

Numerical Analysis of Optimality-System POD for Constrained Optimal Control

Eva Grimm, Martin Gubisch, and Stefan Volkwein

Abstract In this work linear-quadratic optimal control problems for parabolic equations with control and state constraints are considered. Utilizing a Lavrentiev regularization we obtain a linear-quadratic optimal control problem with mixed control-state constraints. For the numerical solution a Galerkin discretization is applied utilizing proper orthogonal decomposition (POD). Based on a perturbation method it is determined by a-posteriori error analysis how far the suboptimal control, computed on the basis of the POD method, is from the (unknown) exact one. POD basis updates are computed by optimality-system POD. Numerical examples illustrate the theoretical results for control and state constrained optimal control problems.

1 Introduction

In this paper we consider a certain class of linear-quadratic optimal control problems governed by linear evolution equations together with control and state constraints. Such linear-quadratic problems are especially interesting as they occur for example as subproblems in each step of sequential quadratic programming (SQP) methods for solving nonlinear problems. For the numerical solution we apply a Galerkin approximation, which is based on proper orthogonal decomposition (POD), a method for deriving reduced-order models of dynamical systems; see [7, 11, 19], for instance. In order to ensure that the POD suboptimal solutions are sufficiently accurate, we derive an a-posteriori error estimate for the difference between the exact (unknown) optimal control and its suboptimal POD approximations. The proof relies on a perturbation argument [5] and extends the results of [8, 22, 25].

However, to obtain the state data underlying the POD reduced order model, it is necessary to solve once the full state system and consequently the POD approximations depend on the chosen parameters for this solve. To be more precise,

E. Grimm • M. Gubisch • S. Volkwein (✉)

Department of Mathematics and Statistics, University of Konstanz, Universitätsstraße 10, 78457 Konstanz, Germany

e-mail: Stefan.Volkwein@uni-konstanz.de; martin.gubisch@uni-konstanz.de

© Springer International Publishing Switzerland 2015

M. Mehl et al. (eds.), *Recent Trends in Computational Engineering - CE2014*,

Lecture Notes in Computational Science and Engineering 105,

DOI 10.1007/978-3-319-22997-3_18

the choice of an initial control turned out to be essential. When using an arbitrary control, the obtained accuracy was not at all satisfying even when using a huge number of basis functions whereas an optimal POD basis (computed from the FE optimally controlled state) led to far better results. To overcome this problem different techniques for improving the POD basis have been proposed. Here, we will apply the so called optimality system POD (OS-POD) introduced in [17]. The idea of OS-POD is straightforward: include the equations determining the POD basis in the optimization process. A thereby obtained basis would be optimal for the considered problem. We follow the ideas in [6, 26], where OS-POD is combined efficiently with an a-posteriori error estimation to compute a better initializing control. The POD basis is then determined from this control and the a-posteriori error estimate ensures that the optimal control problem is solved up to a desired accuracy. Let us refer to [1] where the trust-region POD method is introduced as a different update strategy for the POD basis.

The paper is organized in the following manner: In Sect. 2 we introduce our optimal control problem with control and state constraints. To deal numerically with the state constraints a Lavrentiev regularization is utilized in Sect. 3. The POD method is explained briefly in Sect. 4. In Sect. 5 the existing a-posteriori error analysis is extended to our state-constrained control problem. The combination of the a-posteriori error estimation and OS-POD is explained in Sect. 6. In Sect. 7 we propose two algorithms to solve the reduced optimal control problem. Numerical examples are presented in Sect. 8.

2 The State-Constrained Optimal Control Problem

Suppose that $\Omega \subset \mathbb{R}^d$, $d \in \{1, 2, 3\}$, is an open and bounded domain with Lipschitz-continuous boundary $\Gamma = \partial\Omega$. Let V be a Hilbert space with $H_0^1(\Omega) \subset V \subset H^1(\Omega)$. We endow the Hilbert spaces $H = L^2(\Omega)$ and V with the usual inner products

$$\langle \varphi, \psi \rangle_H = \int_{\Omega} \varphi \psi \, \mathbf{d}\mathbf{x}, \quad \langle \varphi, \psi \rangle_V = \int_{\Omega} \varphi \psi + \nabla \varphi \cdot \nabla \psi \, \mathbf{d}\mathbf{x}$$

Let $T > 0$ be the final time. We introduce a continuous bilinear form $a(\cdot, \cdot) : V \times V \rightarrow \mathbb{R}$ satisfying

$$a(\varphi, \varphi) \geq \alpha_1 \|\varphi\|_V^2 - \alpha_2 \|\varphi\|_H^2 \quad \text{for all } \varphi \in V$$

for constants $\alpha_1 > 0$ and $\alpha_2 \geq 0$. Let us mention that the results can be extended easily to time-dependent bilinear forms in a straightforward way. Recall the Hilbert space $W(0, T) = \{\varphi \in L^2(0, T; V) \mid \varphi_t \in L^2(0, T; V')\}$ endowed with the common inner product [4, pp. 472–479]. Let \mathcal{D} be a bounded subset of \mathbb{R}^d with $d \in \mathbb{N}$. Then the control space is given by $\mathcal{U} = L^2(\mathcal{D}; \mathbb{R}^m)$ for $m \in \mathbb{N}$. By $\mathcal{U}_{\text{ad}} \subset \mathcal{U}$ we define

the closed, convex and bounded subset $\mathcal{U}_{\text{ad}} = \{u \in \mathcal{U} \mid u_a \leq u \leq u_b \text{ in } \mathcal{U}\}$, where $u_a, u_b \in \mathcal{U}$ holds with $u_a \leq u_b$. In particular, we identify \mathcal{U} with its dual space \mathcal{U}' . For $u \in \mathcal{U}_{\text{ad}}, y_\circ \in H$ and $f \in L^2(0, T; V')$ we consider the linear evolution problem

$$\begin{aligned} \frac{d}{dt} \langle y(t), \varphi \rangle_H + a(y(t), \varphi) &= \langle (f + \mathcal{B}u)(t), \varphi \rangle_{V',V} \quad \forall \varphi \in V \text{ in } (0, T], \\ y(0) &= y_\circ \quad \text{in } H, \end{aligned} \tag{1}$$

where $\langle \cdot, \cdot \rangle_{V',V}$ stands for the dual pairing between V and its dual space V' and $\mathcal{B} : \mathcal{U} \rightarrow L^2(0, T; V')$ is a continuous, linear operator. It is known that for every $f \in L^2(0, T; V'), u \in \mathcal{U}$ and $y_\circ \in H$ there is a unique weak solution $y \in W(0, T)$ satisfying (1) and

$$\|y\|_{W(0,T)} \leq C (\|y_\circ\|_H + \|f\|_{L^2(0,T;V')} + \|u\|_{\mathcal{U}}) \tag{2}$$

for a constant $C > 0$ which is independent of y_\circ, f and u . For a proof of the existence of a unique solution we refer to [4, pp. 512–520]. The a-priori error estimate follows from standard variational techniques and energy estimates.

