

# **Hamburger Beiträge** zur Angewandten Mathematik

## **Discretization of optimal control problems**

Michael Hinze and Arnd Rösch

Nr. 2010-07  
August 2010



# Discretization of optimal control problems

Michael Hinze and Arnd Rösch

**Abstract.** Solutions to optimization problems with pde constraints inherit special properties; the associated state solves the pde which in the optimization problem takes the role of a equality constraint, and this state together with the associated control solves an optimization problem, i.e. together with multipliers satisfies first and second order necessary optimality conditions. In this note we review the state of the art in designing discrete concepts for optimization problems with pde constraints with emphasis on structure conservation of solutions on the discrete level, and on error analysis for the discrete variables involved. As model problem for the state we consider an elliptic pde which is well understood from the analytical point of view. This allows to focus on structural aspects in discretization. We discuss the approaches *First discretize, then optimize* and *First optimize, then discretize*, and consider in detail two variants of the *First discretize, then optimize* approach, namely variational discretization, a discrete concept which avoids explicit discretization of the controls, and piecewise constant control approximations. We consider general constraints on the control, and also consider pointwise bounds on the state. We outline the basic ideas for providing optimal error analysis and accomplish our analytical findings with numerical examples which confirm our analytical results. Furthermore we present a brief review on recent literature which appeared in the field of discrete techniques for optimization problems with pde constraints.

**Mathematics Subject Classification (2000).** 49J20, 49K20, 35B37.

**Keywords.** elliptic optimal control problem, state & control constraints, error analysis.

## 1. Introduction

In PDE-constrained optimization, we have usually a pde as state equation and constraints on control and/or state. Let us write the pde for the state  $y \in Y$  with the control  $u \in U$  in the form  $e(y, u) = 0$  in  $Z$ . Assuming smoothness, we are then lead to optimization problems of the form

$$\min_{(y,u) \in Y \times U} J(y, u) \text{ s.t. } e(y, u) = 0, \quad c(y) \in \mathcal{K}, \quad u \in U_{ad}, \quad (1)$$

where  $e : Y \times U \rightarrow Z$  and  $c : Y \rightarrow \mathcal{R}$  are continuously Fréchet differentiable,  $\mathcal{K} \subset \mathcal{R}$  is a closed convex cone representing the state constraints, and  $U_{ad} \subset U$  is a closed convex set representing the control constraints.

Let us give two examples.

**Example 1.1.** Consider the distributed optimal control of a semilinear elliptic PDE:

$$\begin{aligned} \min J(y, u) &:= \frac{1}{2} \|y - y_d\|_{L^2(\Omega)}^2 + \frac{\alpha}{2} \|u\|_{L^2(\Omega)}^2 \\ \text{subject to} & \\ -\Delta y + y^3 &= \gamma u \quad \text{on } \Omega, \\ y &= 0 \quad \text{on } \partial\Omega, \\ a \leq u \leq b &\text{ on } \Omega, \text{ and } y \leq c \text{ on } D, \end{aligned} \tag{2}$$

where  $\gamma \in L^\infty(\Omega) \setminus \{0\}$ ,  $a, b \in L^\infty(\Omega)$ , and  $a \leq b$ . We require  $D \subset\subset \Omega$  to avoid restrictions on the bound  $c$  which would have to be imposed in the case  $D \equiv \Omega$  due to homogeneous boundary conditions required for  $y$ . Let  $n \leq 3$ . By the theory of monotone operators one can show that there exists a unique bounded solution operator of the state equation

$$u \in U_{ad} := \{v \in L^2(\Omega); a \leq v \leq b \text{ a.e.}\} \rightarrow y \in Y := H_0^1(\Omega).$$

Here  $c(y) = y$  with  $\mathcal{K} = \{y \in Y; y \leq c \text{ on } D \subset\subset \Omega\} \subset \mathcal{R} := C^0(\bar{\Omega})$ . Let  $A : H_0^1(\Omega) \rightarrow H_0^1(\Omega)^*$  be the operator associated with the bilinear form  $a(y, v) = \int_\Omega \nabla y \cdot \nabla v \, dx$  for the Laplace operator  $-\Delta y$  and let  $N : y \rightarrow y^3$ . Then the weak formulation of the state equation can be written in the form

$$e(y, u) := Ay + N(y) - \gamma u = 0.$$

**Example 1.2.** We consider optimal control of the time-dependent incompressible Navier-Stokes system for the velocity field  $y \in \mathbb{R}^d$  and the pressure  $p$ . Let  $\Omega \subset \mathbb{R}^d$  denote the flow domain, let  $f : [0, T] \times \Omega \rightarrow \mathbb{R}^d$  be the force per unit mass acting on the fluid and denote by  $y_0 : \Omega \rightarrow \mathbb{R}^d$  the initial velocity of the fluid at  $t = 0$ . Then the Navier Stokes equations can be written in the form

$$\begin{aligned} y_t - \nu \Delta y + (y \cdot \nabla) y + \nabla p &= f \quad \text{on } \Omega_T := (0, T) \times \Omega, \\ \nabla \cdot y &= 0 \quad \text{on } \Omega_T, \\ y(0, \cdot) &= y_0 \quad \text{on } \Omega, \end{aligned} \tag{3}$$

and have to be accomplished by appropriate boundary conditions. For the functional analytic setting let us define the Hilbert spaces

$$V := \text{cl}_{H_0^1(\Omega)^2} \{y \in C_c^\infty(\Omega)^2; \nabla \cdot y = 0\}, \quad H := \text{cl}_{L^2(\Omega)^2} \{y \in C_c^\infty(\Omega)^2; \nabla \cdot y = 0\},$$

and the associated parabolic solution space

$$Y := W(I) \equiv W(I; H, V) = \{y \in L^2(I; V); y_t \in L^2(I; V^*)\}.$$

We say that

$$\begin{aligned} y_t + (y \cdot \nabla) y - \nu \Delta y &= f \text{ in } L^2(I; V^*) =: Z \iff \\ \langle y_t, v \rangle_{V^*, V} + \nu \langle \nabla y, \nabla v \rangle_{L^2(\Omega)^{2 \times 2}} + \langle (y \cdot \nabla) y, v \rangle_{V^*, V} &= \langle f, v \rangle_{V^*, V} \quad \forall v \in V. \end{aligned}$$

Let  $U_{ad} \subset U$  be nonempty, convex and closed with  $U$  denoting a Hilbert space, and  $B : U \rightarrow L^2(I, V^*)$  a linear, bounded control operator. Furthermore we set  $\mathcal{K} := Y$ , i.e. we omit state constraints. Let finally  $J : Y \times U \rightarrow \mathbb{R}$  be a Fréchet differentiable functional. Then we may consider the following optimal control problem

$$\min_{u \in U_{ad}, y \in Y} J(y, u) \text{ s.t. } e(y, u) = 0 \text{ in } Z,$$

where the state equation  $e(y, u) = 0$  is the weak Navier-Stokes equation, i.e.,

$$e : Y \times U \rightarrow Z \times H, \quad e(y, u) = \begin{pmatrix} y_t + (y \cdot \nabla) y - \nu \Delta y - Bu \\ y(0, \cdot) - y_0 \end{pmatrix}.$$

We are interested in illuminating discrete approaches to problem (1), where we place particular emphasis on structure preservation on the discrete level, and also on analysing the contributions to the total error of the discretization errors in the variables and multipliers involved. To approach an optimal control problem of the form (1) numerically one may either discretize this problem by substituting all appearing function spaces by finite dimensional

spaces, and all appearing operators by suitable approximate counterparts which allow their numerical evaluation on a computer, say. Denoting by  $h$  the discretization parameter, one ends up with the problem

$$\min_{(y_h, u_h) \in Y_h \times U_h} J_h(y_h, u_h) \text{ s.t } e_h(y_h, u_h) = 0 \text{ and } c_h(y_h) \in \mathcal{K}_h, u_h \in U_{ad}^h, \quad (4)$$

where  $J_h : Y_h \times U_h \rightarrow \mathbb{R}$ ,  $e_h : Y_h \times U_h \rightarrow Z$ , and  $c_h : Y_h \rightarrow R$  with  $\mathcal{K}_h \subset R$ . For the finite dimensional subspaces one may require  $Y_h \subset Y, U_h \subset U$ , say, and  $\mathcal{K}_h \subseteq R$  a closed and convex cone,  $U_{ad}^h \subseteq U_h$  closed and convex. This approach in general is referred to as first discretize, then optimize. On the other hand one may switch to the Karush-Kuhn-Tucker system associated to (1)

$$e(\bar{y}, \bar{u}) = 0, \quad c(\bar{y}) \in \mathcal{K}, \quad (5)$$

$$\bar{\lambda} \in \mathcal{K}^\circ, \quad \langle \bar{\lambda}, c(\bar{y}) \rangle_{\mathcal{R}^*, \mathcal{R}} = 0, \quad (6)$$

$$L_y(\bar{y}, \bar{u}, \bar{p}) + c'(\bar{y})^* \bar{\lambda} = 0, \quad (7)$$

$$\bar{u} \in U_{ad}, \quad \langle L_u(\bar{y}, \bar{u}, \bar{p}), u - \bar{u} \rangle_{U^*, U} \geq 0 \quad \forall u \in U_{ad}. \quad (8)$$

and substitute all appearing function spaces and operators accordingly, where  $L(y, u, p) := J(y, u) - \langle p, e(y, u) \rangle_{Z^*, Z}$  denotes the Lagrangian associated to (1). This leads to solving

$$e_h(y_h, u_h) = 0, \quad c_h(y_h) \in \mathcal{K}_h, \quad (9)$$

$$\lambda_h \in \mathcal{K}_h^\circ, \quad \langle \lambda_h, c_h(y_h) \rangle_{\mathcal{R}^*, \mathcal{R}} = 0, \quad (10)$$

$$L_{h_y}(y_h, u_h, p_h) + c'_h(y_h)^* \lambda_h = 0, \quad (11)$$

$$\bar{u}_h \in U_{ad}^h, \quad \langle L_{h_u}(y_h, u_h, p_h), u - u_h \rangle_{U^*, U} \geq 0 \quad \forall u \in U_{ad}^h \quad (12)$$

for  $\bar{y}_h, \bar{u}_h, \bar{p}_h, \bar{\lambda}_h$ , where  $L_h$  denotes a discretized version of  $L$ . Of course,  $L \equiv L_h$  is possible. This approach in general is referred to as first optimize, then discretize, since it builds the discretization upon the first order necessary optimality conditions.

Instead of applying discrete concepts to problem (1) or (5)-(8) directly we may first apply an SQP approach on the continuous level and then apply first discretize, then optimize to the related linear quadratic constrained subproblems, or first optimize, then discretize to the SQP systems appearing in each iteration of the Newton method on the infinite dimensional level. This motivates us to illustrate all discrete concepts at hand of linear model pdes which are well understood w.r.t. analysis and discretization concepts and to focus the presentation on structural aspects inherent to optimal control problems with pde constraints. However, error analysis for optimization problems with nonlinear state equations in the presence of constraints on controls and/or state is not straightforward and requires special techniques such as extensions of Newton-Kontorovich-type theorems, and second order sufficient optimality conditions. This complex of questions also will be discussed.

The outline of this work is as follows. In Section 2, we consider an elliptic model optimal control problem containing many relevant features which need to be resolved by a numerical approach. In Section 3 we preview the case of nonlinear state equations and highlight the solution approaches taken so far, as well as the analytical difficulties one is faced with in this situation. Section 4 is devoted to the finite element method for the discretization of the state equation in our model problem. We propose two different approximation approaches of the *First discretize, then optimize*-type to the optimal control problem, including detailed numerical analysis. In Section 5 we present an introduction to relaxation approaches used in presence of state constraints. For Lavrentiev regularization applied to the model problem we present numerical analysis which allows to adapt the finite element discretization error to the regularization error. Finally, we present in Section 6 a brief review on recent literature which appeared in the field of discrete techniques for optimization problems with pde constraints.

## 2. A model problem

As model problem with pointwise bounds on the state we take the Neumann problem

$$(\mathbb{S}) \quad \left\{ \begin{array}{l} \min_{(y,u) \in Y \times U_{ad}} J(y, u) := \frac{1}{2} \int_{\Omega} |y - y_0|^2 + \frac{\alpha}{2} \|u\|_U^2 \\ \text{s.t.} \\ Ay = Bu \quad \text{in } \Omega, \\ \partial_{\eta} y = 0 \quad \text{on } \Gamma, \\ \text{and} \\ y \in Y_{ad} := \{y \in Y, y(x) \leq b(x) \text{ a.e. in } \Omega\}. \end{array} \right\} : \iff y = \mathcal{G}(Bu) \quad (13)$$

Here,  $Y := H^1(\Omega)$ ,  $A$  denotes an uniformly elliptic operator, for example  $Ay = -\Delta y + y$ , and  $\Omega \subset \mathbb{R}^d$  ( $d = 2, 3$ ) denotes an open, bounded sufficiently smooth (or polyhedral) domain. Furthermore, we suppose that  $\alpha > 0$  and that  $y_0 \in L^2(\Omega)$ , and  $b \in W^{2,\infty}(\Omega)$  are given.  $(U, (\cdot, \cdot)_U)$  denotes a Hilbert space and  $B : U \rightarrow L^2(\Omega) \subset H^1(\Omega)^*$  the linear, continuous control operator. By  $R : U^* \rightarrow U$  we denote the inverse of the Riesz isomorphism. Furthermore, we associate to  $A$  the continuous, coercive bilinear form  $a(\cdot, \cdot)$ .

**Example 2.1.** There are several examples for the choice of  $B$  and  $U$ .

- (i) Distributed control:  $U = L^2(\Omega)$ ,  $B = Id : L^2(\Omega) \rightarrow Y'$ .
- (ii) Boundary control:  $U = L^2(\partial\Omega)$ ,  $Bu(\cdot) = \int_{\partial\Omega} u \gamma_0(\cdot) dx \in Y'$ , where  $\gamma_0$  is the boundary trace operator defined on  $Y$ .
- (iii) Linear combinations of input fields:  $U = \mathbb{R}^n$ ,  $Bu = \sum_{i=1}^n u_i f_i$ ,  $f_i \in Y'$ .

If not stated otherwise we from here onwards consider the situation (i) of the previous example. In view of  $\alpha > 0$ , it is standard to prove that problem (13) admits a unique solution  $(y, u) \in Y_{ad} \times U_{ad}$ . In pde constrained optimization, the pde for given data frequently is uniquely solvable. In equation (13) this is also the case, so that for every control  $u \in U_{ad}$  we have a unique state  $y = \mathcal{G}(Bu) \in H^1(\Omega) \cap C^0(\bar{\Omega})$ . We need  $y \in C^0(\bar{\Omega})$  to satisfy the Slater condition required below. Problem (13) therefore is equivalent to the so called reduced optimization problem

$$\min_{v \in U_{ad}} \hat{J}(v) := J(\mathcal{G}(Bv), v) \text{ s.t. } \mathcal{G}(Bv) \in Y_{ad}. \quad (14)$$

The key to the proper numerical treatment of problems (13) and (14) can be found in the first order necessary optimality conditions associated to these control problems. To formulate them properly we require the following constraint qualification, often referred to as *Slater condition*. It requires the existence of a state in the interior of the set  $Y_{ad}$  considered as a subset of  $C^0(\bar{\Omega})$  and ensures the existence of a Lagrange multiplier in the associated dual space. Moreover, it is useful for deriving error estimates.

