

1 **ERROR ESTIMATES FOR THE FEM APPROXIMATION OF**
2 **STATE-CONSTRAINED ELLIPTIC OPTIMAL CONTROL**
3 **PROBLEMS WITH SPARSITY IN A FINITE-DIMENSIONAL**
4 **CONTROL SPACE***

5 PEDRO MERINO[†] AND ALEXANDER NENJER [†]

6 **Abstract.** In this work, we derive an a-priori error estimate of order $h^2|\log(h)|$ for the finite
7 element approximation of a sparse optimal control problem governed by an elliptic equation, which
8 is controlled in a finite dimensional space. Furthermore, box-constraints on the control are considered
9 and finitely many pointwise state-constraints are imposed in specific points in the domain. With this
10 choice for the control space, the achieved order of approximation for the optimal control is optimal, in
11 the sense that the order of the error for the optimal control is of the same order of the approximation
12 for the state equation.

13 **Key words.** Optimal control, elliptic partial differential equations, error estimates, finite ele-
14 ment method, sparsity.

15 **AMS subject classifications.** 80M10, 41A25, 49J20, 49K20.

16 **1. Introduction.** It is well known that sparse optimal controls have very at-
17 tractive properties for their practical implementation c.f. [17]. When the decision
18 variables are functions, sparse optimal control problems produce solutions with small
19 supports on the domain, which can be viewed as a localized action of the optimal con-
20 trol c.f. [32]. However, when the control is a vector, sparsity takes the form of many
21 null components in the solution vector. This is interesting since a simpler solution is
22 not only easier to understand and implement, but it can also be used as an strategy
23 to detect and select the most relevant quantities taking part in the control process,
24 aiming to save scarce resources.

25 In this paper, we consider finite dimensional sparse optimal control problems and
26 its approximation by the standard finite element method. We study the estimation
27 of the convergence order (with respect to the size of the mesh h) of a state con-
28 strained linear-quadratic optimal control problem. Here, the control space is finite
29 dimensional; therefore, in order to promote sparsity in the solution an 1-norm pe-
30 nalization term is included in the cost functional. The choice of finite-dimensional
31 control space as well as pointwise state-constraints are motivated by technological
32 requirements since, in practice, we often choose from a set of finite quantities in order
33 to control certain processes in real applications. For instance, this is the case of digital
34 technology, where finitely many quantities are used rather than functions.

35 In turn, pointwise constraints on the state are also important for applications
36 where it is critical that the associated state remains constrained on specific points
37 on the domain. Several applications in medicine have these kind of situations. For
38 example, in the treatment of cancer with hypertermia or cryosurgery techniques,
39 state-constraints prevent damage on healthy tissue and organs exposed to extreme
40 temperatures, see for example [22] and [31].

*Submitted to the editors DATE.

Funding: This research has been partially supported by project PIJ-15-26 granted by Escuela Politécnica Nacional. Also, we acknowledge partial support of MATHAmSud project SOCDE “Sparse Optimal Control of Differential Equations”.

[†]Department of Mathematics, *Ladrón de Guevara E11-253, Escuela Politécnica Nacional (Quito-Ecuador)*, Research center on mathematical modeling MODEMAT, Quito-Ecuador (pedro.merino@epn.edu.ec, hernan.alex.nenjer@hotmail.com).

41 As mentioned in [7, 8], there are mainly two ways for inducing sparsity on the
 42 solutions: by considering a sparsity inducing term or, alternatively, by choosing an
 43 appropriate control space; for instance, the space of regular Borel measures. In recent
 44 years a considerable amount of literature has been published on sparse optimal control
 45 problems and many authors have studied several questions arising in its analysis,
 46 numerical approximation and numerical methods. See for example [7, 8, 9, 32] for
 47 a complete review of the theoretical results for sparse optimal control of elliptic and
 48 parabolic equations.

49 The question of error estimates by the finite element method for sparse optimal
 50 control problems is rather a recent subject of research. In the recent contributions [6]
 51 and [7], the linear and semilinear elliptic cases have been considered and an error of
 52 order h was derived. The case where the control space consists of the regular Borel
 53 measures was studied in [5]. There, an the error of order h^κ was deduced, where $\kappa = 1$
 54 and $\kappa = 1/2$ for two and three dimensions respectively.

55 From the finite-dimensional perspective, error estimates for optimal control prob-
 56 lems, where the control space is finite dimensional, have also been considered in [25].
 57 However, there sparsity terms were not considered.

58 In this article, we derive an error of convergence of order $h^2|\log(h)|$, which is
 59 optimal in the sense that that the order of the error for the optimal control is of the
 60 same order of the approximation of the state equation. This is not highly surprising
 61 since the control variables are vectors and they do not need to be approximated. The
 62 only source of error comes from the approximation of the state equation. However,
 63 unlike the case of functional controls, the derivation of the error estimates follows
 64 a different analysis. Even though the results obtained by [7] can be applied to the
 65 finite-dimensional case, they do not take into account the finite-dimensional nature
 66 of the control space. In this regard, our results improve the existing estimations by
 67 taking advance of this particularity.

68 We consider an analogous problem as in [24] but the novelty consists in including
 69 the non-differentiable penalization by means of the 1-norm. The lack of differen-
 70 tiability involves new numerical challenges in the approximation and its numerical
 71 solution. At our best knowledge, this is the first contribution of sparse optimal con-
 72 trol in a finite-dimensional control space. Our strategy for deriving error estimates
 73 consists in circumventing the non differentiability of the optimization problem by for-
 74 mulating an equivalent differentiable optimal control problem, obtained by splitting
 75 the control variable in its positive and negative parts. This splitting procedure is
 76 known in mathematical programming in finite dimensions but, due to the duplication
 77 of the number of variables, it is not useful for numerical algorithms. On the other
 78 hand, in context of optimal control, where a state equation has to be approximated
 79 by finite elements, this trick helps to analyze the numerical approximation. Indeed,
 80 having reformulated the original problem we proceed as in [25], by transforming the
 81 original optimal control problem into an optimization problem in finite dimensions
 82 for which we carry out a similar stability analysis and, subsequently to establish opti-
 83 mal error estimates. Our stability analysis developed for the equivalent differentiable
 84 problem differs from [25] in that we exploit convexity, and the so called *strong regu-*
 85 *larity* is weakened by the concept of *regularity*, allowing us to only require a Slater
 86 type condition as constraint qualification.

87 In addition, under further considerations we show that these results also cover the
 88 case when the Tikhonov regularization is not present, which an additional contribution
 89 of this paper.

90 Our paper is organized as follows: in the first section we introduce the optimal

91 control problem and the associated terminology and definitions. Then, in Section 2
 92 we derive optimality conditions characterizing sparse solutions. In the third section
 93 we describe the numerical approximation by the finite element method of the optimal
 94 control problem and state the main result of this paper in Theorem 3.3. The fourth
 95 section focuses on an equivalent differentiable formulation of our problem which is the
 96 bridge to apply the stability analysis. The fifth section is devoted to the derivation
 97 of the error estimate and the sixth and final section, numerical tests confirming the
 98 theory are presented.

99 **2. The sparse optimal control problem and optimality conditions.** Let
 100 $\Omega \subset \mathbb{R}^2$ be an open bounded convex set with Lipschitz boundary Γ and $M \in \mathbb{N} \setminus \{0\}$.
 101 We are interested in the finite–element approximation of the following linear–quadratic
 102 sparse optimal control problem in \mathbb{R}^M :

$$103 \quad (\mathbf{P}_M) \quad \left\{ \begin{array}{l} \min_{(y,u)} J(y, u) = \frac{1}{2} \int_{\Omega} (y - y_d)^2 dx + \frac{\alpha}{2} \|u\|_2^2 + \beta \|u\|_1 \\ \text{subject to:} \\ Ay(x) = \sum_{i=1}^M u_i e_i(x), \quad \text{in } \Omega, \\ y(x) = 0, \quad \text{on } \Gamma, \\ y(x_j) \leq b_j, \quad \forall j = 1, \dots, \ell, \\ u \in \mathcal{U}_{ad}, \end{array} \right.$$

104 In this problem, we assume that $y_d \in L^2(\Omega)$ is a desired state, $\alpha > 0$, $\beta > 0$, as well
 105 as fixed basis functions e_i , for $i = 1, 2, \dots, M$, belong to $C^{0,\gamma}(\Omega)$, for some $0 < \gamma < 1$.
 106 The admissible set of controls is defined by

$$107 \quad \mathcal{U}_{ad} = \{u \in \mathbb{R}^M : u_a \leq u \leq u_b\},$$

108 with $u_a, u_b \in \mathbb{R}^M$ and $u_a < u_b$. For vectors, the symbols \leq and $<$ are understood
 109 component–wise.

110 The differential equation is based upon the uniform ellipticity and symmetry of
 111 the differential operator A , defined by:

$$112 \quad Ay(x) := - \sum_{i,j=1}^2 \partial_j (a_{ij}(x) \partial_i y(x)) + a_0(x) y(x),$$

113 with coefficients $a_{ij} \in C^{1+\delta}(\Omega)$, $0 < \delta < 1$, and $a_0 \in C^{0,\gamma}(\Omega)$ such that $a_0(x) \geq 0$ in
 114 Ω .

115 The constraints on the state, whose bound is given by the vector $b \in \mathbb{R}^\ell$, are
 116 imposed in $\ell \in \mathbb{N}$ a priori selected points $x_i \in \Omega_0$, $i = 1, \dots, \ell$, where Ω_0 is a
 117 subdomain of Ω .

118 Let us summarize some fundamental results on elliptic partial differential equa-
 119 tions which will be applied to the state and adjoint equations associated to our opti-
 120 mal control problem (\mathbf{P}_M) . In what follows, we will use the symbols (\cdot, \cdot) and $\|\cdot\|$ to
 121 denote the inner product and its corresponding norm in $L^2(\Omega)$.

122 Let f be a given function in $L^2(\Omega)$. We say that $y \in H_0^1(\Omega)$ is a weak solution of
 123 the equation

$$124 \quad (2.1) \quad \begin{cases} Ay(x) = f(x), & \text{in } \Omega, \\ y(x) = 0, & \text{on } \Gamma, \end{cases}$$

125 if y fulfills the variational equation:

$$126 \quad \sum_{i,j=1}^2 (a_{ij} \partial_i y, \partial_j \phi) + (a_0 y, \phi) = (f, \phi),$$

127 for all $\phi \in H_0^1(\Omega)$.

128 **PROPOSITION 2.1.** *For each function $f \in L^2(\Omega)$, there exists a unique weak so-*
 129 *lution $y \in H_0^1(\Omega) \cap H^2(\Omega)$ of (2.1) and the mapping $G : L^2(\Omega) \rightarrow H_0^1(\Omega)$ defined*
 130 *by*

$$131 \quad G(f) = y, \quad \forall f \in L^2(\Omega),$$

132 *is linear and continuous. Moreover, if $f \in C^{0,\gamma}(\Omega)$, then $y \in C^{2,\gamma}(\Omega)$.*

133 *Proof.* For the existence and uniqueness we refer to [16] [Theorem 3, p. 301] and
 134 for the regularity result see [18],[19], since Ω is convex. \square

135 **PROPOSITION 2.2.** *For each $u \in \mathbb{R}^M$, the state equation has a unique weak solu-*
 136 *tion $y_u \in H_0^1(\Omega) \cap H^2(\Omega) \cap C^{2,\gamma}(\Omega)$, that is y_u satisfies*

$$137 \quad (2.2) \quad \sum_{i,j=1}^2 (a_{ij} \partial_i y_u, \partial_j \phi) + (a_0 y_u, \phi) = \left(\sum_{i=1}^M u_i e_i, \phi \right),$$

138 for all $\phi \in H_0^1(\Omega)$.

139 For the analysis of the existence of an optimal solution, we reformulate the optimal
 140 control problem (\mathbf{P}_M) as an optimization problem in terms of u only. Let y_i be the
 141 weak solution of (2.1) associated to the right-hand side $f = e_i$. By utilizing the
 142 superposition principle, we define the control-to-state mapping S , which associates to
 143 each vector $u \in \mathbb{R}^M$ the corresponding state y_u defined by

$$144 \quad (2.3) \quad \begin{aligned} S : \mathbb{R}^M &\rightarrow C(\bar{\Omega}) \\ u &\mapsto y_u = Su = \sum_{i=1}^M u_i y_i. \end{aligned}$$

145 Clearly, by the Proposition 2.2, Su is the solution to the state equation for $u \in \mathbb{R}^M$.

146 **Remark 2.3.** The mapping S is linear and continuous.

147 Now, by replacing $y = Su$ in the cost functional J we introduce the reduced
 148 optimal control problem:

$$149 \quad (\mathcal{P}) \quad \left\{ \begin{array}{l} \min_{u \in \mathcal{U}_{ad}} f(u) = \frac{1}{2} \int_{\Omega} \left(\sum_{i=1}^M u_i y_i - y_d \right)^2 dx + \frac{\alpha}{2} \|u\|_2^2 + \beta \|u\|_1 \\ \text{subject to:} \\ \sum_{i=1}^M u_i y_i(x_j) \leq b_j, \quad \forall j = 1, \dots, \ell. \end{array} \right.$$

150 We denote the feasible set of controls by

$$151 \quad (2.4) \quad \mathcal{U}_{feas} = \left\{ u \in \mathcal{U}_{ad} : \sum_{i=1}^M u_i y_i(x_j) \leq b_j, \quad \forall j = 1, \dots, \ell \right\}.$$

152 Due to the presence of the state constraints, depending on the selection of b the
 153 last set has no reason to be non empty. Moreover, since we are looking for sparse
 154 solutions, control constraints should contain 0 as feasible vector. Thus, we make the
 155 following hypothesis.

156 *Hypothesis 2.4.* The upper and lower bounds u_b and u_a satisfy $u_a < 0 < u_b$
 157 pointwise. Moreover, we assume that \mathcal{U}_{feas} is not empty.

158 **THEOREM 2.5.** *Let the Tikhonov parameter α be positive, then there exists a*
 159 *unique solution for the problem (P).*

160 *Proof.* It is clear from the definition of f that it is a continuous and strictly convex
 161 function because the Tikhonov term $\frac{\alpha}{2}\|u\|_2^2$. In this case, the existence of a unique
 162 solution is obtained even if \mathcal{U}_{ad} is unbounded since \mathcal{U}_{feas} remains convex and closed. \square

163 *Remark 2.6.* Let \bar{u} be the optimal control of (P), then its associated state is given
 164 by

$$165 \quad \bar{y} = \sum_{i=1}^M \bar{u}_i y_i.$$

166 We will refer to \bar{y} as the *optimal state*, and to (\bar{y}, \bar{u}) as the *optimal pair*. Next,
 167 we are going to derive the necessary and sufficient conditions for problem (P). For
 168 convenience, in the subsequent formulation the description of inequality constraints
 169 for problem (P), will be represented using the function $g : \mathbb{R}^M \rightarrow \mathbb{R}^{\ell+2M}$ defined by

$$170 \quad (2.5) \quad \begin{aligned} g_j(u) &= \sum_{i=1}^M u_i y_i(x_j) - b_j, & \text{for } j = 1, \dots, \ell, \\ g_{\ell+i}(u) &= u_{a,i} - u_i, & \text{for } i = 1, \dots, M, \\ g_{\ell+M+i}(u) &= u_i - u_{b,i}, & \text{for } i = 1, \dots, M. \end{aligned}$$

171 With this definition, we are able to formulate our optimal control problem as an
 172 optimization problem in finite dimensions as follows:

$$173 \quad (\mathcal{P}') \quad \begin{cases} \min_{u \in \mathbb{R}^M} f(u) = \frac{1}{2} \int_{\Omega} \left(\sum_{i=1}^M u_i y_i - \bar{y}_d \right)^2 dx + \frac{\alpha}{2} \|u\|_2^2 + \beta \|u\|_1 \\ \text{subject to:} \\ g(u) \leq 0. \end{cases}$$

174

175 *Hypothesis 2.7* (Slater condition). There exist a control vector $u^\circ \in \text{int } \mathcal{U}_{ad}$, such
 176 that

$$177 \quad (2.6) \quad \sum_{i=1}^M u_i^\circ y_i(x_j) < b_j, \quad \forall j = 1, \dots, \ell.$$

178 *Remark 2.8.* In the absence of equality constraints, we shall notice that Hypoth-
 179 esis 2.7 implies Mangasarian–Fromovitz constraint qualification.