Remark 1 Let $\hat{y} \in W(0, T)$ be the unique solution to the problem

$$\frac{d}{dt} \langle y(t), \varphi \rangle_H + a(y(t), \varphi) = \langle f(t), \varphi \rangle_{V',V} \quad \forall \varphi \in V \text{ in } (0, T], \quad y(0) = y_\circ \text{ in } H.$$

We introduce the bounded, linear solution operator $\mathcal{S} : L^2(0, T; V') \rightarrow W(0, T)$: for $g \in L^2(0, T; V')$ the function $\mathcal{S}g \in W(0, T)$ is the unique solution to

$$\frac{d}{dt} \langle y(t), \varphi \rangle_H + a(y(t), \varphi) = \langle g(t), \varphi \rangle_{V',V} \quad \forall \varphi \in V \text{ in } (0, T], \quad y(0) = 0 \text{ in } H.$$

Then, the unique solution to (1) is given by $y = \hat{y} + \mathcal{S}\mathcal{B}u$. ◇

We set $\mathcal{W} = L^2(0, T; \mathbb{R}^n)$. Let us introduce the set of admissible states

$$\tilde{\mathcal{Y}}_{\text{ad}} = \{y \in W(0, T) \mid y_a \leq \mathcal{S}y \leq y_b \text{ in } \mathcal{W}\},$$

where $\mathcal{S} : L^2(0, T; V) \rightarrow \mathcal{W}$ is a bounded, linear operator with $n \in \mathbb{N}, y_a, y_b \in \mathcal{W}$ with $y_a \leq y_b$. It follows that $\tilde{\mathcal{Y}}_{\text{ad}}$ is closed and convex in $W(0, T)$. We introduce the Hilbert space $\tilde{\mathcal{X}} = W(0, T) \times \mathcal{U}$ endowed with the natural product topology. Moreover, we define the closed and convex subset $\tilde{\mathcal{X}}_{\text{ad}} = \tilde{\mathcal{Y}}_{\text{ad}} \times \mathcal{U}_{\text{ad}} \subset \tilde{\mathcal{X}}$. The cost function $\tilde{J} : \tilde{\mathcal{X}} \rightarrow \mathbb{R}$ is given by

$$\tilde{J}(y, u) = \frac{\sigma_\Omega}{2} \|y(T) - y_\Omega\|_H^2 + \frac{\sigma_Q}{2} \int_0^T \|y(t) - y_Q(t)\|_H^2 dt + \frac{\sigma_u}{2} \|u\|_{\mathcal{U}}^2 \tag{3}$$

for $x = (y, u) \in \tilde{\mathcal{X}}$, where σ_Q, σ_Ω are nonnegative weighting parameters, $\sigma_u > 0$ is a regularization parameter and $y_Q \in L^2(0, T; H)$, $y_\Omega \in H$ are given desired states. Then, we consider the following convex optimal control problem

$$\min \tilde{J}(x) \quad \text{subject to (s.t.) } x \in \mathcal{F}(\mathbf{P}) \tag{P}$$

with the set $\mathcal{F}(\mathbf{P}) = \{(\hat{y} + \mathcal{S}Bu, u) \in \tilde{\mathcal{X}}_{\text{ad}}\}$ of feasible solutions. By (2) the cost functional is radially unbounded. Since J is weakly lower semicontinuous, (P) admits a global optimal solution $\bar{x} = (\bar{y}, \bar{u})$ provided $\mathcal{F}(\mathbf{P})$ is nonempty. Since $\sigma_u > 0$ holds, \bar{x} is uniquely determined. Uniqueness follows from the strict convexity properties of the objective functional on $\tilde{\mathcal{X}}_{\text{ad}}$. For a proof we refer to [14, Sect. 1.5.2] or [24], for instance.

Example 1 (Boundary Control Without State Constraints) For $T > 0$ we set $Q = (0, T) \times \Omega$ and $\Sigma = (0, T) \times \Gamma$. Let $V = H^1(\Omega)$. For the control space we choose $\mathcal{D} = \Sigma$ and $m = 1$, i.e., $\mathcal{U} = L^2(\Sigma)$. Then, for given control $u \in \mathcal{U}$ and initial condition $y_\circ \in H$ we consider

$$c_p y_t(t, \mathbf{x}) - \Delta y(t, \mathbf{x}) = \tilde{f}(t, \mathbf{x}) \quad \text{in } Q, \tag{4a}$$

$$\frac{\partial y}{\partial n}(t, \mathbf{x}) + qy(t, \mathbf{x}) = u(t, \mathbf{x}) \quad \text{on } \Sigma, \tag{4b}$$

$$y(0, \mathbf{x}) = y_\circ(\mathbf{x}) \quad \text{in } \Omega. \tag{4c}$$

In (4) we suppose $c_p > 0$, $q \geq 0$ and $\tilde{f} \in L^2(0, T; H)$. Setting $f = \tilde{f}/c_p$, introducing the bounded (symmetric) bilinear form $a : V \times V \rightarrow \mathbb{R}$ by

$$a(\varphi, \psi) = \frac{1}{c_p} \int_\Omega \nabla \varphi(\mathbf{x}) \cdot \nabla \psi(\mathbf{x}) \, d\mathbf{x} + \frac{q}{c_p} \int_\Gamma \varphi(\mathbf{x}) \psi(\mathbf{x}) \, d\mathbf{x} \quad \text{for } \varphi, \psi \in V$$

and the linear, bounded operator $\mathcal{B} : U \rightarrow L^2(0, T; V')$ by

$$\langle (\mathcal{B}u)(t), \varphi \rangle_{V', V} = \frac{1}{c_p} \int_\Gamma u(t, \mathbf{x}) \varphi(\mathbf{x}) \, d\mathbf{x} \quad \text{for } \varphi \in V, t \in [0, T]$$

then the weak formulation of (4) can be expressed in the form (1). More details on this example one can found in [6]. ◇

Example 2 (Distributed Control with State Constraints) Let $\Omega, \Gamma, T, Q, \Sigma$ as in Example 1. Let $\chi_i \in H$, $1 \leq i \leq m$, denote given control shape functions. For the control space we choose $\mathcal{D} = (0, T)$ and set $\mathcal{U} = L^2(0, T; \mathbb{R}^m)$. Then, for given control $u \in \mathcal{U}$, initial condition $y_\circ \in H$ and inhomogeneity $f \in L^2(0, T; H)$ we

consider the linear heat equation

$$\begin{aligned}
 y_t(t, \mathbf{x}) - \nu \Delta y(t, \mathbf{x}) + \beta \cdot \nabla y(t, \mathbf{x}) &= f(t, \mathbf{x}) + \sum_{i=1}^m u_i(t) \chi_i(\mathbf{x}), & \text{in } Q, \\
 y(t, \mathbf{x}) &= 0 & \text{on } \Sigma, \\
 y(0, \mathbf{x}) &= y_0(\mathbf{x}) & \text{in } \Omega.
 \end{aligned} \tag{5}$$

with $\nu > 0$ and $\beta \in \mathbb{R}^d$. We introduce the bounded form

$$a(\varphi, \psi) = \nu \int_{\Omega} \nabla \varphi \cdot \nabla \psi \, d\mathbf{x} + \int_{\Omega} \beta \cdot \nabla \varphi \, d\mathbf{x} \quad \text{for } \varphi, \psi \in V$$

and the bounded, linear operator $\mathcal{B} : \mathcal{U} \rightarrow L^2(0, T; H) \hookrightarrow L^2(0, T; V')$ as

$$(\mathcal{B}u)(t, \mathbf{x}) = \sum_{i=1}^m u_i(t) \chi_i(\mathbf{x}) \quad \text{for } (t, \mathbf{x}) \in Q \text{ and } u \in \mathcal{U}.$$