**Assumption 2.2.**  $\exists \bar{u} \in U_{ad} \quad \mathcal{G}(B\bar{u})(x) < b(x)$  for all  $x \in \bar{\Omega}$ .

Following Casas [9, Theorem 5.2] for the problem under consideration we now have the following theorem, which specifies the KKT system (5)-(8) for the setting of problem (13).

**Theorem 2.3.** *Let  $u \in U_{ad}$  denote the unique solution to (13). Then there exist  $\mu \in \mathcal{M}(\bar{\Omega})$  and  $p \in L^2(\Omega)$  such that with  $y = \mathcal{G}(Bu)$  there holds*

$$\int_{\Omega} pAv = \int_{\Omega} (y - y_0)v + \int_{\bar{\Omega}} v d\mu \quad \forall v \in H^2(\Omega) \text{ with } \partial_{\eta} v = 0 \text{ on } \partial\Omega, \quad (15)$$

$$(RB^*p + \alpha u, v - u)_U \geq 0 \quad \forall v \in U_{ad}, \quad (16)$$

$$\mu \geq 0, \quad y(x) \leq b(x) \text{ in } \Omega \text{ and } \int_{\bar{\Omega}} (b - y) d\mu = 0. \quad (17)$$

Here,  $\mathcal{M}(\bar{\Omega})$  denotes the space of Radon measures which is defined as the dual space of  $C^0(\bar{\Omega})$  and endowed with the norm

$$\|\mu\|_{\mathcal{M}(\bar{\Omega})} = \sup_{f \in C^0(\bar{\Omega}), |f| \leq 1} \int_{\bar{\Omega}} f d\mu.$$

Since  $\hat{J}'(v) = B^*p + \alpha(\cdot, u)_U$  a short calculation shows that the variational inequality (16) is equivalent to

$$u = P_{U_{ad}}(u - \sigma R \hat{J}'(u)) \quad (\sigma > 0),$$

where  $P_{U_{ad}}$  denotes the orthogonal projection in  $U$  onto  $U_{ad}$ . This nonsmooth operator equation constitutes a relation between the optimal control  $u$  and its associated adjoint state  $p$ . In the present situation, when we consider the special case  $U \equiv L^2(\Omega)$  with  $B$  denoting the injection from  $L^2(\Omega)$  into  $H^1(\Omega)^*$ , and without control constraints, i.e.  $U_{ad} \equiv L^2(\Omega)$ , this relation boils down to

$$\alpha u + p = 0 \text{ in } L^2(\Omega),$$

since  $\sigma > 0$ . This relation already gives a hint to the discretization of the state  $y$  and the control  $u$  in problem (13), if one wishes to conserve the structure of this algebraic relation also on the discrete level.

### 3. Nonlinear state equations

Practical applications are usually characterized by nonlinear partial differential equations, see Example 1.1 and Example 1.2. We will here only focus on optimization and discretization aspects. Let us assume that there is a solution operator  $S$ :

$$e(y, u) = 0 \iff y = S(u) \quad (18)$$

which maps  $U$  in  $Y$ . The nonlinearity of  $S$  results in a nonconvex optimization problem (1). Therefore, we have to replace the Slater condition (Assumption 2.2) by a Mangasarian-Fromovitz constraint qualification to get the necessary optimality condition. Of course, one needs differentiability properties of the solution operator  $S$ . Let us assume that the operator  $S$  is two times Fréchet differentiable. All these assumptions are satisfied for both examples.

Moreover, let us assume there is a discrete solution operator  $S_h : U \rightarrow Y_h$  with

$$e_h(y_h, u) = 0 \iff y_h = S_h(u). \quad (19)$$

Consequently, the discretization error of the PDE is described by  $\|S(u) - S_h(u)\|_Z$ . Usually, such a priori error estimates are known for a lot of discretizations and for different spaces  $Z$ . However, this is only a first small step in estimating the discretization error for the optimization problem.

Since the optimal control problem is not convex, we have to work with local minima and local convexity properties. We can only expect that a numerical optimization method generates a sequence of locally optimal (discrete) solutions converging to a local optimal solution of the undiscretized problem. To get error estimates one has to deal with local convexity properties. These local convexity properties are described by second-order sufficient optimality conditions. However, these properties lead only to a priori error estimates if one already knows that the discretized solution is sufficiently close to the undiscretized one. This complicated situation requires innovative techniques to obtain the desired results.

Until now such techniques are known only for control constrained optimal control problems. Let us first sketch a technique which was presented in [16]. Here, an auxiliary problem is introduced with the additional constraint

$$\|u - \bar{u}\|_{U_1} \leq r. \quad (20)$$

The choice of the space  $U_1$  is connected to the differentiability properties of the operator  $S$ . In general the  $U_1$ -norm is a stronger norm than the  $U$ -norm. For instance,  $U = L^2(\Omega)$  and  $U_1 = L^\infty(\Omega)$  is a typical choice. Let us mention that the two-norm discrepancy can be avoided for certain elliptic optimal control problems. The radius  $r$  is chosen in such a way that the auxiliary problem is now a strictly convex problem. Consequently, the solution  $\bar{u}$  is the unique solution of the auxiliary optimal control problem.

Using the second-order sufficient optimality condition one shows in a next step that

$$\|\bar{u}_h^r - \bar{u}\|_U + \|\bar{y}_h^r - \bar{y}\|_Y \leq ch^{\hat{\kappa}} \quad (21)$$

where  $\bar{u}_h^r$  is the solution of a discretized version of the auxiliary problem and  $y_h^r$  the corresponding state. Note, that the convergence order does not depend on  $r$ . A similar estimate is obtained for the adjoint state. A projection formula is used to derive an error estimate

$$\|\bar{u}_h^r - \bar{u}\|_{U_1} \leq ch^{\hat{\kappa}}.$$

If  $h$  is sufficiently small, then the additional inequality (20) cannot be active. Consequently  $\bar{u}_h^r$  is also a local minimizer of the discretized problem without this inequality and the error estimate (21) is valid. Let us mention the practical drawback of that result. Since the solution  $\bar{u}$  is unknown, we have no information about what  $h$  is *small enough* means.

A second approach was used in [68, 67]. Only information on the numerical solution are used in that approach. The main idea is to construct a ball around the numerical solution where the objective value of the undiscretized problem on the surface of that ball is greater than the objective value of the discretized control.

Let us mention that the available techniques cannot be applied to state constrained problems. Low regularity properties, instability of dual variables, and missing smoothing properties are some of the reasons that a priori error estimates for nonlinear state constrained problems are challenging.

Let us sketch a third approach which is based on the first order necessary optimality conditions (5)-(8) and (9)-(12), respectively, and which does not use second order sufficient optimality conditions. To begin with we consider

$$\min_{u \in U_{ad}} \hat{J}(u) \equiv J(S(u), u), \quad (22)$$

with  $J$  as in problem (13). In this situation (7) reduces to

$$\langle \alpha(u, \cdot) + B^*p, v - u \rangle_{U^*U} \geq 0 \text{ for all } v \in U_{ad}. \quad (23)$$

Then this variational inequality is equivalent to the semi-smooth operator equation

$$G(u) := u - P_{U_{ad}}\left(-\frac{1}{\alpha}B^*p\right) = 0 \text{ in } U, \quad (24)$$

where  $P_{U_{ad}}$  denotes the orthogonal projection onto  $U_{ad}$  in  $U$ , and where we assume that the Riesz isomorphism is the identity map. Analogously, for the variational discrete approach and its numerical solutions  $u_h \in U_{ad}$  (see next Section),

$$G_h(u_h) := u_h - P_{U_{ad}}\left(-\frac{1}{\alpha}B^*p_h\right) = 0 \text{ in } U. \quad (25)$$

We now pose the following two question; 1. *Given a solution  $u \in U_{ad}$  to (22), does there exist a solution  $u_h \in U_{ad}$  of (25) in a neighborhood of  $u$ ?* 2. *If yes, is this solution unique?* It is clear that solutions to (22) might not be local solutions to the optimization problem, and that every local solution is a solution to (22). In this respect the following exposition generalizes the classical Newton-Kantorovich concept.

To provide positive answers to these questions we have to pose appropriate assumptions on a solution  $u \in U_{ad}$  of (22).



**Definition 3.1.** A solution  $u \in U_{ad}$  of (22) is called regular, if  $M \in \partial G(u)$  exists with  $\|G'(v) - G'(u) - M(v - u)\|_U = o(1)$  for  $v \rightarrow u$ , and  $M$  is invertible with bounded inverse  $M^{-1}$ .

We note that in the case of box constraints with  $U := L^2(\Omega)$  and  $U_{ad} = \{v \in U; a \leq u \leq b\}$  this regularity requirement is satisfied if the gradient of the adjoint state associated to  $u$  admits a non-vanishing gradient on the boarder of the active set, see [28].

In the following we write  $G'(u) := M$ . Now let  $u \in U_{ad}$  denote a regular solution to (22) and consider the operator

$$\Phi(v) := v - G'(u)^{-1}G_h(v). \quad (26)$$

We now, under certain assumptions, show that  $\Phi$  has a fixed point  $u_h$  in a neighborhood of  $u$  which we then consider as discrete approximation to the solution  $u \in U_{ad}$  of (22). A positive answer to question 1 is given in

**Theorem 3.2.** *Let  $u \in U_{ad}$  denote a regular solution to (22). Furthermore let for  $v \in U_{ad}$  the error estimate*

$$\|G_h(v) - G(v)\| \leq ch^\kappa \text{ for } h \rightarrow 0$$

*be satisfied and let  $\Phi$  be compact. Then a neighborhood  $B_r(u) \subset U$  exists such that  $\Phi$  admits a least one fixed point  $u_h \in U_{ad} \cap B_r(u)$ . For  $u_h$  the error estimate*

$$\|u - u_h\|_U \leq Ch^\kappa \text{ for all } 0 < h \leq h_0$$

*holds.*

If we strengthen the regularity requirement on  $u$  by requiring strict differentiability of  $G$  at  $u$ , also uniqueness can be argued. Details are given in [28], where also error estimates for approximation schemes related to Example (1.2) are presented.

## 4. Finite element discretization

For the convenience of the reader we recall the finite element setting. To begin with let  $\mathcal{T}_h$  be a triangulation of  $\Omega$  with maximum mesh size  $h := \max_{T \in \mathcal{T}_h} \text{diam}(T)$  and vertices  $x_1, \dots, x_m$ . We suppose that  $\bar{\Omega}$  is the union of the elements of  $\mathcal{T}_h$  so that element edges lying on the boundary are curved. In addition, we assume that the triangulation is quasi-uniform in the sense that there exists a constant  $\kappa > 0$  (independent of  $h$ ) such that each  $T \in \mathcal{T}_h$  is contained in a ball of radius  $\kappa^{-1}h$  and contains a ball of radius  $\kappa h$ . Let us define the space of linear finite elements,

$$X_h := \{v_h \in C^0(\bar{\Omega}) \mid v_h \text{ is a linear polynomial on each } T \in \mathcal{T}_h\}$$

with the appropriate modification for boundary elements. In what follows it is convenient to introduce a discrete approximation of the operator  $\mathcal{G}$ . For a given function  $v \in L^2(\Omega)$  we denote by  $z_h = \mathcal{G}_h(v) \in X_h$  the solution of the discrete Neumann problem

$$a(z_h, v_h) = \int_{\Omega} v v_h \quad \text{for all } v_h \in X_h.$$

It is well-known that for all  $v \in L^2(\Omega)$

$$\|\mathcal{G}(v) - \mathcal{G}_h(v)\| \leq Ch^2 \|v\|, \quad (27)$$

$$\|\mathcal{G}(v) - \mathcal{G}_h(v)\|_{L^\infty} \leq Ch^{2-\frac{d}{2}} \|v\|. \quad (28)$$

The estimate (28) can be improved provided one strengthens the assumption on  $v$ .

**Lemma 4.1.**

(a) Suppose that  $v \in W^{1,s}(\Omega)$  for some  $1 < s < \frac{d}{d-1}$ . Then

$$\|\mathcal{G}(v) - \mathcal{G}_h(v)\|_{L^\infty} \leq Ch^{3-\frac{d}{s}} |\log h| \|v\|_{W^{1,s}}.$$

(b) Suppose that  $v \in L^\infty(\Omega)$ . Then

$$\|\mathcal{G}(v) - \mathcal{G}_h(v)\|_{L^\infty} \leq Ch^2 |\log h|^2 \|v\|_{L^\infty}.$$

*Proof.* (a): Let  $z = \mathcal{G}(v)$ ,  $z_h = \mathcal{G}_h(v)$ . Elliptic regularity theory implies that  $z \in W^{3,s}(\Omega)$  from which we infer that  $z \in W^{2,q}(\Omega)$  with  $q = \frac{ds}{d-s}$  using a well-known embedding theorem. Furthermore, we have

$$\|z\|_{W^{2,q}} \leq c\|z\|_{W^{3,s}} \leq c\|v\|_{W^{1,s}}. \quad (29)$$

Using Theorem 2.2 and the following estimate from [69] we have

$$\|z - z_h\|_{L^\infty} \leq c |\log h| \inf_{\chi \in X_h} \|z - \chi\|_{L^\infty}, \quad (30)$$

which, combined with a well-known interpolation estimate, yields

$$\|z - z_h\|_{L^\infty} \leq ch^{2-\frac{d}{q}} |\log h| \|z\|_{W^{2,q}} \leq ch^{3-\frac{d}{s}} |\log h| \|v\|_{W^{1,s}}$$

in view (29) and the relation between  $s$  and  $q$ .