180 The following proposition establishes first order optimality conditions for (\mathbf{P}_M) .
 181 The following conditions are obtained by applying Theorem 2.2 in [2].

182 **THEOREM 2.9.** *Let $\bar{u} \in \mathcal{U}_{feas}$ be the optimal solution for the problem (\mathbf{P}_M) and \bar{y}*
 183 *be the optimal state. Then, under Hypothesis 2.7 there exists $(\nu, \mu_a, \mu_b) \in \mathbb{R}^\ell \times \mathbb{R}^M \times$*

184 \mathbb{R}^M such that for all $i = 1, \dots, M$, it follows that

$$185 \quad (2.7a) \quad \int_{\Omega} (\bar{y} - y_d) y_i dx + \alpha \bar{u}_i + \sum_{j=1}^{\ell} \nu_j y_i(x_j) + \mu_{b_i} - \mu_{a_i} = -\beta, \quad \text{if } \bar{u}_i > 0,$$

$$186 \quad (2.7b) \quad \int_{\Omega} (\bar{y} - y_d) y_i dx + \alpha \bar{u}_i + \sum_{j=1}^{\ell} \nu_j y_i(x_j) + \mu_{b_i} - \mu_{a_i} = \beta, \quad \text{if } \bar{u}_i < 0,$$

$$187 \quad (2.7c) \quad \left| \int_{\Omega} (\bar{y} - y_d) y_i dx + \sum_{j=1}^{\ell} \nu_j y_i(x_j) + \mu_{b_i} - \mu_{a_i} \right| \leq \beta, \quad \text{if } \bar{u}_i = 0,$$

188
189

$$190 \quad (2.8a) \quad \nu \geq 0, \quad \text{and} \quad \nu_j (\bar{y}(x_j) - b_j) = 0,$$

$$191 \quad (2.8b) \quad \mu_a \geq 0, \quad \text{and} \quad \mu_{a_i} (u_{a,i} - \bar{u}_i) = 0,$$

$$192 \quad (2.8c) \quad \mu_b \geq 0, \quad \text{and} \quad \mu_{b_i} (\bar{u}_i - u_{b,i}) = 0.$$

194 *Proof.* By composition, $f := J(S \cdot, \cdot)$ is a Lipschitz function. Moreover, g is
195 a continuously differentiable affine function. By Hypothesis 2, the condition (i) of
196 Theorem 2 in [2] is satisfied, therefore we conclude the existence of $(\nu, \mu_a, \mu_b) \in$
197 $\mathbb{R}^{\ell} \times \mathbb{R}^M \times \mathbb{R}^M$ satisfying:

$$198 \quad (2.9a) \quad \sum_{j=1}^{\ell} \nu_j (v_j - g_j(\bar{u})) \leq 0, \quad \forall v \in \mathbb{R}^{\ell}, \text{ with } v \leq 0$$

$$199 \quad (2.9b) \quad \sum_{i=1}^M \mu_{a_i} (u_i - g_{\ell+i}(\bar{u})) \leq 0, \quad \forall u \in \mathbb{R}^M, \text{ with } u \leq 0$$

$$200 \quad (2.9c) \quad \sum_{i=1}^M \mu_{b_i} (u_i - g_{\ell+M+i}(\bar{u})) \leq 0, \quad \forall u \in \mathbb{R}^M, \text{ with } u \leq 0.$$

201

202 and

$$203 \quad (2.10) \quad 0 \in \partial f(\bar{u}) + (\nu, \mu_a, \mu_b)^{\top} \nabla g(\bar{u})$$

205 Since the previous inequalities (2.9) fulfill for any $v \leq 0$ and $u \leq 0$ respectively, by
206 appropriately testing v and u , we deduce that $(\nu, \mu_a, \mu_b) \geq 0$ and we also get the
207 relations:

$$208 \quad \nu_j g_j(\bar{u}) = \nu_j \left(\sum_{i=1}^M \bar{u}_i y_i(x_j) - b_j \right) = 0, \quad \forall j = 1, \dots, \ell,$$

$$209 \quad \mu_{a_i} g_{\ell+i}(\bar{u}) = \mu_{a_i} (u_{a,i} - \bar{u}_i) = 0, \quad \forall i = 1, \dots, M$$

$$210 \quad \mu_{b_i} g_{\ell+M+i}(\bar{u}) = \mu_{b_i} (\bar{u}_i - u_{b,i}) = 0, \quad \forall i = 1, \dots, M.$$

212 On the other hand, from (2.10) and since f is the sum of a differentiable term and
213 the 1-norm term, the usual rules of subdifferential calculus (see [12]) imply

$$214 \quad 0 \in [(\bar{y} - y_d, y_i)]_{i=1}^M + \alpha \bar{u} + \beta \partial \|\cdot\|_1(\bar{u}) + \left[\sum_{j=1}^{\ell} \nu_j y_i(x_j) \right]_{i=1}^M - \mu_a + \mu_b,$$

215 which is equivalent to:

$$216 \quad [(\bar{y} - y_d, y_i)]_{i=1}^M + \alpha \bar{u} + \left[\sum_{j=1}^{\ell} \nu_j y_i(x_j) \right]_{i=1}^M - \mu_a + \mu_b \in -\beta \partial \|\cdot\|_1(\bar{u}).$$

217 Finally, (2.7) is a consequence the subdifferential of the 1-norm:

$$218 \quad \partial \|\cdot\|_1(\bar{u}) = \left\{ z \in \mathbb{R}^M : z_i \in \begin{cases} \{\text{sign}(\bar{u}_i)\}, & \text{if } \bar{u}_i \neq 0, \\ [-1, 1], & \text{if } \bar{u}_i = 0. \end{cases} \right\}.$$

219 □

220 **3. Approximation by the finite element method.** In the previous section
 221 our problem was transformed into a finite-dimensional programming problem. How-
 222 ever, its numerical solution still requires the approximation of the state equation. Let
 223 us describe the numerical approximation of the state equation by finite elements. The
 224 associated definitions and classical results of the approximation can be found in [10]
 225 and [26].

226 **3.1. Discretization of the state equation .** We consider a family of meshes
 227 $(\mathcal{T}_h)_{h>0}$ of $\bar{\Omega}$, which consist of triangles $T \in \mathcal{T}_h$ such that the following conditions are
 228 satisfied:

- 229 • $\bigcup_{T \in \mathcal{T}_h} T = \bar{\Omega}$ and
- 230 • For two triangular elements T_i and T_j , $i \neq j$ they share a vertex, a side or
 231 are disjoint.

232 For each triangle $T \in \mathcal{T}_h$, we denote $\rho(T)$ the diameter of T , and $\sigma(T)$ the diameter
 233 of the largest ball contained in T . The mesh size h is defined by

$$234 \quad h = \max_{T \in \mathcal{T}_h} \rho(T).$$

235 In addition, we impose in addition the following regularity assumption on the grid:

236 *Hypothesis 3.1.* There exist two positive constants ρ y σ such that

$$237 \quad \frac{\rho(T)}{\sigma(T)} \leq \sigma \quad \text{y} \quad \frac{h}{\rho(T)} \leq \rho, \quad \forall T \in \mathcal{T}_h,$$

238 for all $h > 0$.

239 Associated with the triangulation \mathcal{T}_h , we define the set of piecewise linear and con-
 240 tinuous functions:

$$241 \quad Y_h = \{y_h \in C(\bar{\Omega}) : y_h|_T \in P_1(T), \forall T \in \mathcal{T}_h, y_h = 0 \text{ on } \Gamma\},$$

242 where $P_1(T)$ denotes the set of affine real-valued functions defined on T . Additionally,
 243 for each $i = 1, \dots, M$, we define the discrete state y_i^h , as the unique function of Y_h
 244 that satisfies

$$245 \quad (3.1) \quad \sum_{j,k=1}^2 (a_{jk} \partial_j y_i^h, \partial_k \phi_h) + (a_0 y_i^h, \phi_h) = (e_i, \phi_h),$$

246 for all $\phi_h \in Y_h$. The the error of the approximation of the solution of (2.1) by the
 247 solution of equation (3.1) can be estimated in the several norms, summarized in the
 248 next result.

249 PROPOSITION 3.2. *There exists a constant $c > 0$ independent of h such that*

250 (3.2a) $\|y_i - y_i^h\| \leq ch^2,$

251 (3.2b) $\|y_i - y_i^h\|_{L^\infty(\Omega)} \leq ch,$

252 (3.2c) $\|y_i - y_i^h\|_{L^\infty(\Omega_0)} \leq ch^2 |\log h|,$

254 for $i = 1, \dots, M$.

255 For the proof of (3.2a) and (3.2c) see [24, Proposition 3.3]. The estimate (3.2b) can
256 be found in [25, Theorem 3.1].

257 The *discrete state* associated to $u \in \mathcal{U}_{ad}$ is denoted by y_u^h and it is defined as the
258 unique element y_u^h of Y_h that satisfies the equation

259 (3.3)
$$\sum_{j,k=1}^2 (a_{jk} \partial_j y_u^h, \partial_k \phi_h) + (a_0 y_u^h, \phi_h) = \left(\sum_{i=1}^M u_i e_i, \phi_h \right), \quad \forall \phi \in Y_h.$$

260 Analogously, using the superposition principle the discrete state y_u^h associated to
261 $u \in \mathcal{U}_{ad}$ can be expressed as:

262 (3.4)
$$y_u^h = \sum_{i=1}^M u_i y_i^h.$$

263 Observe that by the triangular inequality it follows that the Proposition 3.2 holds
264 for y_u and y_u^h . Furthermore, we can also define the *discrete control-to-state operator*
265 $S^h : \mathbb{R}^M \rightarrow Y_h$ by

266 (3.5)
$$u \mapsto \sum_{i=1}^M u_i y_i^h.$$

267 **3.2. Discretization of the optimal control problem .** We are in place to
268 define the following discrete counterpart of problem (\mathbf{P}_M) by replacing the discrete
269 state y_u^h defined in (3.4) in (\mathbf{P}_M) instead of y . Therefore, the discrete optimal control
270 problem reads:

271 (\mathbf{P}_M^h)
$$\left\{ \begin{array}{l} \min_{(y^h, u)} J(y^h, u) := \frac{1}{2} \int_{\Omega} (y^h - y_d)^2 dx + \frac{\alpha}{2} \|u\|_2^2 + \beta \|u\|_1 \\ \text{subject to:} \\ \text{discrete equation (3.3),} \\ y^h(x_j) \leq b_j, \quad \forall j = 1, \dots, \ell, \\ u \in \mathcal{U}_{ad}. \end{array} \right.$$

272 The existence of a unique optimal control \bar{u}^h for (\mathbf{P}_M^h) is obtained by the same finite-
273 dimensional arguments as in the case of problem (\mathbf{P}_M) . The main result of this paper
274 is given in the following theorem.

275 THEOREM 3.3. *Let \bar{u} and \bar{u}^h be the optimal controls of problems (\mathbf{P}_M) and (\mathbf{P}_M^h)
276 respectively. Then, the estimate*

277 (3.6) $\|\bar{u} - \bar{u}^h\| \leq Ch^2 |\log h|,$

278 holds for some constant $C > 0$, independent of mesh parameter h .

279 We claim that this estimate is optimal in the sense that it has the same order of error
 280 than for the state equation given in Proposition 3.2. These kind of estimates are
 281 important for a-priori knowledge of how fast we can improve the expected quality of
 282 the solution when the number of elements in the discretization is increased. Moreover,
 283 error estimates can be used to design strategies for the reduction of the dimensionality
 284 of the optimal control problem, see [23].

285 The error estimate (3.6) is analogous to the estimate obtained in [24] for the
 286 smooth case. However, its derivation requires to cope with the non differentiability
 287 of the 1-norm. In order to derive the estimate (3.6) we proceed to reformulate prob-
 288 lem (\mathcal{P}) as a equivalent differentiable problem by splitting the control variable in its
 289 positive and negative parts. Then, a stability analysis is carried out which allows us
 290 to establish the desired estimate. We reserve the proof of this theorem for Section 5.

291 **4. Reformulation of (\mathbf{P}_M) as a smooth problem.** In order to reformulate
 292 our problem, we introduce the following notation. The positive and negative parts of
 293 a vector $u \in \mathbb{R}^M$ are denoted by the vectors $u^+ \in \mathbb{R}^M$ and $u^- \in \mathbb{R}^M$, respectively
 294 and whose components are defined by

$$295 \quad (4.1) \quad u_i^+ = \begin{cases} u_i & \text{if } u_i \geq 0, \\ 0 & \text{otherwise,} \end{cases} \quad \text{and} \quad u_i^- = \begin{cases} -u_i & \text{if } u_i \leq 0, \\ 0 & \text{otherwise.} \end{cases}$$

296 We recall following properties of the positive and negative parts of $u \in \mathbb{R}^M$:

$$297 \quad (4.2a) \quad u = u^+ - u^-,$$

$$298 \quad (4.2b) \quad u^{+\top} u^- = 0,$$

$$299 \quad (4.2c) \quad \|u\|_1 = \sum_{i=1}^M (u_i^+ + u_i^-),$$

$$300 \quad (4.2d) \quad \|u\|_2^2 = \sum_{i=1}^M [u_i^{+2} + u_i^{-2}].$$

302 Using relations (4.2) we reformulate problem (\mathbf{P}_M) as a smooth problem, however
 303 the number of variables is duplicated. In the following we will use bold characters
 304 to denote variables associated with (\mathbf{P}_M). The reformulated smooth optimal control
 305 problem reads:

$$306 \quad (\mathbf{P}_{2M}) \quad \left\{ \begin{array}{l} \min_{(y, \mathbf{u})} \mathbf{J}(y, \mathbf{u}) = \frac{1}{2} \int_{\Omega} (y - y_d)^2 dx + \frac{\alpha}{2} \|\mathbf{u}\|_2^2 + \beta \sum_{i=1}^{2M} \mathbf{u}_i \\ \text{subject to:} \\ Ay(x) = \sum_{i=1}^M (\mathbf{u}_i - \mathbf{u}_{i+M}) e_i(x), \quad \text{in } \Omega, \\ y(x) = 0, \quad \text{on } \Gamma, \\ y(x_j) \leq b_j, \quad \forall j = 1, \dots, \ell, \\ [\mathbf{u}_i - \mathbf{u}_{i+M}]_{i=1}^M \in \mathcal{U}_{ad}, \\ \mathbf{u}_i \geq 0, \quad \forall i = 1, \dots, 2M. \end{array} \right.$$

307 Note that problem (\mathbf{P}_{2M}) has a unique solution if $\alpha > 0$. Let us denote its solution
 308 by $\bar{\mathbf{u}}$. Before we proof the equivalence of problems (\mathbf{P}_M) and (\mathbf{P}_{2M}) and the relation

309 of the positive and negative parts of a vector, the following orthogonality result is
 310 expected from the decomposition (4.2a) and its proof can be found in [33].