It follows that the weak formulation of (5) can be expressed in the form (1). We choose certain shape functions $\pi_1, \dots, \pi_n \in H$ and introduce the operator $\mathcal{I} : L^2(0, T; V) \rightarrow \mathcal{W}$ by

$$(\mathcal{I}\varphi)(t) = \begin{pmatrix} (\mathcal{I}_1\varphi)(t) \\ \vdots \\ (\mathcal{I}_n\varphi)(t) \end{pmatrix} \quad \text{with} \quad (\mathcal{I}_i\varphi)(t) = \int_{\Omega} \pi_i(\mathbf{x})\varphi(t, \mathbf{x}) \, d\mathbf{x}$$

for $\varphi \in L^2(0, T; V)$. Then, the state constraints have the form

$$y_{ai}(t) \leq \int_{\Omega} \pi_i(\mathbf{x})y(t, \mathbf{x}) \, d\mathbf{x} \leq y_{bi}(t) \quad \text{in } [0, T] \text{ and for } 1 \leq i \leq n,$$

where $(y, w) \in W(0, T) \times \mathcal{W}$ holds; see also [7]. ◇

3 The Lavrentiev Regularization

It is well-known that the (sufficient) first-order optimality conditions for (P) involve a measure-valued Lagrange multiplier associated with the state constraint $\bar{y} \in \tilde{\mathcal{Y}}_{\text{ad}}$; see [14, Sect. 1.7.3]. To develop a fast numerical solution methods (by combining semismooth Newton techniques with reduced-order modelling) we apply a Lavrentiev regularization of the state constraints. For that purpose we introduce

an additional (artificial) control variable and approximate the pure state by mixed control-state constraints, which enjoy L^2 -regularity; see [23].

Instead of $\tilde{\mathcal{X}}$ we consider the Hilbert space $\mathcal{X} = W(0, T) \times \mathcal{U} \times \mathcal{W}$, again supplied with the product topology. For given $\varepsilon > 0$ the subset $\tilde{\mathcal{X}}_{\text{ad}}$ is replaced by the closed and convex subset

$$\mathcal{X}_{\text{ad}}^\varepsilon = \{(y, u, w) \in \mathcal{X} \mid y_a \leq \varepsilon w + \mathcal{I}y \leq y_b \text{ in } \mathcal{W}, u \in \mathcal{U}_{\text{ad}}\}. \quad (6)$$

For a chosen weight $\sigma_w > 0$ we also extend the cost functional \tilde{J} by defining $J : \mathcal{X} \rightarrow \mathbb{R}$ with

$$J(y, u, w) = \tilde{J}(y, u) + \frac{\sigma_w}{2} \|w\|_{\mathcal{W}}^2, \quad x = (y, u, w) \in \mathcal{X}.$$

Now the regularized optimal control problem has the following form

$$\min J(x) \quad \text{s.t.} \quad x \in \mathcal{F}(\mathbf{P}^\varepsilon) \quad (\mathbf{P}^\varepsilon)$$

with the feasible set $\mathcal{F}(\mathbf{P}^\varepsilon) = \{\hat{y} + \mathcal{S}\mathcal{B}u, u, w) \in \mathcal{X}_{\text{ad}}^\varepsilon\}$. If $\mathcal{F}(\mathbf{P}^\varepsilon) \neq \emptyset$ holds, it follows by similar arguments as above that (\mathbf{P}^ε) possesses a unique global optimal solution \bar{x} .

Let us define the control space $\mathcal{V} = \mathcal{U} \times \mathcal{W}$. We introduce the reduced cost functional \hat{J} by $\hat{J}(v) = J(\hat{y} + \mathcal{S}\mathcal{B}u, u, w)$ for $v = (u, w) \in \mathcal{V}$. By Remark 1 the solution to (1) can be expressed as $y = \hat{y} + \mathcal{S}\mathcal{B}u$. Thus, the set of admissible controls is given by

$$\mathcal{V}_{\text{ad}}^\varepsilon = \{v = (u, w) \in \mathcal{V} \mid u \in \mathcal{U}_{\text{ad}} \text{ and } \hat{y}_a \leq \varepsilon w + \mathcal{I}\mathcal{S}\mathcal{B}u \leq \hat{y}_b \text{ in } \mathcal{W}\}$$

with $\hat{y}_a = y_a - \mathcal{I}\hat{y}$ and $\hat{y}_b = y_b - \mathcal{I}\hat{y}$. Now, (\mathbf{P}^ε) is equivalent to the reduced problem

$$\min \hat{J}(v) \quad \text{s.t.} \quad v \in \mathcal{V}_{\text{ad}}^\varepsilon. \quad (\hat{\mathbf{P}}^\varepsilon)$$

The control $\bar{v} = (\bar{u}, \bar{w})$ is the unique solution to $(\hat{\mathbf{P}}^\varepsilon)$ if and only if $\bar{x} = (\hat{y} + \mathcal{S}\mathcal{B}\bar{u}, \bar{v})$ is the unique solution to (\mathbf{P}^ε) .

Next we formulate first-order sufficient optimality conditions for (\mathbf{P}^ε) (see [24], for instance):

Theorem 1 *Suppose that the feasible set $\mathcal{F}(\mathbf{P}^\varepsilon)$ is nonempty. The point $\bar{x} = (\bar{y}, \bar{u}, \bar{w}) \in \mathcal{X}_{\text{ad}}^\varepsilon$ is a (global) optimal solution to (\mathbf{P}^ε) if and only if there are unique Lagrange multipliers $(\bar{p}, \bar{\lambda}_u, \bar{\lambda}_y) \in \mathcal{X}$ satisfying the dual equations*

$$-\frac{d}{dt} \langle \bar{p}(t), \varphi \rangle_H + a(\varphi, \bar{p}(t)) + \langle (\mathcal{I}^* \bar{\lambda}_y)(t), \varphi \rangle_{V', V} = \sigma_Q \langle (y_Q - \bar{y})(t), \varphi \rangle_H \quad (7)$$

$$\forall \varphi \in V \text{ in } [0, T), \quad \bar{p}(T) = \sigma_\Omega (y_\Omega - \bar{y}(T)) \text{ in } H,$$

and the optimality conditions

$$\sigma_u \bar{u} - \mathcal{B}^* \bar{p} + \bar{\lambda}_u = 0 \text{ in } \mathcal{U}, \quad \sigma_w \bar{w} + \varepsilon \bar{\lambda}_y = 0 \text{ in } \mathcal{W},$$

where $\mathcal{I}^* : \mathcal{W} \rightarrow L^2(0, T; V')$ and $\mathcal{B}^* : L^2(0, T; V) \rightarrow \mathcal{U}$ denote the adjoint operators of \mathcal{I} and \mathcal{B} , respectively. For the Lagrange multipliers $\bar{\lambda}_u$ and $\bar{\lambda}_y$ we have

$$\begin{aligned} \bar{\lambda}_u &= \max(0, \bar{\lambda}_u + \gamma_u(\bar{u} - u_b)) + \min(0, \bar{\lambda}_u + \gamma_u(\bar{u} - u_a)) && \text{in } \mathcal{U}, \\ \bar{\lambda}_y &= \max(0, \bar{\lambda}_y + \gamma_w(\varepsilon \bar{w} + \mathcal{I} \bar{y} - y_b)) + \min(0, \bar{\lambda}_y + \gamma_w(\varepsilon \bar{w} + \mathcal{I} \bar{y} - y_a)) && \text{in } \mathcal{W}, \end{aligned}$$

where $\gamma_u, \gamma_w > 0$ are arbitrarily chosen.

