

A NOTE ON A BILEVEL PROBLEM FOR PARAMETER LEARNING FOR OPTIMAL CONTROL PROBLEMS WITH THE WAVE EQUATION

WIEBKE GÜNTHER AND AXEL KRÖNER

ABSTRACT. In this paper we consider a bilevel problem for determining the optimal regularization parameter in an optimal control problem with the linear wave equation transferring results from [Holler, Kunisch, and Barnard, *A bilevel approach for parameter learning in inverse problems, Inverse Problems* 34 (2018) 115012] where a general function space setting and applications to (bilinear) elliptic problems have been addressed. We analyze the well-posedness and derive the optimality conditions for the bilevel problem for the wave equation. Moreover, for given noisy data the numerical performance of the approach to find the regularization parameter is compared for different choices of priors in the Tikhonov regularization term of the lower level problem.

1. INTRODUCTION

In this paper we consider bilevel problems for determining the optimal regularization parameter in optimal control problems with the linear wave equation. Let $\Omega \subset \mathbb{R}^d$, $d = 1, 2, 3$, be an open and bounded subset, $T > 0$, $I := (0, T)$, $Q := \Sigma$, $\Sigma := I \times \partial\Omega$. For initial data $(y_0, y_1) \in H_0^1(\Omega) \times L^2(\Omega)$ and distributed control $u \in L^2(Q)$ the linear wave equation is given by

$$(1.1) \quad y_{tt} - \Delta y = u \quad \text{in } Q, \quad y(0, \cdot) = y_0 \quad \text{in } \Omega, \quad y_t(0, \cdot) = y_1 \quad \text{in } \Omega, \quad y = 0 \quad \text{on } \Sigma.$$

We denote the unique solution of (1.1) associated with control u by $y = y[u]$. For given ground truth $u^\dagger \in L^2(Q)$ and noisy measurements of the exact state $y_{\delta_j} \in L^2(Q)$, $1 \leq j \leq m$, $m \in \mathbb{N}$, the bilevel problem contains the lower level problem

$$(LP) \quad S := \underset{(y, u) \in Y \times U}{\operatorname{argmin}} J(y, u)$$

with

$$(1.2) \quad J(y, u) := \left\{ \alpha \cdot \Psi(u) + \frac{1}{2m} \sum_{j=1}^m \|y[u] - y_{\delta_j}\|_{L^2(Q)}^2 \right\}$$

and is given as

$$(UP) \quad \min_{\substack{\hat{\alpha} \leq \alpha \leq \hat{\alpha}, \\ (y_\alpha, u_\alpha) \in Y \times U}} \|u_\alpha - u^\dagger\|_{L^2(Q)}^2, \quad \text{s.t. } (y_\alpha, u_\alpha) \in S$$

2010 *Mathematics Subject Classification.* 49K20 and 35L05.

Key words and phrases. bilevel optimization and wave equation and multi-penalty regularization and optimal control.

with control space $U = L^2(I; H^k(\Omega))$, $k \in \mathbb{N}$ and state space

$$(1.4) \quad Y := L^2(I; H_0^1(\Omega)) \cap H^1(I; L^2(\Omega)) \cap H^2(I; H^{-1}(\Omega)).$$

The Tikhonov regularization is given by $\Psi(u) := \sum_{i=1}^r \Psi_i(u)$ for penalty functionals

$$(1.5) \quad \Psi_i(u) := \frac{1}{2} \|K_i u\|_{L^2(Q)}^2$$

with so-called priors $K_i \in L(H^k(\Omega), L^2(\Omega))$, $1 \leq i \leq r$, $r \in \mathbb{N}$, and lower and upper bounds in \mathbb{R}^r of the regularization parameter $\alpha \in \mathbb{R}^r$ given by

$$(1.6) \quad 0 < \tilde{\alpha} \leq \hat{\alpha} < \infty$$

where the inequalities are understood component-wise. This bilevel problem consists of a lower level problem (**LP**) determining for given regularization parameter α the optimal control u , and an upper level problem determining the optimal regularization parameter with respect to the ground truth u^\dagger .

There are many publications on the topic of casting parameter learning in regularized optimization problems as a bilevel optimization problem and investigating its solvability and optimality conditions. Examples for treatment of the finite dimensional case can be found in, e.g., Kunisch and Pock [14], and De los Reyes, Schönlieb, Valkonen [4] having applications in imaging. Publications in the infinite case are often concerned with inverse optimal control problems with partial differential equations, see, e.g., Harder and Wachsmuth [8], Holler, Kunisch, and Barnard [11]. Those works deal with inverse optimal control problems governed by (bi-)linear elliptic partial differential equations. For learning nonlocal regularization operators see Holler and Kunisch [10] and algorithms for a bilevel optimal control problems with a non-smooth lower level problem [3]. For problems with varying regularization in a weighted total variation model see Hintermüller et al. [9]. Bilevel problems with non-smooth lower level problems combined with convolutional neural networks are considered, e.g., in Ochs, Ranftl, Brox, and Pock [18].

The contribution of this note is a proof of concept of the ideas developed in [11] for optimal control problems governed by the linear wave equation. The well-posedness of the bilevel problem (**UP**) is shown and necessary optimality conditions are derived. For the latter one we transfer the problem to a single level one. We present several numerical examples considering the effect of different priors in the problem setting.

Finally, we remark that single level optimal control problems for the wave equation have been studied with respect to various aspects, see, e.g., [7, 13, 15, 16, 17].

The paper is organized as follows: In Section 2 we consider the multiple prior case, in Section 3 the problem is discretized and the algorithm is formulated, and in Section 4 numerical examples are presented.

Notation: Let H be Hilbert space. By $L(H)$ we denote the linear bounded mappings on H . The associated duality product is denoted by $\langle \cdot, \cdot \rangle_H$. Throughout the paper we use the standard notation for Lebesgue, Sobolev, and Bochner spaces. The inner product in $L^2(\Omega)$ is denoted by (\cdot, \cdot) and the inner product in $L^2(I; L^2(\Omega))$ by $(\cdot, \cdot)_I$. For Banach space X the associated norm is denoted by $\|\cdot\|_X$. The absolute value is denoted by $|\cdot|$.