(b): Elliptic regularity theory in the present case implies that  $z \in W^{2,q}(\Omega)$  for all  $1 \leq q < \infty$  with

$$\|z\|_{W^{2,q}} \leq Cq \|v\|_{L^q}$$

where the constant  $C$  is independent of  $q$ . For the dependence on  $q$  in this estimate we refer to the work of Agmon, Douglis and Nirenberg [1], see also [31] and [33, Chapter 9]. Proceeding as in (a) we have

$$\|z - z_h\|_{L^\infty} \leq Ch^{2-\frac{d}{q}} |\log h| \|z\|_{W^{2,q}} \leq Cqh^{2-\frac{d}{q}} |\log h| \|v\|_{L^q} \leq Cqh^{2-\frac{d}{q}} |\log h| \|v\|_{L^\infty},$$

so that choosing  $q = |\log h|$  gives the result.  $\square$

An important ingredient in our analysis is an error bound for a solution of a Neumann problem with a measure valued right hand side. Let  $A$  be as above and consider

$$\begin{aligned} A^*q &= \tilde{\mu}_\Omega & \text{in } \Omega \\ \sum_{i=1}^d (\sum_{j=1}^d a_{ij}q_{x_j} + b_i q)\nu_i &= \tilde{\mu}_{\partial\Omega} & \text{on } \partial\Omega. \end{aligned} \quad (31)$$

**Theorem 4.2.** Let  $\tilde{\mu} \in \mathcal{M}(\bar{\Omega})$ . Then there exists a unique weak solution  $q \in L^2(\Omega)$  of (31), i.e.

$$\int_\Omega qAv = \int_\Omega v d\tilde{\mu} \quad \forall v \in H^2(\Omega) \text{ with } \sum_{i,j=1}^d a_{ij}v_{x_i}\nu_j = 0 \text{ on } \partial\Omega.$$

Furthermore,  $q$  belongs to  $W^{1,s}(\Omega)$  for all  $s \in (1, \frac{d}{d-1})$ . For the finite element approximation  $q_h \in X_h$  of  $q$  defined by

$$a(v_h, q_h) = \int_\Omega v_h d\tilde{\mu} \quad \text{for all } v_h \in X_h$$

the following error estimate holds:

$$\|q - q_h\| \leq Ch^{2-\frac{d}{2}} \|\tilde{\mu}\|_{\mathcal{M}(\bar{\Omega})}. \quad (32)$$

*Proof.* A corresponding result is proved in [7] for the case of an operator  $A$  without transport term subject to Dirichlet conditions, but the arguments can be adapted to our situation. We omit the details.  $\square$

#### 4.1. Variational discretization

The discretization of the partial differential equations induces a natural discretization of the control via the optimality condition. Every a priori discretization of the control introduces a significant additional error which may reduce the approximation rates. Therefore, only the partial differential equations are discretized in the variational discretization concept. Problem (13) is now approximated by the following sequence of so called *variational discrete* control problems [44] depending on the mesh parameter  $h$ :

$$\begin{aligned} \min_{u \in U_{ad}} \hat{J}_h(u) &:= \frac{1}{2} \int_{\Omega} |y_h - y_0|^2 + \frac{\alpha}{2} \|u\|_U^2 \\ \text{subject to } y_h &= \mathcal{G}_h(Bu) \text{ and } y_h(x_j) \leq b(x_j) \text{ for } j = 1, \dots, m. \end{aligned} \quad (33)$$

Notice that the integer  $m$  is not fixed and tends to infinity as  $h \rightarrow 0$ , so that the number of state constraints in this optimal control problem increases with decreasing mesh size of underlying finite element approximation of the state space. This discretization approach can be understood as a generalization of the *First discretize, then optimize* approach in that it avoids discretization of the control space  $U$ . Problem (33) represents a convex infinite-dimensional optimization problem of similar structure as problem (13), but with only finitely many equality and inequality constraints for the state, which form a convex admissible set. So we are again in the setting of (1) with  $Y$  replaced by the finite element space  $X_h$  (compare also the analysis of Casas presented in [10])

**Lemma 4.3.** *Problem (33) has a unique solution  $u_h \in U_{ad}$ . There exist  $\mu_1, \dots, \mu_m \in \mathbb{R}$  and  $p_h \in X_h$  such that with  $y_h = \mathcal{G}_h(Bu_h)$  and  $\mu_h = \sum_{j=1}^m \mu_j \delta_{x_j}$  we have*

$$a(v_h, p_h) = \int_{\Omega} (y_h - y_0)v_h + \int_{\bar{\Omega}} v_h d\mu_h \quad \forall v_h \in X_h, \quad (34)$$

$$(RB^*p_h + \alpha u_h, v - u_h)_U \geq 0 \quad \forall v \in U_{ad}, \quad (35)$$

$$\mu_j \geq 0, y_h(x_j) \leq b(x_j), j = 1, \dots, m \text{ and } \int_{\bar{\Omega}} (I_h b - y_h) d\mu_h = 0. \quad (36)$$

Here,  $\delta_x$  denotes the Dirac measure concentrated at  $x$  and  $I_h$  is the usual Lagrange interpolation operator. We have  $\hat{J}'_h(v) = B^*p_h + \alpha(\cdot, u_h)_U$ , so that the considerations after Theorem 2.3 also apply in the present setting, but with  $p$  replaced by the discrete function  $p_h$ , i.e. there holds

$$u_h = P_{U_{ad}}(u_h - \sigma R \hat{J}'_h(u_h)) \quad (\sigma > 0).$$

For  $\sigma = \frac{1}{\alpha}$  we obtain

$$u = P_{U_{ad}} \left( -\frac{1}{\alpha} RB^*p \right) \text{ and } u_h = P_{U_{ad}} \left( -\frac{1}{\alpha} RB^*p_h \right). \quad (37)$$

Due to the presence of  $P_{U_{ad}}$  in variational discretization the function  $u_h \in U_{ad}$  will in general not belong to  $X_h$  even in the case  $U = L^2(\Omega), B = Id$ . This is different for the purely state constrained problem, for which  $P_{U_{ad}} \equiv Id$ , so that in this specific setting  $u_h = -\frac{1}{\alpha} p_h \in X_h$  by (37). In that case the space  $U = L^2(\Omega)$  in (33) may be replaced by  $X_h$  to obtain the same discrete solution  $u_h$ , which results in a finite-dimensional discrete optimization problem instead. However, we emphasize, that the infinite-dimensional formulation of (33) is very useful in numerical analysis [46, Chap. 3].

As a first result for (33) it is proved in e.g. [46, Chap. 3] that the sequence of optimal controls, states and the measures  $\mu_h$  are uniformly bounded.

**Lemma 4.4.** *Let  $u_h \in U_{ad}$  be the optimal solution of (33) with corresponding state  $y_h \in X_h$  and adjoint variables  $p_h \in X_h$  and  $\mu_h \in \mathcal{M}(\bar{\Omega})$ . Then there exists  $\bar{h} > 0$  so that*

$$\|y_h\|, \|u_h\|_U, \|\mu_h\|_{\mathcal{M}(\bar{\Omega})} \leq C \quad \text{for all } 0 < h \leq \bar{h}.$$

*Proof.* Since  $\mathcal{G}(B\tilde{u})$  is continuous, Assumption 2.2 implies that there exists  $\delta > 0$  such that

$$\mathcal{G}(B\tilde{u}) \leq b - \delta \quad \text{in } \bar{\Omega}. \quad (38)$$

It follows from (28) that there is  $h_0 > 0$  with

$$\mathcal{G}_h(B\tilde{u}) \leq b \quad \text{in } \bar{\Omega} \text{ for all } 0 < h \leq h_0$$

so that  $J_h(u_h) \leq J_h(\tilde{u}) \leq C$  uniformly in  $h$  giving

$$\|u_h\|_U, \|y_h\| \leq C \quad \text{for all } h \leq h_0. \quad (39)$$

Next, let  $u$  denote the unique solution to problem (13). We infer from (38) and (28) that  $v := \frac{1}{2}u + \frac{1}{2}\tilde{u}$  satisfies

$$\begin{aligned} \mathcal{G}_h(Bv) &\leq \frac{1}{2}\mathcal{G}(Bu) + \frac{1}{2}\mathcal{G}(B\tilde{u}) + Ch^{2-\frac{d}{2}}(\|Bu\| + \|B\tilde{u}\|) \\ &\leq b - \frac{\delta}{2} + Ch^{2-\frac{d}{2}}(\|u\|_U + \|\tilde{u}\|_U) \leq b - \frac{\delta}{4} \quad \text{in } \bar{\Omega} \end{aligned} \quad (40)$$

provided that  $h \leq \bar{h}, \bar{h} \leq h_0$ . Since  $v \in U_{ad}$ , (35), (34), (39) and (40) imply

$$\begin{aligned} 0 &\leq (RB^*p_h + \alpha(u_h - u_{0,h}), v - u_h)_U = \int_{\Omega} B(v - u_h)p_h + \alpha(u_h - u_{0,h}, v - u_h)_U \\ &= a(\mathcal{G}_h(Bv) - y_h, p_h) + \alpha(u_h - u_{0,h}, v - u_h)_U \\ &= \int_{\Omega} (\mathcal{G}_h(Bv) - y_h)(y_h - y_0) + \int_{\bar{\Omega}} (\mathcal{G}_h(Bv) - y_h)d\mu_h + \alpha(u_h - u_{0,h}, v - u_h)_U \\ &\leq C + \sum_{j=1}^m \mu_j (b(x_j) - \frac{\delta}{4} - y_h(x_j)) = C - \frac{\delta}{4} \sum_{j=1}^m \mu_j \end{aligned}$$

where the last equality is a consequence of (36). It follows that

$$\|\mu_h\|_{\mathcal{M}(\bar{\Omega})} \leq C$$

and the lemma is proved.  $\square$

## 4.2. Piecewise constant controls

We consider the special case  $U = L^2(\Omega)$ , so that  $B$  denotes the injection of  $L^2(\Omega)$  into  $H^1(\Omega)^*$  with box constraints  $a_l \leq u \leq a_r$  on the control. Controls are approximated by element-wise piecewise constant functions. For details we refer to [23]. We define the space of piecewise constant functions

$$Y_h := \{v_h \in L^2(\Omega) \mid v_h \text{ is constant on each } T \in \mathcal{T}_h\}.$$

and denote by  $Q_h : L^2(\Omega) \rightarrow Y_h$  the orthogonal projection onto  $Y_h$  so that

$$(Q_h v)(x) := \int_T v, \quad x \in T, T \in \mathcal{T}_h,$$

where  $\int_T v$  denotes the average of  $v$  over  $T$ . In order to approximate (13) we introduce a discrete counterpart of  $U_{ad}$ ,

$$U_{ad}^h := \{v_h \in Y_h \mid a_l \leq v_h \leq a_u \text{ in } \Omega\}.$$

Note that  $U_{ad}^h \subset U_{ad}$  and that  $Q_h v \in U_{ad}^h$  for  $v \in U_{ad}$ . Since  $Q_h v \rightarrow v$  in  $L^2(\Omega)$  as  $h \rightarrow 0$  we infer from the continuous embedding  $H^2(\Omega) \hookrightarrow C^0(\bar{\Omega})$  and Lemma 4.1 that

$$\mathcal{G}_h(Q_h v) \rightarrow \mathcal{G}(v) \text{ in } L^\infty(\Omega) \text{ for all } v \in U_{ad}. \quad (41)$$

Problem (13) here now is approximated by the following sequence of control problems depending on the mesh parameter  $h$ :

$$\begin{aligned} \min_{u \in U_{ad}^h} J_h(u) &:= \frac{1}{2} \int_{\Omega} |y_h - y_0|^2 + \frac{\alpha}{2} \int_{\Omega} |u|^2 \\ \text{subject to } y_h &= \mathcal{G}_h(u) \text{ and } y_h(x_j) \leq b(x_j) \text{ for } j = 1, \dots, m. \end{aligned} \quad (42)$$

Problem (42), as problem (33), represents a convex finite-dimensional optimization problem of similar structure as problem (13), but with only finitely many equality and inequality constraints for state and control, which form a convex admissible set. The following optimality conditions can be argued as those given in (4.3) for problem (33).

**Lemma 4.5.** *Problem (42) has a unique solution  $u_h \in U_{ad}^h$ . There exist  $\mu_1, \dots, \mu_m \in \mathbb{R}$  and  $p_h \in X_h$  such that with  $y_h = \mathcal{G}_h(u_h)$  and  $\mu_h = \sum_{j=1}^m \mu_j \delta_{x_j}$  we have*

$$a(v_h, p_h) = \int_{\Omega} (y_h - y_0)v_h + \int_{\bar{\Omega}} v_h d\mu_h \quad \forall v_h \in X_h, \quad (43)$$

$$\int_{\Omega} (p_h + \alpha u_h)(v_h - u_h) \geq 0 \quad \forall v_h \in U_{ad}^h, \quad (44)$$

$$\mu_j \geq 0, y_h(x_j) \leq b(x_j), j = 1, \dots, m \text{ and } \int_{\bar{\Omega}} (I_h b - y_h) d\mu_h = 0. \quad (45)$$

Here,  $\delta_x$  denotes the Dirac measure concentrated at  $x$  and  $I_h$  is the usual Lagrange interpolation operator.

For (42) we now prove bounds on the discrete states and the discrete multipliers. Similar to Lemma 4.4 we have

**Lemma 4.6.** *Let  $u_h \in U_{ad}^h$  be the optimal solution of (42) with corresponding state  $y_h \in X_h$  and adjoint variables  $p_h \in X_h$  and  $\mu_h \in \mathcal{M}(\bar{\Omega})$ . Then there exists  $\bar{h} > 0$  such that*

$$\|y_h\|, \|\mu_h\|_{\mathcal{M}(\bar{\Omega})} \leq C, \quad \|p_h\|_{H^1} \leq C\gamma(d, h) \quad \text{for all } 0 < h \leq \bar{h},$$

where  $\gamma(2, h) = \sqrt{|\log h|}$  and  $\gamma(3, h) = h^{-\frac{1}{2}}$ .

*Proof.* Since  $\mathcal{G}(\tilde{u}) \in C^0(\bar{\Omega})$ , Assumption 2.2 implies that there exists  $\delta > 0$  such that

$$\mathcal{G}(\tilde{u}) \leq b - \delta \quad \text{in } \bar{\Omega}. \quad (46)$$

It follows from (41) that there is  $\bar{h} > 0$  with

$$\mathcal{G}_h(Q_h \tilde{u}) \leq b - \frac{\delta}{2} \quad \text{in } \bar{\Omega} \text{ for all } 0 < h \leq \bar{h}. \quad (47)$$

Since  $Q_h \tilde{u} \in U_{ad}^h$ , (45), (44) and (47) imply

$$\begin{aligned} 0 &\leq \int_{\Omega} (p_h + \alpha u_h)(Q_h \tilde{u} - u_h) = \int_{\Omega} p_h(Q_h \tilde{u} - u_h) + \alpha \int_{\Omega} u_h(Q_h \tilde{u} - u_h) \\ &= a(\mathcal{G}_h(Q_h \tilde{u}) - y_h, p_h) + \alpha \int_{\Omega} u_h(Q_h \tilde{u} - u_h) \\ &= \int_{\Omega} (\mathcal{G}_h(Q_h \tilde{u}) - y_h)(y_h - y_0) + \int_{\bar{\Omega}} (\mathcal{G}_h(Q_h \tilde{u}) - y_h) d\mu_h + \alpha \int_{\Omega} u_h(Q_h \tilde{u} - u_h) \\ &\leq C - \frac{1}{2} \|y_h\|^2 + \sum_{j=1}^m \mu_j (b(x_j) - \frac{\delta}{2} - y_h(x_j)) = C - \frac{1}{2} \|y_h\|^2 - \frac{\delta}{2} \sum_{j=1}^m \mu_j \end{aligned}$$

where the last equality is a consequence of (45). It follows that  $\|y_h\|, \|\mu_h\|_{\mathcal{M}(\bar{\Omega})} \leq C$ . In order to bound  $\|p_h\|_{H^1}$  we insert  $v_h = p_h$  into (44) and deduce with the help of the coercivity of  $A$ , a well-known inverse estimate and the bounds we have already obtained that

$$\begin{aligned} c_1 \|p_h\|_{H^1}^2 &\leq a(p_h, p_h) = \int_{\Omega} (y_h - y_0)p_h + \int_{\bar{\Omega}} p_h d\mu_h \\ &\leq \|y_h - y_0\| \|p_h\| + \|p_h\|_{L^\infty} \|\mu_h\|_{\mathcal{M}(\bar{\Omega})} \leq C \|p_h\| + C\gamma(d, h) \|p_h\|_{H^1}. \end{aligned}$$

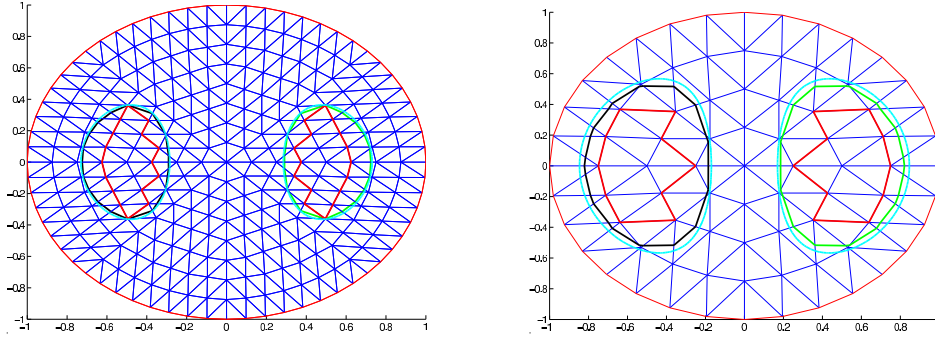


FIGURE 1. Numerical comparison of active sets obtained by variational discretization, and those obtained by a conventional approach with piecewise linear, continuous controls:  $h = \frac{1}{8}$  and  $\alpha = 0.1$  (left),  $h = \frac{1}{4}$  and  $\alpha = 0.01$  (right). The red line depicts the boarder of the active set in the conventional approach, the cyan line the exact boarder, the black and green lines, respectively the boarders of the active set in variational discretization.