Remark 2

- (1) Analogous to Remark 1 we split the adjoint variable into one part depending on the fixed desired states and into two other parts, which depend linearly on the control variable and on the multiplier λ . Recall that we have defined \hat{y} as well as the operator \mathcal{S} in Remark 1. For given $y_Q \in L^2(0, T; H)$ and $y_\Omega \in H$ let $\hat{p} \in W(0, T)$ denote the unique solution to the adjoint equation

$$\begin{aligned} -\frac{d}{dt} \langle \hat{p}(t), \varphi \rangle_H + a(\varphi, \hat{p}(t)) &= \sigma_Q \langle (y_Q - \hat{y})(t), \varphi \rangle_H \quad \forall \varphi \in V \text{ in } [0, T), \\ \hat{p}(T) &= \sigma_\Omega (y_\Omega - \hat{y}(T)) \quad \text{in } H. \end{aligned}$$

Further, we define the linear, bounded operators $\mathcal{A}_1 : \mathcal{U} \rightarrow W(0, T)$ and $\mathcal{A}_2 : \mathcal{W} \rightarrow W(0, T)$ as follows: for any $u \in \mathcal{U}$ the function $p = \mathcal{A}_1 u$ is the unique solution to

$$\begin{aligned} -\frac{d}{dt} \langle p(t), \varphi \rangle_H + a(\varphi, p(t)) &= -\sigma_Q \langle (\mathcal{S} \mathcal{B} u)(t), \varphi \rangle_H \quad \forall \varphi \in V \text{ in } [0, T), \\ p(T) &= -\sigma_\Omega (\mathcal{S} \mathcal{B} u)(T) \quad \text{in } H \end{aligned}$$

and for given $\lambda \in \mathcal{W}$ the function $p = \mathcal{A}_2 \lambda$ uniquely solves $p(T) = 0$ in H and

$$-\frac{d}{dt} \langle p(t), \varphi \rangle_H + a(\varphi, p(t)) + \langle (\mathcal{I}^* \lambda_y)(t), \varphi \rangle_{V', V} = 0 \quad \forall \varphi \in V \text{ in } [0, T).$$

Then, the solution to (7) can be expressed as $\bar{p} = \hat{p} + \mathcal{A}_1 \bar{u} + \mathcal{A}_2 \bar{\lambda}_y$.

- (2) To solve (\mathbf{P}^ε) numerically for fixed $\varepsilon > 0$ we use a primal-dual active set strategy. This method is equivalent to a locally superlinearly convergent semi-smooth Newton algorithm applied to the first-order optimality conditions [8–10]. \diamond

4 The POD Method

Let \mathcal{Z} be either the space H or the space V . In \mathcal{Z} we denote by $\langle \cdot, \cdot \rangle_{\mathcal{Z}}$ and $\| \cdot \|_{\mathcal{Z}} = \langle \cdot, \cdot \rangle_{\mathcal{Z}}^{1/2}$ the inner product and the associated norm, respectively. For fixed $\wp \in \mathbb{N}$ let the so-called *snapshots* $z^k(t) \in \mathcal{Z}$ be given for $t \in [0, T]$ and $1 \leq k \leq \wp$. To avoid a trivial case we suppose that at least one of the z^k 's is nonzero. Then, we introduce the linear subspace

$$\mathcal{Z}^{\wp} = \text{span} \left\{ z^k(t) \mid t \in [0, T] \text{ and } 1 \leq k \leq \wp \right\} \subset \mathcal{Z} \tag{8}$$

with dimension $\mathfrak{d} \geq 1$. We call the set \mathcal{Z}^{\wp} *snapshot subspace*. The method of POD consists in choosing a complete orthonormal basis $\{\psi_i\}_{i=1}^{\infty}$ in \mathcal{Z} such that for every $\ell \leq \mathfrak{d}$ the mean square error between the \wp elements z^k and their corresponding ℓ th partial Fourier sum is minimized:

$$\left\{ \begin{array}{l} \min \sum_{k=1}^{\wp} \int_0^T \left\| z^k(t) - \sum_{i=1}^{\ell} \langle z^k(t), \psi_i \rangle_{\mathcal{Z}} \psi_i \right\|_{\mathcal{Z}}^2 dt \\ \text{s.t. } \{\psi_i\}_{i=1}^{\ell} \subset \mathcal{Z} \text{ and } \langle \psi_i, \psi_j \rangle_{\mathcal{Z}} = \delta_{ij}, \quad 1 \leq i, j \leq \ell. \end{array} \right. \tag{9}$$

In (9) the symbol δ_{ij} denotes the Kronecker symbol satisfying $\delta_{ii} = 1$ and $\delta_{ij} = 0$ for $i \neq j$. An optimal solution $\{\bar{\psi}_i\}_{i=1}^{\ell}$ to (9) is called a *POD basis of rank ℓ* .

Remark 3 In real computations, we do not have the whole trajectories $z^k(t)$ at hand for all $t \in [0, T]$ and $1 \leq k \leq \wp$. Here we apply a discrete variant of the POD method; see [7, 16] for more details. \diamond

To solve (9) we define the linear operator $\mathcal{R} : \mathcal{Z} \rightarrow \mathcal{Z}^{\wp}$ as follows:

$$\mathcal{R}\psi = \sum_{k=1}^{\wp} \int_0^T \langle \psi, z^k(t) \rangle_{\mathcal{Z}} z^k(t) dt \quad \text{for } \psi \in \mathcal{Z}. \tag{10}$$

Then, \mathcal{R} is a compact, nonnegative and selfadjoint operator. Suppose that $\{\bar{\lambda}_i\}_{i=1}^{\infty}$ and $\{\bar{\psi}_i\}_{i=1}^{\infty}$ denote the nonnegative eigenvalues and associated orthonormal eigenfunctions of \mathcal{R} satisfying

$$\mathcal{R}\bar{\psi}_i = \bar{\lambda}_i \bar{\psi}_i, \quad \bar{\lambda}_1 \geq \dots \geq \bar{\lambda}_{\mathfrak{d}} > \bar{\lambda}_{\mathfrak{d}+1} = \dots = 0. \tag{11}$$

Then, for every $\ell \leq \mathfrak{d}$ the first ℓ eigenfunctions $\{\bar{\psi}_i\}_{i=1}^{\ell}$ solve (9) and

$$\sum_{k=1}^{\wp} \int_0^T \left\| z^k(t) - \sum_{i=1}^{\ell} \langle z^k(t), \bar{\psi}_i \rangle_{\mathcal{Z}} \bar{\psi}_i \right\|_{\mathcal{Z}}^2 dt = \sum_{i=\ell+1}^{\mathfrak{d}} \bar{\lambda}_i.$$

For more details we refer the reader to [11, 12] and [7, Chap. 2], for instance.