2. THE LOWER-LEVEL PROBLEM AND BILEVEL PROBLEM

In this section we state existence results and derive the optimality systems for the lower level and the bilevel problem. Let

$$(2.1) \quad Z := L^2(I, H^{-1}(\Omega)) \times H^{-1}(\Omega) \times L^2(\Omega)$$

identifying $L^2(\Omega)$ with its dual. We introduce $e : Y \times U \rightarrow Z$ given by

$$(2.2) \quad e(y, u) := \int_0^T \langle y_{tt}(t, \cdot), \cdot \rangle_{H_0^1(\Omega)} dt + (\nabla y(t, \cdot), \nabla \cdot)_I - (u(t, \cdot), \cdot)_I \\ + \langle y_0 - y(0, \cdot), \cdot \rangle_{H_0^1(\Omega)} - \langle y_1 - y_t(0, \cdot), \cdot \rangle_{L^2(\Omega)}.$$

For $(y_0, y_1) \in H_0^1(\Omega) \times L^2(\Omega)$ the weak formulation of (1.1) is given by

$$(2.3) \quad e(y, u)(w) = 0 \quad \text{for all } w := (v, q_1, q_2) \in Z^*.$$

Proposition 2.1. Equation (2.3) has for control $u \in L^2(I; L^2(\Omega))$, initial data $(y_0, y_1) \in H_0^1(\Omega) \times L^2(\Omega)$ a unique weak solution $y \in Y$ of (2.3) satisfying the stability estimate

$$(2.4) \quad \|y\|_Y \leq C(\|u\|_{L^2(I; L^2(\Omega))} + \|y_0\|_{H_0^1(\Omega)} + \|y_1\|_{L^2(\Omega)}).$$

Proof. We refer to [5, Theorem 3 and 4, p. 384f]. \square

The solution mapping of (1.1) has the structure (neglecting arguments of G_Q being zero)

$$(2.5) \quad y = G_Q[u] + G_Q[y_0] + G_Q[y_1],$$

with $G_Q[u, y_0, y_1] : L^2(Q) \times H_0^1(\Omega) \times L^2(\Omega) \rightarrow Y$. Each term in (2.5) is linear continuous with respect to u , y_0 , and y_1 , respectively. Since by Aubin–Lions (see [2]) $Y \subset L^2(Q)$ continuous we can view $\tilde{G}_Q := G_Q[\cdot, 0, 0]$ as a linear continuous operator with range in $L^2(Q)$. In the following, we will consider $G := E_Y \tilde{G}_Q$ instead of \tilde{G}_Q where $E_Y : Y \rightarrow L^2(Q)$ denotes the embedding operator. We thus have the operator $G : L^2(Q) \rightarrow L^2(Q)$ defined by $u \mapsto y[u] =: Gu$. Using this embedding has the advantage that the adjoint operator G^* also acts on the space $L^2(Q)$.

We can formulate problem (LP) equivalently as

$$(2.6) \quad \min_{u \in U} \left(\frac{1}{2m} \sum_{j=1}^m \|Gu - w_j\|_{L^2(Q)}^2 + \frac{1}{2} \sum_{i=1}^r \alpha_i \|K_i u\|_{L^2(Q)}^2 \right)$$

with the abbreviation $w_j := -G_Q[y_0] - G_Q[y_1] + y_{\delta_j}$.

2.1. Lower level problem. Existence and the optimality system for the lower level problem is derived in this section. In the following we make the following hypothesis.

Hypothesis 2.2. There exists $c > 0$ such that

$$(2.7) \quad \sum_{i=1}^r \|K_i u\|_{L^2(Q)}^2 \geq c \|u\|_U^2 \quad \text{for all } u \in U.$$

Then, obviously for sequences $(u_n) \subset U$ with $\|u_n\|_U \rightarrow \infty$ ($n \rightarrow \infty$) we have $\sum_{i=1}^r \|K_i u_n\|_{L^2(Q)} \rightarrow \infty$. Thus, we have in particular

$$(2.8) \quad \begin{cases} \Psi_i : U \rightarrow [0, \infty) \text{ for } 1 \leq i \leq r \text{ are weakly lower semi-continuous on } U, \\ \sum_{i=1}^r \Psi_i \text{ is coercive on } U \text{ and} \\ \text{proper on the set of feasible points of the (LP)}. \end{cases}$$

Furthermore, we observe that for $d = 2$ Hypothesis 2.2 is satisfied with equality for

$$(2.9) \quad U = L^2(I; H^1(\Omega)), \quad K_1 = \text{id}, \quad K_2 = \partial_{x_1}, \quad K_3 = \partial_{x_2}$$

which will be considered in the numerical examples later.

Theorem 2.3. *Assume that Hypothesis 2.2 is satisfied. Then the (LP) has a unique solution $(\bar{y}, \bar{u}) \in Y \times U$.*

Proof. This follows by classical arguments using convexity of (LP), see, e.g., [12]. \square

Furthermore, by convexity, the necessary and sufficient optimality conditions for a pair (\bar{y}, \bar{u}) to be optimal for (LP) are given by

$$(2.10) \quad \frac{1}{m} \sum_{j=1}^m G^*(G\bar{u} - w_j) + \sum_{i=1}^r \alpha_i \mathcal{K}_i \bar{u} = 0$$

with the abbreviation $\mathcal{K}_i := K_i^* K_i$, $1 \leq i \leq r$.

We introduce the costate equation for given state $\bar{y} \in Y$ by

$$(2.11) \quad \lambda_{tt} - \Delta \lambda = -\frac{1}{m} \sum_{i=1}^m (\bar{y} - y_{\delta_i}) \text{ in } Q, \quad \lambda = 0 \text{ on } \Sigma, \quad \lambda(T, \cdot) = \lambda_t(T, \cdot) = 0 \text{ in } \Omega$$

whose solution is understood in a weak sense as follows

$$(2.12) \quad \begin{aligned} & \int_0^T (\langle \lambda_{tt}(t, \cdot), v \rangle_{H_0^1(\Omega)} + (\nabla \lambda(t, \cdot), \nabla v(t, \cdot))_{L^2(\Omega)}) dt - \langle \lambda(T, \cdot), q_1 \rangle_{H_0^1(\Omega)} \\ & + \langle \lambda_t(T, \cdot), q_2 \rangle_{L^2(\Omega)} = -\frac{1}{m} \sum_{i=1}^m (\bar{y} - y_{\delta_i}, v)_{L^2(Q)} \quad \text{for all } (v, q_1, q_2) \in Z. \end{aligned}$$

Using the reversibility of the wave equation and Proposition 2.1 we obtain a unique solution $\lambda \in Y$.

Theorem 2.4. *Let $\bar{u} \in U$ be a control with associated state \bar{y} such that (\bar{y}, \bar{u}) solves problem (LP). Then, there exists a $\lambda \in Y$ such that the state equation (1.1), the costate equation (2.11), and the optimality condition*

$$(2.13) \quad \sum_{i=1}^r \alpha_i \mathcal{K}_i \bar{u} - \lambda = 0$$

are satisfied.

Proof. Observe that by Proposition 2.1 the operator $D_y e(\bar{y}, \bar{u}) \in \mathcal{L}(Y, Z)$ has a bounded inverse. Hence, there exists a unique $(\bar{\lambda}, \mu_0, \mu_1) \in Z^*$ such that by Theorem A.4 we have

$$(2.14) \quad \int_0^T (\langle \delta y_{tt}(t, \cdot), \bar{\lambda}(t, \cdot) \rangle_{H_0^1(\Omega)} + (\nabla \delta y(t, \cdot), \nabla \bar{\lambda}(t, \cdot))_{L^2(\Omega)}) dt + \frac{1}{m} \sum_{i=1}^m (\bar{y} - y_{\delta_i}, \delta y)_{L^2(Q)} - \langle \delta y(0, \cdot), \mu_0 \rangle_{H_0^1(\Omega)} + (\delta y_t(0, \cdot), \mu_1)_{L^2(\Omega)} = 0 \quad \text{for all } \delta y \in Y.$$

Using $C_0^\infty(I; V) \subset Y$ is dense in $L^2(I; V)$ we get that assuming $\bar{\lambda} \in Y$, the adjoint equation (2.14) is equivalent to (2.12) which has a unique solution $\bar{\lambda} \in Y$ and which is together with (μ_0, μ_1) the unique adjoint state.