311 LEMMA 4.1. *If $\bar{\mathbf{u}}$ is the optimal control of (\mathbf{P}_{2M}) then $\bar{\mathbf{u}}_i \bar{\mathbf{u}}_{i+M} = 0$ for all $i =$
 312 $1, \dots, M$.*

313 The equivalence between problems (\mathbf{P}_M) and (\mathbf{P}_{2M}) is stated in the following
 314 lemma.

315 LEMMA 4.2. *Let \bar{u} and $\bar{\mathbf{u}}$ be the solutions of the problems (\mathbf{P}_M) and (\mathbf{P}_{2M}) re-
 316 spectively, then*

$$317 \quad [\bar{\mathbf{u}}_i]_{i=1}^M = \bar{u}^+ \quad \text{and} \quad [\bar{\mathbf{u}}_{i+M}]_{i=1}^M = \bar{u}^-,$$

$$318 \quad \text{i.e. } \bar{u} = [\bar{\mathbf{u}}_i]_{i=1}^M - [\bar{\mathbf{u}}_{i+M}]_{i=1}^M.$$

319 *Proof.* Let $\bar{\mathbf{u}}$ be the solution for problem (\mathbf{P}_{2M}) and \bar{u} the solution of (\mathbf{P}_M) with
 320 \bar{u}^+ and \bar{u}^- their positive and negative parts respectively. Then, let us define the
 321 vectors:

$$322 \quad \bar{\mathbf{w}} = \begin{bmatrix} \bar{u}^+ \\ \bar{u}^- \end{bmatrix} \in \mathbb{R}^{2M}, \quad \text{and} \quad \bar{w} = [\bar{\mathbf{u}}_i - \bar{\mathbf{u}}_{i+M}]_{i=1}^M \in \mathbb{R}^M.$$

323 Since $[\bar{\mathbf{w}}_i - \bar{\mathbf{w}}_{i+M}]_{i=1}^M = \bar{u}^+ - \bar{u}^- = \bar{u} \in \mathcal{U}_{ad}$ and

$$324 \quad \sum_{i=1}^M (\bar{\mathbf{w}}_i - \bar{\mathbf{w}}_{i+M}) y_i(x_j) - b_j = \sum_{i=1}^M \bar{u}_i y_i(x_j) - b_j \leq 0, \quad \forall j = 1, \dots, \ell,$$

325 it follows that $\bar{\mathbf{w}}$ is an admissible control for (\mathbf{P}_{2M}) with associated state $y_{\bar{\mathbf{w}}} =$
 326 $S[\bar{\mathbf{w}}_i - \bar{\mathbf{w}}_{i+M}]_{i=1}^M$. Moreover, taking into account (4.2b) we have that

$$\begin{aligned} \mathbf{J}(y_{\bar{\mathbf{w}}}, \bar{\mathbf{w}}) &= \frac{1}{2} \|y_{\bar{\mathbf{w}}} - y_d\|^2 + \frac{\alpha}{2} \|\bar{\mathbf{w}}\|_2^2 + \beta \sum_{i=1}^{2M} \bar{\mathbf{w}}_i \\ &= \frac{1}{2} \left\| \sum_{i=1}^M (\bar{\mathbf{w}}_i - \bar{\mathbf{w}}_{i+M}) y_i - y_d \right\|^2 + \frac{\alpha}{2} \sum_{i=1}^{2M} \bar{\mathbf{w}}_i^2 + \beta \sum_{i=1}^{2M} \bar{\mathbf{w}}_i \\ 327 \quad (4.3) \quad &= \frac{1}{2} \left\| \sum_{i=1}^M \bar{u}_i y_i - y_d \right\|^2 + \frac{\alpha}{2} \sum_{i=1}^M (\bar{u}_i^+ + \bar{u}_i^-)^2 + \beta \sum_{i=1}^M \bar{u}_i^+ + \bar{u}_i^- \\ &= \frac{1}{2} \|\bar{y} - y_d\|^2 + \frac{\alpha}{2} \sum_{i=1}^M (\bar{u}_i)^2 + \beta \sum_{i=1}^M |\bar{u}_i| = J(\bar{y}, \bar{u}). \end{aligned}$$

328 Analogously, \bar{w} is admissible for problem (\mathbf{P}_M) and $J(y_{\bar{w}}, \bar{w}) = \mathbf{J}(y_{\bar{\mathbf{u}}}, \bar{\mathbf{u}})$. Then, by
 329 (4.3) and the optimality of (\bar{y}, \bar{u}) we have that

$$330 \quad \mathbf{J}(y_{\bar{\mathbf{w}}}, \bar{\mathbf{w}}) = J(\bar{y}, \bar{u}) \leq J(y_{\bar{w}}, \bar{w}) = \mathbf{J}(y_{\bar{\mathbf{u}}}, \bar{\mathbf{u}}).$$

331 Therefore, uniqueness of the solution of the optimal control problem (\mathbf{P}_{2M}) allows
 332 us to conclude that $\bar{\mathbf{w}} = \bar{\mathbf{u}}$. Similarly, we deduce that $\bar{u} = \bar{w}$, therefore the proof is
 333 complete. \square

334 Let us rewrite problem (\mathbf{P}_{2M}) as a finite dimensional programming problem.
 335 First, by replacing the control-to-state operator S , consider the reduced cost function

$$336 \quad (4.4) \quad \mathbf{f}(\mathbf{u}) = \frac{1}{2} \int_{\Omega} (S([\mathbf{u}_i - \mathbf{u}_{i+M}]_{i=1}^M) - y_d)^2 dx + \frac{\alpha}{2} \|\mathbf{u}\|_2^2 + \beta \sum_{i=1}^{2M} \mathbf{u}_i.$$

337 For convenience, we introduce the symmetric and positive semi-definite matrix $\mathcal{Y} \in$
 338 $\mathbb{R}^{2M \times 2M}$ which has the form

$$339 \quad (4.5) \quad \mathcal{Y} = \begin{bmatrix} Y & -Y \\ -Y & Y \end{bmatrix},$$

340 where the entries of $Y \in \mathbb{R}^{M \times M}$ are given by

$$341 \quad (4.6) \quad Y_{ij} = \int_{\Omega} y_i(x)y_j(x) dx.$$

342 Moreover, the vector $\mathbf{b} \in \mathbb{R}^{2M}$ is defined by its components $\mathbf{b}_i = \int_{\Omega} y_i y_d dx$ and
 343 $\mathbf{b}_{i+M} = -\int_{\Omega} y_i y_d dx$, for $i = 1, \dots, M$. In this way, we can rewrite

$$344 \quad (4.7) \quad \mathbf{f}(\mathbf{u}) = \frac{1}{2} \mathbf{u}^{\top} (\mathcal{Y} + \alpha I) \mathbf{u} + \mathbf{b}^{\top} \mathbf{u} + \beta \sum_{i=1}^{2M} \mathbf{u}_i + c,$$

345 with corresponding constant c term depending on y_d . In addition, we collect all
 346 inequality constraints by defining the constraint function $\mathbf{g} : \mathbb{R}^{2M} \rightarrow \mathbb{R}^{\ell+4M}$ as follows:
 347

$$348 \quad \mathbf{g}_j(\mathbf{u}) = \sum_{i=1}^M (\mathbf{u}_i - \mathbf{u}_{i+M}) y_i(x_j) - b_j$$

$$349 \quad (4.8a) \quad = [\mathcal{Z} \mathbf{u} - b]_j, \quad \text{for } j = 1, \dots, \ell,$$

$$350 \quad (4.8b) \quad \mathbf{g}_{\ell+i}(\mathbf{u}) = u_{a,i} - \mathbf{u}_i + \mathbf{u}_{i+M}, \quad \text{for } i = 1, \dots, M,$$

$$351 \quad (4.8c) \quad \mathbf{g}_{i+\ell+M}(\mathbf{u}) = \mathbf{u}_i - \mathbf{u}_{i+M} - u_{b,i}, \quad \text{for } i = 1, \dots, M,$$

$$352 \quad (4.8d) \quad \mathbf{g}_{i+\ell+2M}(\mathbf{u}) = -\mathbf{u}_i, \quad \text{for } i = 1, \dots, 2M.$$

354 Here, the matrix $\mathcal{Z} \in \mathbb{R}^{\ell \times 2M}$ is such that $\mathcal{Z}_{ji} = y_i(x_j)$ for and $\mathcal{Z}_{j\ell+i} = -y_i(x_j)$
 355 $i = 1, \dots, M$ and $j = 1, \dots, \ell$. Hence, using the above definitions we can write
 356 problem (\mathbf{P}_{2M}) in the compact form: we can write problem (\mathbf{P}_{2M}) in the compact
 357 form:

$$358 \quad (4.9) \quad \begin{cases} \min_{\mathbf{u} \in \mathbb{R}^{2M}} \mathbf{f}(\mathbf{u}) \\ \text{subject to:} \\ \mathbf{g}(\mathbf{u}) \leq 0. \end{cases}$$

359 The *Lagrangian* associated to problem (4.9) is given by:

$$360 \quad (4.10) \quad \mathcal{L}(\mathbf{u}, \boldsymbol{\nu}, \boldsymbol{\mu}_a, \boldsymbol{\mu}_b, \boldsymbol{\eta}) = \mathbf{f}(\mathbf{u}) + \sum_{j=1}^{\ell} \nu_j \mathbf{g}_j(\mathbf{u}) + \sum_{i=1}^M \mu_{a,i} \mathbf{g}_{i+\ell}(\mathbf{u}) + \sum_{i=1}^M \mu_{b,i} \mathbf{g}_{i+\ell+M}(\mathbf{u})$$

$$361 \quad + \sum_{i=1}^{2M} \eta_i \mathbf{g}_{i+\ell+2M}(\mathbf{u}).$$

362

363 In the following lemma we establish *first-order necessary optimality conditions*
 364 c.f. [3].

365 LEMMA 4.3. *Under Hypothesis 2.7, it follows that the constraints given by \mathbf{g} are*
 366 *regular in the sense of [27] and [28]. That is,*

$$367 \quad (4.11) \quad 0 \in \text{int}\{\mathbf{g}(\bar{\mathbf{u}}) + \mathbf{g}'(\bar{\mathbf{u}})(\mathbb{R}^{2M} - \bar{\mathbf{u}}) + \mathbb{R}_+^{4M+\ell}\}.$$

368 *Proof.* Since \mathbf{g} is defined in terms of affine functions c.f. (4.8), then there are a
 369 matrix $G \in \mathbb{R}^{\ell+4M \times 2M}$ and a vector $d \in \mathbb{R}^{\ell+4M}$ such that the constraint $\mathbf{g}(\mathbf{u}) \leq 0$
 370 equivalently expressed as $G\mathbf{u} \leq d$ for all $\mathbf{u} \in \mathbb{R}^{2M}$. In this case, condition (4.11) is
 371 equivalent to

$$372 \quad (4.12) \quad d \in \text{int}\{G(\mathbb{R}^{2M}) + \mathbb{R}_+^{\ell+4M}\},$$

373 which is verified with the help of \mathbf{u}° . Indeed, by its definition we have that $0 > \mathbf{g}(\mathbf{u}^\circ) =$
 374 $G\mathbf{u}^\circ - d$. This fact implies that $d \in \text{int}\{\mathbb{R}_+^{\ell+4M}\} + G\mathbf{u}^\circ \subset \text{int}\{G(\mathbb{R}^{2M}) + \mathbb{R}_+^{\ell+4M}\}$ \square

375 THEOREM 4.4. *Let $\bar{\mathbf{u}} \in \mathbb{R}^{2M}$ the solution of problem (4.9) with associated optimal*
 376 *state $\bar{\mathbf{y}} = \sum_{i=1}^M (\bar{\mathbf{u}}_i - \bar{\mathbf{u}}_{i+M})y_i$. Then, there exists $(\bar{\nu}, \bar{\mu}_a, \bar{\mu}_b, \bar{\eta}) \in \mathbb{R}_+^\ell \times \mathbb{R}_+^M \times \mathbb{R}_+^M \times$*
 377 *$\mathbb{R}_+^{2M} \cup \{0\}$ such that the conditions below are satisfied,*

$$378 \quad (4.13a) \quad (\bar{\mathbf{y}} - y_d, y_i) + \alpha \bar{\mathbf{u}}_i + \beta + \sum_{j=1}^{\ell} \nu_j y_i(x_j) + \mu_{b_i} - \mu_{a_i} - \eta_i = 0,$$

$$379 \quad (4.13b) \quad -(\bar{\mathbf{y}} - y_d, y_i) + \alpha \bar{\mathbf{u}}_{i+M} + \beta - \sum_{j=1}^{\ell} \nu_j y_i(x_j) - \mu_{b_i} + \mu_{a_i} - \eta_{i+M} = 0.$$

380

381 *for $i = 1, \dots, M$, holds, together with the following conditions:*

$$382 \quad (4.14a) \quad \bar{\mathbf{u}} \geq 0,$$

$$383 \quad (4.14b) \quad u_{a_i} \leq \bar{\mathbf{u}}_i - \bar{\mathbf{u}}_{i+M} \leq u_{b_i}, \quad \forall i = 1, \dots, 2M.$$

$$384 \quad (4.14c) \quad \nu_j (\bar{\mathbf{y}}(x_j) - b_j) = 0, \quad \forall j = 1, \dots, \ell,$$

$$385 \quad (4.14d) \quad \mu_{a_i} (u_{a,i} - \bar{\mathbf{u}}_i + \bar{\mathbf{u}}_{i+M}) = 0, \quad \forall i = 1, \dots, M,$$

$$386 \quad (4.14e) \quad \mu_{b_i} (\bar{\mathbf{u}}_i - \bar{\mathbf{u}}_{i+M} - u_{b,i}) = 0, \quad \forall i = 1, \dots, M,$$

$$387 \quad (4.14f) \quad \eta_i \bar{\mathbf{u}}_i = 0, \quad \forall i = 1, \dots, 2M.$$

389 *Proof.* Hypothesis 2.7 and Lemma 4.2 imply that $\bar{\mathbf{u}}$ is a regular point for problem
 390 (\mathbf{P}_{2M}) . Therefore, the result follows by applying the standard theory of constrained
 391 optimization c.f. [3]. \square

392 Similarly, we consider the reduced discretized counterpart of problem (\mathbf{P}_{2M}) by
 393 replacing the state equation by the finite element approximation (3.3). The problem
 394 is formulated as follows:

$$395 \quad (\mathbf{P}_{2M}^h) \left\{ \begin{array}{l} \min_{(y^h, \mathbf{u})} \mathbf{J}(y^h, \mathbf{u}) = \frac{1}{2} \int_{\Omega} (y^h - y_d)^2 dx + \frac{\alpha}{2} \|\mathbf{u}\|_2^2 + \beta \sum_{i=1}^{2M} \mathbf{u}_i \\ \text{subject to:} \\ \text{discrete equation (3.3) with } u_i = \mathbf{u}_i - \mathbf{u}_{i+M}, i = 1, \dots, M. \\ y^h(x_j) \leq b_j, \quad \forall j = 1, \dots, \ell, \\ [\mathbf{u}_i - \mathbf{u}_{i+M}]_{i=1}^M \in \mathcal{U}_{ad}, \\ \mathbf{u}_i \geq 0, \quad \forall i = 1, \dots, 2M. \end{array} \right.$$

396 Observe that $\bar{\mathbf{u}}^h$ is a unique solution for (\mathbf{P}_{2M}^h) . This fact can be deduced by the
 397 same arguments as in the non discretized case. Analogously to Lemma 4.2 for problem
 398 (\mathbf{P}_{2M}) , it is also true that (\mathbf{P}_{2M}^h) is equivalent to problem (\mathbf{P}_M^h) and this equivalence
 399 is expressed by the relations:

$$400 \quad (4.15) \quad [\bar{\mathbf{u}}_i^h]_{i=1}^M = \bar{u}^{h+} \quad \text{and} \quad [\bar{\mathbf{u}}_{i+M}^h]_{i=1}^M = \bar{u}^{h-}$$

$$401 \quad (4.16) \quad \text{i.e., } \bar{u}^h = [\bar{\mathbf{u}}_i^h]_{i=1}^M - [\bar{\mathbf{u}}_{i+M}^h]_{i=1}^M.$$

403 Because of the equivalence of problems (\mathbf{P}_M) and (\mathbf{P}_{2M}) (also between (\mathbf{P}_M^h) and
 404 (\mathbf{P}_{2M}^h)) and the relation of the positive and negative parts of a vector, the following
 405 orthogonality result is expected and its proof can be found in [33].