Remark 4

- (a) In the context of the optimal control problem (\mathbf{P}^ε) a reasonable choice for the snapshots is $z^1 = y$ and $z^2 = p$. Utilizing new POD error estimates for evolution problems [3, 20] and optimal control problems [13, 25] convergence and rate of convergence results are derived for linear-quadratic control constrained problems in [7] for the choices $\mathcal{Z} = H$ and $\mathcal{Z} = V$.
- (b) For the numerical realization the space \mathcal{Z} has to be discretized by, e.g., finite element discretizations. In this case the Hilbert space \mathcal{Z} has to be replaced by an Euclidean space \mathbb{R}^l endowed with a weighted inner product; see [7].

If a POD basis $\{\psi_i\}_{i=1}^\ell$ of rank ℓ is computed, we set $V^\ell = \text{span}\{\psi_1, \dots, \psi_\ell\}$. Then, one can derive a reduced-order model (ROM) for (1): for any $g \in L^2(0, T; V')$ the function $q^\ell = \mathcal{S}^\ell g$ is given by $q^\ell(0) = 0$ in H and

$$\frac{d}{dt} \langle q^\ell(t), \psi \rangle_H + a(q^\ell(t), \psi) = \langle g(t), \psi \rangle_{V', V} \quad \forall \psi \in V^\ell \text{ in } (0, T].$$

For any $u \in \mathcal{U}_{\text{ad}}$ the POD approximation y^ℓ for the state solution is $y^\ell = \hat{y} + \mathcal{S}^\ell \mathcal{B}u$. Analogously, a ROM can be derived for the adjoint equation; see, e.g., [7]. The POD Galerkin approximation of $(\hat{\mathbf{P}}^\varepsilon)$ is given by

$$\min J^\ell(v) = J(\hat{y} + \mathcal{S}^\ell \mathcal{B}u, v) \quad \text{s.t.} \quad v = (u, w) \in \mathcal{V}_{\text{ad}}^{\varepsilon, \ell} \quad (\hat{\mathbf{P}}^{\varepsilon, \ell})$$

where the set of admissible controls is

$$\mathcal{V}_{\text{ad}}^{\varepsilon, \ell} = \{v = (u, w) \in \mathcal{V} \mid u \in \mathcal{U}_{\text{ad}} \text{ and } \hat{y}_a \leq \varepsilon w + \mathcal{I} \mathcal{S}^\ell \mathcal{B}u \leq \hat{y}_b \text{ in } \mathcal{W}\}.$$

5 A-Posteriori Error Analysis

Let us consider (\mathbf{P}) with control, but no state constraints. Based on a perturbation argument [5] it is derived in [25] how far the suboptimal POD control \bar{u}^ℓ , computed on the basis of the POD model, is from the (unknown) exact \bar{u} . Then, the error estimate reads as follows:

$$\|\bar{u}^\ell - \bar{u}\|_{\mathcal{U}} \leq \frac{1}{\sigma_u} \|\zeta^\ell\|_{\mathcal{U}}, \quad (12)$$

where the computable perturbation function $\zeta^\ell \in \mathcal{U}$ is given by

$$\zeta^\ell = \begin{cases} -\min(0, \sigma_u \bar{u}^\ell - \mathcal{B}^* \tilde{p}^\ell) & \text{in } \mathcal{A}_a^\ell = \{s \in \mathcal{D} \mid \bar{u}^\ell(s) = u_a(s)\}, \\ -\max(0, \sigma_u \bar{u}^\ell - \mathcal{B}^* \tilde{p}^\ell) & \text{in } \mathcal{A}_b^\ell = \{s \in \mathcal{D} \mid \bar{u}^\ell(s) = u_b(s)\}, \\ -(\sigma_u \bar{u}^\ell - \mathcal{B}^* \tilde{p}^\ell) & \text{in } \mathcal{D} \setminus (\mathcal{A}_a^\ell \cup \mathcal{A}_b^\ell), \end{cases}$$

with $\tilde{p}^\ell = \hat{p} + \mathcal{A}_1 \tilde{u}^\ell$. It is shown in [7, 25] that $\|\zeta^\ell\|_{\mathcal{U}}$ tends to zero as ℓ tends to infinity. Hence, increasing the number of POD ansatz functions leads to more accurate POD suboptimal controls.

Estimate (12) can be generalized for the mixed control-state constraints. First-order sufficient optimality conditions for $(\hat{\mathbf{P}}^\varepsilon)$ are of the form

$$\langle \hat{J}'(\bar{v}), v - \bar{v} \rangle_{\mathcal{V}} \geq 0 \quad \text{for all } v \in \mathcal{V}_{\text{ad}}^\varepsilon, \quad (13)$$

where the gradient at a point $v = (u, w) \in \mathcal{V}$ is given by Tröltzsch [24]

$$\langle \hat{J}'(v), v_\delta \rangle_{\mathcal{V}} = \langle \sigma_u u - \mathcal{B}^*(\hat{p} + \mathcal{A}_1 u), u_\delta \rangle_{\mathcal{U}} + \langle \sigma_w w_\delta, w_\delta \rangle_{\mathcal{W}} \quad \forall v_\delta = (u_\delta, w_\delta) \in \mathcal{V}.$$

Let us introduce the bounded, linear transformation $\mathcal{T} : \mathcal{V} \rightarrow \mathcal{V}$ as

$$\mathcal{T}(v) = (u, \varepsilon w + \mathcal{I} \mathcal{S} \mathcal{B} u) \quad \text{for } v = (u, w) \in \mathcal{V}. \quad (14)$$

We assume that \mathcal{T} is continuously invertible. For sufficient conditions we refer to [8, Lemma 2.1]. Then, $v = (u, w)$ belongs to $\mathcal{V}_{\text{ad}}^\varepsilon$ if and only if $\mathbf{v} = (u, \mathfrak{w}) = \mathcal{T}(v)$ satisfies

$$u_a \leq u \leq u_b \text{ in } \mathcal{U} \quad \text{and} \quad \hat{y}_a \leq \mathfrak{w} \leq \hat{y}_b \text{ in } \mathcal{W}. \quad (15)$$