Furthermore, we have

$$(2.15) \quad \sum_{i=1}^r \alpha_i (K_i \bar{u}, K_i u)_{L^2(Q)} - (u, \lambda)_{L^2(Q)} = 0 \quad \text{for all } u \in U = L^2(Q).$$

□

2.2. Upper level problem. In this section the existence of a solution of the bilevel problem is stated and the optimality system for the problem is derived.

Lemma 2.5. *Hypotheses A.1 and A.2 are satisfied.*

Proof. *Hypotheses A.1:* (H1) is obvious. Since

$$(2.16) \quad (0, 0) \in F_{\text{ad}} := \{(y, u) \in Y \times U \mid e(y, u) = 0\}$$

the admissible set is non-empty implying (H2). (H3) follows from linearity of the state equation. (H4) follows directly from the estimate in Proposition 2.3. (H5) and (H6) follow from (2.8). (H7) follows from Hypothesis 2.2 and that norm together with weak convergence is equivalent to strong convergence. (H8) is given by Proposition 2.1.

Hypothesis A.2: (B1) and (B2) are satisfied trivially. (B3) follows from the first part of the proof of Theorem 2.4 and Proposition 2.1. □

Consequently, we derive by Theorem A.3 the following

Corollary 2.6. The bilevel problem (UP) has a solution.

The optimality condition (2.10) for the (LP) is necessary and sufficient. Thus, the bilevel problem (UP) is equivalent to

$$(2.17) \quad \min_{\substack{\bar{\alpha} \leq \alpha \leq \hat{\alpha}, \\ (y_\alpha, u_\alpha) \in Y \times U}} \|u_\alpha - u^\dagger\|_{L^2(Q)}^2, \quad \text{s.t.} \quad (2.10) \text{ holds for } (y_\alpha, u_\alpha).$$

or equivalently,

$$(2.18) \quad \left\{ \begin{array}{l} \min_{\substack{\bar{\alpha} \leq \alpha \leq \hat{\alpha}, \\ (y_\alpha, u_\alpha) \in Y \times U}} \|u_\alpha - u^\dagger\|_{L^2(Q)}^2 \quad \text{s.t.} \\ \lambda_{tt} - \Delta \lambda = -\frac{1}{m} \sum_{i=1}^m (y - y_{\delta_i}), \quad \lambda(T, \cdot) = \lambda_t(T, \cdot) = 0, \quad \lambda|_\Sigma = 0, \\ \sum_{i=1}^r \alpha_i \mathcal{K}_i u_\alpha - \lambda = 0, \\ y_{tt} - \Delta y = u_\alpha, \quad y(0, \cdot) = y_0, \quad y_t(0, \cdot) = y_1, \quad y|_\Sigma = 0. \end{array} \right.$$

Defining the lagrangian $L: [\bar{\alpha}, \hat{\alpha}] \times U \times Y \times Z^* \rightarrow \mathbb{R}$ as

$$(2.19) \quad L(\alpha, u, y, w) = \langle w, e(u, y) \rangle_Z + J_\alpha(y, u).$$

by the linear-quadratic structure of the lower level problem we have immediately the second order condition for some $\eta > 0$ given as

$$(2.20) \quad D_{(y,u)}^2 L(\alpha, u, y, w)[(\delta y, \delta u), (\delta y, \delta u)] \geq \eta \|(y, u)\|_{Y \times U}^2.$$

Now, we can formulate the specific version of Lemma A.5 for this example.

Lemma 2.7. *Let $(\bar{\alpha}, \bar{y}, \bar{u})$ be a solution to (UP) with uniquely determined Lagrange multiplier $w = (\lambda, \mu_0, \mu_1) \in Z^*$ for the (LP). Then, there exists a unique $(p, q, z) \in Y \times U \times Y$ such that*

$$(2.21) \quad (q, \Psi_u(\bar{u}))_{L^2(Q)}(\alpha - \bar{\alpha}) \geq 0, \text{ for all } \alpha \in [\check{\alpha}, \hat{\alpha}],$$

$$(2.22) \quad z_{tt} - \Delta z = -p, \quad z(T, \cdot) = z_t(T, \cdot) = 0,$$

$$(2.23) \quad \bar{u} - u^\dagger + \sum_{i=1}^r \bar{\alpha}_i \mathcal{K}_i q + z = 0,$$

$$(2.24) \quad p_{tt} - \Delta p = -q, \quad p(0, \cdot) = p_t(0, \cdot) = 0.$$

Proof. This is a direct consequence of Lemma A.5 using regularity results from Proposition 2.1 for the solution of the linear wave equation. \square

Remark 2.8. Let $(\bar{\alpha}, \bar{y}, \bar{u})$ be a solution to (UP). Then, we can interpret (2.22)–(2.24) as the optimality system of the following minimization problem

$$(2.25) \quad \begin{cases} \min_{(p,q) \in Y \times U} J(p, q) = \frac{1}{2} \sum_{i=1}^r \alpha_i \|K_i q\|_{L^2(Q)}^2 + \frac{1}{2} \|p\|_{L^2(Q)}^2 + (\bar{u} - u^\dagger, q)_{L^2(Q)} & \text{s.t.} \\ p_{tt} - \Delta p = -q \text{ in } Q, \quad p(0, \cdot) = p_t(0, \cdot) = 0 \text{ in } \Omega, \quad p = 0 \text{ on } \Sigma. \end{cases}$$

This reformulation will be used for the implementation of the upper level problem in the software package *dolphin-adjoint* [6].

3. DISCRETIZATION

The problem is discretized following ideas from [13, Section 5]. The wave equation is written as a first order system in time. We introduce a partition of the time interval I as

$$(3.1) \quad I = 0 \cup I_1 \cup \dots \cup I_M$$

with $I_m = (t_{m-1}, t_m]$ of size k_m and time points

$$(3.2) \quad 0 = t_0 < t_1 < \dots < t_{M-1} < t_M = T.$$

For the spatial discretization we consider two dimensional shape regular meshes; see, e.g., [6]. A mesh consists of triangles K , which constitute a nonoverlapping cover of the computational domain Ω . The corresponding mesh is denoted by $\mathcal{T} = \{K\}$, where we define the discretization parameter h as a cellwise function by setting $h|_K = h_K$ with the diameter h_K of the cell K . On the mesh \mathcal{T}_h we construct conforming finite element spaces $V_h \subset H^1(\Omega)$ and $V_h^0 \subset H_0^1(\Omega)$ in the following standard way:

$$(3.3) \quad V_h := \{v \in H^1(\Omega) | v|_K \in \mathcal{P}^1(K) \text{ for } K \in \mathcal{T}_h\},$$

$$(3.4) \quad V_h^0 := \{v \in H_0^1(\Omega) | v|_K \in \mathcal{P}^1(K) \text{ for } K \in \mathcal{T}_h\}.$$