Hence  $\|p_h\|_{H^1} \leq C\gamma(d, h)$  and the lemma is proved.  $\square$

Similar considerations hold for control approximations by continuous, piecewise polynomial functions. Discrete approaches to problem (13) relying on control approximations directly lead to fully discrete optimization problems like (42). These approaches lead to large-scale finite-dimensional optimization problems, since the discretization of the pde in general introduces a large number of degrees of freedom. Numerical implementation then is easy, which certainly is an important advantage of control approximations over variational discretization, whose numerical implementation is more involved. The use of classical NLP solvers for the numerical solution of the underlying discretized problems only is feasible, if the solver allows to exploit the underlying problem structure e.g. by providing user interfaces for first- and second-order derivatives.

On the other hand, the numerical implementation of variational discretization is not straightforward. The great advantage of variational discretization however is its property of optimal approximation accuracy, which is completely determined by that of the related state and adjoint state. Fig. 1 compares active sets obtained by variational discretization and piecewise linear control approximations in the presence of box constraints. One clearly observes that the active sets are resolved much more accurate when using variational discretization. In particular, the boundary of the active set is in general different from finite element edges.

The error analysis for problem (13) relies on the regularity of the involved variables, which is reflected by the optimality system presented in (15)-(17). If only control constraints are present, neither the multiplier  $\mu$  in (15) nor the complementarity condition (17) appear. Then the variational inequality (16) restricts the regularity of the control  $u$ , and thus also that of the state  $y$ . If the desired state  $y_0$  is regular enough, the adjoint variable  $p$  admits the highest regularity properties among all variables involved in the optimality system. Error analysis in this case then should involve the adjoint variable  $p$  and exploit its regularity properties.

If pointwise state constraints, are present, the situation is completely different. Now the adjoint variable only admits low regularity due to the presence of the multiplier  $\mu$ , which in general is only a measure. The state now admits the highest regularity in the optimality system. This fact should be exploited in the error analysis. However, the presence of the complementarity system (17) requires  $L^\infty$ -error estimates for the state.

### 4.3. Error bounds

For the approximation error of variational discretization we have the following theorem, whose proof can be found in [46, Chap. 3].

**Theorem 4.7.** *Let  $u$  and  $u_h$  be the solutions of (13) and (33) respectively. Then*

$$\|u - u_h\|_U + \|y - y_h\|_{H^1} \leq Ch^{1-\frac{d}{4}}.$$

*If in addition  $Bu \in W^{1,s}(\Omega)$  for some  $s \in (1, \frac{d}{d-1})$  then*

$$\|u - u_h\|_U + \|y - y_h\|_{H^1} \leq Ch^{\frac{3}{2}-\frac{d}{2s}} \sqrt{|\log h|}.$$

*If  $Bu, Bu_h \in L^\infty(\Omega)$  with  $(Bu_h)_h$  uniformly bounded in  $L^\infty(\Omega)$  also*

$$\|u - u_h\|_U, \|y - y_h\|_{H^1} \leq Ch |\log h|,$$

*where the latter estimate is valid for  $d = 2, 3$ .*

*Proof.* We test (16) with  $u_h$ , (35) with  $u$  and add the resulting inequalities. This gives

$$(RB^*(p - p_h) - \alpha(u_0 - u_{0,h}) + \alpha(u - u_h), u_h - u)_U \geq 0,$$

which in turn yields

$$\alpha \|u - u_h\|_U^2 \leq \int_{\Omega} B(u_h - u)(p - p_h) - \alpha(u_0 - u_{0,h}, u_h - u)_U. \quad (48)$$

Let  $y^h := \mathcal{G}_h(Bu) \in X_h$  and denote by  $p^h \in X_h$  the unique solution of

$$a(w_h, p^h) = \int_{\Omega} (y - y_0)w_h + \int_{\bar{\Omega}} w_h d\mu \quad \text{for all } w_h \in X_h.$$

Applying Theorem 4.2 with  $\tilde{\mu} = (y - y_0)dx + \mu$  we infer

$$\|p - p^h\| \leq Ch^{2-\frac{d}{2}} (\|y - y_0\| + \|\mu\|_{\mathcal{M}(\bar{\Omega})}). \quad (49)$$

Recalling that  $y_h = \mathcal{G}_h(Bu_h)$ ,  $y^h = \mathcal{G}_h(Bu)$  and observing (34) as well as the definition of  $p^h$  we can rewrite the first term in (48)

$$\begin{aligned} \int_{\Omega} B(u_h - u)(p - p_h) &= \int_{\Omega} B(u_h - u)(p - p^h) + \int_{\Omega} B(u_h - u)(p^h - p_h) \\ &= \int_{\Omega} B(u_h - u)(p - p^h) + a(y_h - y^h, p^h - p_h) \\ &= \int_{\Omega} B(u_h - u)(p - p^h) + \int_{\Omega} (y - y_h)(y_h - y^h) + \int_{\bar{\Omega}} (y_h - y^h)d\mu - \int_{\bar{\Omega}} (y_h - y^h)d\mu_h \\ &= \int_{\Omega} B(u_h - u)(p - p^h) - \|y - y_h\|^2 + \int_{\Omega} (y - y_h)(y - y^h) \\ &\quad + \int_{\bar{\Omega}} (y_h - y^h)d\mu + \int_{\bar{\Omega}} (y^h - y_h)d\mu_h. \end{aligned} \quad (50)$$

After inserting (50) into (48) and using Young's inequality we obtain in view of (49), (27) and the properties of the  $L^2$ -projection

$$\begin{aligned} &\frac{\alpha}{2} \|u - u_h\|_U^2 + \frac{1}{2} \|y - y_h\|^2 \\ &\leq C(\|p - p^h\|^2 + \|y - y^h\|^2 + \|u_0 - u_{0,h}\|^2) + \int_{\bar{\Omega}} (y_h - y^h)d\mu + \int_{\bar{\Omega}} (y^h - y_h)d\mu_h \\ &\leq Ch^{4-d} + \int_{\bar{\Omega}} (y_h - y^h)d\mu + \int_{\bar{\Omega}} (y^h - y_h)d\mu_h. \end{aligned} \quad (51)$$

It remains to estimate the integrals involving the measures  $\mu$  and  $\mu_h$ . Since

$$y_h - y^h \leq (I_h b - b) + (b - y) + (y - y^h) \quad \text{in } \bar{\Omega}$$

we deduce with the help of (17)

$$\int_{\bar{\Omega}} (y_h - y^h) d\mu \leq \|\mu\|_{\mathcal{M}(\bar{\Omega})} (\|I_h b - b\|_{\infty} + \|y - y^h\|_{\infty}).$$

Similarly, (36) implies

$$\int_{\bar{\Omega}} (y^h - y_h) d\mu_h \leq \|\mu_h\|_{\mathcal{M}(\bar{\Omega})} (\|b - I_h b\|_{\infty} + \|y - y^h\|_{\infty}).$$

Inserting the above estimates into (51) and using Lemma 4.4 as well as an interpolation estimate we infer

$$\|u - u_h\|_U^2 + \|y - y_h\|^2 \leq Ch^{4-d} + C\|y - y^h\|_{L^\infty}. \quad (52)$$

The estimates on  $\|u - u_h\|_U$  now follow from (28) and Lemma 4.1 respectively. Finally, in order to bound  $\|y - y_h\|_{H^1}$  we note that

$$a(y - y_h, v_h) = \int_{\Omega} B(u - u_h)v_h$$

for all  $v_h \in X_h$ , from which one derives the desired estimates using standard finite element techniques and the bounds on  $\|u - u_h\|_U$ . In order to avoid the dependence on the dimension we should avoid finite element approximations of the adjoint variable  $p$ , which due to its low regularity only allows error estimates in the  $L^2$  norm. We therefore provide a proof technique which completely avoids the use of finite element approximations of the adjoint variable. To begin with we start with the basic estimate (48)

$$\alpha\|u - u_h\|_U^2 \leq \int_{\Omega} B(u_h - u)(p - p_h) - \alpha(u_0 - u_{0,h}, u_h - u)_U$$

and write

$$\begin{aligned} \int_{\Omega} B(u_h - u)(p - p_h) &= \int_{\Omega} pA(\tilde{y} - y) - a(y_h - y^h, p_h) = \\ &= \int_{\Omega} (y - y_0)(\tilde{y} - y) + \int_{\Omega} \tilde{y} - y d\mu - \int_{\Omega} (y_h - y_0)(y_h - y^h) + \int_{\Omega} y_h - y^h d\mu_h, \end{aligned}$$

where  $\tilde{y} := \mathcal{G}(u_h)$ . Proceeding similar as in the proof of the previous theorem we obtain

$$\begin{aligned} \int_{\Omega} (y - y_0)(\tilde{y} - y) + \int_{\Omega} \tilde{y} - y d\mu - \int_{\Omega} (y_h - y_0)(y_h - y^h) + \int_{\Omega} y_h - y^h d\mu_h &\leq \\ &\leq C\{\|\mu\|_{\mathcal{M}(\bar{\Omega})} + \|\mu_h\|_{\mathcal{M}(\bar{\Omega})}\}\{\|b - I_h b\|_{\infty} + \|y - y^h\|_{\infty} + \|\tilde{y} - y_h\|_{\infty}\} - \\ &\quad - \|y - y_h\|^2 + C\{\|y - y^h\| + \|\tilde{y} - y_h\|\}. \end{aligned}$$

Using Lemma 4.4 together with Lemma 4.1 then yields

$$\alpha\|u - u_h\|_U^2 + \|y - y_h\|^2 \leq C\{h^2 + h^2|\log h|^2\},$$

so that the claim follows as in the proof of the previous theorem.  $\square$

**Remark 4.8.** Let us note that the approximation order of the controls and states in the presence of control and state constraints is the same as in the purely state constrained case, if  $Bu \in W^{1,s}(\Omega)$ . This assumption holds for the important example  $U = L^2(\Omega)$ ,  $B = Id$  and  $u_{0,h} = P_h u_0$ , with  $u_0 \in H^1(\Omega)$  and  $P_h : L^2(\Omega) \rightarrow X_h$  denoting the  $L^2$ -projection, and subsets of the form

$$U_{ad} = \{v \in L^2(\Omega), a_l \leq v \leq a_u \text{ a.e. in } \Omega\},$$

with bounds  $a_l, a_u \in W^{1,s}(\Omega)$ , since  $u_0 \in H^1(\Omega)$ , and  $p \in W^{1,s}(\Omega)$ . Moreover,  $u, u_h \in L^\infty(\Omega)$  with  $\|u_h\|_{\infty} \leq C$  uniformly in  $h$  holds if for example  $a_l, a_u \in L^\infty(\Omega)$ .



For piecewise constant control approximations and the setting of Section 4.2 the following theorem is proved in [23].

**Theorem 4.9.** *Let  $u$  and  $u_h$  be the solutions of (13) and (42) respectively with  $(u_h)_h \subset L^\infty(\Omega)$  uniformly bounded. Then we have for  $0 < h \leq \bar{h}$*

$$\|u - u_h\| + \|y - y_h\|_{H^1} \leq \begin{cases} Ch|\log h|, & \text{if } d = 2 \\ C\sqrt{h}, & \text{if } d = 3. \end{cases}$$

*Proof.* We test (16) with  $u_h$ , (45) with  $Q_h u$  and add the resulting inequalities. Keeping in mind that  $u - Q_h u \perp Y_h$  we obtain

$$\begin{aligned} & \int_{\Omega} (p - p_h + \alpha(u - u_h))(u_h - u) \\ & \geq \int_{\Omega} (p_h + \alpha u_h)(u - Q_h u) = \int_{\Omega} (p_h - Q_h p_h)(u - Q_h u). \end{aligned}$$

As a consequence,

$$\alpha \|u - u_h\|^2 \leq \int_{\Omega} (u_h - u)(p - p_h) - \int_{\Omega} (p_h - Q_h p_h)(u - Q_h u) \equiv I + II. \quad (53)$$

Let  $y^h := \mathcal{G}_h(u) \in X_h$  and denote by  $p^h \in X_h$  the unique solution of

$$a(w_h, p^h) = \int_{\Omega} (y - y_0)w_h + \int_{\bar{\Omega}} w_h d\mu \quad \text{for all } w_h \in X_h.$$

Applying Theorem 4.2 with  $\tilde{\mu} = (y - y_0)dx + \mu$  we infer

$$\|p - p^h\| \leq Ch^{2-\frac{d}{2}} (\|y - y_0\| + \|\mu\|_{\mathcal{M}(\bar{\Omega})}). \quad (54)$$

Recalling that  $y_h = \mathcal{G}_h(u_h)$ ,  $y^h = \mathcal{G}_h(u)$  and observing (44) as well as the definition of  $p^h$  we can rewrite the first term in (53)

$$\begin{aligned} I &= \int_{\Omega} (u_h - u)(p - p^h) + \int_{\Omega} (u_h - u)(p^h - p_h) \\ &= \int_{\Omega} (u_h - u)(p - p^h) + a(y_h - y^h, p^h - p_h) \\ &= \int_{\Omega} (u_h - u)(p - p^h) + \int_{\Omega} (y - y_h)(y_h - y^h) + \int_{\bar{\Omega}} (y_h - y^h)d\mu - \int_{\bar{\Omega}} (y_h - y^h)d\mu_h \\ &= \int_{\Omega} (u_h - u)(p - p^h) - \|y - y_h\|^2 + \int_{\Omega} (y - y_h)(y - y^h) \\ &\quad + \int_{\bar{\Omega}} (y_h - y^h)d\mu + \int_{\bar{\Omega}} (y^h - y_h)d\mu_h. \end{aligned} \quad (55)$$

Applying Young's inequality we deduce

$$\begin{aligned} |I| &\leq \frac{\alpha}{4} \|u - u_h\|^2 - \frac{1}{2} \|y - y_h\|^2 + C(\|p - p^h\|^2 + \|y - y^h\|^2) \\ &\quad + \int_{\bar{\Omega}} (y_h - y^h)d\mu + \int_{\bar{\Omega}} (y^h - y_h)d\mu_h. \end{aligned} \quad (56)$$

Let us estimate the integrals involving the measures  $\mu$  and  $\mu_h$ . Since  $y_h - y^h \leq (I_h b - b) + (b - y) + (y - y^h)$  in  $\bar{\Omega}$  we deduce with the help of (17), Lemma 4.1 and an interpolation estimate

$$\int_{\bar{\Omega}} (y_h - y^h)d\mu \leq \|\mu\|_{\mathcal{M}(\bar{\Omega})} (\|I_h b - b\|_{\infty} + \|y - y^h\|_{\infty}) \leq Ch^2 |\log h|^2.$$

On the other hand  $y^h - y_h \leq (y^h - y) + (b - I_h b) + (I_h b - y_h)$ , so that (45), Lemma 4.1 and Lemma 4.6 yield

$$\int_{\bar{\Omega}} (y^h - y_h)d\mu_h \leq \|\mu_h\|_{\mathcal{M}(\bar{\Omega})} (\|b - I_h b\|_{\infty} + \|y - y^h\|_{\infty}) \leq Ch^2 |\log h|^2.$$

Inserting these estimates into (56) and recalling (27) as well as (32) we obtain

$$|I| \leq \frac{\alpha}{4} \|u - u_h\|^2 - \frac{1}{2} \|y - y_h\|^2 + Ch^{4-d} + Ch^2 |\log h|^2. \quad (57)$$

Let us next examine the second term in (53). Since  $u_h = Q_h u_h$  and  $Q_h$  is stable in  $L^2(\Omega)$  we have

$$\begin{aligned} |II| &\leq 2 \|u - u_h\| \|p_h - Q_h p_h\| \leq \frac{\alpha}{4} \|u - u_h\|^2 + Ch^2 \|p_h\|_{H^1}^2 \\ &\leq \frac{\alpha}{4} \|u - u_h\|^2 + Ch^2 \gamma(d, h)^2 \end{aligned}$$

using an interpolation estimate for  $Q_h$  and Lemma 4.6. Combining this estimate with (57) and (53) we finally obtain

$$\|u - u_h\|^2 + \|y - y_h\|^2 \leq Ch^{4-d} + Ch^2 |\log h|^2 + Ch^2 \gamma(d, h)^2$$

which implies the estimate on  $\|u - u_h\|$ . In order to bound  $\|y - y_h\|_{H^1}$  we note that

$$a(y - y_h, v_h) = \int_{\Omega} (u - u_h) v_h$$

for all  $v_h \in X_h$ , from which one derives the desired estimate using standard finite element techniques and the bound on  $\|u - u_h\|$ .  $\square$

**Remark 4.10.** An inspection of the proof of Theorem 4.7 shows that we also could avoid to use error estimates for the auxiliary function  $p^h$  if we would use a technique for the term I similar to that used in the proof of the third part of Theorem 4.7. However, our approach to estimate II is based on inverse estimates which finally lead to the dimension dependent error estimate presented in Theorem 4.7.