Notice that (13) can be expressed equivalently as

$$\langle \mathcal{T}^{-*} \hat{J}'(\mathcal{T}^{-1} \bar{\mathbf{v}}), \mathbf{v} - \bar{\mathbf{v}} \rangle_{\mathcal{V}} \geq 0 \quad \text{for all } \mathbf{v} \in \mathcal{V} \text{ satisfying (15)}, \quad (16)$$

where \mathcal{T}^{-*} denotes the inverse of the operator \mathcal{T}^* . Suppose that $\bar{v}^\ell = (\bar{u}^\ell, \bar{w}^\ell) \in \mathcal{V}_{\text{ad}}^{\varepsilon, \ell}$ is the solution to $(\hat{\mathbf{P}}^{\varepsilon, \ell})$. Our goal is to estimate the norm

$$\|\bar{v} - \bar{v}^\ell\|_{\mathcal{V}}$$

without the knowledge of the optimal solution $\bar{v} = \mathcal{T}^{-1} \bar{\mathbf{v}}$. We set $\bar{\mathbf{v}}^\ell = \mathcal{T} \bar{v}^\ell = (\bar{u}^\ell, \varepsilon \bar{w}^\ell + \mathcal{I} \mathcal{S} \mathcal{B} \bar{u}^\ell)$. If $\bar{v}^\ell \neq \bar{v}$ holds, then $\bar{\mathbf{v}}^\ell \neq \bar{\mathbf{v}}$. In particular, \bar{v}^ℓ does not satisfy the sufficient optimality condition (13). However, there exists a function $\zeta^\ell \in \mathcal{V}$ such that

$$\langle \mathcal{T}^{-*} \hat{J}'(\mathcal{T}^{-1} \bar{\mathbf{v}}^\ell) + \zeta^\ell, \mathbf{v} - \bar{\mathbf{v}}^\ell \rangle_{\mathcal{V}} \geq 0 \quad \text{for all } \mathbf{v} \in \mathcal{V} \text{ satisfying (15)}. \quad (17)$$

Choosing $\mathbf{v} = \bar{\mathbf{v}}^\ell$ in (16), $\mathbf{v} = \bar{\mathbf{v}}$ in (17) and adding both inequality we infer that

$$\begin{aligned} 0 &\leq \left\langle \mathcal{T}^{-*} (\hat{J}'(\mathcal{T}^{-1} \bar{\mathbf{v}}^\ell) + \mathcal{T}^* \zeta^\ell - \hat{J}'(\mathcal{T}^{-1} \bar{\mathbf{v}})), \bar{\mathbf{v}} - \bar{\mathbf{v}}^\ell \right\rangle_{\mathcal{V}} \\ &= \left\langle \hat{J}'(\bar{v}^\ell) - \hat{J}'(\bar{v}) + \mathcal{T}^* \zeta^\ell, \mathcal{T}^{-1}(\bar{\mathbf{v}} - \bar{\mathbf{v}}^\ell) \right\rangle_{\mathcal{V}} \end{aligned}$$

$$\begin{aligned}
 &= \langle (\sigma_u(\bar{u}^\ell - \bar{u}) - \mathcal{B}^* \mathcal{A}_1(\bar{u}^\ell - \bar{u}), \sigma_w(\bar{w}^\ell - \bar{w})) + \mathcal{T}^* \zeta^\ell, \bar{v} - \bar{v}^\ell \rangle_{\mathcal{V}} \\
 &\leq -\sigma \|\bar{v} - \bar{v}^\ell\|_{\mathcal{V}}^2 + \langle \mathcal{B}^* \mathcal{A}_1(\bar{u} - \bar{u}^\ell), \bar{u} - \bar{u}^\ell \rangle_{\mathcal{U}} + \langle \mathcal{T}^* \zeta^\ell, \bar{v} - \bar{v}^\ell \rangle_{\mathcal{V}}
 \end{aligned}$$

with $\sigma = \min(\sigma_u, \sigma_w) > 0$. In [8, Lemma 2.2] it is shown that $\langle \mathcal{B}^* \mathcal{A}_1(\bar{u} - \bar{u}^\ell), \bar{u} - \bar{u}^\ell \rangle_{\mathcal{U}} \leq 0$ holds. Consequently,

$$0 \leq -\sigma \|\bar{v} - \bar{v}^\ell\|_{\mathcal{V}}^2 + \langle \mathcal{T}^* \zeta^\ell, \bar{v} - \bar{v}^\ell \rangle_{\mathcal{V}} \leq -\sigma \|\bar{v} - \bar{v}^\ell\|_{\mathcal{V}}^2 + \|\mathcal{T}^* \zeta\|_{\mathcal{V}} \|\bar{v} - \bar{v}^\ell\|_{\mathcal{V}}$$

which implies the following proposition.

Proposition 1 *Let the operator \mathcal{T} —introduced in (14)—possess a bounded inverse. Suppose that \bar{v} and \bar{v}^ℓ are the optimal solution to $(\hat{\mathbf{P}}^\varepsilon)$ and $(\hat{\mathbf{P}}^{\varepsilon, \ell})$, respectively, satisfying $\bar{v}^\ell = \mathcal{T} \bar{v}^\ell \in \mathcal{V}_{\text{ad}}^\varepsilon$. Then, there is a perturbation $\zeta^\ell = (\zeta_u^\ell, \zeta_w^\ell) \in \mathcal{V}$ satisfying*

$$\|\bar{v} - \bar{v}^\ell\|_{\mathcal{V}} \leq \frac{1}{\sigma} \|\mathcal{T}^* \zeta^\ell\|_{\mathcal{V}} \quad \text{with } \sigma = \min(\sigma_u, \sigma_w) > 0. \quad (18)$$

The perturbation $\zeta^\ell = (\zeta_u^\ell, \zeta_w^\ell)$ can be computed as follows: Let $\xi^\ell = (\xi_u^\ell, \xi_w^\ell) = \mathcal{T}^{-*} \hat{\mathcal{J}}'(\bar{v}^\ell) \in \mathcal{V}$. Then, ξ^ℓ solves the linear system

$$\begin{pmatrix} \text{id}_{\mathcal{U}} & 0 \\ \mathcal{B}^* \mathcal{I}^* \mathcal{I}^* & \varepsilon \text{id}_{\mathcal{W}} \end{pmatrix} \begin{pmatrix} \xi_u^\ell \\ \xi_w^\ell \end{pmatrix} = \begin{pmatrix} \sigma_u \bar{u}^\ell - \mathcal{B}^*(\hat{p} + \mathcal{A}_1 \bar{u}^\ell) \\ \sigma_w \bar{w}^\ell \end{pmatrix},$$

where, e.g., $\text{id}_{\mathcal{U}} : \mathcal{U} \rightarrow \mathcal{U}$ stands for the identity operator. Note that (17) can be written as $\langle \xi + \zeta, \mathbf{v} - \bar{v}^\ell \rangle_{\mathcal{V}} \geq 0$ for all $\mathbf{v} \in \mathcal{V}$ satisfying (15). We find

$$\zeta_u^\ell = \begin{cases} -\min(0, \xi_u^\ell) & \text{in } \mathcal{A}_a^u = \{\bar{u}^\ell = u_a\} \subset \mathcal{U}, \\ -\max(0, \xi_u^\ell) & \text{in } \mathcal{A}_b^u = \{\bar{u}^\ell = u_b\} \subset \mathcal{U}, \\ -\xi_u^\ell & \text{in } \mathcal{U} \setminus (\mathcal{A}_a^u \cup \mathcal{A}_b^u) \end{cases}$$

and

$$\zeta_w^\ell = \begin{cases} -\min(0, \xi_w^\ell) & \text{in } \mathcal{A}_a^w = \{\varepsilon \bar{w}^\ell + \mathcal{I} \mathcal{I} \mathcal{B} \bar{u}^\ell = \hat{y}_a\} \subset \mathcal{W}, \\ -\max(0, \xi_w^\ell) & \text{in } \mathcal{A}_b^w = \{\varepsilon \bar{w}^\ell + \mathcal{I} \mathcal{I} \mathcal{B} \bar{u}^\ell = \hat{y}_b\} \subset \mathcal{W} \\ -\xi_w^\ell & \text{in } \mathcal{W} \setminus (\mathcal{A}_a^w \cup \mathcal{A}_b^w). \end{cases}$$