We define the following space-time finite element ansatz and test spaces:

$$(3.5) \quad X_{kh} := \{v_{kh} \in C(I, V_h) | v_{kh}|_{I_m} \in \mathcal{P}^1(I_m, V_h)\},$$

$$(3.6) \quad X_{kh}^0 := \{v_{kh} \in C(I, V_h^0) | v_{kh}|_{I_m} \in \mathcal{P}^1(I_m, V_h^0)\},$$

$$(3.7) \quad \tilde{X}_{kh} := \{v_{kh} \in L^2(I, V_h) | v_{kh}|_{I_m} \in P_0(I_m, V_h) \text{ and } v_{kh}(0, \cdot) \in V_h\},$$

$$(3.8) \quad \tilde{X}_{kh}^0 := \{v_{kh} \in L^2(I, V_h^0) | v_{kh}|_{I_m} \in \mathcal{P}_0(I_m, V_h^0) \text{ and } v_{kh}(0, \cdot) \in V_h\},$$

where $P^r(I_m, V_h)$ denotes the space of polynomials up to degree r on I_m with values in V_h . Thus, the spaces X_{kh} and X_{kh}^0 consist of piecewise linear and continuous functions in time with values in the usual spatial finite element space, whereas the functions in \tilde{X}_{kh} and \tilde{X}_{kh}^0 are piecewise constant in time and therefore discontinuous. Based on the equivalent formulation of the state equations as first-order systems we introduce the Galerkin finite element formulation of the state equations. We introduce the discrete control and state space

$$(3.9) \quad Y_{kh} := X_{kh} \times X_{kh}, \quad U_{kh} = X_{kh}$$

and the bilinear form $a: X_{kh} \times X_{kh} \times \tilde{X}_{kh} \times \tilde{X}_{kh} \rightarrow \mathbb{R}$ by

$$(3.10) \quad a(y, v) := a(y_1, y_2, v_1, v_2) = (\partial_t y_2, v_1)_I + (\nabla y_1, \nabla v_1)_I \\ + (\partial_t y_1, v_2)_I - (y_2, v_2)_I + (y_2(0, \cdot), v_1(0, \cdot)) - (y_1(0, \cdot), v_2(0, \cdot))$$

with $y = (y_1, y_2)$ and $\xi = (\xi_1, \xi_2)$. The discrete state equation is given as

$$(3.11) \quad a(u_{kh}, v_{kh}) = (u_{kh}, v_{kh})_I + (y_1, v_{kh}^1(0, \cdot)) - (y_0, v_{kh}^2(0, \cdot)) \text{ for all } v_{kh} \in \tilde{X}_{kh}^0 \times \tilde{X}_{kh}$$

defining (by classical arguments) the linear and continuous mapping $U_{kh} \rightarrow Y_{kh}$, $u_{kh} \mapsto y_{kh}$. That means, the first and second component of the state are discretized by continuous piecewise linear finite elements where we impose zero boundary conditions for the first component; the control is discretized by continuous piecewise linear finite elements. The discrete bilevel problem contains the lower level problem

(LP_{discr})

$$S_{kh} := \operatorname{argmin}_{(y_{kh}, u_{kh}) \in Y_{kh} \times U_{kh}} \left\{ \alpha \cdot \Psi(u_{kh}) + \frac{1}{2m} \sum_{j=1}^m \|y_{kh}[u_{kh}] - y_{\delta_j}\|_{L^2(Q)}^2 \right\}$$

and is given as

$$(UP_{discr}) \quad \min_{\tilde{\alpha} \leq \alpha \leq \hat{\alpha}, (y_{kh}^\alpha, u_{kh}^\alpha) \in Y_{kh} \times U_{kh}} \|u_{kh}^\alpha - u^\dagger\|_{L^2(Q)}^2, \quad \text{s.t. } (y_{kh}^\alpha, u_{kh}^\alpha) \in S_{kh}$$

As a time stepping scheme a Crank-Nicolson scheme is applied. The problem is implemented in *dolfin-adjoint* [6] based on *FEniCS* [1].

For solving problem (UP_{discr}) we use a gradient descent scheme with dynamic stepsize adaption. By avoiding an Armijo backtracking line search the computational effort can be strongly reduced which is an important issue for this bilevel problem involving time-depending partial differential equations. For details see Algorithm 1.

Data: Initialize $\check{\alpha}$ and $\hat{\alpha}$ in \mathbb{R} , $\alpha_0 \in [\check{\alpha}, \hat{\alpha}]$, the step size $\tau := \bar{\tau} > 0$, $\text{tol}_\tau \ll 1$ tolerance $\text{tol} > 0$, and $k := 0$.

Compute u_0 by solving the lower-level problem (LP) with $\alpha = \alpha_0$.

Calculate the corresponding Lagrange multiplier (p_0, q_0, z_0) by solving the optimality system (2.22)–(2.24) (i.e. by solving (2.25)).

Set

$$(3.13) \quad g_k^{(i)} := \langle K_i q_k, K_i u_k \rangle_{L^2(Q)} \quad \text{for } i = 1, \dots, r.$$

while $|g_k|^2 > \text{tol}$ and $\tau > \text{tol}_\tau$ **do**

Set

$$(3.14) \quad \alpha := \alpha - \tau \frac{g_k}{|g_k|}.$$

Compute u_α by solving the lower-level problem (LP) for current α .

if $\|u_\alpha - u^\dagger\|_{L^2(Q)}^2 < \|u_k - u^\dagger\|_{L^2(Q)}^2$ **then**

 Compute the corresponding Lagrange multiplier (p_k, q_k, z_k) via (2.22)–(2.24) and the gradient g_k by (3.13).

 Set $\alpha_k = \alpha$ and $u_k := u_\alpha$.

 Set $k = k + 1$.

else

 Set $\tau = 0.5\tau$.

end

end

Algorithm 1: Solver for bilevel problem.

4. NUMERICAL RESULTS

Set $\Omega = (0, 1) \times (0, 1)$, $T = 1.5$, i.e. $Q = (0, 1.5) \times (0, 1)^2$, $m = 1$, and let the state y^\dagger corresponding to the exact control be given for $(t, x_1, x_2) \in Q$ by

$$u^\dagger(t, x_1, x_2) := \begin{cases} 2.5t(\sin^4(4\pi x_1)) + x_2^2 & \text{if } |x_1 - 0.5| < 0.2 \text{ and } |x_2 - 0.5| < 0.2, \\ 2.5 & \text{else.} \end{cases}$$

We choose this control for analyzing the performance of the parameter learning algorithm, since it exhibits a lot of distinct structure. Furthermore, we set $\bar{\tau} = 10^{-4}$, $\text{tol}_\tau = 10^{-10}$, $\check{\alpha}_i = 10^{-13}$, and $\hat{\alpha}_i = 10^{-2}$ for $i = 1$ (one prior) and $i = 1, 2, 3$ (three priors), respectively. The exact control and corresponding state (with $y_0 = 0$) are shown in Figure 1 and Figure 2. The problem is discretized on a uniform 64×64 spatial mesh and a temporal mesh with 50 timesteps. Noisy data is generated in each time step via pointwise setting

$$(4.1) \quad y_\delta(t, \cdot) := y^\dagger(t, \cdot) + \varepsilon \xi,$$

where $\xi \in \mathcal{N}(0, 1)$ is a standard normal distributed random variable and we set

$$(4.2) \quad \varepsilon := \epsilon \|y^\dagger\|_{L^\infty(Q)}$$

with the relative noise level $\epsilon > 0$. The noisy data for a noise level of $\epsilon = 0.1$ is shown in Figure 3. In the experiments we use a fixed seed for random number generation to ensure comparability. We consider the case with one prior $K_1 := \text{id}$ as well as the case of three priors with additional

$$K_2 := \partial_{x_1}, \quad K_3 := \partial_{x_2}.$$

For the case of only one prior K_1 the value of the cost functional of the higher-level problem, which is the squared L^2 -distance between the exact control and the recovered control, in dependence on the regularization parameter α is shown in Figure 4. The algorithm is initialized with $\alpha_0 = 2 \cdot 10^{-6}$.