The theorems above have in common that a control error estimate is only available for  $\alpha > 0$ . However, the appearance of  $\alpha$  in these estimates indicates that in the *bang-bang*-case  $\alpha = 0$  an error estimate for  $\|y - y_h\|_{L^2}$  still is available, whereas no information for the control error  $\|u - u_h\|_U$  seems to remain. In [25] a refined analysis of bang-bang controls without state constraints also provides estimates for the control error on inactive regions in the  $L^1$ -norm. We further observe that piecewise constant control approximations in 2 space dimensions deliver the same approximation quality as variational discrete controls. Only in 3 space dimensions variational discretization provides a better error estimate. This is caused by the fact that state constraints limit the regularity of the adjoint state, so that optimal error estimates can be expected by techniques which avoid its use. Currently the analysis for piecewise constant control approximations involves an inverse estimate for  $\|p_h\|_{H^1}$ , which explains the lower approximation order in the case  $d = 3$ .

Let us mention that the bottleneck in the analysis here is not formed by control constraints, but by the state constraints. In fact, if one uses  $U_{ad} = U$ , then variational discretization (33) delivers the same numerical solution as the approach (42) with piecewise linear, continuous control approximations. Variational discretization really pays off if only control constraints are present and the adjoint variable is smooth, compare [44],[46, Chap.. 3].

For the numerical solution of problem (33), (42) several approaches exist in the literature. Common are so called regularization methods which relax the state constraints in (13) by either substituting it by a mixed control-state constraint (Lavrentiev relaxation [62]), or by adding suitable penalty terms to the cost functional instead requiring the state constraints (barrier methods [47, 70], penalty methods [38, 40]). These approaches will be discussed in the following section.

## 5. Regularization and Discretization

### 5.1. Motivation

In this section we will focus on optimal control problems with pointwise state constraints. Let us consider

$$\min F(y, u) = \frac{1}{2} \|y - y_d\|_{L^2(\Omega)}^2 + \frac{\nu}{2} \|u\|_{L^2(\Omega)}^2 \quad (58)$$

subject to the state equation

$$y = S(u), \quad (59)$$

and pointwise state constraints

$$y \geq y_c \quad \text{a.e. on } D, \quad (60)$$

and  $u \in U_{ad}$  which may be the whole space  $U = L^2(\Omega)$  or contain additional control constraints. Moreover,  $D \subset \Omega$  denotes the set where the state constraints are given. The operator  $S$  plays the role of a solution operator of a linear or nonlinear partial differential equation. Example 1.1 represents the nonlinear case.

Let us now focus on linear partial differential equations. For the numerical solution of problem (33), (42) several approaches exist in the literature. Common are so called regularization methods which relax the state constraints in (13) by either substituting it by a mixed control-state constraint (Lavrentiev relaxation [62]), or by adding suitable penalty terms to the cost functional instead requiring the state constraints (barrier methods [47, 70], penalty methods [38, 40]).

The analysis of unregularized and regularized optimal control problems is quite similar. Regularization leads often to more smooth solution. However, this effect disappears when the regularization parameter tends to zero. Consequently solving methods and numerical aspects are the main reason for regularization.

### Numerical approaches for discretized problems

Discretized optimal control problems with pointwise state constraints can be attacked by different techniques. Projected or conditional gradient methods are very robust, but slow.

Active set strategies become very popular in the recent years, see [5, 6, 52]. Active and inactive sets are fixed in every iteration. In contrast to the most classical techniques in nonlinear optimization, whole sets can change from one iteration to the next. In each iteration a problem has to be solved without inequality constraints. In many cases such methods can be interpreted as semismooth Newton methods, see [39]. Thus, active set strategies are fast convergent and mesh independent solving methods [42]. However, a direct application of active set strategies to the discretized is often impossible: If one starts completely inactive, then all violated inequalities leads to elements of active sets. It is easy to construct situations where the number of free optimization variables is smaller then the number of new active constraints. Moreover, subproblems in the active set algorithm may be ill-posed or ill-conditioned. Therefore, a direct application of active set methods to the discretized problem cannot be recommended.

Interior point methods and barrier techniques are an alternative approach. The objective functional is modified by a penalty term in such a way that the feasible iterates stay away from the bounds in the inequalities. In contrast to the active set strategies the new problem is smooth, but nonlinear. The penalty term modifies the problem. The obtained solution is not the solution of the original problem. In general the solution of the original problem is obtained by tending the penalty parameter to zero or infinity.

### Regularization techniques

We have introduced numerical techniques for solving optimal control problems with inequality constraints. The most techniques introduce some parameters and modify the (discretized) optimal control problem. One can have different views on the whole solution process. Our first approach was to discretize the optimal control problem first and then to find a solving method. Another view is to fix a discretization and a regularization or penalty parameter and to look for over all error.

Exactly this is the issue of this section. The tuning of discretization and regularization can significantly reduce the computational effort without loss on accuracy. Before we start to find error estimates for this combined approach, we will shortly explain the different techniques. For simplicity we chose a linear operator  $S$ .

### 1. Moreau-Yosida-regularization

The Moreau-Yosida-regularization (see [40]) uses a quadratic penalty term

$$\min F_\gamma(y, u) = \frac{1}{2}\|y - y_d\|_Y^2 + \frac{\nu}{2}\|u\|_U^2 + \frac{\gamma}{2}\|(y_c - y)_+\|_{L^2(D)}^2. \quad (61)$$

Moreover the inequality constraint (60) is dropped. The quantity  $\gamma$  plays the role of a regularization parameter. The original problem is obtained in the limit  $\gamma \rightarrow \infty$ . Combined error estimates should help to find a reasonable size of  $\gamma$  for a given discretization. Let us mention that the nonsmooth penalty term can be treated by a semismooth Newton approach which can be reinterpreted as an active set approach.

### 2. Lavrentiev regularization

The Lavrentiev regularization (see [59, 62]) modifies only the inequality constraint

$$y + \lambda u \geq y_c \quad \text{a.e. on } D. \quad (62)$$

This approach is only possible if the control  $u$  acts on the whole set  $D$ . This modification overcomes the ill-posedness effect in active set strategies, since

$$y_c = Su + \lambda u \text{ on } D, \quad (63)$$

requires no longer the inversion of a compact operator. In contrast to the original problem, the Lagrange multipliers associated with the mixed constraints are regular functions. This technique works well for problems without additional control constraints. However, there are difficulties for problems with additional control constraints if the control constraints and the mixed constraints are active simultaneously. Then the dual variables are not uniquely determined. Consequently, the corresponding active set strategy is not well defined. In this approach, the Lavrentiev parameter  $\lambda$  has to tend to zero to obtain the solution of the original problem. A generalization of the Lavrentiev regularization is used in the source representation method, see [72, 73].

### 3. Virtual control approach

The virtual control approach (see [50, 51]) modifies the complete problem. A new control  $v$  is introduced on the domain  $D$ :

$$\min F(y, u) = \frac{1}{2}\|y - y_d\|_Y^2 + \frac{\nu}{2}\|u\|_U^2 + \frac{f(\varepsilon)}{2}\|v\|_{L^2(D)}^2 \quad (64)$$

subject to

$$y = Su + g(\varepsilon)Tv, \quad (65)$$

and

$$y + h(\varepsilon)v \geq y_c \quad \text{a.e. on } D, \quad (66)$$

with suitable chosen functions  $f$ ,  $g$ , and  $h$ . The operator  $T$  represents a solution operator of a partial differential equation for the source term  $g(\varepsilon)v$ . Let us mention that this technique can be interpreted as Moreau-Yosida approach for the choice  $g \equiv 0$ . The original

problem is obtained for  $\varepsilon \rightarrow 0$ . In contrast to the Lavrentiev regularization, the dual variables are unique. Moreover, this approach guarantees well defined subproblems in active set strategies in contrast to the Lavrentiev regularization.

#### 4. Barrier methods - Interior point methods

Interior point methods or barrier methods deliver regularized solution which are feasible for the original problem. The objective is modified to

$$\min F_\mu(y, u) = \frac{1}{2} \|y - y_d\|_Y^2 + \frac{\nu}{2} \|u\|_U^2 + \mu \varphi(y - y_c) \quad (67)$$

where  $\varphi$  denotes a suitable smooth barrier function, see [70]. A typical choice would be a logarithmic function, i.e.,

$$\varphi(y - y_c) = - \int_D \log(y - y_c) dx.$$

The inequality constraints are dropped. The regularization parameter  $\mu$  tends to zero or infinity to obtain the original problem. Interior point methods have the advantage that no nonsmooth terms occurs. However, the barrier function  $\varphi$  generates a new nonlinearity.

Before we start with the presentation of the main ideas, we give an overview on combined regularization and discretization error estimates. Error estimates for the Lavrentiev regularization can be found in [18, 45]. The Moreau-Yosida approach is analyzed in [38]. Discretization error estimates for virtual control concept are derived in [48]. Results for the interior point approach are published in [47].

### 5.2. Error estimates for variational discretization

In this subsection we will demonstrate the technique to obtain error estimates for the variational discretization concept. We will focus on linear solution operators  $S$  here. The case of a nonlinear solution operator is discussed later.

#### Regularization error

The regularization techniques presented above modify the original problem in the objective, the state equation or in the inequality constraints. To deal with all these concepts in a general framework would lead to a very technical presentation. Therefore we pick a specific approach.

Let us explain the main issues for the Lavrentiev regularization. Here, the inequality constraint was changed to

$$y + \lambda u \geq y_c \quad \text{a.e. on } D. \quad (68)$$

Remember, that this regularization is only possible, if the set  $D$  is a subset of the set where the control acts. The inequality (68) changes the set of admissible controls. Therefore, one has to ensure that the admissible set of the regularized problem is nonempty. This is one motivation to require a Slater type condition:

**Assumption 5.1.** There exists a control  $\hat{u} \in U_{ad}$  and a real number  $\tau > 0$  with  $\hat{y} = S\hat{u} \geq y_c + \tau$  and  $\|\hat{u}\|_{L^\infty(D)} \leq c$ .

Assumption 2.2 ensure the existence of at least one feasible point if  $\lambda \leq \tau/c$ . Now one has at least three possibilities to derive regularization error estimates:

1. Work with the complete optimality system including Lagrange multipliers and use the uniform boundedness of the Lagrange multipliers. That technique was used in section 4.

2. Use the optimality conditions for a multiplier free formulation. Test the corresponding variational inequalities with suitable functions. We will demonstrate this technique in this section.
3. Work again multiplier free. The definition of admissible sets is the same as in approach 2. Now, the aim is to construct estimates of the form

$$|J(\bar{y}, \bar{u}) - J(\bar{y}_\lambda, \bar{u}_\lambda)| \leq \psi(\lambda).$$

This can be used to obtain the desired estimates for the regularization error.

All three approaches are used in literature to derive estimates. We will focus on the second approach. Let us define the sets:

$$U^0 = \{u \in U_{ad} : Su \geq y_c\},$$

$$U^\lambda = \{u \in U_{ad} : \lambda u + Su \geq y_c\}.$$

Then the first-order optimal conditions for the solution  $\bar{u}$  of the original problem and the solution  $\bar{u}_\lambda$  for the regularized problem read

$$(S^*(S\bar{u} - y_d) + \nu\bar{u}, u - \bar{u}) \geq 0 \text{ for all } u \in U^0 \quad (69)$$

$$(S^*(S\bar{u}_\lambda - y_d) + \nu\bar{u}_\lambda, u - \bar{u}_\lambda) \geq 0 \text{ for all } u \in U^\lambda \quad (70)$$

Now one has to look for suitable test functions  $u$ . Suitable test functions should be feasible for one of these problems and close to the solutions of the other problem. We assume  $u_0 \in U^0$  and  $u_\lambda \in U^\lambda$  and obtain

$$(S^*(S\bar{u} - y_d) + \nu\bar{u}, u_0 - \bar{u}) \geq 0 \quad (71)$$

$$(S^*(S\bar{u}_\lambda - y_d) + \nu\bar{u}_\lambda, u_\lambda - \bar{u}_\lambda) \geq 0 \quad (72)$$

Adding these inequalities yields

$$\begin{aligned} & (S^*(S\bar{u} - y_d) + \nu\bar{u}, u_0 - \bar{u}_\lambda) + (S^*(S\bar{u} - y_d) + \nu\bar{u}, \bar{u}_\lambda - \bar{u}) + \\ & (S^*(S\bar{u}_\lambda - y_d) + \nu\bar{u}_\lambda, u_\lambda - \bar{u}) + (S^*(S\bar{u}_\lambda - y_d) + \nu\bar{u}_\lambda, \bar{u} - \bar{u}_\lambda) \geq 0. \end{aligned} \quad (73)$$

We obtain for the first and the third term

$$(S^*(S\bar{u} - y_d) + \nu\bar{u}, u_0 - \bar{u}_\lambda) \leq c\|u_0 - \bar{u}_\lambda\|_U \quad (74)$$

$$(S^*(S\bar{u}_\lambda - y_d) + \nu\bar{u}_\lambda, u_\lambda - \bar{u}) \leq c\|u_\lambda - \bar{u}\|_U \quad (75)$$

For the sum of the second and the fourth term we find

$$\begin{aligned} & (S^*(S\bar{u} - y_d) + \nu\bar{u}, \bar{u}_\lambda - \bar{u}) + \\ & (S^*(S\bar{u}_\lambda - y_d) + \nu\bar{u}_\lambda, u_\lambda - \bar{u}) = -\nu\|\bar{u} - \bar{u}_\lambda\|_U^2 - \|S(\bar{u} - \bar{u}_\lambda)\|_Y^2 \end{aligned} \quad (76)$$

Consequently, we end up with

$$\nu\|\bar{u} - \bar{u}_\lambda\|_U^2 + \|S(\bar{u} - \bar{u}_\lambda)\|_Y^2 \leq c\|u_0 - \bar{u}_\lambda\|_U + c\|u_\lambda - \bar{u}\|_U. \quad (77)$$

Thus, the choice of the test functions  $u_0, u_\lambda$  is an important issue. We will investigate two cases.