406 LEMMA 4.5. *If $\bar{\mathbf{u}}$ (respectively $\bar{\mathbf{u}}^h$) is the optimal control of (\mathbf{P}_{2M}) (respectively
 407 (\mathbf{P}_{2M}^h)) then $\bar{\mathbf{u}}_i \bar{\mathbf{u}}_{i+M} = 0$ (respectively $\bar{\mathbf{u}}_i^h \bar{\mathbf{u}}_{i+M}^h = 0$) for all $i = 1, \dots, M$.*

408 Moreover, we denote by \mathbf{g}^h and \mathcal{Y}^h the discrete analogues of \mathbf{g} and \mathcal{Y} defined in (4.8)
 409 and in (4.5), respectively, which are obtained by replacing y_i by y_i^h correspondingly.
 410 \mathbf{b}^h is defined analogously. The discrete objective function \mathbf{f}^h is also defined in this
 411 manner,

$$412 \quad (4.17) \quad \mathbf{f}^h(\mathbf{u}) = \frac{1}{2} \mathbf{u}^\top (\mathcal{Y}^h + \alpha I) \mathbf{u} + \mathbf{b}^{h\top} \mathbf{u} + \beta \sum_{i=1}^{2M} \mathbf{u}_i + c^h,$$

413 where c^h is an appropriate constant depending on y_d and h . Now we are able to write
 414 problem (\mathbf{P}_{2M}^h) compactly, as the following finite-dimensional optimization problem.

$$415 \quad (4.18) \quad \begin{cases} \min_{\mathbf{u} \in \mathbb{R}^{2M}} \mathbf{f}^h(\mathbf{u}) \\ \text{subject to:} \\ \mathbf{g}^h(\mathbf{u}) \leq 0. \end{cases}$$

416

417 *Remark 4.6.* From (3.2b) it follows that $|y_i(x_j) - y_i^h(x_j)| \leq ch^2 |\log(h)|$ for all
 418 $i = 1, \dots, M$ and all $j = 1, \dots, \ell$. This fact, together with Hypothesis 2.7 imply
 419 that u° is also a Slater point for (\mathbf{P}_M^h) . Therefore, $\bar{\mathbf{u}}^h$ is a regular point for (\mathbf{P}_{2M}^h)
 420 if h is sufficiently small. Applying the Lagrange principle again, we obtain a similar
 421 optimality system as in Theorem 4.4. The *discrete optimality system* for the discrete
 422 optimal pair $(\bar{\mathbf{u}}^h, \bar{y}^h)$ and the corresponding associated multipliers $(\bar{\nu}^h, \bar{\mu}_\alpha^h, \bar{\mu}_b^h, \bar{\eta}^h) \in$
 423 $\mathbb{R}_+^\ell \times \mathbb{R}_+^M \times \mathbb{R}_+^M \times \mathbb{R}_+^{2M} \cup \{0\}$, is given below.

$$424 \quad (4.19a) \quad (\bar{y}^h - y_d, y_i^h) + \alpha \bar{\mathbf{u}}_i^h + \beta + \sum_{j=1}^{\ell} \nu_j^h y_i^h(x_j) + \mu_{b_i}^h - \mu_{a_i}^h - \eta_i^h = 0,$$

$$425 \quad (4.19b) \quad -(\bar{y}^h - y_d, y_i^h) + \alpha \bar{\mathbf{u}}_{i+M}^h + \beta - \sum_{j=1}^{\ell} \nu_j^h y_i^h(x_j) - \mu_{b_i}^h + \mu_{a_i}^h - \eta_{i+M}^h = 0.$$

426

427 for $i = 1, \dots, M$. Moreover, the following conditions are satisfied:

$$428 \quad (4.20a) \quad \bar{\mathbf{u}}^h \geq 0,$$

$$429 \quad (4.20b) \quad u_{ai} \leq \bar{\mathbf{u}}_i^h - \bar{\mathbf{u}}_{i+M}^h \leq u_{bi}, \quad \forall i = 1, \dots, 2M.$$

$$430 \quad (4.20c) \quad \boldsymbol{\nu}_j^h (\bar{\mathbf{y}}^h(x_j) - b_j) = 0, \quad \forall j = 1, \dots, \ell,$$

$$431 \quad (4.20d) \quad \boldsymbol{\mu}_{a,i}^h (u_{a,i} - \bar{\mathbf{u}}_i^h + \bar{\mathbf{u}}_{i+M}^h) = 0, \quad \forall i = 1, \dots, M,$$

$$432 \quad (4.20e) \quad \boldsymbol{\mu}_{b,i}^h (\bar{\mathbf{u}}_i^h - \bar{\mathbf{u}}_{i+M}^h - u_{b,i}) = 0, \quad \forall i = 1, \dots, M,$$

$$433 \quad (4.20f) \quad \boldsymbol{\eta}_i^h \bar{\mathbf{u}}_i^h = 0, \quad \forall i = 1, \dots, 2M.$$

435 **5. Derivation of the error of convergence.** This section is devoted to the
 436 derivation of the order of convergence for the finite element approximation (\mathbf{P}_{2M}^h)
 437 (\mathbf{P}_{2M}) which is the intermediate step in order to derive an estimation of the error for
 438 the solution of (\mathbf{P}_M^h) . Therefore, we develop the necessary tool for the error estimation
 439 by regarding problem (4.18) as a perturbation of (4.9).

440 Continuing with our analysis, we define the cone $K = \mathbb{R}_+^\ell \times \mathbb{R}_+^M \times \mathbb{R}_+^M \times \mathbb{R}_+^{2M} \cup \{0\}$.
 441 Then, we are interested in the stability analysis, with respect to h , of the problem

$$442 \quad \begin{cases} \min_{\mathbf{u}} \mathbf{f}(\mathbf{u}) \\ \mathbf{g}(\mathbf{u}) \leq_K 0, \end{cases} \quad \text{and its perturbation} \quad \begin{cases} \min_{\mathbf{u}} \mathbf{f}^h(\mathbf{u}) \\ \mathbf{g}^h(\mathbf{u}) \leq_K 0. \end{cases}$$

443 For the comparison of the solutions of the preceding problems we first obtain the
 444 following estimation of the objective and the constraint functions.

445 **LEMMA 5.1.** *There exist $C > 0$, independent of h , such that for h small enough*
 446 *it follows that*

$$447 \quad \begin{aligned} & \|\mathbf{f}(\mathbf{u}) - \mathbf{f}^h(\mathbf{w})\| + \|\nabla \mathbf{f}(\mathbf{u}) - \nabla \mathbf{f}^h(\mathbf{w})\| + \|\nabla^2 \mathbf{f}(\mathbf{u}) - \nabla^2 \mathbf{f}^h(\mathbf{w})\| \\ & + \sum_{j=1}^{\ell} (|\mathbf{g}_j(\mathbf{u}) - \mathbf{g}_j^h(\mathbf{w})| + \|\nabla \mathbf{g}_j(\mathbf{u}) - \nabla \mathbf{g}_j^h(\mathbf{w})\|) \leq C(\|\mathbf{u} - \mathbf{w}\| + h^2 |\log h|) \end{aligned}$$

448 for all $\mathbf{u}, \mathbf{w} \in \mathbb{R}^{2M}$.

449 *Proof.* The proof of this Lemma can be found in [25], and follows by taking into
 450 account the estimate (3.2c) within the Taylor expansion of \mathbf{f} . The estimates for \mathbf{g} are
 451 direct since \mathbf{g} is affine linear with respect to \mathbf{u} . \square

452 **LEMMA 5.2.** *There exists a constant $C > 0$, independent of h , such that for h*
 453 *small enough, we have the estimate*

$$454 \quad (5.1) \quad \|\mathcal{Y} - \mathcal{Y}^h\| \leq Ch^2,$$

455 in the Frobenius norm. Also the estimate,

$$456 \quad (5.2) \quad \|\mathbf{b} - \mathbf{b}^h\| \leq Ch^2.$$

457 holds.

458 *Proof.* In order to estimate the Frobenius norm of $\mathcal{Y} - \mathcal{Y}^h$, we consider the dif-

459 fference of their entries. By applying Cauchy–Schwarz inequality we get:

$$\begin{aligned}
 460 \quad |Y_{ij} - Y_{ij}^h| &= \left| \int_{\Omega} y_i(x)y_j(x) dx - \int_{\Omega} y_i^h(x)y_j^h(x) dx \right| \\
 461 \quad &\leq \left| \int_{\Omega} (y_i - y_i^h)y_j dx \right| + \left| \int_{\Omega} (y_j^h - y_j)y_i^h dx \right| \\
 462 \quad &\leq \|y_i - y_i^h\| \|y_j\| + \|y_j^h - y_j\| \|y_i^h\|,
 \end{aligned}$$

for all $i, j = 1, \dots, M$. Therefore, recalling that $\|y_i^h\|$ is bounded for $i = 1, \dots, M$, and in view of (3.2a), we conclude that there exist a constant $C > 0$ such that

$$|Y_{ij} - Y_{ij}^h| \leq Ch^2.$$

464 From this estimate we compute (5.1). The estimate (5.2) follows by analogous argu-
465 ments. \square

466 Among other properties of the approximated problem (\mathbf{P}_{2M}^h) , boundedness of its
467 multipliers is an important feature in the derivation of error estimates. We show this
468 property in the following Lemma.

469 LEMMA 5.3. *For sufficiently small h , the corresponding multipliers $\boldsymbol{\nu}^h, \boldsymbol{\mu}_a^h, \boldsymbol{\mu}_b^h, \boldsymbol{\eta}^h$*
470 *associated with optimality conditions for problem (\mathbf{P}_{2M}^h) are uniformly bounded.*

471 *Proof.* Let us multiply equations (4.19) by $\mathbf{u}_i^\circ - \bar{\mathbf{u}}_i^h$. By adding them up we obtain:

$$\begin{aligned}
 472 \quad 0 &= \nabla \mathbf{f}^h(\bar{\mathbf{u}}^h)^\top (\mathbf{u}^\circ - \bar{\mathbf{u}}^h) + \sum_{i=1}^M \sum_{j=1}^{\ell} \boldsymbol{\nu}_j^h (\mathbf{u}_i^\circ - \bar{\mathbf{u}}_i^h - \mathbf{u}_{i+M}^\circ + \bar{\mathbf{u}}_{i+M}^h) y_i^h(x_j) \\
 473 \quad &+ \sum_{i=1}^M \boldsymbol{\mu}_{b_i}^h (\mathbf{u}_i^\circ - \bar{\mathbf{u}}_i^h - \mathbf{u}_{i+M}^\circ + \bar{\mathbf{u}}_{i+M}^h) - \sum_{i=1}^M \boldsymbol{\mu}_{a_i}^h (\mathbf{u}_i^\circ - \bar{\mathbf{u}}_i^h - \mathbf{u}_{i+M}^\circ + \bar{\mathbf{u}}_{i+M}^h) \\
 474 \quad &- \sum_{i=1}^{2M} \boldsymbol{\eta}_i^h (\mathbf{u}_i^\circ - \bar{\mathbf{u}}_i^h) \\
 475 \quad &= \nabla \mathbf{f}^h(\bar{\mathbf{u}}^h)^\top (\mathbf{u}^\circ - \bar{\mathbf{u}}^h) + \sum_{j=1}^{\ell} \boldsymbol{\nu}_j^h (\mathbf{y}^{\circ h}(x_j) - b_j) \\
 476 \quad &+ \sum_{i=1}^M \boldsymbol{\mu}_{b_i}^h (\mathbf{u}_i^\circ - \mathbf{u}_{i+M}^\circ - u_{b_i}) - \sum_{i=1}^M \boldsymbol{\mu}_{a_i}^h (\mathbf{u}_i^\circ - \mathbf{u}_{i+M}^\circ - u_{a_i}) - \sum_{i=1}^{2M} \boldsymbol{\eta}_i^h \mathbf{u}_i^\circ, \\
 477 \quad &
 \end{aligned}$$

478 which results in the relation

$$\begin{aligned}
 479 \quad \sum_{j=1}^{\ell} \boldsymbol{\nu}_j^h (b_j - \mathbf{y}^{\circ h}(x_j)) &+ \sum_{i=1}^M \boldsymbol{\mu}_{b_i}^h (u_{b_i} - \mathbf{u}_i^\circ + \mathbf{u}_{i+M}^\circ) \\
 480 \quad &+ \sum_{i=1}^M \boldsymbol{\mu}_{a_i}^h (\mathbf{u}_i^\circ - \mathbf{u}_{i+M}^\circ - u_{a_i}) + \sum_{i=1}^{2M} \boldsymbol{\eta}_i^h \mathbf{u}_i^\circ = \nabla \mathbf{f}^h(\bar{\mathbf{u}}^h)^\top (\mathbf{u}^\circ - \bar{\mathbf{u}}^h). \\
 481 \quad &
 \end{aligned}$$

482 Notice that $\min_j \{b_j - \mathbf{y}^{\circ h}(x_j)\} > \min_j \{b_j - \mathbf{y}^\circ(x_j)\} - ch^2 |\log(h)| > \gamma_1 > 0$
483 for some positive γ_1 and h sufficiently small. Moreover, taking into account that
484 $\gamma_2 := \min_i (u_{b_i} - \mathbf{u}_i^\circ + \mathbf{u}_{i+M}^\circ) > 0$, $\gamma_3 := \min_i (\mathbf{u}_i^\circ - \mathbf{u}_{i+M}^\circ - u_{a_i}) > 0$ and that $|\boldsymbol{\eta}^h| \leq \beta$,

485 together with the nonnegativity of the multipliers, it follows that

$$486 \quad \gamma_1 \sum_{j=1}^{\ell} |\boldsymbol{\nu}_j^h| + \gamma_2 \sum_{i=1}^M |\boldsymbol{\mu}_{b_i}^h| + \gamma_3 \sum_{i=1}^M |\boldsymbol{\mu}_{a_i}^h| \leq \nabla \mathbf{f}^h(\bar{\mathbf{u}}^h)^\top (\mathbf{u}^o - \bar{\mathbf{u}}^h).$$

487 The right-hand side in the last inequality is bounded because of the boundedness of
488 $\bar{\mathbf{u}}^h$ and its associated state. Thus, the last inequality just implies that the multipliers
489 are bounded. \square

490 *Remark 5.4.* The optimality systems (4.18) and (4.9) are obtained by the stan-
491 dard Lagrange principle, i.e. the corresponding optimal vectors satisfy:

$$492 \quad \begin{cases} \nabla_u \mathcal{L}(\bar{\mathbf{u}}, \boldsymbol{\nu}, \boldsymbol{\mu}_a, \boldsymbol{\mu}_b, \boldsymbol{\eta}) = 0, \\ \mathbf{g}(\bar{\mathbf{u}}) \leq_K 0, \\ (\boldsymbol{\nu}, \boldsymbol{\mu}_a, \boldsymbol{\mu}_b, \boldsymbol{\eta})^\top \mathbf{g}(\bar{\mathbf{u}}) = 0, \end{cases} \quad \text{and} \quad \begin{cases} \nabla_u \mathcal{L}^h(\bar{\mathbf{u}}^h, \boldsymbol{\nu}^h, \boldsymbol{\mu}_a^h, \boldsymbol{\mu}_b^h, \boldsymbol{\eta}^h) = 0, \\ \mathbf{g}^h(\bar{\mathbf{u}}^h) \leq_K 0, \\ (\boldsymbol{\nu}^h, \boldsymbol{\mu}_a^h, \boldsymbol{\mu}_b^h, \boldsymbol{\eta}^h)^\top \mathbf{g}^h(\bar{\mathbf{u}}^h) = 0. \end{cases}$$