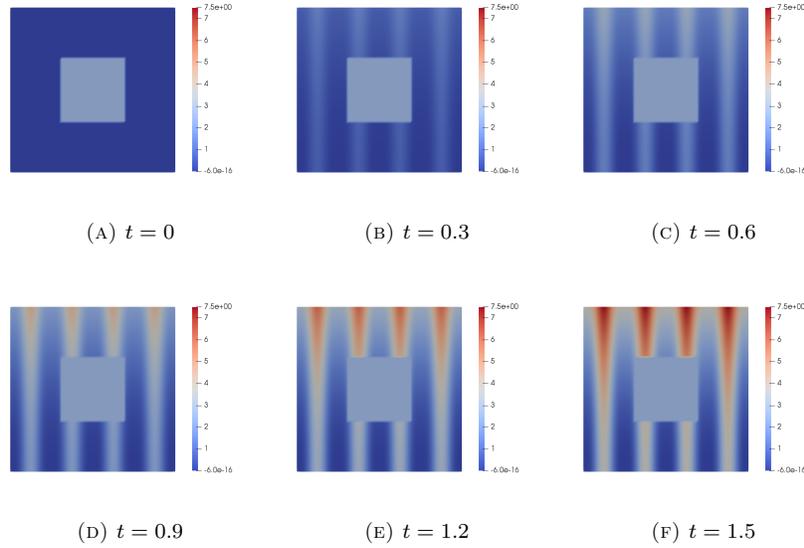


FIGURE 1. Exact control u^\dagger plotted at several timepoints.

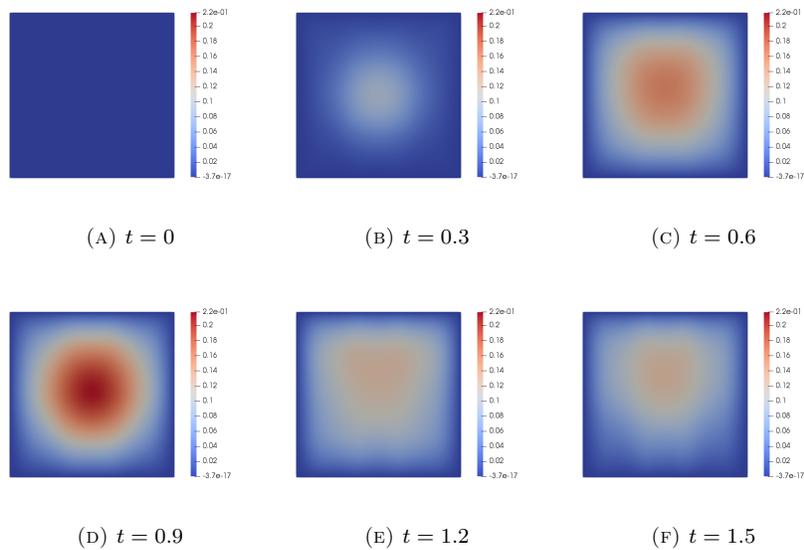


FIGURE 2. Exact state y^\dagger plotted at several timepoints.

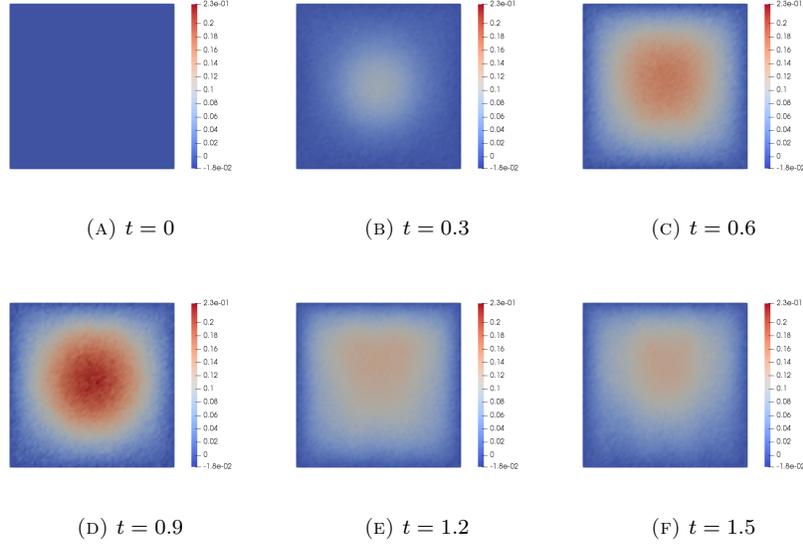


FIGURE 3. Noisy data y_s plotted at several timepoints with a noise level of 10%.

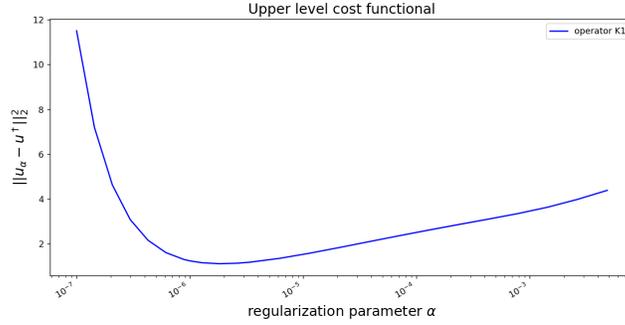


FIGURE 4. Optimal cost of the bilinear problem in dependence on the regularization parameter α with one operator $K_1 = \text{id}$.

Used operators	(Locally) optimal α^*	Error $\ u_{\alpha^*} - u^\dagger\ _{L^2(Q)}^2$
K_1	$1.91 \cdot 10^{-6}$	1.105
(K_1, K_2, K_3)	$(9.81 \cdot 10^{-9}, 1.67 \cdot 10^{-10}, 9.99 \cdot 10^{-7})$	0.586

TABLE 1. Optimal α^* computed for different combinations of regularization operators.