6 Optimality-System POD

The accuracy of the reduced-order model can be controlled by the a-posteriori error analysis presented in Sect. 5. However, if the POD basis is created from a reference trajectory containing features which are quite different from those of

the optimally controlled trajectory, a rather huge number of POD ansatz functions has to be included in the reduced-order model. This fact may lead to non-efficient reduced-order models and numerical instabilities. To avoid these problems the POD basis is generated in an initialization step utilizing *optimality system POD* (OS-POD) introduced in [17]. In OS-POD the POD basis is updated in the direction of the minimum of the cost. Recall that the POD basis is computed from the state $y = \hat{y} + \mathcal{S} \mathcal{B} u$ with some control $u^0 \in \mathcal{U}_{\text{ad}}$. Thus, the reduced-order Galerkin projection depends on the state variable and hence on the control u at which the eigenvalue $\mathcal{R} \psi_i = \lambda_i \psi_i$ for $i = 1, \dots, \ell$ is solved for the basis $\{\psi_i\}_{i=1}^{\ell}$. If the optimal control \bar{u} differs significantly from the initially chosen control u^0 , the POD basis does not reflect the dynamics of the system in a sufficiently accurate manner. Therefore, we consider the extended problem:

$$\min \hat{J}^{\ell}(v) \text{ s.t. } \begin{cases} v = (u, w) \in \mathcal{V}_{\text{ad}}^{\varepsilon, \ell}, \\ (\psi, \lambda) \text{ satisfies (11) for } \wp = 1 \text{ and } z^1 = \hat{y} + \mathcal{S} \mathcal{B} u. \end{cases} \quad (\hat{\mathbf{P}}_{\text{os}}^{\varepsilon, \ell})$$

Notice that the first line of the constraints in $(\hat{\mathbf{P}}_{\text{os}}^{\varepsilon, \ell})$ coincides with the constraints in $(\hat{\mathbf{P}}^{\varepsilon, \ell})$, whereas the second line of the constraints in $(\hat{\mathbf{P}}_{\text{os}}^{\varepsilon, \ell})$ are the infinite-dimensional eigenvalue problem defining the POD basis. For the optimal solution the problem formulation $(\hat{\mathbf{P}}_{\text{os}}^{\varepsilon, \ell})$ has the property that the associated POD reduced system is computed from the trajectory corresponding to the optimal control and thus, differently from $(\hat{\mathbf{P}}^{\varepsilon, \ell})$, the problem of unmodelled dynamics is removed. Of course, $(\hat{\mathbf{P}}_{\text{os}}^{\varepsilon, \ell})$ is more complicated than $(\hat{\mathbf{P}}^{\varepsilon, \ell})$. For practical realization an operator splitting approach is used in [17], where also sufficient conditions are given so that $(\hat{\mathbf{P}}_{\text{os}}^{\varepsilon, \ell})$ possesses a unique optimal solution, which can be characterized by first-order necessary optimality conditions; compare [17] for more details. Convergence results for OS-POD are studied in [18]. The combination of OS-POD and a-posteriori error analysis is suggested in [26] and tested successfully in [6]. The resulting strategy is presented in the next section.

7 Algorithms

For *pure control constraints*, i.e., \hat{J}^{ℓ} depends only on u , a variable splitting is proposed, where a good POD basis is initialized by applying a few projected gradient steps [15]. Then, the POD basis is kept fixed and $(\hat{\mathbf{P}}^{\varepsilon, \ell})$ is solved. If the a-posteriori error estimator $\|\zeta^{\ell}\|_{\mathcal{U}}/\sigma_u$ is too large [compare (12)], the number ℓ of POD basis elements is increased and a new solution to $(\hat{\mathbf{P}}^{\varepsilon, \ell})$ is computed. This process is repeated until we obtain convergence; see Algorithm 1. Let us mention that we also utilize snapshots of the adjoint variable in order to compute a POD basis as described in Remark 4(a).

For the *mixed constraints*, this iteration does not turn out to be efficient enough. The gradient steps do not lead to a satisfactorily fast and accurate POD basis.

Algorithm 1 (OS-POD with a-posteriori error estimation for control constraints)

- Require:** Maximal number ℓ_{\max} of POD basis elements, $\ell_{\min} < \ell_{\max}$, initial control u^0 , and a-posteriori error tolerance $\varepsilon_{\text{apo}} > 0$;
- 1: Determine the state $y = \hat{y} + \mathcal{S}\mathcal{B}u^0$ and adjoint $p = \hat{p} + \mathcal{A}_1 u^0$;
 - 2: Compute a POD basis $\{\psi_i(u)\}_{i=1}^{\ell}$ as described in Remark 4(a);
 - 3: Perform $k \geq 0$ projected gradient steps (PGS) with an Armijo line search for $(\hat{\mathbf{P}}_{\text{os}}^{\varepsilon, \ell})$ in order to get u^k and associated POD basis $\{\psi_i(u^k)\}_{i=1}^{\ell_{\max}}$; set $\ell = \ell_{\min}$;
 - 4: Solve $(\hat{\mathbf{P}}^{\varepsilon, \ell})$ for $\bar{u}^{\ell} \in \mathcal{U}_{\text{ad}}$ by the primal-dual active set strategy;
 - 5: Compute the perturbation $\zeta^{\ell} = \zeta^{\ell}(\bar{u}^{\ell})$ as explained in Sect. 5;
 - 6: **if** $\|\zeta^{\ell}\|_{\mathcal{U}}/\sigma_u > \varepsilon_{\text{apo}}$ **and** $\ell < \ell_{\max}$ **then**
 - 7: Enlarge ℓ and go back to step 4;
 - 8:
 - 9: **end if**
-

Therefore, we invest more effort in the gradient steps by interacting between the projected gradient method and the primal-dual active set strategy (PDASS). In contrast to the situation of pure control constraints, we can provide basis updates based on the more accurate PDASS controls. The strategy is explained in Algorithm 2.