#### First case

In the first case we assume that

$$u \in U_{ad} \quad \Rightarrow \quad \|u\|_{L^\infty(D)} \leq K.$$

Let us start with the construction of the function  $u_0$ . A reasonable choice is

$$u_0 = (1 - \delta)\bar{u}_\lambda + \delta\hat{u}.$$

Of course we have  $u_0 \in U_{ad}$ . Therefore, we have only to check the state constraints

$$\begin{aligned} Su_0 &= (1 - \delta)S\bar{u}_\lambda + \delta S\hat{u} \\ &\geq (1 - \delta)(y_c - \lambda\bar{u}) + \delta(y_c + \tau) \\ &\geq y_c + \delta\tau - (1 - \delta)\lambda K \end{aligned}$$

Consequently, we have to choose  $\delta \sim \lambda$  to satisfy the state constraints. In the same way we get for

$$u_\lambda = (1 - \sigma)\bar{u} + \sigma\hat{u}$$

the relation

$$Su_\lambda + \lambda u_\lambda \geq y_c + \sigma(\tau - \lambda\|\hat{u}\|_{L^\infty(D)}) - (1 - \sigma)\lambda K$$

and we need for feasibility  $\sigma \sim \lambda$ . This leads to the final regularization error estimate

$$\nu\|\bar{u} - \bar{u}_\lambda\|_U^2 + \|S(\bar{u} - \bar{u}_\lambda)\|_Y^2 \leq c\lambda. \quad (78)$$

### Second case

In the second case we assume  $U_{ad} = U$ . Here we have no bounds for the supremum norm of the control  $u$  on the set  $D$ . Therefore, the way to derive the two inequalities for  $u_0, u_\lambda$  is not longer possible.

To get error estimates one has to require that the operator  $S$  is sufficiently smoothing, self-adjoint and  $S + \lambda I$  is continuously invertible. Let us assume that we have

$$\|Su\|_{H^s(\Omega)} \leq c\|u\|_U$$

with  $s > d/2$ .

Since the control  $\bar{u}$  may be unbounded in a point where the state constraint is active, we have to find a new construction of  $u_\lambda$ . We choose

$$u_\lambda = (\lambda I + S)^{-1}S\bar{u} \quad (79)$$

and a simple computation shows that this function satisfies the regularized state constraints. Moreover, we get

$$u_\lambda - \bar{u} = -\lambda(\lambda I + S)^{-1}\bar{u}.$$

Note that the operator  $(\lambda I + S)^{-1}$  becomes unbounded for  $\lambda \rightarrow 0$ . The optimality condition with Lagrange multipliers yields a representation

$$\bar{u} = S^*w = Sw$$

with some  $w \in (L^\infty(D))^*$ . This can be written as

$$\bar{u} = S^k(S^{1-k}w)$$

and  $S^{1-k}w$  is an  $L^2$ -function. An easy computation yields  $0 < k < 1 - \frac{d}{2s}$ . Such a property is called source representation in the theory of inverse problems. Applying standard spectral methods from that theory, we obtain

$$\|u_\lambda - \bar{u}\|_U \leq c\lambda^k \quad (80)$$

with  $0 < k < 1 - \frac{d}{2s}$ , see [17]. It remains to construct  $u_0$ . Here we can choose again the construction

$$u_0 = (1 - \delta)\bar{u}_\lambda + \delta\hat{u},$$

but we have to change the estimation technique since the constant  $K$  appears in the estimates in the first case. We have to avoid the term  $\|\bar{u}_\lambda\|_{L^\infty(D)}$  to get an error estimate. Our aim is now to replace the  $L^\infty(D)$ -norm by the  $L^2(D)$ -norm.

Due to the form of the objective, the controls  $u_\lambda$  are uniformly bounded in  $U$ . Next we use

$$\|Su\|_{H^s(\Omega)} \leq c\|u\|_U$$

and obtain a uniform bound of  $Su_0$  in  $H^s(\Omega)$ . This space is continuously embedded in  $C^{0,s-d/2}(\bar{\Omega})$  if  $s - d/2 < 1$  and we get

$$\|Su_0\|_{C^{0,\gamma}(\bar{\Omega})} \leq c$$

with  $\gamma = s - d/2$ . Moreover, one needs the estimate

$$\|f\|_{L^\infty(D)} \leq c \|f\|_{L^2(D)}^{\frac{\gamma}{\gamma+d/2}}$$

for the specific function  $f := (y_c - Su_0)_+$ , see [50], Lemma 3.2. A short computation yields

$$\frac{\gamma}{\gamma + d/2} = 1 - \frac{d}{2s}.$$

Combining these inequalities, we find

$$\|(y_c - Su_0)_+\|_{L^\infty(D)} \leq c \|(y_c - Su_0)_+\|_{L^2(D)}^{1 - \frac{d}{2s}}.$$

This is essential ingredient for the error estimate. Now one can proceed like in the first case.

The final estimate is given by

$$\|u_0 - \bar{u}_\lambda\|_U \leq c\lambda^{1 - \frac{d}{2s}}. \quad (81)$$

Together with (77) and (80) we end up by

$$\nu \|\bar{u} - \bar{u}_\lambda\|_U^2 + \|S(\bar{u} - \bar{u}_\lambda)\|_Y^2 \leq c\lambda^k \quad (82)$$

with  $0 < k < 1 - \frac{d}{2s}$ .

### Discretization error

Only the partial differential equations are discretized in the variational discretization concept. This is reflected by the discretized state equation

$$y_h = S_h u_h. \quad (83)$$

Let us define the set of admissible controls for the discretized and regularized problem

$$U_h^\lambda = \{u_h \in U_{ad} : \lambda u_h^\lambda + S_h u_h^\lambda \geq y_c\}.$$

We obtain for the optimal solution  $\bar{u}_h^\lambda$  the following necessary and sufficient optimality condition

$$(S_h^*(S_h \bar{u}_h^\lambda - y_d) + \nu \bar{u}_h^\lambda, u_h - \bar{u}_h^\lambda) \geq 0 \text{ for all } u_h \in U_h^\lambda \quad (84)$$

We estimate the total regularization and discretization error by means of the triangle inequality

$$\|\bar{u} - \bar{u}_h^\lambda\|_U \leq \|\bar{u} - \bar{u}_\lambda\|_U + \|\bar{u}_\lambda - \bar{u}_h^\lambda\|_U. \quad (85)$$

The first term was already estimated in (78) and (82). Let us mention that a direct estimate of the total error will lead to the same result, see [45]. The estimation of the second term can be done in a similar manner as for the regularization error. Let us define

$$u_0^\lambda = (1 - \delta)\bar{u}_h^\lambda + \delta\hat{u}.$$

Then we obtain

$$\begin{aligned} \lambda u_0^\lambda + S u_0^\lambda &= (1 - \delta)(\lambda \bar{u}_h^\lambda + S u_h^\lambda) + \delta(\lambda \hat{u} + S \hat{u}) \\ &\geq y_c + (1 - \delta)(S \bar{u}_h^\lambda - S_h \bar{u}_h^\lambda) + \delta\tau - \delta\lambda \|\hat{u}\|_{L^\infty(D)} \end{aligned}$$

For  $\lambda < \frac{\tau}{\|\hat{u}\|_{L^\infty(D)}}$  we obtain

$$\lambda u_0^\lambda + S u_0^\lambda \geq y_c + \frac{\delta\tau}{2} - (1 - \delta) \|S \bar{u}_h^\lambda - S_h \bar{u}_h^\lambda\|_{L^\infty(D)}$$

which allows a choice  $\delta \sim \|S \bar{u}_h^\lambda - S_h \bar{u}_h^\lambda\|_{L^\infty(D)}$ . The same technique can be applied to

$$u_h^\lambda = (1 - \sigma)\bar{u}_\lambda + \sigma\hat{u}.$$



Again, we find for  $\lambda < \frac{\tau}{\|\bar{u}\|_{L^\infty(D)}}$

$$\lambda u_h^\lambda + S_h u_h^\lambda \geq y_c + \frac{\sigma\tau}{2} - (1 - \sigma)\|S\bar{u}_\lambda - S_h\bar{u}_\lambda\|_{L^\infty(D)} - \sigma\|S\hat{u} - S_h\hat{u}\|_{L^\infty(D)}.$$

Consequently, we can choose  $\sigma \sim \|S\bar{u}_\lambda - S_h\bar{u}_\lambda\|_{L^\infty(D)} + \|S\hat{u} - S_h\hat{u}\|_{L^\infty(D)}$ .

The derivation of the error estimate is similar to that one of the regularization error. We start with

$$\begin{aligned} (S^*(S\bar{u}_\lambda - y_d) + \nu\bar{u}_\lambda, u_0^\lambda - \bar{u}_\lambda) &\geq 0 \\ (S_h^*(S_h\bar{u}_h^\lambda - y_d) + \nu\bar{u}_h^\lambda, u_h^\lambda - \bar{u}_h^\lambda) &\geq 0 \end{aligned}$$

and add these two inequalities. However, the different operators  $S$  and  $S_h$  leads to modifications. Let us estimate the term

$$\begin{aligned} &(S^*(S\bar{u}_\lambda - y_d) + \nu\bar{u}_\lambda, \bar{u}_h^\lambda - \bar{u}_\lambda) + (S_h^*(S_h\bar{u}_h^\lambda - y_d) + \nu\bar{u}_h^\lambda, \bar{u}^\lambda - \bar{u}_h^\lambda) \\ &= -\nu\|\bar{u}^\lambda - \bar{u}_h^\lambda\|_U^2 + (S^*(S\bar{u}_\lambda - y_d) - S_h^*(S_h\bar{u}_h^\lambda - y_d), \bar{u}_h^\lambda - \bar{u}_\lambda) \\ &\leq -\nu\|\bar{u}^\lambda - \bar{u}_h^\lambda\|_U^2 + \|S^*y_d - S_h^*y_d\|_U\|\bar{u}^\lambda - \bar{u}_h^\lambda\|_U - \|S_h(\bar{u}^\lambda - \bar{u}_h^\lambda)\|_Y^2 \\ &\quad + \|S^*S\bar{u}^\lambda - S_h^*S_h\bar{u}^\lambda\|_U\|\bar{u}^\lambda - \bar{u}_h^\lambda\|_U \end{aligned}$$

The final estimate is obtained by means of Young's inequality

$$\begin{aligned} \frac{\nu}{2}\|\bar{u}^\lambda - \bar{u}_h^\lambda\|_U^2 + \|S_h(\bar{u}^\lambda - \bar{u}_h^\lambda)\|_Y^2 &\leq c(\|S^*S\bar{u}^\lambda - S_h^*S_h\bar{u}^\lambda\|_U^2 + \|S^*y_d - S_h^*y_d\|_U^2 \\ &\quad + \|S\bar{u}_h^\lambda - S_h\bar{u}_h^\lambda\|_{L^\infty(D)} + \|S\bar{u}_\lambda - S_h\bar{u}_\lambda\|_{L^\infty(D)} \\ &\quad + \|S\hat{u} - S_h\hat{u}\|_{L^\infty(D)}) \end{aligned} \quad (86)$$

Let us specify the quantities for the problems (13),(33) for the elliptic equation

$$-\Delta y = u \text{ in } \Omega, \quad y = 0 \text{ on } \Gamma$$

where the domain  $\Omega$  is a polygonal (polyhedral) domain or has smooth boundary. Because of the Dirchlet boundary condition we require  $D \subset\subset \Omega$ . Moreover, we assume a standard quasiuniform finite element discretization. The first two terms of (86) are of higher order

$$\begin{aligned} \|S^*S\bar{u}^\lambda - S_h^*S_h\bar{u}^\lambda\|_U &\leq ch^2\|u^\lambda\|_{L^2(\Omega)}, \\ \|S^*y_d - S_h^*y_d\|_U &\leq ch^2\|y_d\|_{L^2(\Omega)}. \end{aligned}$$

The three remaining terms of (86) are responsible for the approximation rate

$$\begin{aligned} \|S\bar{u}_h^\lambda - S_h\bar{u}_h^\lambda\|_{L^\infty(D)} &\leq c|\log h|^2h^2\|\bar{u}_h^\lambda\|_{L^\infty(\Omega)} \\ \|S\bar{u}_\lambda - S_h\bar{u}_\lambda\|_{L^\infty(D)} &\leq c|\log h|^2h^2\|\bar{u}^\lambda\|_{L^\infty(\Omega)} \\ \|S\hat{u} - S_h\hat{u}\|_{L^\infty(D)} &\leq c|\log h|^2h^2\|\hat{u}\|_{L^\infty(\Omega)} \end{aligned}$$

Combining all results we end up with

$$\|\bar{u} - \bar{u}_h^\lambda\|_U \leq c(\sqrt{\lambda} + |\log h|/h)$$

for the case of additional control constraints. Consequently a choice  $\lambda \sim |\log h|^2h^2$  leads to a balanced error contribution in this case. The problem without additional control constraints is more difficult. Then, norm like  $\|\bar{u}^\lambda\|_{L^\infty(\Omega)}$  are not uniformly bounded with respect to  $\lambda$ . In that case one has to deal with weaker error estimates where the corresponding norms of  $\bar{u}^\lambda$  are uniformly bounded with respect to  $\lambda$ .

### 5.3. Full discretization

Now, we will discuss a full discretization. In an abstract setting, we replace the control space  $U$  by an arbitrary finite dimensional control subspace  $U_h$ . The admissible discrete control set is defined by  $U_{ad}^h = U_{ad} \cap U_h$ .

Let us directly estimate the norm  $\|\bar{u} - \bar{u}_h^\lambda\|_U$ . We will only emphasize the key points in the estimation process. Again, one test function can be constructed by

$$u_0 = (1 - \delta)\bar{u}_h^\lambda + \delta\hat{u}. \quad (87)$$

This term can be analyzed similar to the variational discretization concept. The construction of the other test function depends again on the presence of additional control constraints. If additional control constraints are given, then a choice

$$u_h^\lambda = (1 - \sigma)P_h\bar{u} + \sigma P_h\hat{u} \quad (88)$$

is reasonable. The test function  $u_h^\lambda$  has to belong to  $U_h$  and has to satisfy the control constraints. This is reflected by the choice of a suitable projection or interpolation operator  $P_h$ . For piecewise constant controls one can choose  $P_h$  as the  $L^2$ -projection operator to  $U_h$ . A quasi-interpolation operator can be used as  $P_h$  for piecewise linear controls. Similar to the variational discretization, a choice of  $\lambda$  in the size of  $L^\infty(D)$ -error  $\lambda \sim \|S\bar{u} - S_h\bar{u}\|_{L^\infty(D)}$  leads again to a balanced error contribution. For problems without control constraints we need a construction similar to (79).

Next, we will point out a specific feature of the derivation process. Let us recall the inequality (75)

$$(S^*(S\bar{u}_\lambda - y_d) + \nu\bar{u}_\lambda, u_\lambda - \bar{u}) \leq c\|u_\lambda - \bar{u}\|_U$$

for the variational discretization. Proceeding the same way for the full discretization we would get

$$(S_h^*(S_h\bar{u}_h^\lambda - y_d) + \nu\bar{u}_h^\lambda, u_h^\lambda - \bar{u}) \leq c\|u_h^\lambda - \bar{u}\|_U.$$

Using (88), we find

$$\|u_h^\lambda - \bar{u}\|_U \leq \sigma\|P_h\hat{u} - \bar{u}\|_U + (1 - \sigma)\|P_h\bar{u} - \bar{u}\|_U.$$

The first term becomes small because of the factor  $\sigma$ . The second term has bad approximation properties because of the low regularity properties of the control.