Table 1 shows the (locally) optimal α^* computed for the single operator K_1 and the set of three operators K_i , $I = 1, 2, 3$ for a noise level of $\epsilon = 0.1$. In the latter case the algorithm is initialized with $\alpha_0 = (10^{-8}, 10^{-8}, 10^{-6})$. In case of only one prior

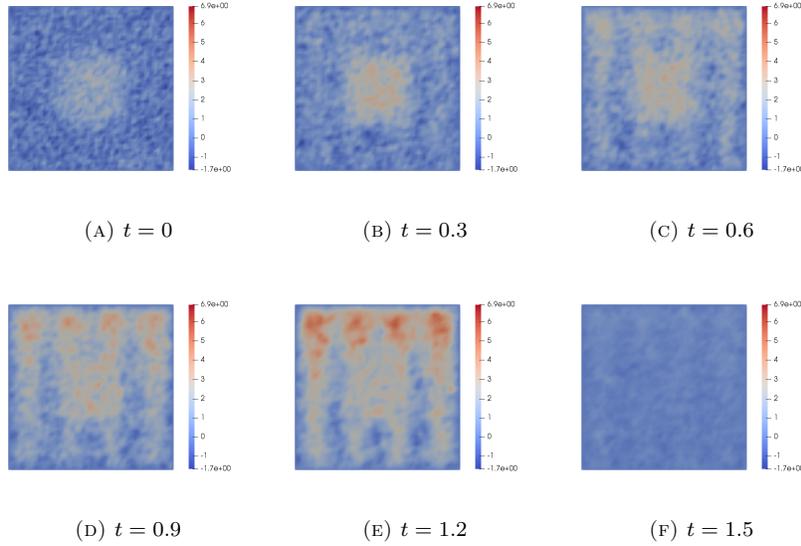


FIGURE 5. Reconstructed control u_α plotted at several timepoints using the optimal (as determined by the algorithm) parameter $\alpha = 1.91 \cdot 10^{-6}$ with regularization operator $K_1 = \text{id}$, $\|u_{\alpha^*} - u^\dagger\|_{L^2(Q)}^2 = 1.105$.

K_1 the reconstructed control for the optimal numerically determined regularization parameter α is shown in Figure 5, for comparison in Figure 6 the control for a suboptimal regularization parameter is shown. The numerical results confirm the expected behaviour. For the case of three priors K_1 , K_2 , and K_3 the control for the optimal numerically determined parameter α is given in Figure 7. In comparison to only one prior it shows a better reconstruction of the vertical elements but a slightly worse one of the square in the middle. Although not covered by the theory (Hypothesis 2.2 is no longer satisfied) we consider numerically the case of only one single prior K_2 with here $\alpha_0 = 10^{-3}$. The resulting sub-optimal and optimal controls can be seen in Figures 8–10 showing again, that the numerically optimal regularization parameter α leads to the best reconstruction of the control.

APPENDIX A. GENERAL SETTING

We recall the general setting from Holler et al. [11]: The lower level problem is given by

$$(P_{\alpha, y_\delta}) \quad \begin{cases} \min_{(y, u) \in Y \times U} \mathcal{I}_{\alpha, y_\delta}(y, u) = \frac{1}{2m} \sum_{j=1}^m \|y - y_{\delta_j}\|_{\tilde{Y}}^2 + \alpha \cdot \Psi(u) \text{ subject to} \\ e(y, u) = 0, \end{cases}$$

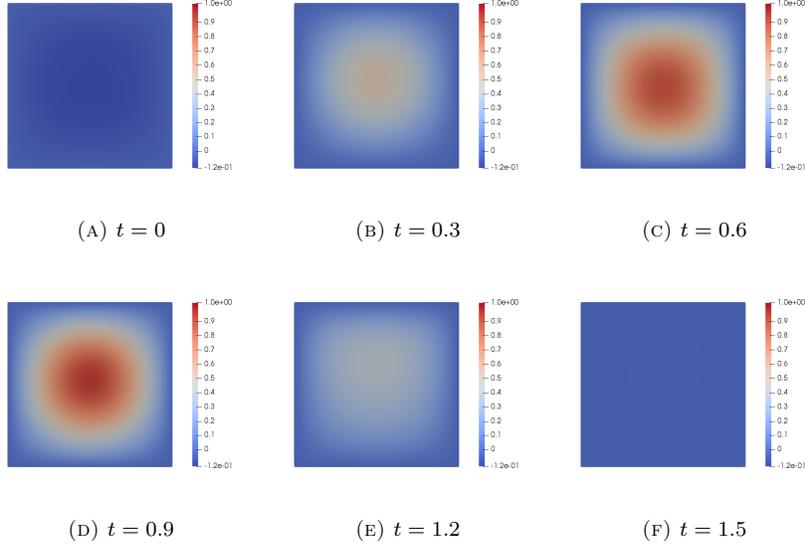


FIGURE 6. Reconstructed control u_α plotted at several timepoints using the sub-optimal (too big) parameter $\alpha = 10^{-2}$ with regularization operator $K_1 = \text{id}$, $\|u_{\alpha^*} - u^\dagger\|_{L^2(Q)}^2 = 4.864$.

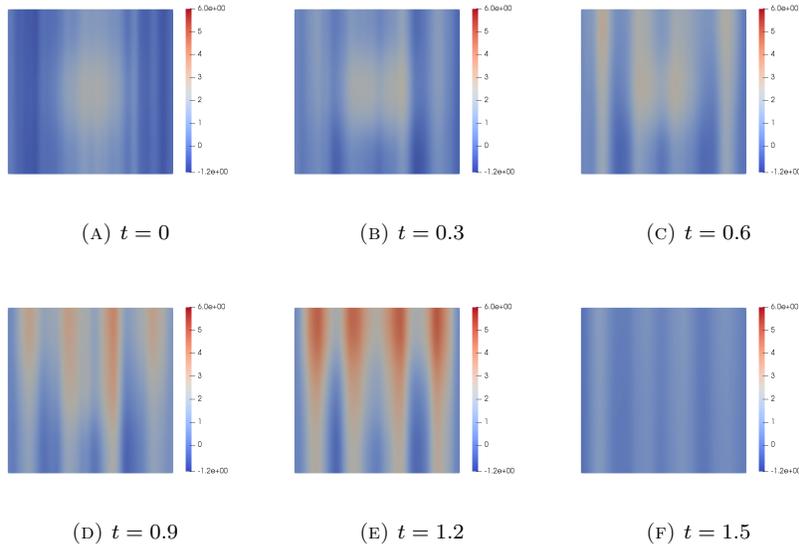


FIGURE 7. Reconstructed control u_α plotted at several timepoints using the parameters $\alpha = (9.81 \cdot 10^{-9}, 1.67 \cdot 10^{-10}, 9.99 \cdot 10^{-7})$ with regularization operators K_1, K_2, K_3 , $\|u_{\alpha^*} - u^\dagger\|_{L^2(Q)}^2 = 0.586$.

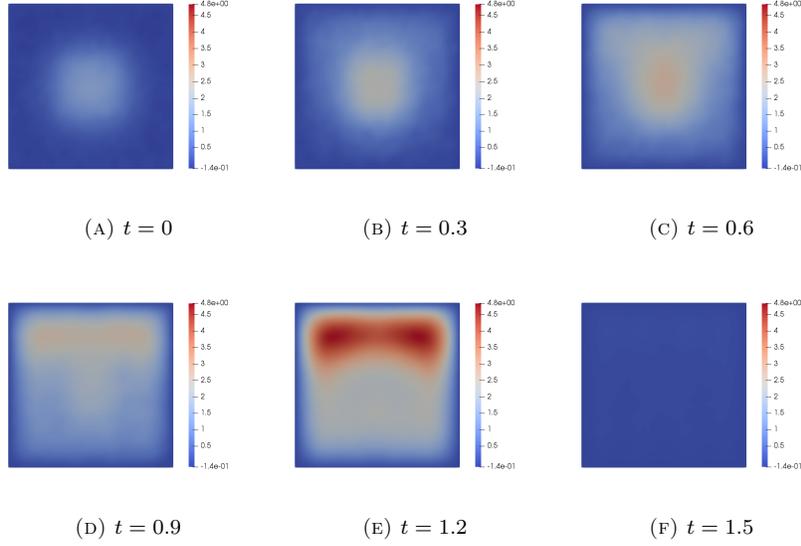


FIGURE 8. Reconstructed control u_α plotted at several timepoints using the sub-optimal (too small) parameter $\alpha = 1 \cdot 10^{-8}$ with regularization operator $K_3 = \partial_{x_2}$, $\|u_{\alpha^*} - u^\dagger\|_{L^2(Q)}^2 = 2.15$.