493 By introducing the *normal cone* of K :

$$494 \quad \partial \psi_K(x) = \begin{cases} z \in \mathbb{R}^\ell \times \mathbb{R}^M \times \mathbb{R}^M \times \mathbb{R}^{2M} \text{ with } z^T(w - x) \leq 0 \quad \forall w \in K, & \text{if } x \in K, \\ \emptyset, & \text{if } x \notin K, \end{cases}$$

495 these optimality systems can be rewritten as the generalized equations. Indeed, fol-
496 lowing [29] we have:

$$497 \quad (5.3) \quad 0 \in \begin{bmatrix} \nabla_u \mathcal{L}(\bar{\mathbf{u}}, \boldsymbol{\nu}, \boldsymbol{\mu}_a, \boldsymbol{\mu}_b, \boldsymbol{\eta}) \\ -\mathbf{g}(\bar{\mathbf{u}}) \end{bmatrix} + \begin{bmatrix} \{0\} \\ \partial \psi_K(\boldsymbol{\nu}, \boldsymbol{\mu}_a, \boldsymbol{\mu}_b, \boldsymbol{\eta}) \end{bmatrix}$$

498 and

$$499 \quad (5.4) \quad 0 \in \begin{bmatrix} \nabla_u \mathcal{L}^h(\bar{\mathbf{u}}^h, \boldsymbol{\nu}^h, \boldsymbol{\mu}_a^h, \boldsymbol{\mu}_b^h, \boldsymbol{\eta}^h) \\ -\mathbf{g}^h(\bar{\mathbf{u}}^h) \end{bmatrix} + \begin{bmatrix} \{0\} \\ \partial \psi_K(\boldsymbol{\nu}^h, \boldsymbol{\mu}_a^h, \boldsymbol{\mu}_b^h, \boldsymbol{\eta}^h) \end{bmatrix}$$

501 respectively. Moreover, according to [25][Section 5.3] the solutions of (5.3) and (5.4)
502 can be viewed as solutions of a perturbed problem. Therefore, we consider the per-
503 turbed problem

$$504 \quad (5.5) \quad 0 \in \begin{bmatrix} \nabla_u \mathcal{L}(\mathbf{u}, \boldsymbol{\nu}, \boldsymbol{\mu}_a, \boldsymbol{\mu}_b, \boldsymbol{\eta}) \\ -\mathbf{g}(\mathbf{u}) \end{bmatrix} + \begin{bmatrix} \{0\} \\ \partial \psi_K(\boldsymbol{\nu}, \boldsymbol{\mu}_a, \boldsymbol{\mu}_b, \boldsymbol{\eta}) \end{bmatrix} - \delta.$$

506 It is clear that we can express (5.3) in terms of (5.5) by taking $(\mathbf{u}, \boldsymbol{\nu}, \boldsymbol{\mu}_a, \boldsymbol{\mu}_b, \boldsymbol{\eta}) =$
507 $(\bar{\mathbf{u}}, \bar{\boldsymbol{\nu}}, \bar{\boldsymbol{\mu}}_a, \bar{\boldsymbol{\mu}}_b, \bar{\boldsymbol{\eta}})$ and $\delta = 0$. For (5.4) we take $(\mathbf{u}, \boldsymbol{\nu}, \boldsymbol{\mu}_a, \boldsymbol{\mu}_b, \boldsymbol{\eta}) = (\bar{\mathbf{u}}^h, \bar{\boldsymbol{\nu}}^h, \bar{\boldsymbol{\mu}}_a^h, \bar{\boldsymbol{\mu}}_b^h, \bar{\boldsymbol{\eta}}^h)$
508 with associated perturbation $\delta(h)$ defined by:

$$509 \quad (5.6) \quad \delta(h) = \begin{cases} 0 & \text{if } h = 0 \\ \left[\begin{array}{c} \nabla_u \mathcal{L}(\bar{\mathbf{u}}^h, \bar{\boldsymbol{\nu}}^h, \bar{\boldsymbol{\mu}}_a^h, \bar{\boldsymbol{\mu}}_b^h, \bar{\boldsymbol{\eta}}^h) - \nabla_u \mathcal{L}^h(\bar{\mathbf{u}}^h, \bar{\boldsymbol{\nu}}^h, \bar{\boldsymbol{\mu}}_a^h, \bar{\boldsymbol{\mu}}_b^h, \bar{\boldsymbol{\eta}}^h) \\ \mathbf{g}(\bar{\mathbf{u}}^h) - \mathbf{g}^h(\bar{\mathbf{u}}^h) \end{array} \right], & \text{if } h \neq 0. \end{cases}$$

510 At this point, we have all the required elements to proof the error estimate for the
511 approximation of the optimal control. Observe that by Lemma 4.3 we do not need
512 strong regularity as in [25].

513 **Proof of Theorem 3.3 .**

514 **Non optimal estimate.** Consider the Slater point defined in (2.6). In view of
 515 the state constraints, we have that $\bar{\mathbf{u}}$ does not need to be feasible for (\mathbf{P}_{2M}^h) neither
 516 $\bar{\mathbf{u}}^h$ need to be feasible for (\mathbf{P}_{2M}) . Therefore, we define the intermediate vectors

$$\begin{aligned} 517 \quad \tilde{\mathbf{u}} &= (1-t)\bar{\mathbf{u}}^h + t\mathbf{u}^\circ, \quad \text{for } t \in (0, 1), \text{ and} \\ 518 \quad \hat{\mathbf{u}} &= (1-t)\bar{\mathbf{u}} + t\mathbf{u}^\circ, \quad \text{for } t \in (0, 1), \end{aligned}$$

520 where \mathbf{u}° is defined by $\mathbf{u}_i^\circ = u_i^{\circ+}$ and $\mathbf{u}_{i+M}^\circ = u_i^{\circ-}$, for $i = 1, \dots, M$. By convexity, $\tilde{\mathbf{u}}$
 521 is feasible for (\mathbf{P}_{2M}^h) and so does $\hat{\mathbf{u}}$ for problem (\mathbf{P}_{2M}) .

522 By defining $\varrho := \max_{j=1, \dots, \ell} \sum_{i=1}^M u_i^\circ y_i(x_j) < b_j$, then $\varrho - b_j < 0$. Moreover,
 523 since $\bar{\mathbf{u}}^h$ satisfies the state constraints for the discrete problem and also satisfies the
 524 estimate (3.2b) then, for h sufficiently small, $\tilde{\mathbf{u}}$ fulfills

$$\begin{aligned} 525 \quad \sum_{i=1}^M (\tilde{\mathbf{u}}_i - \tilde{\mathbf{u}}_{i+M}) y_i(x_j) &= (1-t) \sum_{i=1}^M (\bar{\mathbf{u}}_i^h - \bar{\mathbf{u}}_{i+M}^h) y_i(x_j) + t \sum_{i=1}^M (\mathbf{u}_i^\circ - \mathbf{u}_{i+M}^\circ) y_i(x_j) \\ 526 \quad &< (1-t) \sum_{i=1}^M (\bar{\mathbf{u}}_i^h - \bar{\mathbf{u}}_{i+M}^h) (y_i(x_j) - y_i^h(x_j)) + (1-t)b_j + t\varrho \\ 527 \quad &< ch^2 |\log(h)| \sum_{i=1}^M |\bar{\mathbf{u}}_i^h - \bar{\mathbf{u}}_{i+M}^h| + t(\varrho - b_j) + b_j \\ 528 \quad &< ch^2 |\log(h)| C_{ad} + t(\varrho - b_j) + b_j, \end{aligned}$$

530 where C_{ad} is a positive bound for the admissible set U_{ad} not depending on h . Here,
 531 we realize that the right-hand side is less than b_j by selecting

$$532 \quad (5.7) \quad t = \frac{-c}{\varrho - b_j} h^2 |\log(h)| C_{ad} > 0.$$

533 Therefore, $\tilde{\mathbf{u}}$ is feasible for problems (\mathbf{P}_{2M}) and (\mathbf{P}_{2M}^h) . The same conclusion is
 534 obtained for the control for $\hat{\mathbf{u}}$ with the same arguments. In view of feasibility of $\tilde{\mathbf{u}}$
 535 and $\hat{\mathbf{u}}$, and the optimality of $\bar{\mathbf{u}}$ and $\bar{\mathbf{u}}^h$, we have

$$536 \quad (5.8) \quad \nabla \mathbf{f}(\bar{\mathbf{u}})^\top (\tilde{\mathbf{u}} - \bar{\mathbf{u}}) \geq 0 \quad \text{and} \quad \nabla \mathbf{f}^h(\bar{\mathbf{u}}^h)^\top (\hat{\mathbf{u}} - \bar{\mathbf{u}}^h) \geq 0$$

538 By adding both inequalities in (5.8), and by replacing the corresponding gradients of
 539 \mathbf{f} and \mathbf{f}^h as well as the definition of $\hat{\mathbf{u}}$ and $\tilde{\mathbf{u}}$, we get that

$$\begin{aligned} 540 \quad 0 &\leq \bar{\mathbf{u}}^\top (\mathcal{Y} + \alpha I)(t\mathbf{u}^\circ - \bar{\mathbf{u}}) + (1-t)(\bar{\mathbf{u}}^h - \bar{\mathbf{u}}) + \bar{\mathbf{u}}^{h\top} (\mathcal{Y}^h + \alpha I)(t\mathbf{u}^\circ - \bar{\mathbf{u}}^h) + (1-t)(\bar{\mathbf{u}} - \bar{\mathbf{u}}^h) \\ 541 \quad &+ \mathbf{b}^\top (t\mathbf{u}^\circ - \bar{\mathbf{u}}) + (1-t)(\bar{\mathbf{u}}^h - \bar{\mathbf{u}}) + \mathbf{b}^{h\top} (t\mathbf{u}^\circ - \bar{\mathbf{u}}^h) + (1-t)(\bar{\mathbf{u}} - \bar{\mathbf{u}}^h) \\ 542 \quad &+ \beta \sum_{i=1}^{2M} (t\mathbf{u}_i^\circ - \bar{\mathbf{u}}_i) + (1-t)(\bar{\mathbf{u}}_i^h - \bar{\mathbf{u}}_i) + \beta \sum_{i=1}^{2M} (t\mathbf{u}_i^\circ - \bar{\mathbf{u}}_i^h) + (1-t)(\bar{\mathbf{u}}_i - \bar{\mathbf{u}}_i^h) \\ 543 \quad &= t \left[\bar{\mathbf{u}}^\top (\mathcal{Y} + \alpha I)(\mathbf{u}^\circ - \bar{\mathbf{u}}) + \bar{\mathbf{u}}^{h\top} (\mathcal{Y}^h + \alpha I)(\mathbf{u}^\circ - \bar{\mathbf{u}}^h) + \mathbf{b}^\top (\mathbf{u}^\circ - \bar{\mathbf{u}}) + \mathbf{b}^{h\top} (\mathbf{u}^\circ - \bar{\mathbf{u}}^h) \right. \\ 544 \quad &\quad \left. + \beta \sum_{i=1}^{2M} (\mathbf{u}_i^\circ - \bar{\mathbf{u}}_i) + \beta \sum_{i=1}^{2M} (\mathbf{u}_i^\circ - \bar{\mathbf{u}}_i^h) \right] \\ 545 \quad &+ (1-t) \left[\bar{\mathbf{u}}^\top (\mathcal{Y} + \alpha I)(\bar{\mathbf{u}}^h - \bar{\mathbf{u}}) + \bar{\mathbf{u}}^{h\top} (\mathcal{Y}^h + \alpha I)(\bar{\mathbf{u}} - \bar{\mathbf{u}}^h) + (\mathbf{b} - \mathbf{b}^h)(\bar{\mathbf{u}}^h - \bar{\mathbf{u}}) \right]. \end{aligned}$$

546 By rearranging terms, the last inequality is equivalent to

$$\begin{aligned}
547 & (1-t)(\bar{\mathbf{u}}^h - \bar{\mathbf{u}})^\top (\mathcal{Y} + \alpha I)(\bar{\mathbf{u}}^h - \bar{\mathbf{u}}) \\
548 & \leq t \left[\bar{\mathbf{u}}^\top (\mathcal{Y} + \alpha I)(\mathbf{u}^\circ - \bar{\mathbf{u}}) + \bar{\mathbf{u}}^{h\top} (\mathcal{Y}^h + \alpha I)(\mathbf{u}^\circ - \bar{\mathbf{u}}^h) + \mathbf{b}^\top (\mathbf{u}^\circ - \bar{\mathbf{u}}) + \mathbf{b}^{h\top} (\mathbf{u}^\circ - \bar{\mathbf{u}}^h) \right. \\
549 & \left. + \beta \sum_{i=1}^{2M} (\mathbf{u}_i^\circ - \bar{\mathbf{u}}_i) + (\mathbf{u}_i^\circ - \bar{\mathbf{u}}_i^h) \right] + (1-t) \left[\bar{\mathbf{u}}^{h\top} (\mathcal{Y}^h - \mathcal{Y})(\bar{\mathbf{u}} - \bar{\mathbf{u}}^h) + (\mathbf{b} - \mathbf{b}^h)(\bar{\mathbf{u}}^h - \bar{\mathbf{u}}) \right] \\
550 & = t \left[\nabla \mathbf{f}(\bar{\mathbf{u}})^\top (\mathbf{u}^\circ - \bar{\mathbf{u}}) - \nabla \mathbf{f}^h(\bar{\mathbf{u}}^h)^\top (\bar{\mathbf{u}}^h - \mathbf{u}^\circ) \right] \\
551 & + (1-t) \left[\bar{\mathbf{u}}^{h\top} (\mathcal{Y}^h - \mathcal{Y})(\bar{\mathbf{u}} - \bar{\mathbf{u}}^h) + (\mathbf{b} - \mathbf{b}^h)(\bar{\mathbf{u}}^h - \bar{\mathbf{u}}) \right].
\end{aligned}$$

553 Now, by applying Lemma 5.1, we estimate the right-hand side as follows

$$\begin{aligned}
554 & (1-t)(\bar{\mathbf{u}}^h - \bar{\mathbf{u}})^\top (\mathcal{Y} + \alpha I)(\bar{\mathbf{u}}^h - \bar{\mathbf{u}}) \\
555 & \leq t \left[\nabla \mathbf{f}(\bar{\mathbf{u}})^\top (\mathbf{u}^\circ - \bar{\mathbf{u}}) - \nabla \mathbf{f}^h(\bar{\mathbf{u}}^h)^\top (\bar{\mathbf{u}}^h - \mathbf{u}^\circ) \right] + (1-t) \left[\bar{\mathbf{u}}^{h\top} (\mathcal{Y}^h - \mathcal{Y})(\bar{\mathbf{u}} - \bar{\mathbf{u}}^h) \right. \\
556 & \left. + (\mathbf{b} - \mathbf{b}^h)(\bar{\mathbf{u}}^h - \bar{\mathbf{u}}) \right] \\
557 & \leq tC(\|2\mathbf{u}^\circ - \bar{\mathbf{u}} - \bar{\mathbf{u}}^h\| + h^2|\log(h)|) + (1-t)(\|\bar{\mathbf{u}}^h\| \|\mathcal{Y}^h - \mathcal{Y}\| + \|\mathbf{b} - \mathbf{b}^h\|) \|\bar{\mathbf{u}}^h - \bar{\mathbf{u}}\|.
\end{aligned}$$

559 From Lemma 5.2, one gets

$$\begin{aligned}
560 & (1-t)(\bar{\mathbf{u}}^h - \bar{\mathbf{u}})^\top (\mathcal{Y} + \alpha I)(\bar{\mathbf{u}}^h - \bar{\mathbf{u}}) \\
561 & \leq t(\tilde{C}_{ad} + h^2|\log(h)|) + (\|\bar{\mathbf{u}}^h\| \|\mathcal{Y}^h - \mathcal{Y}\| + \|\mathbf{b} - \mathbf{b}^h\|) \|\bar{\mathbf{u}}^h - \bar{\mathbf{u}}\| \\
562 & \leq t(\tilde{C}_{ad} + h^2|\log(h)|) + \tilde{C}_{ad}h^2,
\end{aligned}$$

564 where $\tilde{C}_{ad} > 0$ is a generic constant independent of h . Therefore, thanks to Lemma
565 5.3, the estimate (5.7) and the positive definiteness of the matrix $\mathcal{Y} + \alpha I$ we can
566 conclude the existence of a positive constant C , such that

$$567 \quad (5.9) \quad \|\bar{\mathbf{u}}^h - \bar{\mathbf{u}}\| \leq C\sqrt{h^2|\log(h)|} \rightarrow 0, \quad \text{whenever } h \rightarrow 0.$$

568 Also, in view of (5.9), we notice that $\bar{\mathbf{y}}^h \rightarrow \bar{\mathbf{y}}$ as $h \rightarrow 0$.

569 **Improved estimate.** From our discussion in Remark 5.4 we know that opti-
570 mality conditions can be represented as perturbed generalized equation depending
571 on the perturbation $\delta(h)$ given in (5.6). We do not apply directly the results in [25]
572 since there, the LICQ constraint qualifications and strong regularity were required.
573 However, we follow their general ideas.