A modification of the estimation process yields

$$(S_h^*(S_h\bar{u}_h^\lambda - y_d) + \nu\bar{u}_h^\lambda, P_h\bar{u} - \bar{u}) = (S_h^*(S_h\bar{u}_h^\lambda - y_d), P_h\bar{u} - \bar{u}) + \nu(\bar{u}_h^\lambda, P_h\bar{u} - \bar{u}).$$

We find for the first term

$$(S_h^*(S_h\bar{u}_h^\lambda - y_d), P_h\bar{u} - \bar{u}) \leq c(\|SP_h\bar{u} - S\bar{u}\|_Y + \|S^*S\bar{u}_h^\lambda - S_h^*S_h\bar{u}_h^\lambda\|_U + \|S^*y_d - S_h^*y_d\|_U)$$

and all terms have good approximation properties. Let us assume that  $P_h$  is the  $L^2$ -projection to  $U_h$ . By orthogonality we get

$$\nu(\bar{u}_h^\lambda, P_h\bar{u} - \bar{u}) = 0.$$

and the all problems are solved. However, this choice is possible only for spaces of piecewise constant functions if control constraints are given. For piecewise linear functions a quasiinterpolation operator yields the desired results, see [29]. In the final result the  $L^\infty(D)$ -error dominates again the approximation behavior. For the complete derivation of the results we refer to [18].

#### 5.4. A short note to nonlinear state equations

We already addressed the main difficulties in Section 3. Let us mention that the approach of the last subsections cannot be used, since the admissible sets  $U^0$ ,  $U^\lambda$ , and  $U_h^\lambda$  are not convex.

There are two techniques available to tackle this problem. A first approach works mainly with objective values. Again feasible points were constructed. Then, the difference of objective values is estimated. The desired error estimate can be derived by means of a local quadratic growth condition. This technique was used in [49].

Local quadratic growth is usually shown by a second-order sufficient optimality condition. If the dual variables of the unregularized problems are unique, then the second-order sufficient optimality conditions are also satisfied for regularized problems. This is true for the Moreau-Yosida approach and for the virtual control concept. For both techniques one has uniqueness of dual variables for the regularized problems by construction. This is not the case for the Lavrentiev regularization. Dual variables are not unique for the Lavrentiev regularization if control constraints and mixed constraints are active simultaneously. Thus, separation of strongly active sets is needed to get the corresponding local uniqueness result, see [63]. These papers contain regularization error results of the form (78) for nonlinear problems.

Another possible technique would be to work with the complete optimality system with Lagrange multipliers. However, this approach was used only for state constrained linear-quadratic problems until now.

Our motivation for the regularization of state constrained problems was that the resulting problems can be solved efficiently. This statement is also correct for nonlinear problems. Main issues for a good performance are local convexity properties of the regularized problems and local uniqueness of stationary points (including dual variables). This can be guaranteed by second-order sufficient optimality conditions.

## 6. A brief discussion of further literature

### 6.1. Literature related to control constraints

There are many contributions to finite element analysis for elliptic control problems with constraints on the controls. For an introduction to the basic techniques we refer to the book [71] of Tröltzsch. Falk [30], and Geveci [32] present finite element analysis for piecewise constant approximations of the controls. For semilinear state equations Arada, Casas, and Tröltzsch in [4] present a finite element analysis for piecewise constant discrete controls. Among other things they prove that the sequence  $(u_h)_h$  of discrete controls contains a subsequence converging to a solution  $u$  of the continuous optimal control problem. Assuming certain second order sufficient conditions for  $u$  they are also able to prove optimal error estimates of the form

$$\|u - u_h\| = \mathcal{O}(h) \text{ and } \|u - u_h\|_\infty = \mathcal{O}(h^\lambda),$$

with  $\lambda = 1$  for triangulations of non-negative type, and  $\lambda = 1/2$  in the general case. In [15] these results are extended in that Casas and Tröltzsch prove that every nonsingular local solution  $u$  (i.e. a solution satisfying a second order sufficient condition) locally can be approximated by a sequence  $(u_h)_h$  of discrete controls, also satisfying these error estimates. There are only few results considering uniform estimates. For piecewise linear controls in

the presence of control constraints Meyer and Rösch in [61] for two-dimensional bounded domains with  $C^{1,1}$ -boundary prove the estimate

$$\|u - u_h\|_\infty = \mathcal{O}(h),$$

which seems to be optimal with regard to numerical results reported in [46, Chap. 3], and which is one order less than the approximation order obtained with variational discretization. The same authors in [60] propose post processing for elliptic optimal control problems which in a preliminary step computes a piecewise constant optimal control  $\bar{u}$  and with its help a projected control  $u^P$  through  $u^P = P_{U_{ad}}(-\frac{1}{\alpha}B^*p_h(\bar{u}))$  which then satisfies

$$\|u - u^P\| = \mathcal{O}(h^2).$$

Casas, Mateos and Tröltzsch in [13] present numerical analysis for Neumann boundary control of semilinear elliptic equations and prove the estimate

$$\|u - u_h\|_{L^2(\Gamma)} = \mathcal{O}(h)$$

for piecewise constant control approximations. In [12] Casas and Mateos extend these investigations to piecewise linear, continuous control approximations, and also to variational discrete controls. Requiring a second order sufficient conditions at the continuous solution  $u$  they are able to prove the estimates

$$\|u - u_h\|_{L^2(\Gamma)} = o(h), \text{ and } \|u - \bar{u}_h\|_{L^\infty(\Gamma)} = o(h^{\frac{1}{2}}),$$

for a general class of control problems, where  $u_h$  denotes the piecewise linear, continuous approximation to  $u$ . For variational discrete controls  $u_h^v$  they show the better estimate

$$\|u - u_h^v\|_{L^2(\Gamma)} = \mathcal{O}(h^{\frac{3}{2}-\epsilon}) \quad (\epsilon > 0).$$

Furthermore, they improve their results for objectives which are quadratic w.r.t. the control and obtain

$$\|u - u_h\|_{L^2(\Gamma)} = \mathcal{O}(h^{\frac{3}{2}}), \text{ and } \|u - u_h\|_{L^\infty(\Gamma)} = \mathcal{O}(h).$$

The dependence of the approximation with respect to the largest angle  $\omega$  of a polygonal domain is studied in Mateos and Rösch [54]. This allows to obtain error estimates of the form

$$\|u - u_h\|_{L^2(\Gamma)} = \mathcal{O}(h^\kappa)$$

with  $\kappa > 3/2$  for convex domains ( $\omega < \pi$ ) and  $\kappa > 1$  for concave domains ( $\omega > \pi$ ).

Let us finally recall the contribution [14] of Casas and Raymond to numerical analysis of Dirichlet boundary control, who for two-dimensional convex polygonal domains prove the optimal estimate

$$\|u - u_h\|_{L^2(\Gamma)} \leq Ch^{1-1/q},$$

where  $u_h$  denotes the optimal discrete boundary control which they sought in the space of piecewise linear, continuous finite elements on  $\Gamma$ . Here  $q \geq 2$  depends on the smallest angle of the boundary polygon. May, Rannacher and Vexler study Dirichlet boundary control without control constraints in [55]. They also consider two dimensional convex polygonal domains and among other things provide optimal error estimates in weaker norms. In particular they address

$$\|u - u_h\|_{H^{-1}(\Gamma)} + \|y - y_h\|_{H^{-1/2}(\Omega)} \sim h^{2-2/q}.$$

Vexler in [74] for  $U_{ad} = \{u \in \mathbb{R}^n; a \leq u \leq b\}$  and  $Bu := \sum_{i=1}^n u_i f_i$  with  $f_i \in H^{5/2}(\Gamma)$  provides finite element analysis for Dirichlet boundary control in bounded, two-dimensional polygonal domains. Among other things he in [74, Theorem 3.4] shows that

$$|u - u_h| \leq Ch^2.$$

Error analysis for general two- and three-dimensional curved domains is presented by Deckelnick, Günther and Hinze in [27]. They prove the error bound

$$\|u - \tilde{u}_h\|_{0,\Gamma} + \|y - \tilde{y}_h\|_{0,\Omega} \leq Ch\sqrt{|\log h|},$$

and for piecewise  $O(h^2)$  regular triangulations of two-dimensional domains the superconvergence result

$$\|u - \tilde{u}_h\|_{0,\Gamma} + \|y - \tilde{y}_h\|_{0,\Omega} \leq Ch^{\frac{3}{2}}.$$

Let us shortly comment on a priori error estimates for non-uniform grids. Mesh grading for reentrant corners was investigated by Apel, Rösch, and Winkler [3] and Apel Rösch, and Sirch [2] for optimal approximation error in the  $L^2$ -norm and in the  $L^\infty$ -norm, respectively. A detailed overview on results with non-uniform grids can be found in the paper of Apel and Sirch inside this book.

## 6.2. Literature for (control and) state constraints

To the authors knowledge only few attempts have been made to develop a finite element analysis for elliptic control problems in the presence of control and state constraints. In [10] Casas proves convergence of finite element approximations to optimal control problems for semi-linear elliptic equations with finitely many state constraints. Casas and Mateos extend these results in [11] to a less regular setting for the states and prove convergence of finite element approximations to semi-linear distributed and boundary control problems. In [58] Meyer considers a fully discrete strategy to approximate an elliptic control problem with pointwise state and control constraints. He obtains the approximation order

$$\|\bar{u} - \bar{u}_h\| + \|\bar{y} - \bar{y}_h\|_{H^1} = \mathcal{O}(h^{2-d/2-\epsilon}) \quad (\epsilon > 0),$$

where  $d$  denotes the spatial dimension. His results confirm those obtained by the Deckelnick and Hinze in [21] for the purely state constrained case, and are in accordance with Theorem 4.7. Meyer also considers variational discretization and in the presence of  $L^\infty$  bounds on the controls shows

$$\|\bar{u} - \bar{u}_h\| + \|\bar{y} - \bar{y}_h\|_{H^1} = \mathcal{O}(h^{1-\epsilon} |\log h|) \quad (\epsilon > 0),$$

which is a result of a similar quality as that given in the third part of Theorem 4.7.

Let us comment also on further approaches that tackle optimization problems for pdes with control and state constraints. A *Lavrentiev-type regularization* of problem (13) is investigated by Meyer, Rösch and Tröltzsch in [62]. In this approach the state constraint  $y \leq b$  in (13) is replaced by the mixed constraint  $\epsilon u + y \leq b$ , with  $\epsilon > 0$  denoting a regularization parameter. It turns out that the associated Lagrange multiplier  $\mu_\epsilon$  belongs to  $L^2(\Omega)$ . Numerical analysis for this approach with emphasis on the coupling of gridsize and regularization parameter  $\epsilon$  is presented by Hinze and Meyer in [45]. The resulting optimization problems are solved either by interior-point methods or primal-dual active set strategies, compare the work [59] by Meyer, Prüfert and Tröltzsch.

Hintermüller and Kunisch in [40, 41] consider the Moreau-Yosida relaxation approach to problem classes containing (13). In this approach the state constraint is relaxed in that it is dropped and a  $L^2$  regularization term of the form  $\frac{1}{2\gamma} \int_\Omega |\max(0, \gamma \mathcal{G}(Bu))|^2$  is added to the cost functional instead, where  $\gamma$  denotes the relaxation parameter. Numerical analysis for this approach with emphasis on the coupling of gridsize and relaxation parameter  $\gamma$  is presented by Hintermüller and Hinze in [38].

Schiela in [70] chooses a different way to relax state constraints in considering barrier functionals of the form  $-\mu \int_\Omega \log(-\mathcal{G}(Bu)) dx$  which penalize the state constraints. In [47] he together with Hinze presents numerical analysis for this approach with emphasis on the coupling of gridsize and barrier parameter  $\mu$ .

### 6.3. Gradient constraints

In many practical applications pointwise constraints on the gradient of the state are required, for example if one aims on avoiding large von Mises stresses, see [26] for a discussion. For elliptic optimal control problems with these kind of constraints Deckelnick, Günter and Hinze in [26] propose a mixed finite element approximation for the state combined with variational discretization and prove the error estimate

$$\|u - u_h\| + \|y - y_h\| \leq Ch^{\frac{1}{2}} |\log h|^{\frac{1}{2}},$$

which is valid for two- and three-dimensional spatial domains. The classical finite element approach using piecewise linear, continuous approximations for the states is investigated by Günter and Hinze in [35]. They are able to show the estimates

$$\|y - y_h\| \leq Ch^{\frac{1}{2}(1-\frac{d}{r})}, \text{ and } \|u - u_h\|_{L^r} \leq Ch^{\frac{1}{r}(1-\frac{d}{r})},$$

which are valid for variational discretization as well as for piecewise constant control approximations. Here,  $d = 2, 3$  denotes the space dimension, and  $r > d$  the integration order of the  $L^r$ -control penalization term in the cost functional. Ortner and Wollner in [65] for the same discretization approach obtain similar results adapting the proof technique of [21] to investigate the numerical approximation of elliptic optimal control problems with pointwise bounds on the gradient of the state.

### 6.4. Literature on control of time-dependent problems

In the literature only few contributions to numerical analysis for control problems with time dependent pdes can be found. For unconstrained linear quadratic control problems with the time dependent Stokes equation in two- and three-dimensional domains Deckelnick and Hinze in [19] prove the error bound

$$\|u - u_{h,\sigma}\|_{L^2((0,T)\times\Omega)} = \mathcal{O}(\sigma + h^2).$$

Here and below  $\sigma$  denotes the discretization parameter for the controls. They use a fully implicit variant of Eulers method for the time discretization which is equivalent to the  $dG(0)$  approximation. In space they use Taylor-Hood finite elements. Using [19, (3.1),(3.6)] this estimate directly extends also to the control constrained case.

Boundary control for the heat equation in one spatial dimension is considered by Malanowski in [53] with piecewise constant, and by Rösch in [66] with piecewise linear, continuous control approximations. Requiring strict complementarity for the continuous solution Rösch is able to prove the estimate

$$\|u - u_\sigma\| = \mathcal{O}(\sigma^{\frac{3}{2}}).$$

Malanowski proves the estimate

$$\|u - u_{h,\sigma}\|_{L^2((0,T)\times\Omega)} = \mathcal{O}(\sigma + h),$$

where  $h$  denotes the discretization parameter for the space discretizations.

In a recent work [56, 57] Meidner and Vexler present extensive research on control problems governed by parabolic equations and their discrete approximation based on  $dG(0)$  in time and finite element in space, where they consider the heat equation as mathematical model on a two- or three-dimensional convex polygonal domain. For variational discretization of [44] they prove the estimate

$$\|u - u_{h,\sigma}\|_{L^2((0,T)\times\Omega)} = \mathcal{O}(\sigma + h^2),$$

which under the assumption of strict complementarity of the continuous solution also holds for post-processing [60].

For control problems with nonlinear time dependent equations one only finds few contributions in the literature. In [36, 37] Gunzburger and Manservigi present a numerical approach

to control of the instationary Navier-Stokes equations (3) using the first discretize then optimize approach. The first optimize then discretize approach applied to the same problem class is discussed by Hinze in [43]. Deckelnick and Hinze provide numerical analysis for a general class of control problems with the instationary Navier Stokes system (3) in [20]. Among other things they prove existence and local uniqueness of variational discrete controls in neighborhoods of nonsingular continuous solutions, and for semi-discretization in space with Taylor-Hood finite elements provide the error estimate

$$\int_0^T \|u - u_h\|_U^2 dt \leq Ch^4.$$

Here,  $u, u_h$  denote the continuous and variational discrete optimal control, respectively. This result also carries over to the case of control constraints under the assumptions made in Section 3.