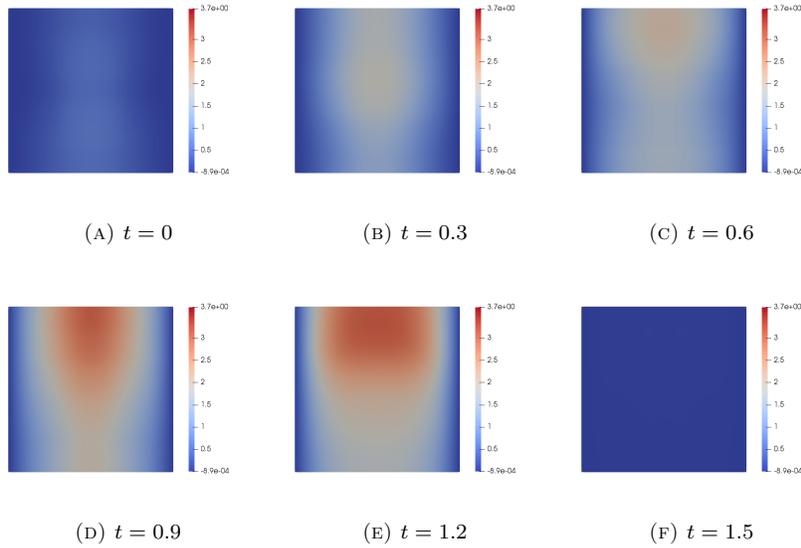


FIGURE 9. Reconstructed control u_α plotted at several timepoints using the sub-optimal (too big) parameter $\alpha = 1.96 \cdot 10^{-5}$ with regularization operator $K_3 = \partial_{x_2}$, $\|u_{\alpha^*} - u^\dagger\|_{L^2(Q)}^2 = 2.494$.

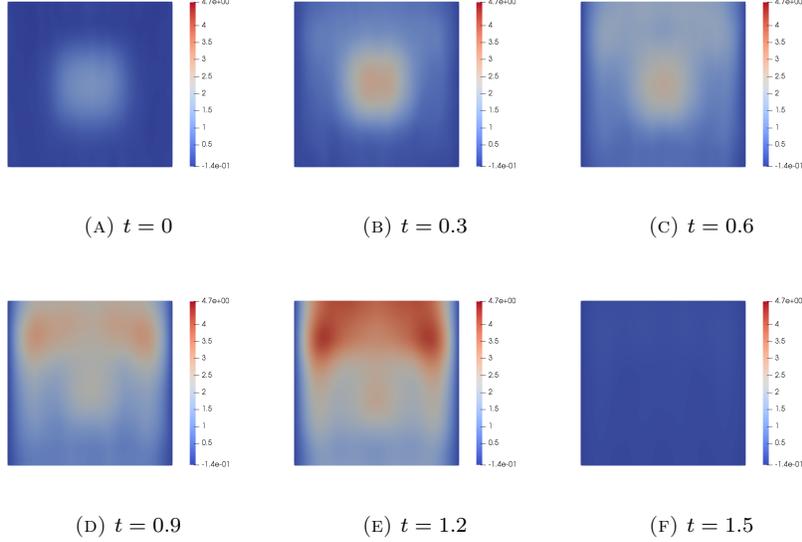


FIGURE 10. Reconstructed control u^* plotted at several timepoints using the learned parameter $\alpha^* = 3.28 \cdot 10^{-7}$ with regularization operator $K_3 = \partial_{x_2}$, $\|u_{\alpha^*} - u^\dagger\|_{L^2(Q)}^2 = 1.46289$.

and the upper level problem

$$(UP) \quad \begin{cases} \min_{\alpha \in [\underline{\alpha}, \bar{\alpha}], (y_\alpha, u_\alpha) \in U} \|u_\alpha - u^\dagger\|_{\tilde{U}}^2 & \text{subject to} \\ u_\alpha \text{ is minimizer of } (P_{\alpha, y_\delta}). \end{cases}$$

with a reflexive Banach space U , Y is a reflexive Banach space, \tilde{U} and \tilde{Y} are Hilbert spaces with $U \subset \tilde{U}$, $Y \subset \tilde{Y}$ continuous, $u^\dagger \in \tilde{U}$ is the ground truth control, and $y_{\delta_j} \in \tilde{Y}$, $1 \leq j \leq m$, are noisy measurements of the ground truth state, $e: Y \times U \rightarrow Z$ represents equality constraints in a reflexive Banach space Z , $\Psi_i: U \rightarrow [0, \infty]$, $1 \leq i \leq r$, are penalty functionals, and α , $\bar{\alpha}$, and $\hat{\alpha}$ as in (1.6).

Hypothesis A.1. (H1) The feasible control set U is closed and convex.

(H2) The feasible set of the lower level problem F_{ad} (defined as in (2.16)) is non-empty.

(H3) For every sequence (y_n, u_n) in $Y \times U$ and $(\bar{y}, \bar{u}) \in Y \times U$ such that $e(y_n, u_n) = 0$ for all $n \in \mathbb{N}$, and $(y_n, u_n) \rightharpoonup (\bar{y}, \bar{u})$ it follows that $e(\bar{y}, \bar{u}) = 0$.

(H4) For every sequence (y_n, u_n) in F_{ad} it holds that if (u_n) is bounded in U , then (y_n) is bounded in Y .

(H5) The function $\sum_{i=1}^r \Psi_i$ is coercive on U and proper on F_{ad} .

(H6) The penalty functionals Ψ_i , $1 \leq i \leq r$, are weakly lower semi-continuous on U .

(H7) For every sequence (u_n) in U and $u \in U$ it holds that, if $u_n \rightharpoonup u$ and $\Psi(u_n) \rightarrow \Psi(u)$, then it follows that $u_n \rightarrow u$.

(H8) For each $u \in U$ there exists a unique $y[u] \in Y$ such that $e(y[u], u) = 0$, and the mapping $u \mapsto y[u]$ is continuous from U to Y .

Hypothesis A.2. (B1) The state equation e is twice continuously Fréchet differentiable on $Y \times U$.

(B2) The penalty function Ψ is twice continuously Fréchet differentiable on U .

(B3) For each $u \in U$ there is a unique $y \in Y$ such that

$$(A.1) \quad e(y, u) = 0.$$

Moreover, $e_y(y, u) \in L(Y, Z)$ is bijective for all $(y, u) \in Y \times U$ satisfying the equation (A.1).

