approach to an inverse interface problem*, Inverse Problems, 17 (2001), p. 1225.

[16] O. A. LADYZHENSKAYA AND N. N. URAL'TSEVA, *Linear and Quasilinear Elliptic Equations*, Translated from the Russian by Scripta Technica, Inc. Translation editor: Leon Ehrenpreis, Academic Press, New York, 1968.

[17] F. MURAT, *Contre-exemples pour divers problèmes où le contrôle intervient dans les coefficients*, Ann. Mat. Pura Appl. (4), 112 (1977), pp. 49–68.

[18] F. MURAT AND L. TARTAR, *H-convergence*, in Topics in the mathematical modelling of composite materials, vol. 31 of Progr. Nonlinear Differential Equations Appl., Birkhäuser Boston, Boston, MA, 1997, pp. 21–43.

[19] P. NEITTAANMAKI, J. SPREKELS, AND D. TIBA, *Optimization of Elliptic Systems*, Springer Monographs in Mathematics, Springer, New York, 2006.

[20] O. PIRONNEAU, *Optimal Shape Design for Elliptic Systems*, Springer Series in Computational Physics, Springer-Verlag, New York, 1984.

[21] W. Schirotzek, *Nonsmooth Analysis*, Universitext, Springer, Berlin, 2007.

[22] J. Sokołowski and A. Żochowski, *On the topological derivative in shape optimization*, SIAM J. Control Optim., 37 (1999), pp. 1251–1272 (electronic).

[23] J. Sokołowski and J.-P. Zolésio, *Introduction to Shape Optimization*, vol. 16 of Springer Series in Computational Mathematics, Springer-Verlag, Berlin, 1992.

[24] M. Ulbrich, *Semismooth Newton Methods for Variational Inequalities and Constrained Optimization Problems in Function Spaces*, vol. 11 of MOS-SIAM Series on Optimization, SIAM, Philadelphia, PA, 2011.

[25] A. Wouk, *A Course of Applied Functional Analysis*, Wiley-Interscience [John Wiley & Sons], New York, 1979.