For problems with state constraints only a few contributions are known. Deckelnick and Hinze in [24] investigate variational discretization for parabolic control problems in the presence of state constraints. Among other things they prove an error bound

$$\alpha \|u - u_h\|^2 + \|y - y_h\|^2 \leq C \begin{cases} h\sqrt{|\log h|}, & (d = 2) \\ \sqrt{h}, & (d = 3). \end{cases}$$

under the natural regularity assumption  $y = \mathcal{G}(Bu) \in W = \{v \in C^0([0, T]; H^2), v_t \in L^2(H^1)\}$  with time stepping  $\delta t \sim h^2$ . Exploiting results of Nochetto and Verdi [64] in the case  $d = 2$  and  $Bu \in L^\infty(\Omega_T)$  it seems possible to us that an error bound of the form

$$\alpha \|u - u_h\|^2 + \|y - y_h\|^2 \lesssim C(h^2 + \tau)$$

can be proved.

Very recently Giles and S. Ulbrich [34] considered the numerical approximation of optimal control problems for scalar conservation laws and provided a detailed numerical analysis for the discrete treatment of control problem.

## References

- [1] Agmon, S., Douglis, A., Nirenberg, L.: Estimates near the boundary for solutions of elliptic partial differential equations satisfying general boundary conditions. *Comm. Pure Appl. Math.*, **12**, 623–727 (1959)
- [2] Apel, T., Rösch, A., Sirch, D.:  $L^\infty$ -Error Estimates on Graded Meshes with Application to Optimal Control. *SIAM Journal Control and Optimization*, 48(3): 1771-1796, 2009.
- [3] Apel, T., Rösch, A., Winkler, G: Optimal control in non-convex domains: a priori discretization error estimates. *Calcolo*, 44(3), 137-158 (2007).
- [4] Arada, N., Casas, E., Tröltzsch, F.: Error estimates for the numerical approximation of a semilinear elliptic control problem. *Computational Optimization and Applications* **23**, 201–229 (2002)
- [5] Maïtine Bergounioux, Kazufumi Ito, and Karl Kunisch. Primal-dual strategy for constrained optimal control problems. *SIAM J. Control and Optimization*, 37:1176–1194, 1999.
- [6] Maïtine Bergounioux and Karl Kunisch. Primal-dual strategy for state-constrained optimal control problems. *Computational Optimization and Applications*, 22:193–224, 2002.
- [7] Casas, E.:  $L^2$  estimates for the finite element method for the Dirichlet problem with singular data. *Numer. Math.* **47**, 627–632 (1985)
- [8] Casas, E.: Control of an elliptic problem with pointwise state constraints. *SIAM J. Cont. Optim.* **4**, 1309–1322 (1986)
- [9] Casas, E.: Boundary control of semilinear elliptic equations with pointwise state constraints. *SIAM J. Cont. Optim.* **31**, 993–1006 (1993)

- [10] Casas, E.: Error Estimates for the Numerical Approximation of Semilinear Elliptic Control Problems with Finitely Many State Constraints. *ESAIM, Control Optim. Calc. Var.* **8**, 345–374 (2002)
- [11] Casas, E., Mateos, M.: Uniform convergence of the FEM. Applications to state constrained control problems. *Comp. Appl. Math.* **21**, (2002)
- [12] Casas, E., Mateos, M.: Error Estimates for the Numerical Approximation of Neumann Control Problems. *Comp. Appl. Math.*, to appear
- [13] Casas, E., Mateos, M., Tröltzsch, F.: Error estimates for the numerical approximation of boundary semilinear elliptic control problems. *Comput. Optim. Appl.* **31**, 193–219 (2005)
- [14] Casas, E., Raymond, J.P.: Error estimates for the numerical approximation of Dirichlet Boundary control for semilinear elliptic equations. *SIAM J. Cont. Optim.* **45**, 1586–1611 (2006)
- [15] Casas, E., Tröltzsch, F.: Error estimates for the finite element approximation of a semilinear elliptic control problems. *Contr. Cybern.* **31**, 695–712 (2005)
- [16] E. Casas, M. Mateos, and F. Tröltzsch. Error estimates for the numerical approximation of boundary semilinear elliptic control problems. *Comput. Optim. Appl.*, 31(2):193–219, 2005.
- [17] Svetlana Cherednichenko, Klaus Krumbiegel, and Arnd Rösch. Error estimates for the Lavrentiev regularization of elliptic optimal control problems. *Inverse Problems*, 24(5):055003, 2008.
- [18] Svetlana Cherednichenko and Arnd Rösch. Error estimates for the discretization of elliptic control problems with pointwise control and state constraints. *Computational Optimization and Applications*, 44:27–55, 2009.
- [19] Deckelnick, K., Hinze, M.: Error estimates in space and time for tracking-type control of the instationary Stokes system. *ISNM* **143**, 87–103 (2002)
- [20] Deckelnick, K., Hinze, M.: Semidiscretization And Error Estimates For Distributed Control Of The Instationary Navier-Stokes Equations. *Numer. Math.* **97**, 297–320 (2004)
- [21] Deckelnick, K., Hinze, M.: Convergence of a finite element approximation to a state constrained elliptic control problem. *SIAM J. Numer. Anal.* **45**, 1937–1953 (2007)
- [22] Deckelnick, K., Hinze, M.: A finite element approximation to elliptic control problems in the presence of control and state constraints. *Hamburger Beiträge zur Angewandten Mathematik HBAM2007-01* (2007)
- [23] Deckelnick, K., Hinze, M.: Numerical analysis of a control and state constrained elliptic control problem with piecewise constant control approximations. *Proceedings of the ENUMATH* (2007).
- [24] Deckelnick, K., Hinze, M.: Variational Discretization of Parabolic Control Problems in the Presence of Pointwise State Constraints. Priority Programme 1253, Preprint-Nr.: SPP1253-08-08, to appear in *JCM* (2009).
- [25] Deckelnick, K., Hinze, M.: A note on the approximation of elliptic control problems with bang-bang controls. Priority Programme 1253, Preprint-Nr.: SPP1253-070 (2009).
- [26] Deckelnick, K., Günther, A., Hinze, M.: Finite element approximations of elliptic control problems with constraints on the gradient. *Numer. Math.* 111:335-350 (2009).
- [27] Deckelnick, K., Günther, A., Hinze, M.: Finite Element Approximation of Dirichlet Boundary Control for Elliptic PDEs on Two- and Three-Dimensional Curved Domains. *SIAM J. Cont. Optim.* 48:2798-2819 (2009).
- [28] Deckelnick, K., Hinze, M., Matthes, U., Schiela, A.: Approximation schemes for constrained optimal control with nonlinear pdes. In preparation (2010).
- [29] J.C. de los Reyes, C. Meyer, and B. Vexler. Finite element error analysis for state-constrained optimal control of the Stokes equations. *Control and Cybernetics*, 37(2):251–284, 2008.
- [30] Falk, R.: Approximation of a class of optimal control problems with order of convergence estimates. *J. Math. Anal. Appl.* **44**, 28–47 (1973)
- [31] Gastaldi, L., Nochetto, R.H.: On  $L^\infty$ -accuracy of mixed finite element methods for second order elliptic problems. *Mat. Apl. Comput.* **7**, 13–39 (1988)
- [32] Geveci, T.: On the approximation of the solution of an optimal control problem governed by an elliptic equation. *Math. Model. Numer. Anal.* **13**, 313–328 (1979)
- [33] Gilbarg, D., Trudinger, N.S.: Elliptic partial differential equations of second order (2nd ed.). Springer (1983)



- [34] Giles, M., Ulbrich, S.: Convergence of linearised and adjoint approximations for discontinuous solutions of conservation laws. Part 1: Siam J. Numer. Anal. 48:882-904, Part 2: Siam J. Numer. Anal. 48:905-921 (2010).
- [35] Günther, A., Hinze, M.: Elliptic Control Problems with Gradient Constraints - Variational Discrete Versus Piecewise Constant Controls, Comput. Optim. Appl., DOI: 10.1007/s10589-009-9308-8 (2009).
- [36] Gunzburger, M.D, Manservigi, S.: Analysis and approximation of the velocity tracking problem for Navier-Stokes flows with distributed control. Siam J, Numer. Anal. **37**, 1481–1512, 2000
- [37] Gunzburger, M.D., Manservigi, S.: The velocity tracking problem for Navier-Stokes flows with boundary controls. Siam J. Control and Optimization **39**, 594–634 (2000)
- [38] Michael Hintermüller and Michael Hinze. Moreau-yosida regularization in state constrained elliptic control problems: error estimates and parameter adjustment. *SIAM J. Numerical Analysis*, 47:1666–1683, 2009.
- [39] Michael Hintermüller, Kazufumi Ito, and Karl Kunisch. The primal-dual active set strategy as a semismooth Newton method. *SIAM J. Optimization*, 13:865–888, 2003.
- [40] Michael Hintermüller and Karl Kunisch. Path-following methods for a class of constrained minimization problems in function space. *SIAM J. Optimization*, 17:159–187, 2006.
- [41] Hintermüller, M., Kunisch, K.: Feasible and non-interior path following in constrained minimization with low multiplier regularity. Report, Universität Graz (2005)
- [42] Hintermüller, M., Ulbrich, M.: A mesh independence result for semismooth Newton methods. *Mathematical Programming* 101:151-184, 2004
- [43] Hinze, M.: Optimal and instantaneous control of the instationary Navier-Stokes equations. Habilitationsschrift, Fachbereich Mathematik, Technische Universität Berlin (2000)
- [44] Hinze, M.: A variational discretization concept in control constrained optimization: the linear-quadratic case. *Computational Optimization and Applications* **30**, 45–63 (2005)
- [45] Michael Hinze and Christian Meyer. Variational discretization of Lavrentiev-regularized state constrained elliptic optimal control problems. *Computational Optimization and Applications*, DOI: 10.1007/s10589-008-9198-1, 2010.
- [46] Michael Hinze, René Pinnau, Michael Ulbrich, Stefan Ulbrich. *Optimization with PDE constraints* MMTA 23, Springer 2009.
- [47] Michael Hinze and Anton Schiela. Discretization of interior point methods for state constrained elliptic optimal control problems: optimal error estimates and parameter adjustment. *Computational Optimization and Applications*, DOI: 10.1007/s10589-009-9278-x, 2010.
- [48] Klaus Krumbiegel, Christian Meyer, and Arnd Rösch. A priori error analysis for neumann boundary control problems. *SIAM J. Control and Optimization*. submitted.
- [49] Klaus Krumbiegel, Ira Neitzel, and Arnd Rösch. Regularization error estimates for semilinear elliptic optimal control problems with pointwise state and control constraints. *Computational Optimization and Applications*. submitted.
- [50] Klaus Krumbiegel and Arnd Rösch. On the regularization error of state constrained Neumann control problems. *Control and Cybernetics*, 37:369–392, 2008.
- [51] Klaus Krumbiegel and Arnd Rösch. A virtual control concept for state constrained optimal control problems. *Computational Optimization and Applications*, 43:213–233, 2009.
- [52] Karl Kunisch and Arnd Rösch. Primal-dual active set strategy for a general class of constrained optimal control problems. *SIAM J. Optimization*, 13(2):321–334, 2002.
- [53] Malanowski, K.: Convergence of approximations vs. regularity of solutions for convex, control constrained optimal-control problems. *Appl. Math. Optim.* **8** 69–95 (1981)
- [54] Mateos, M. and Rösch, A. On saturation effects in the Neumann boundary control of elliptic optimal control problems. *Computational Optimization and Applications*, online published, DOI 10.1007/s10589-009-9299-5.
- [55] May, S., Rannacher, R., Vexler, B.: A priori error analysis for the finite element approximation of elliptic Dirichlet boundary control problems. *Proceedings of ENUMATH 2007, Graz* (2008)
- [56] Meidner, D., Vexler, B.: A Priori Error Estimates for Space-Time Finite Element Discretization of Parabolic Optimal Control Problems. Part I: Problems without Control Constraints. *SIAM Journal on Control and Optimization* **47**, 1150–1177 (2008)

- [57] Meidner, D., Vexler, B.: A Priori Error Estimates for Space-Time Finite Element Discretization of Parabolic Optimal Control Problems. Part II: Problems with Control Constraints. *SIAM Journal on Control and Optimization* **47**, 1301–1329 (2008)
- [58] Meyer, C.: Error estimates for the finite element approximation of an elliptic control problem with pointwise constraints on the state and the control, WIAS Preprint 1159 (2006)
- [59] Christian Meyer, Uwe Prüfert, and Fredi Tröltzsch. On two numerical methods for state-constrained elliptic control problems. *Optimization Methods and Software*, 22:871–899, 2007.
- [60] Meyer, C., Rösch, A.: Superconvergence properties of optimal control problems. *SIAM J. Control Optim.* **43**, 970–985 (2004)
- [61] Meyer, C., Rösch, A.:  $L^\infty$ -estimates for approximated optimal control problems. *SIAM J. Control Optim.* **44**, 1636–1649 (2005)
- [62] Christian Meyer, Arnd Rösch, and Fredi Tröltzsch. Optimal control of PDEs with regularized pointwise state constraints. *Computational Optimization and Applications*, 33(2-3):209–228, 2006.
- [63] I. Neitzel and F. Tröltzsch. On regularization methods for the numerical solution of parabolic control problems with pointwise state constraints. *Control and Cybernetics*, 37:1013–1043, 2008.
- [64] Nochetto, R.H., Verdi, C.: *Convergence past singularities for a fully discrete approximation of curvature-driven interfaces*, *SIAM J. Numer. Anal.* **34**, 490–512 (1997).
- [65] Ortner, C., Wollner, W.: *A priori error estimates for optimal control problems with pointwise constraints on the gradient of the state*. DFG Schwerpunktprogramm 1253, Preprint No. SPP1253-071, (2009).
- [66] Rösch, A.: Error estimates for parabolic optimal control problems with control constraints, *Zeitschrift für Analysis und ihre Anwendungen ZAA* **23**, 353376 (2004)
- [67] A. Rösch, D. Wachsmuth: *How to check numerically the sufficient optimality conditions for infinite-dimensional optimization problems*, in *Optimal control of coupled systems of partial differential equations*, ed. K. Kunisch, G. Leugering, J. Sprekels, F. Tröltzsch, pages 297–318, Birkhäuser 2009.
- [68] A. Rösch, D. Wachsmuth: *Numerical verification of optimality conditions*, *SIAM Journal Control and Optimization*, 47, 5 (2008), 2557–2581.
- [69] Schatz, A.H.: Pointwise error estimates and asymptotic error expansion inequalities for the finite element method on irregular grids. I: Global estimates. *Math. Comput.* **67**(223), 877–899 (1998)
- [70] Anton Schiela. Barrier methods for optimal control problems with state constraints. *SIAM J. Optimization*, 20:1002–1031, 2009.
- [71] Tröltzsch, F.: *Optimale Steuerung mit partiellen Differentialgleichungen*. (2005)
- [72] Fredi Tröltzsch and Irwin Yousept. A regularization method for the numerical solution of elliptic boundary control problems with pointwise state constraints. *Computational Optimization and Applications*, 42:43–66, 2009.
- [73] Fredi Tröltzsch and Irwin Yousept. Source representation strategy for optimal boundary control problems with state constraints. *Journal of Analysis and its Applications*, 28:189–203, 2009.
- [74] Vexler, B.: Finite element approximation of elliptic Dirichlet optimal control problems. *Numer. Funct. Anal. Optim.* **28**, 957–975 (2007)