Theorem A.3. [11, Theorem 4.1] *Let Hypotheses A.1 (H1)–(H6) hold. Then the problem (UP) has a solution.*

Theorem A.4. ([11, Sec. 5.1] and [19]) *Let Hypothesis A.2 be satisfied. Furthermore, let (\bar{y}, \bar{u}) be a solution to (P_{α, y_δ}) such that $e_y(\bar{y}, \bar{u})$ is bijective. Then there exists a unique $\bar{\lambda} \in Z^*$ such that $(\bar{y}, \bar{u}, \bar{\lambda})$ is a KKT point of (P_{α, y_δ}) . In particular, (\bar{y}, \bar{u}) satisfies the first order necessary optimality conditions, i.e. there exists $\bar{\lambda} \in Z^*$ such that*

$$(A.2) \quad \begin{aligned} \bar{y} - \bar{y}_\delta + \bar{\lambda} e_y(\bar{y}, \bar{u}) &= 0, \\ \alpha \cdot \Psi_u(\bar{u}) + \bar{\lambda} e_u(\bar{y}, \bar{u}) &= 0, \\ e(\bar{y}, \bar{u}) &= 0, \end{aligned}$$

where $\bar{y}_\delta := \frac{1}{m} \sum_{j=1}^m y_{\delta_j}$.

Theorem A.5. [11, Lemma 5.1] *Let Hypothesis A.1 (1)–(8) and Hypothesis A.2 be satisfied. Let $(\bar{\alpha}, \bar{y}, \bar{u}, \bar{\lambda})$ be a local solution to (UP), and the second-order condition (2.20) be satisfied in (\bar{y}, \bar{u}) . Then there exists a unique $(p, q, z) \in Y \times U \times Z^*$ such that*

$$(A.3) \quad \begin{aligned} \langle \Psi_u(\bar{u})q, \alpha - \bar{\alpha} \rangle_U &\geq 0, \text{ for all } \alpha \in [\underline{\alpha}, \bar{\alpha}], \\ p + \bar{\lambda} e_{yy}(\bar{y}, \bar{u})p + \bar{\lambda} e_{yu}(\bar{y}, \bar{u})q + z e_y(\bar{y}, \bar{u}) &= 0, \\ \bar{u} - u^\dagger + \bar{\lambda} e_{uy}(\bar{y}, \bar{u})p + \bar{\alpha} \cdot \Psi_{uu}(\bar{u})q + \bar{\lambda} e_{uu}(\bar{y}, \bar{u})q + z e_u(\bar{y}, \bar{u}) &= 0, \\ e_y(\bar{y}, \bar{u})p + e_u(\bar{y}, \bar{u})q &= 0. \end{aligned}$$

Proof. By [11, Cor. 3.1] we have stability with respect to data and regularization parameter, and by [11, Lem. 5.1] we derive the statement of the theorem. Thereby, we use the fact that the necessary condition is also sufficient for optimality for the lower level problem. \square

REFERENCES

1. M. Alnæs, J. Blechta, J. Hake, A. Johansson, B. Kehlet, A. Logg, C. Richardson, J. Ring, M. Rognes, and G. Wells, *The FEniCS Project Version 1.5*, Archive of Numerical Software **3** (2015), no. 100.
2. J.-P. Aubin, *Un théorème de compacité*, C. R. Acad. Sci. Paris **256** (1963), 5042–5044.
3. C. Christof, *Gradient-based solution algorithms for a class of bilevel optimization and optimal control problems with a non-smooth lower level*, SIAM: Journal on Optimization **30** (2020), no. 1, 290–318.
4. J. C. De los Reyes, C.-B. Schönlieb, and T. Valkonen, *Bilevel parameter learning for higher-order total variation regularisation models*, J. Math. Imaging Vision **57** (2017), no. 1, 1–25. MR 3592840
5. L.C. Evans, *Partial differential equations*, Amer. Math Soc., Providence, RI, 1998, Graduate Studies in Mathematics 19.

6. P. E. Farrell, D. A. Ham, S. W. Funke, and M. E. Rognes, *Automated derivation of the adjoint of high-level transient finite element programs*, SIAM Journal on Scientific Computing **35** (2013), no. 4, C369–C393.
7. M. Gugat, A. Keimer, and G. Leugering, *Optimal distributed control of the wave equation subject to state constraints*, ZAMM - Journal of Applied Mathematics and Mechanics / Zeitschrift für Angewandte Mathematik und Mechanik **89** (2009), no. 6, 420–444.
8. F. Harder and G. Wachsmuth, *Optimality conditions for a class of inverse optimal control problems with partial differential equations*, Optimization **68** (2019), no. 2-3, 615–643.
9. M. Hintermüller, C.N. Rautenberg, T. Wu, and A. Langer, *Optimal selection of the regularization function in a weighted total variation model. Part II: Algorithm, its analysis and numerical tests*, J. Math. Imaging Vision **59** (2017), no. 3, 515–533. MR 3712428
10. G. Holler and K. Kunisch, *Learning nonlocal regularization operators*, preprint (2020).
11. G. Holler, K. Kunisch, and R.C Barnard, *A bilevel approach for parameter learning in inverse problems*, Inverse Problems **34** (2018), no. 11, 115012.
12. J. Jahn, *Introduction to the theory of nonlinear optimization*, Mathematical Methods of Operations Research-ZOR **44** (1996), no. 3, 291–291.
13. A. Kröner, K. Kunisch, and B. Vexler, *Semismooth Newton methods for optimal control of the wave equation with control constraints*, SIAM J. Control Optim. **49** (2011), no. 2, 830–858.
14. K. Kunisch and T. Pock, *A bilevel optimization approach for parameter learning in variational models*, SIAM Journal on Imaging Sciences **6** (2013), no. 2, 938–983.
15. K. Kunisch, P. Trautmann, and B. Vexler, *Optimal control of the undamped linear wave equation with measure valued controls*, SIAM Journal on Control and Optimization **54** (2016), no. 3, 1212–1244.
16. K. Kunisch and D. Wachsmuth, *On time optimal control of the wave equation, its regularization and optimality system*, ESAIM: Control, Optimisation and Calculus of Variations **19** (2013), no. 2, 317336.
17. J.-L. Lions, *Optimal control of systems governed by partial differential equations.*, Translated from the French by S. K. Mitter. Die Grundlehren der mathematischen Wissenschaften, Band 170, Springer-Verlag, New York, 1971.
18. P. Ochs, R. Ranftl, T. Brox, and T. Pock, *Techniques for gradient-based bilevel optimization with non-smooth lower level problems*, J. Math. Imaging Vision **56** (2016), no. 2, 175–194. MR 3535017
19. J. Zowe and S. Kurcyusz, *Regularity and stability for the mathematical programming problem in banach spaces*, Applied mathematics and Optimization **5** (1979), no. 1, 49–62.

INSTITUT FÜR MATHEMATIK, HUMBOLDT-UNIVERSITÄT ZU BERLIN, 10099 BERLIN, GERMANY
E-mail address: guenther@math.hu-berlin.de

WEIERSTRASS INSTITUTE FOR APPLIED ANALYSIS AND STOCHASTICS, MOHRENSTR. 39, 10117 BERLIN, GERMANY
E-mail address: axel.kroener@wias-berlin.de