574 Using the fact that the constraints are regular at $\bar{\mathbf{u}}$ (in the sense given in Lemma
575 4.3) and the properties of the objective and constraint functions we are able to satisfy
576 the stability result given in [30][Theorem 4.2]

577 Let us begin by considering $\delta(h)$ defined by (5.6) for $h \neq 0$. It is clear that Lemma
578 5.1 implies that the second component of $\delta(h)$ obeys $\|\mathbf{g}^h(\bar{\mathbf{u}}^h) - \mathbf{g}(\bar{\mathbf{u}}^h)\| \leq Ch^2|\log(h)|$.
579 On the other hand, a direct computation of the first components gives

$$\begin{aligned}
580 & \|\nabla_{\mathbf{u}} \mathcal{L}^h(\bar{\mathbf{u}}^h, \bar{\mathbf{v}}^h, \bar{\boldsymbol{\mu}}_a^h, \bar{\boldsymbol{\mu}}_b^h, \bar{\boldsymbol{\eta}}^h) - \nabla_{\mathbf{u}} \mathcal{L}(\bar{\mathbf{u}}^h, \bar{\mathbf{v}}^h, \bar{\boldsymbol{\mu}}_a^h, \bar{\boldsymbol{\mu}}_b^h, \bar{\boldsymbol{\eta}}^h)\| \\
581 & \leq \|\nabla \mathbf{f}^h(\bar{\mathbf{u}}^h) - \nabla \mathbf{f}(\bar{\mathbf{u}}^h)\| + \left\| \sum_{j=1}^{\ell} \bar{\mathbf{v}}_j^h y_i^h(x_j) - \sum_{j=1}^{\ell} \bar{\mathbf{v}}_j^h y_i(x_j) \right\| \leq Ch^2|\log(h)|,
\end{aligned}$$

582

583 where the last bound follow from Lemma 5.1 and estimate (3.2c). In particular, this
 584 implies that the perturbation $\delta(h)$ is of order $h^2|\log(h)|$.

585 As established in Lemma (4.3), our constraints, considered in the polihedral cone
 586 $\mathbb{R}_+^{\ell+4M}$, are regular. Moreover, the second order sufficient conditions are satisfied in
 587 view of the positive definiteness of the matrix $\mathcal{L}''(\bar{\mathbf{u}}, \bar{\boldsymbol{\nu}}, \bar{\boldsymbol{\mu}}_a, \bar{\boldsymbol{\mu}}_b, \bar{\boldsymbol{\eta}}) = \mathbf{f}''(\bar{\mathbf{u}}) = \mathcal{Y} + \alpha I$.
 588 Hence, we fulfill the hypothesis of Theorem 4.2 in [30] for the function $\mathcal{F} : \mathbb{R}^{2M} \times$
 589 $\mathbb{R}^{\ell+4M} \times (\mathcal{U}_{feas} \times \mathbb{R}^\ell \times \mathbb{R}^M \times \mathbb{R}^M \times \mathbb{R}^{2M}) \rightarrow \mathbb{R}^{2M} \times \mathbb{R}^\ell \times \mathbb{R}^M \times \mathbb{R}^M \times \mathbb{R}^{2M}$ given by

$$590 \quad (5.10) \quad \mathcal{F}(\delta, (\mathbf{u}, \boldsymbol{\nu}, \boldsymbol{\mu}_a, \boldsymbol{\mu}_b, \boldsymbol{\eta})) = \begin{bmatrix} \nabla_{\mathbf{u}} \mathcal{L}(\mathbf{u}, \boldsymbol{\nu}, \boldsymbol{\mu}_a, \boldsymbol{\mu}_b, \boldsymbol{\eta}) \\ -\mathbf{g}(\mathbf{u}) \end{bmatrix} - \delta.$$

591 Using the following definitions from [30]:

$$592 \quad \boldsymbol{\Lambda}(\mathbf{u}, \delta) := \{(\boldsymbol{\nu}, \boldsymbol{\mu}_a, \boldsymbol{\mu}_b, \boldsymbol{\eta}) : (\mathbf{u}, (\boldsymbol{\nu}, \boldsymbol{\mu}_a, \boldsymbol{\mu}_b, \boldsymbol{\eta}), \delta) \text{ satisfying (5.5)}\},$$

$$593 \quad \boldsymbol{\Lambda}_0 := \boldsymbol{\Lambda}(\bar{\mathbf{u}}, 0),$$

$$594 \quad \mathcal{SP}(\delta) := \{\mathbf{u} \in U : (\mathbf{u}, (\boldsymbol{\nu}, \boldsymbol{\mu}_a, \boldsymbol{\mu}_b, \boldsymbol{\eta}), \delta) \text{ satisfying (5.5), for some } (\boldsymbol{\nu}, \boldsymbol{\mu}_a, \boldsymbol{\mu}_b, \boldsymbol{\eta})\},$$

596 Here U is a neighborhood of $\bar{\mathbf{u}}$. We have as a consequence of [30][Theorem 4.2] that
 597 there exists a positive constant c , such that the inequality

$$598 \quad \text{dist}[(\mathbf{u}, (\boldsymbol{\nu}, \boldsymbol{\mu}_a, \boldsymbol{\mu}_b, \boldsymbol{\eta})), \bar{\mathbf{u}} \times \boldsymbol{\Lambda}_0] \leq c \|\mathcal{F}(\delta, (\mathbf{u}, \boldsymbol{\nu}, \boldsymbol{\mu}_a, \boldsymbol{\mu}_b, \boldsymbol{\eta})) - \mathcal{F}(0, (\mathbf{u}, \boldsymbol{\nu}, \boldsymbol{\mu}_a, \boldsymbol{\mu}_b, \boldsymbol{\eta}))\| \\ 599 \quad (5.11) \quad \quad \quad = c \|\delta\|$$

601 is satisfied for δ sufficiently close to 0 and for all $(\mathbf{u}, (\boldsymbol{\nu}, \boldsymbol{\mu}_a, \boldsymbol{\mu}_b, \boldsymbol{\eta})) \in \mathcal{SP}(\delta) \times \boldsymbol{\Lambda}(\mathbf{u}, \delta)$.

Since $\delta(h) \rightarrow 0$ and $\bar{\mathbf{u}}^h \in U$ for $h \rightarrow 0$, and $(\bar{\mathbf{u}}^h, (\boldsymbol{\nu}^h, \boldsymbol{\mu}_a^h, \boldsymbol{\mu}_b^h, \boldsymbol{\eta}^h), \delta(h))$ satisfies (5.5), in particular (5.11) implies that

$$\|\bar{\mathbf{u}}^h - \bar{\mathbf{u}}\| \leq ch^2|\log(h)|.$$

602 Finally, by Lemma 4.2, we know that $\bar{u} = [\bar{\mathbf{u}}_i]_{i=1}^M - [\bar{\mathbf{u}}_{i+M}]_{i=1}^M$ and $\bar{u}^h = [\bar{\mathbf{u}}_i^h]_{i=1}^M -$
 603 $[\bar{\mathbf{u}}_{i+M}^h]_{i=1}^M$. Hence, the estimate

$$\begin{aligned} \|\bar{u} - \bar{u}^h\| &= \left\| \left([\bar{\mathbf{u}}_i]_{i=1}^M - [\bar{\mathbf{u}}_{i+M}]_{i=1}^M \right) - \left([\bar{\mathbf{u}}_i^h]_{i=1}^M - [\bar{\mathbf{u}}_{i+M}^h]_{i=1}^M \right) \right\| \\ 604 \quad &\leq C \|\bar{\mathbf{u}} - \bar{\mathbf{u}}^h\| \\ &\leq Ch^2|\log h|, \end{aligned}$$

605 holds for $h > 0$ small enough.

606 **6. Comments on the case $\alpha = 0$.** We have discussed so far the case when the
 607 Tikhonov regularization term is present in the cost function. Here, we briefly discuss
 608 the case $\alpha = 0$. We will denote this problem by (\mathbf{P}_M^0) . Several results are analogous to
 609 the previous section, therefore, their proofs will not be given. If $\alpha = 0$, the lack of strict
 610 convexity of the cost function prevents us to fulfill second order sufficient conditions
 611 which are needed in our analysis. Therefore, we make the following assumption.

612 *Hypothesis 6.1.* We assume that the matrix Y is positive definite.

613 *Remark 6.2.* Note that this assumption can be satisfied for instance if the set of
 614 states $\{y_i \in H_0^1(\Omega) : i = 1, \dots, M\}$ is orthonormal. For example, we may choose
 615 finitely many normalized eigenfunctions of the Laplace operator (i.e. $A = -\Delta$).

616 We also notice that this assumption implies the uniqueness of the solution for
 617 problem (\mathbf{P}_M^0) .

618 First order necessary optimality conditions for (\mathbf{P}_M^0) are obtained analogously.

619 THEOREM 6.3. *Under assumption 2.7, there exists $(\nu, \mu_a, \mu_b) \in \mathbb{R}_+^\ell \times \mathbb{R}_+^M \times \mathbb{R}_+^M \cup$
620 $\{0\}$ such that for all $i = 1, \dots, M$, it follows that*

$$621 \quad (6.1a) \quad \int_{\Omega} (\bar{y} - y_d) y_i \, dx + \sum_{j=1}^{\ell} \nu_j y_i(x_j) + \mu_{b_i} - \mu_{a_i} = -\beta, \quad \text{if } \bar{u}_i > 0,$$

$$622 \quad (6.1b) \quad \int_{\Omega} (\bar{y} - y_d) y_i \, dx + \sum_{j=1}^{\ell} \nu_j y_i(x_j) + \mu_{b_i} - \mu_{a_i} = \beta, \quad \text{if } \bar{u}_i < 0,$$

$$623 \quad (6.1c) \quad \left| \int_{\Omega} (\bar{y} - y_d) y_i \, dx + \sum_{j=1}^{\ell} \nu_j y_i(x_j) + \mu_{b_i} - \mu_{a_i} \right| \leq \beta, \quad \text{if } \bar{u}_i = 0,$$

624

625 moreover, the following conditions are satisfied:

$$626 \quad (6.2a) \quad \nu \geq 0, \quad \text{and} \quad \nu_j (\bar{y}(x_j) - b_j) = 0,$$

$$627 \quad (6.2b) \quad \mu_a \geq 0, \quad \text{and} \quad \mu_{a_i} (u_{a,i} - \bar{u}_i) = 0,$$

$$628 \quad (6.2c) \quad \mu_b \geq 0, \quad \text{and} \quad \mu_{b_i} (\bar{u}_i - u_{b,i}) = 0.$$

630 Again, by using (4.2) we reformulate problem for the case $\alpha = 0$ in \mathbb{R}^{2M} as follows:

$$631 \quad (\mathbf{P}_{2M}^0) \quad \left\{ \begin{array}{l} \min_{(y, \mathbf{u})} \mathbf{J}(y, \mathbf{u}) = \frac{1}{2} \int_{\Omega} (y - y_d)^2 \, dx + \beta \sum_{i=1}^{2M} \mathbf{u}_i \\ \text{subject to:} \\ Ay(x) = \sum_{i=1}^M (\mathbf{u}_i - \mathbf{u}_{i+M}) e_i(x), \quad \text{in } \Omega, \\ y(x) = 0, \quad \text{on } \Gamma, \\ y(x_j) \leq b_j, \quad \forall j = 1, \dots, \ell, \\ [\mathbf{u}_i - \mathbf{u}_{i+M}]_{i=1}^M \in \mathcal{U}_{ad}, \\ \mathbf{u}_i \geq 0, \quad \forall i = 1, \dots, 2M. \end{array} \right.$$

632 In turn, we deduce the corresponding first order optimality conditions in the same
633 fashion as in Theorem 4.4.

634 THEOREM 6.4. *Under Assumption 2.7, there exists $(\bar{\nu}, \bar{\mu}_a, \bar{\mu}_b, \bar{\eta}) \in \mathbb{R}_+^\ell \times \mathbb{R}_+^M \times$
635 $\mathbb{R}_+^M \times \mathbb{R}_+^{2M} \cup \{0\}$ such that the conditions below are satisfied*

$$636 \quad (6.3a) \quad (\bar{y} - y_d, y_i) + \beta + \sum_{j=1}^{\ell} \nu_j y_i(x_j) + \mu_{b_i} - \mu_{a_i} - \eta_i = 0,$$

$$637 \quad (6.3b) \quad -(\bar{y} - y_d, y_i) + \beta - \sum_{j=1}^{\ell} \nu_j y_i(x_j) - \mu_{b_i} + \mu_{a_i} - \eta_{i+M} = 0.$$

638

639 for $i = 1, \dots, M$, holds, together with the following slackness conditions

$$\begin{aligned}
 640 \quad (6.4a) \quad & \bar{\mathbf{u}} \geq 0, \\
 641 \quad (6.4b) \quad & u_{a_i} \leq \bar{\mathbf{u}}_i - \bar{\mathbf{u}}_{i+M} \leq u_{b_i}, \quad \forall i = 1, \dots, 2M. \\
 642 \quad (6.4c) \quad & \boldsymbol{\nu}_j (\bar{\mathbf{y}}(x_j) - b_j) = 0, \quad \forall j = 1, \dots, \ell, \\
 643 \quad (6.4d) \quad & \boldsymbol{\mu}_{a_i} (u_{a,i} - \bar{\mathbf{u}}_i + \bar{\mathbf{u}}_{i+M}) = 0, \quad \forall i = 1, \dots, M, \\
 644 \quad (6.4e) \quad & \boldsymbol{\mu}_{b_i} (\bar{\mathbf{u}}_i - \bar{\mathbf{u}}_{i+M} - u_{b,i}) = 0, \quad \forall i = 1, \dots, M, \\
 645 \quad (6.4f) \quad & \boldsymbol{\eta}_i \bar{\mathbf{u}}_i = 0, \quad \forall i = 1, \dots, 2M.
 \end{aligned}$$

647 Following the same procedure of discretization using the finite element method we
 648 formulate the associated discretized problem, using the superscript h to indicate the
 649 corresponding discretized variables. Similar to Theorem 4.4, we conclude that there
 650 exists $(\bar{\boldsymbol{\nu}}^h, \bar{\boldsymbol{\mu}}_a^h, \bar{\boldsymbol{\mu}}_b^h, \bar{\boldsymbol{\eta}}^h) \in \mathbb{R}_+^\ell \times \mathbb{R}_+^M \times \mathbb{R}_+^M \times \mathbb{R}_+^{2M} \cup \{0\}$, such that

$$\begin{aligned}
 651 \quad & (\bar{\mathbf{y}}^h - y_d, y_i^h) + \beta + \sum_{j=1}^{\ell} \boldsymbol{\nu}_j^h y_i^h(x_j) + \boldsymbol{\mu}_{b_i}^h - \boldsymbol{\mu}_{a_i}^h - \boldsymbol{\eta}_i^h = 0, \\
 652 \quad & -(\bar{\mathbf{y}}^h - y_d, y_i^h) + \beta - \sum_{j=1}^{\ell} \boldsymbol{\nu}_j^h y_i^h(x_j) - \boldsymbol{\mu}_{b_i}^h + \boldsymbol{\mu}_{a_i}^h - \boldsymbol{\eta}_{i+M}^h = 0. \\
 653
 \end{aligned}$$

654 for $i = 1, \dots, M$, together with the following conditions

$$\begin{aligned}
 655 \quad & \bar{\mathbf{u}}^h \geq 0, \\
 656 \quad & u_{a_i} \leq \bar{\mathbf{u}}_i^h - \bar{\mathbf{u}}_{i+M}^h \leq u_{b_i}, \quad \forall i = 1, \dots, 2M. \\
 657 \quad & \boldsymbol{\nu}_j^h (\bar{\mathbf{y}}^h(x_j) - b_j) = 0, \quad \forall j = 1, \dots, \ell, \\
 658 \quad & \boldsymbol{\mu}_{a_i}^h (u_{a,i} - \bar{\mathbf{u}}_i^h + \bar{\mathbf{u}}_{i+M}^h) = 0, \quad \forall i = 1, \dots, M, \\
 659 \quad & \boldsymbol{\mu}_{b_i}^h (\bar{\mathbf{u}}_i^h - \bar{\mathbf{u}}_{i+M}^h - u_{b,i}) = 0, \quad \forall i = 1, \dots, M, \\
 660 \quad & \boldsymbol{\eta}_i^h \bar{\mathbf{u}}_i^h = 0, \quad \forall i = 1, \dots, 2M.
 \end{aligned}$$

662 As mentioned earlier, second order sufficient conditions must be satisfied at the
 663 optimum. These conditions are given in [29][Definition 2.1] which, after the compu-
 664 tation of the tangent cone at $\mathbf{g}(\bar{\mathbf{u}})$, are formulated as follows:

665 Let $\bar{\mathbf{u}} \in \mathbb{R}^{2M}$ a feasible point for (\mathbf{P}_{2M}^0) and $(\bar{\boldsymbol{\nu}}, \bar{\boldsymbol{\mu}}_a, \bar{\boldsymbol{\mu}}_b, \bar{\boldsymbol{\eta}}) \in \mathbb{R}_+^\ell \times \mathbb{R}_+^M \times \mathbb{R}_+^M \times$
 666 $\mathbb{R}_+^{2M} \cup \{0\}$ satisfying the first order necessary conditions (6.1) and (6.2). The *second-*
 667 *order sufficient conditions* are satisfied at $\bar{\mathbf{u}}$ with multipliers $(\bar{\boldsymbol{\nu}}, \bar{\boldsymbol{\mu}}_a, \bar{\boldsymbol{\mu}}_b, \bar{\boldsymbol{\eta}})$ if for each
 668 $\mathbf{v} \in \mathbb{R}^{2M}$ such that

$$669 \quad (6.7) \quad \nabla \mathbf{g}_j(\bar{\mathbf{u}})^\top \mathbf{v} \leq 0 \quad \text{for } j \in \{1, \dots, \ell + 4M : \mathbf{g}_j(\bar{\mathbf{u}}) = 0\} \quad \text{and} \quad \nabla \mathbf{f}(\bar{\mathbf{u}})^\top \mathbf{v} = 0,$$

670 the quadratic growth condition:

$$671 \quad (6.8) \quad \mathbf{v}^\top \mathcal{L}''(\bar{\mathbf{u}}, \bar{\boldsymbol{\nu}}, \bar{\boldsymbol{\mu}}_a, \bar{\boldsymbol{\mu}}_b, \bar{\boldsymbol{\eta}}) \mathbf{v} \geq c \|\mathbf{v}\|^2,$$

672 holds, for some constant $c > 0$.

673 LEMMA 6.5. *Under Hypothesis 6.1 the unique solution $\bar{\mathbf{u}}$ for problem (\mathbf{P}_M^0) sat-*
 674 *isfies the second-order sufficient conditions given by (6.7) and (6.8).*

675 *Proof.* Let $\mathbf{v} \in \mathbb{R}^{2M}$ be a vector satisfying condition (6.7). By using the first-
676 order optimality conditions (6.3), we have that

$$677 \quad \nabla \mathbf{f}(\bar{\mathbf{u}})^\top \mathbf{v} + \sum_{j=1}^{\ell} \bar{\nu}_j y_{\mathbf{v}}(x_j) + \sum_{i=1}^M \bar{\mu}_{\mathbf{b}_i}(\mathbf{v}_i - \mathbf{v}_{i+M}) - \sum_{i=1}^M \bar{\mu}_{\mathbf{a}_i}(\mathbf{v}_i - \mathbf{v}_{i+M}) - \sum_{i=1}^{2M} \bar{\eta}_i \mathbf{v}_i = 0,$$

678 where $y_{\mathbf{v}} = \sum_{i=1}^M (\mathbf{v}_i - \mathbf{v}_{i+M}) y_i$. Moreover, by taking into account (6.7), we replace
679 the identity $\nabla \mathbf{f}(\bar{\mathbf{u}})^\top \mathbf{v} = 0$ in the previous equation, which leads to

$$680 \quad (6.9) \quad \sum_{j=1}^{\ell} \bar{\nu}_j y_{\mathbf{v}}(x_j) + \sum_{i=1}^M \bar{\mu}_{\mathbf{b}_i}(\mathbf{v}_i - \mathbf{v}_{i+M}) - \sum_{i=1}^M \bar{\mu}_{\mathbf{a}_i}(\mathbf{v}_i - \mathbf{v}_{i+M}) - \sum_{i=1}^{2M} \bar{\eta}_i \mathbf{v}_i = 0.$$

681 From the complementarity slackness conditions (6.4c)-(6.4f) we know that if $\mathbf{g}_j(\bar{\mathbf{u}}) \neq 0$
682 for some $j \in \{1, \dots, \ell + 4M\}$, thus the corresponding multiplier vanishes. Hence, we
683 have that the sums on expression (6.9) are meaningful only on those j 's for which
684 $\mathbf{g}_j(\bar{\mathbf{u}}) = 0$. In view of this observation, we have that (6.9) becomes

$$685 \quad \sum_{\substack{j \in \{1, \dots, \ell\} \\ \mathbf{g}_j(\bar{\mathbf{u}}) = 0}} \bar{\nu}_j y_{\mathbf{v}}(x_j) + \sum_{\substack{i \in \{1, \dots, M\} \\ \mathbf{g}_{\ell+i}(\bar{\mathbf{u}}) = 0}} \bar{\mu}_{\mathbf{b}_i}(\mathbf{v}_i - \mathbf{v}_{i+M}) - \sum_{\substack{i \in \{1, \dots, M\} \\ \mathbf{g}_{\ell+M+i}(\bar{\mathbf{u}}) = 0}} \bar{\mu}_{\mathbf{a}_i}(\mathbf{v}_i - \mathbf{v}_{i+M}) \\ 686 \quad (6.10) \quad - \sum_{\substack{i \in \{1, \dots, 2M\} \\ \bar{\mathbf{u}}_i = 0}} \bar{\eta}_i \mathbf{v}_i = 0,$$

688 On the other hand, since the Lagrange multipliers are nonnegative and (6.7) holds,
689 we have

$$690 \quad \bar{\nu}_j y_{\mathbf{v}}(x_j) \leq 0, \text{ for all } j \in \{j \in \{1, \dots, \ell\} : \mathbf{g}_j(\bar{\mathbf{u}}) = 0\}, \\ 691 \quad \bar{\mu}_{\mathbf{b}_i}(\mathbf{v}_i - \mathbf{v}_{i+M}) \leq 0, \text{ for all } i \in \{i \in \{1, \dots, M\} : \mathbf{g}_{\ell+i}(\bar{\mathbf{u}}) = 0\}, \\ 692 \quad -\bar{\mu}_{\mathbf{b}_i}(\mathbf{v}_i - \mathbf{v}_{i+M}) \leq 0, \text{ for all } i \in \{i \in \{1, \dots, M\} : \mathbf{g}_{\ell+M+i}(\bar{\mathbf{u}}) = 0\}, \\ 693 \quad -\bar{\eta}_i \mathbf{v}_i \leq 0, \text{ for all } i \in \{i \in \{1, \dots, 2M\} : \bar{\mathbf{u}}_i = 0\}.$$

695 From these relations, we conclude that all terms in (6.10) vanish. In particular we get

$$696 \quad (6.11) \quad \bar{\eta}_i \mathbf{v}_i = 0, \text{ for all } i \in \{i \in \{1, \dots, 2M\} : \bar{\mathbf{u}}_i = 0\}.$$

697 The first order optimality conditions (6.1) and (6.2) imply that $0 < \beta = \frac{\bar{\eta}_i + \bar{\eta}_{i+M}}{2}$ for
698 all $i = 1, \dots, M$. We are going to show that $\mathbf{v}_i \mathbf{v}_{i+M} = 0$ for all $i = 1, \dots, M$. Indeed,
699 using the fact that $\bar{\mathbf{u}}_i \bar{\mathbf{u}}_{i+M} = 0$ we analyze the following cases:

- 700 (i) If $\bar{\mathbf{u}}_i = 0$ and $\bar{\mathbf{u}}_{i+M} > 0$, from the complementarity slackness conditions we have
701 $\bar{\eta}_{i+M} \bar{\mathbf{u}}_{i+M} = 0$, therefore $\bar{\eta}_{i+M} = 0$ which implies $\bar{\eta}_i > 0$, by (6.11) it follows
702 that $\mathbf{v}_i = 0$.
- 703 (ii) If $\bar{\mathbf{u}}_i > 0$ and $\bar{\mathbf{u}}_{i+M} = 0$ we get that $\mathbf{v}_{i+M} = 0$. The arguments are analogous
704 to the case (i).
- 705 (iii) If $\bar{\mathbf{u}}_i = 0$ and $\bar{\mathbf{u}}_{i+M} = 0$, again, from the complementarity slackness conditions
706 we have $\bar{\eta}_i \bar{\mathbf{u}}_i = 0$ and $\bar{\eta}_{i+M} \bar{\mathbf{u}}_{i+M} = 0$, and since $\bar{\eta}_i + \bar{\eta}_{i+M} > 0$ we get $\bar{\eta}_i > 0$
707 or $\bar{\eta}_{i+M} > 0$ implying that $\mathbf{v}_i = 0$ or $\mathbf{v}_{i+M} = 0$. Thus $\mathbf{v}_i \mathbf{v}_{i+M} = 0$.

708 We have shown that $\mathbf{v}_i \mathbf{v}_{i+M} = 0$ for all $i = 1, \dots, M$, therefore

$$\begin{aligned}
 709 \quad \mathbf{v}^\top \mathcal{L}''(\bar{\mathbf{u}}, \bar{\mathcal{V}}, \bar{\boldsymbol{\mu}}_a, \bar{\boldsymbol{\mu}}_b, \bar{\boldsymbol{\eta}}) \mathbf{v} &= \mathbf{v}^\top \mathcal{Y} \mathbf{v} \\
 710 &= ([\mathbf{v}_j - \mathbf{v}_{j+M}]_{j=1}^M)^\top Y[\mathbf{v}_i - \mathbf{v}_{i+M}]_{i=1}^M \\
 711 &= ([\mathbf{v}_j]_{j=1}^M)^\top Y[\mathbf{v}_i]_{i=1}^M + ([\mathbf{v}_{j+M}]_{j=1}^M)^\top Y[\mathbf{v}_{i+M}]_{i=1}^M \\
 712 &\quad - 2 \sum_{i=1}^M Y_{ii} \mathbf{v}_i \mathbf{v}_{i+M} \geq c \|\mathbf{v}\|^2, \\
 713
 \end{aligned}$$

714 for a constant $c > 0$. □

715 Finally, having satisfied the second order conditions in Lemma 6.5, we follow the
 716 same analysis as in Theorem 3.3, obtaining the following error estimate.

717 **THEOREM 6.6.** *Let \bar{u} be the optimal control of problems (\mathbf{P}_M^0) and \bar{u}_h the solution*
 718 *of the corresponding discretized problem. Then, the estimate*

$$719 \quad (6.12) \quad \|\bar{u} - \bar{u}^h\| \leq Ch^2 |\log h|,$$

720 holds for some constant $C > 0$, independent of mesh parameter h .

721 **7. Numerical experimentation.** In this section we focus on the numerical con-
 722 firmation of our theoretical estimates by means of examples implemented in Matlab.
 723 In our experiments, we observe that the order of the error corresponds approximately
 724 to our prediction $h^2 \approx h^2 |\log h|$. The $|\log(h)|$ -term is hard to detect in our compu-
 725 tations because the size of the mesh is small.

726 In order to compute the numerical solution in our examples, we transform our
 727 problem in a finite-dimensional problem, after the discretization of the state equation
 728 and the variables represented in the finite-element space. The cost function is also
 729 transformed to a function defined in finite-dimensional spaces using the mid-point
 730 rule to approximate the integrals. The contribution of the integration error is known to
 731 be of order h^2 . Therefore, it should not affect to the observation of the error estimates
 732 of the optimal variables. The associated optimization problems are nonsmooth since
 733 a term with 1-norm is present. In order to solve these problems numerically we apply
 734 the *second order orthant-wise method* developed in [14], which is able to compute
 735 solutions with exact null components where the solution is sparse.

736 The measures of convergence rate with respect to h were calculated by using
 737 decreasing mesh sizes at different values of h . As in [25] we use a very thin mesh
 738 in order to compute a high precision solution subsequently considered as the “exact”
 739 solution, with which we will compare the approximated solutions for every mesh. The
 740 experimental order of convergence is measure with the help of the following formula:

$$741 \quad (7.1) \quad EOC = \frac{\log(\|\bar{u}^{h_1} - \bar{u}_h^*\|) - \log(\|\bar{u}_h^* - \bar{u}^{h_2}\|)}{\log(h_1) - \log(h_2)},$$

742 for two consecutive mesh sizes h_1 y h_2 , and u_h^* the approximate the reference solution
 743 of the problem.

744 **Example 1..** In our first example we consider \mathbb{R}^8 as control space and we chose
 745 5 different points of the domain where the state constraints are imposed.

$$\begin{aligned}
 & \min_{u \in \mathcal{U}_{ad} \subset \mathbb{R}^8} J(y, u) = \frac{1}{2} \int_{\Omega} (y - y_d)^2 dx + \frac{1}{80} \|u\|_2^2 + 5 \|u\|_1 \\
 & \text{subject to} \\
 \text{(E1)} \quad & \left\{ \begin{array}{l} -\Delta y(x) + y(x) = \sum_{i=1}^8 u_i e_i(x), \quad \text{in } \Omega = (0, 1)^2, \\ y(x) = 0, \quad \text{on } \Gamma, \\ y(x_j) \leq -10, \quad \forall j = 1, 2, 3, 4, 5; \end{array} \right.
 \end{aligned}$$

The state constraints are given in the following points of the domain: $x_1 = \begin{pmatrix} 0.08 \\ 0.4 \end{pmatrix}$, $x_2 = \begin{pmatrix} 0.4 \\ 0.4 \end{pmatrix}$, $x_3 = \begin{pmatrix} 0.84 \\ 0.12 \end{pmatrix}$, $x_4 = \begin{pmatrix} 0.12 \\ 0.44 \end{pmatrix}$, $x_5 = \begin{pmatrix} 0.2 \\ 0.24 \end{pmatrix}$ and the functions e_i are defined by

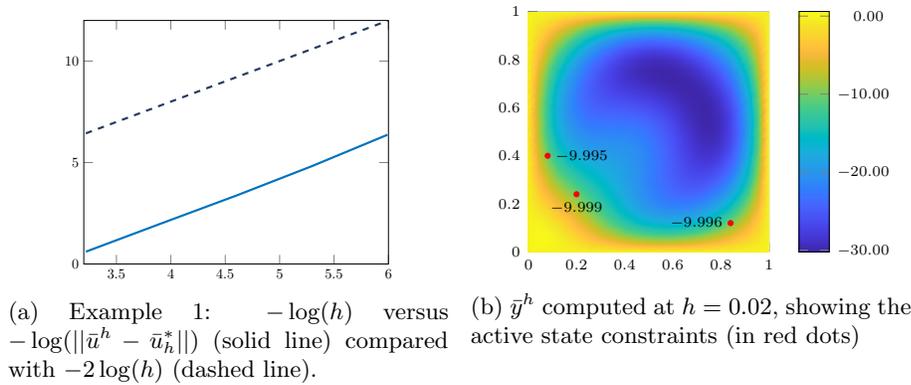
$$\begin{aligned}
 e_1(x) &= x_1 + x_2, & e_2(x) &= 8\pi^2 \sin(2\pi x_1) \sin(2\pi x_2), \\
 e_3(x) &= x_1 - x_2, & e_4(x) &= \cos^3(\pi x_1) \\
 e_5(x) &= 4\pi^2 \cos(2\pi(x_1 + x_2)), & e_6(x) &= 4\pi^2 \cos(2\pi(x_1 - x_2)), \\
 e_7(x) &= x_1^2 + x_2^2, & e_8(x) &= (x_1^2 - 1)(x_2^2 - 1)(x_1^2 + x_2^2).
 \end{aligned}$$

Additionally, the set of admissible controls is determined by the constants $u_a = -500$ and $u_b = 500$. The chosen desired state is $y_d(x) = -2 \sin(2\pi x_1) \sin(2\pi x_2) + 5$.

The exact solution of this problem is not known. Therefore, we compute a reference solution $\bar{u}_h^* = \bar{u}^h$ for $h = 0.00125$. The order of error computed with (7.1) gives the following results:

h	$\ \bar{u}^h - \bar{u}_h^*\ $	EOC
0.04	0.54515935	-
0.02	0.13564030	2.0069
0.01	0.03410037	1.9919
0.005	0.00802171	2.0878
0.0025	0.00171625	2.2247

Table 1: Experimental order of error for Example 1.



759 We notice that the error estimate is approximately h^2 . Figure 1a, shows similar
 760 slopes of the numerical error compared with h^2 in logarithmic scale.

761 The approximated optimal control for (E1) with $h = 0.04$, is given by

$$762 \quad \bar{u}^h = \begin{pmatrix} -455.9386 \\ 0 \\ 0 \\ 0 \\ 9.7017 \\ 13.1094 \\ -160.0708 \\ 0 \end{pmatrix};$$

763 were sparsity property is satisfied in several of its entries.

764 In the Figure 1b, we observe that the approximated optimal state \bar{y}^h for $h = 0.04$
 765 satisfies (approximately) the state constraints at the prescribed points of the domain.

$$766 \quad [\bar{y}^h(x_j)]_{j=1}^5 = \begin{pmatrix} -9.9948 \\ -19.7645 \\ -9.9964 \\ -14.5492 \\ -9.9994 \end{pmatrix},$$

767 i.e. in the points x_1, x_3 and x_5 the constraints are active.

768 **Example 2.** Here, we are interested in solving the optimal control problem with-
 769 out the Tikhonov regularization ($\alpha = 0$). We also consider five disjoint points on the
 770 domain where state constraints must be satisfied. Our problem reads

$$771 \quad (\mathbf{E2}) \quad \begin{cases} \min_{u \in \mathcal{U}_{ad} \subset \mathbb{R}^5} J(y, u) = \frac{1}{2} \|y - y_d\|^2 + 2\|u\|_1 \\ \text{subject to} \\ -\Delta y(x) = \sum_{i=1}^5 u_i e_i(x), \quad \text{in } \Omega = (0, 1)^2, \\ y(x) = 0, \quad \text{on } \Gamma, \\ y(x_j) \leq b_j, \quad \forall j = 1, 2, 3, 4, 5; \end{cases}$$

772 with $x_1 = \begin{pmatrix} 0.2 \\ 0.5 \end{pmatrix}$, $x_2 = \begin{pmatrix} 0.7 \\ 0.5 \end{pmatrix}$, $x_3 = \begin{pmatrix} 0.8 \\ 0.5 \end{pmatrix}$, $x_4 = \begin{pmatrix} 0.2 \\ 0.1 \end{pmatrix}$, $x_5 = \begin{pmatrix} 0.5 \\ 0.1 \end{pmatrix}$ and $b =$
 773 $(-8, -10, -12, -12, 9)^\top$ Functions e_i are defined by

$$774 \quad \begin{aligned} e_1(x) &= \sin(\pi(4x_1 + 5x_2)), & e_2(x) &= \sin(\pi(2x_1 + 3x_2)), \\ 775 \quad e_3(x) &= -\cos(\pi(-3x_1 + 17x_2)), & e_4(x) &= \cos(\pi(x_1 + x_2)), \\ 776 \quad e_5(x) &= -\sin(\pi(x_1 - x_2)). & & z \end{aligned}$$

778 On set of admissible controls is given by

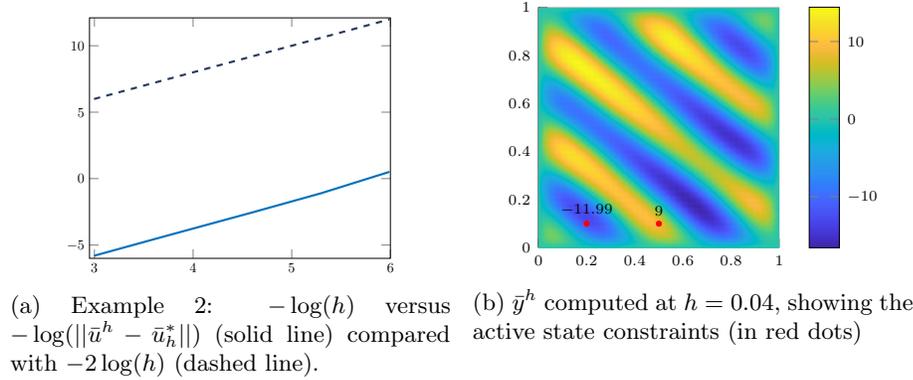
$$779 \quad \mathcal{U}_{ad} = \{u \in \mathbb{R}^8 : -5000 \leq u_i \leq 5000, \quad \forall i = 1, 2, \dots, 8\},$$

780 and the desired state is $y_d(x) = -2 \sin(2\pi x_1) \cos(2\pi x_2) + 5$.

781 We proceed as in the first example since exact solution of this problem is also not
 782 known. The reference solution $\bar{u}_h^* = \bar{u}^h$ is computed at $h = 0.00125$. The experimental
 783 order of convergence can be observed in Table 2.

h	$\ \bar{u}^h - \bar{u}_h^*\ $	EOC
0.05	332.1845	-
0.025	80.5232	2.0445
0.020	51.2134	2.028
0.010	12.6040	2.0226
0.005	3.0017	2.07
0.0025	0.6	2.3227

Table 2: Experimental order of error for Example 2.



(a) Example 2: $-\log(h)$ versus $-\log(\|\bar{u}^h - \bar{u}_h^*\|)$ (solid line) compared with $-2\log(h)$ (dashed line). (b) \bar{y}^h computed at $h = 0.04$, showing the active state constraints (in red dots)

784 Note that the error estimator is close to 2. Therefore, the order of the error is
 785 approximately h^2 . Figure 2a, shows the numerical error compared with the function
 786 h^2 in logarithmic scale. The optimal control calculated for $h = 0.02$, is given by:

$$787 \quad \bar{u}^h = \begin{pmatrix} 4657.8 \\ 0 \\ 0 \\ 29 \\ 150.3 \end{pmatrix}.$$

788 Sparsity is satisfied in components \bar{u}_2^h and \bar{u}_3^h . In Figure 2a the approximated optimal
 789 state \bar{y}^h for $h = 0.02$ is depicted. According to the definition of the bound b , the values
 790 of the state at the constrained points are given by

$$791 \quad [\bar{y}^h(x_j)]_{j=1}^5 = \begin{pmatrix} -8 \\ -12.7414 \\ -13.116 \\ -11.992 \\ 9.000 \end{pmatrix},$$

792 where activity at the points x_4 and x_5 is observed.

793

REFERENCES

794 [1] Bartle, R.G. *The elements of integration and Lebesgue Measure*. EU : Wiley Classics Library
 795 (1995).

- 796 [2] Bonnans, J., Casas, E. . Contrôle de systèmes elliptiques similineaires comportant des con-
797 traintes sur l'état. (French. English summary) [Control of semilinear elliptic systems with
798 state constraints] *Nonlinear partial differential equations and their applications. Collège*
799 *de France seminar, VIII* (1988), 69-86.
- 800 [3] Bonnans, J. F., Gilbert, J. C., Lemaréchal C., Sagastizábal C. A. *Numerical Optimization.*
801 Paris, Grenoble, Rio de Janeiro, Springer (2006).
- 802 [4] Brezis, H. *Functional analysis, sobolev spaces and partial differential equations.* EU: Springer
803 (2010).
- 804 [5] Casas, E., Clason, C., Kunisch, K. Approximation of elliptic control problems in measure spaces
805 with sparse solutions. *SIAM J. Control Optim.*, 50 (2012), 1735–1752.
- 806 [6] Casas, E., Herzog, R., Wachsmuth, G. Approximation of sparse controls in semilinear equations
807 by piecewise linear functions. *Numer. Math.* 122 (2012), no. 4, 645–669.
- 808 [7] Casas, E., Herzog, R., Wachsmuth, G. Optimality conditions and error analysis of semilinear
809 elliptic control problems with L^1 cost functional. *SIAM J. Optim.*, 22 (2012), 795–820.
- 810 [8] Casas, E., Kunisch, K. Optimal control of semilinear elliptic equations in measure spaces. *SIAM*
811 *J. Control Optim.*, 52 (2014), 339-364.
- 812 [9] Casas, E., Tröltzsch, F. Second-order and stability analysis for state-constrained elliptic optimal
813 control problems with sparse controls. *SIAM J. Control Optim.*, 52 (2014), 1010–1033.
- 814 [10] Ciarlet, P. *The finite element method for elliptic problems.* Amsterdam, NL: North-Holland
815 (1978).
- 816 [11] Clarke, F. *Functional Analysis, Calculus of Variations and Optimal Control.* FR: Springer
817 (2013).
- 818 [12] Clarke, F. H. *Optimization and Nonsmooth Analysis.* EU: John Wiley & Sons (1983).
- 819 [13] De Los Reyes, J. C. *Numerical PDE-Constrained Optimization.* EC: Springer (2015).
- 820 [14] De Los Reyes, J. C., Loayza, E., Merino, P. Second-order orthant-based methods with enriched
821 Hessian information for sparse ℓ_1 -optimization. *Comput. Optim. Appl.*, 67 (2017), 225-258.
- 822 [15] De Los Reyes, J. C., Merino, P., Rehberg, J., Tröltzsch, F. Optimality conditions for state-
823 constrained PDE control problems with time-dependent controls. *Control Cybernet*, 37
824 (2008), 5–38.
- 825 [16] Evans, L. *Partial Differential Equations.* EU: AMS (1998).
- 826 [17] Fu, H., Ng, M.K., Nikolova, M., Barlow, J.L. Efficient minimization methods of mixed $\ell_2 - \ell_1$
827 and $\ell_1 - \ell_1$ norms for image restoration. *SIAM J. Sci. Comput.*, 27 (2006), 1881–1902.
- 828 [18] Gilbarg, D & Trudinger N. *Elliptic Partial Differential Equations of Second Order.* DE:
829 Springer (1998).
- 830 [19] Grisvard, P. *Elliptic Problems in Nonsmooth Domains.* Philadelphia, EU: Pitman (1985).
- 831 [20] Henrion, R. On constraint qualifications. *Optimization Theory and Applications*, 72, (1992).
832 187-197.
- 833 [21] Kunisch, K., Wachsmuth, D. On time optimal control of the wave equation and its numerical
834 realization as parametric optimization problem. *SIAM J. Control Optim.*, 51 (2013), 1232-
835 1262.
- 836 [22] Kok, H., Wust, P., Stauffer, P., Bardati, F., van Rhooon, G., & Crezee, J. Current state of the
837 art of regional hyperthermia treatment planning: a review. *Radiation Oncology*, 10, 196.
838 (2015). <http://doi.org/10.1186/s13014-015-0503-8>
- 839 [23] Merino, P., Neitzel I., Tröltzsch F. An adaptive numerical method for semi-infinite elliptic control
840 problems based on error estimates. *Optimization Methods and Software*, 30 (2015), x492–
841 515.
- 842 [24] Merino, P., Neitzel I., Tröltzsch F. On linear-quadratic elliptic control problems of semi-infinite
843 type. *Applicable Analysis*, 90 (2011), 1047–1074.
- 844 [25] Merino, P., Tröltzsch, F., Vexler, B. Error Estimates for the Finite Element Approximation
845 of a Semilinear Elliptic Control Problem with State Constraints and Finite Dimensional
846 Control Space. *Mathematical Modelling and Numerical Analysis*, 44 (2010), 167–188.
- 847 [26] Rannacher, R., Vexler, B. A priori error estimates for the finite element discretization of elliptic
848 parameter identification problems with pointwise measurements. *SIAM J. Control Optim.*,
849 44 (2005), 1844 –1863.
- 850 [27] Robinson, S. M. Stability Theory for Systems of Inequalities, part I: Linear Systems*. *SIAM*
851 *J. Numer. Anal.*, 12 (1975), 754–769.
- 852 [28] Robinson, S. M. Stability Theory for Systems of Inequalities, part II: Differentiable Nonlinear
853 Systems*. *SIAM J. Numer. Anal.*, 13 (1976), 497–513.
- 854 [29] Robinson, S. M. Strongly Regular Generalized Equations. *Mathematics of Operations Research*,
855 5 (1980), 43–62.
- 856 [30] Robinson, S. M. Generalized Equations and Their Solutions, Part II Applications to Nonlinear
857 Programming. *Mathematical Programming Study*, 19 (1982), 200–221.

- 858 [31] Roubíček, T. Optimal control of a Stefan problem with state-space constraints. *Numerische*
859 *Mathematik*, 50 (1987), 723–744.
- 860 [32] Stadler, G. Elliptic optimal control problems with L^1 -control cost and applications for the
861 placement of control devices. *Comput. Optim. Appl.*, 44 (2009), 159–181.
- 862 [33] Tianhong, H. *Lasso and General ℓ_1 Regularized Regression under Linear Equality and Inequality*
863 *Constraints*. Purdue University, EU: PhD thesis (2011).
- 864 [34] Tröltzsch, F. *Optimal Control of Partial Differential Equations Theory, Methods and Applica-*
865 *tions*. EU: AMS (2010).