Algorithm 2 (OS-POD with a-posteriori error estimation for state constraints)

- Require:** Maximal number ℓ_{\max} of POD basis elements, $\ell < \ell_{\max}$, initial control u^0 , and a-posteriori error tolerance $\varepsilon_{\text{apo}} > 0$;
- 1: Determine the state $y = \hat{y} + \mathcal{S}\mathcal{B}u^0$ and adjoint $p = \hat{p} + \mathcal{A}_1 u^0$;
 - 2: Compute a POD basis $\{\psi_i(u)\}_{i=1}^{\ell}$ as described in Remark 4(a);
 - 3: Solve $(\hat{\mathbf{P}}^{\varepsilon, \ell})$ for $\bar{v}^{\ell} = (\bar{u}^{\ell}, \bar{w}^{\ell}) \in \mathcal{V}_{\text{ad}}^{\varepsilon, \ell}$ by the primal-dual active set strategy;
 - 4: Perform $k \geq 0$ projected gradient steps with an Armijo line search for $(\hat{\mathbf{P}}_{\text{os}}^{\varepsilon, \ell})$ in order to get u^k and associated POD basis $\{\psi_i(u^k)\}_{i=1}^{\ell}$;
 - 5: Compute the perturbation $\zeta = \zeta(\bar{v}^{\ell})$ as explained in Sect. 5;
 - 6: **if** $\|\mathcal{T}^* \zeta\|_{\mathcal{V}}/\sigma > \varepsilon_{\text{apo}}$ **and** $\ell < \ell_{\max}$ **then**
 - 7: Enlarge ℓ and go back to step 3;
 - 8: **else**
 - 9: Set $\ell = \ell_{\min}$ and go back to step 1;
 - 10: **end if**
-

8 Numerical Experiments

In this section we carry out numerical test examples illustrating the efficiency of the combination of OS-POD and a-posteriori error estimation. The evolution problems are approximated by a standard finite element (FE) method with piecewise linear finite elements for the spatial discretization. The time integration is done by the

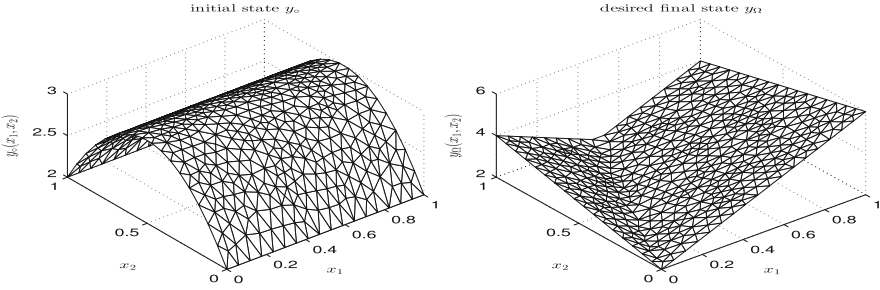


Fig. 1 Run 1: the initial condition y_0 (left) and the desired terminal state y_Ω (right)

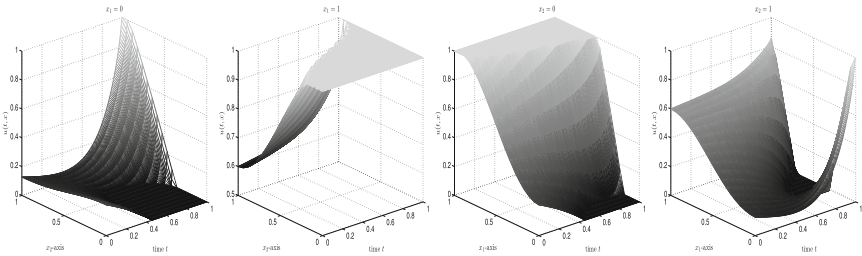


Fig. 2 Run 1: FE optimal control along the boundary parts $x_1 = 0$, $x_1 = 1$, $x_2 = 0$, and $x_2 = 1$

implicit Euler method. All programs are written in MATLAB utilizing the PARTIAL DIFFERENTIAL EQUATION TOOLBOX for the FE discretization.

Run 1 (Example 1) We choose $d = 2$ and consider the unit square $\Omega = (0, 1) \times (0, 1) \subset \mathbb{R}^2$ as spatial domain with time interval $[0, T] = [0, 1]$. The FE triangulation with maximal edge length $h = 0.06$ leads to 498 degrees of freedom. For the time integration we choose an equidistant time grid $t_j = j\Delta t$ for $j = 0, \dots, 250$ with $\Delta t = 0.004$. Motivated by the discretization error we set $\varepsilon_{\text{apo}} = \max(h^2, \Delta t) = \Delta t$. In (4) we choose the data $c_p = 10$, $q = 0.01$, $\tilde{f} = 0$ and $y_0(x_1, x_2) = 3 - 4(x_2 - 0.5)^2$; see left plot in Fig. 1. We use $\sigma_\Omega = 0$, $\sigma_\Omega = 1$ and the regularization $\sigma_u = 0.1$ in the cost function (3) to approximate the desired terminal state $y_\Omega(x_1, x_2) = 2 + 2|2x_1 - x_2|$; see right plot in Fig. 1. The control constraints are chosen to be $u_a = 0$ and $u_b = 1$. The FE primal-dual active set strategy needs five iterations and 860.75 s. The optimal FE control is presented in Fig. 2. We apply Algorithm 1 with $\ell_{\text{max}} = 40$, $\ell = 10$ and initial control $u^0 = 0$. First we do not perform any OS-POD strategy (i.e., $k = 0$ in Algorithm 1). The method stops in 110.77 s with $\ell = 35 < \ell_{\text{max}}$ ansatz functions with $\|\zeta^\ell\|_U/\sigma_u \approx 0.0034 < \varepsilon_{\text{apo}}$. Each solve of $(\hat{\mathbf{P}}^{\varepsilon, \ell})$ needs four or five iterations to determine the suboptimal POD solutions. If we initialize Algorithm 1 with the optimal FE control \bar{u}^{FE} as initial control and perform no OS-POD strategy, only $\ell = 13$ POD basis functions are required. We get $\|\zeta^\ell\|_U/\sigma_u \approx 0.0019 < \varepsilon_{\text{apo}}$ and the CPU time is 11.48 s, which is ten times faster than with the initial control $u^0 = 0$. With one OS-POD gradient

Table 1 Run 1: performance of Algorithm 1

	$k = 0$	$k = 1$	$k = 2$	With \bar{u}^{FE}
Required ℓ	35	40	13	13
CPU time	110.77 s	147.14 s	18.39 s	11.48 s
$\ \zeta^\ell\ _{\mathcal{U}}/\sigma_u$	$3.43 \cdot 10^{-3}$	$1.14 \cdot 10^{-2}$	$2.82 \cdot 10^{-3}$	$1.94 \cdot 10^{-3}$
$\ \bar{u}^\ell - \bar{u}^{FE}\ _{\mathcal{U}}$	$3.15 \cdot 10^{-3}$	$9.53 \cdot 10^{-3}$	$2.62 \cdot 10^{-3}$	$1.93 \cdot 10^{-3}$

Table 2 Run 1: comparison of POD suboptimal solutions for $\ell = 15$ and k OS-POD steps

	$k = 0$	$k = 1$	$k = 2$	With \bar{u}^{FE}
$\ \zeta^\ell\ _{\mathcal{V}}/\sigma_u$	$2.50 \cdot 10^{-2}$	$1.45 \cdot 10^{-2}$	$2.27 \cdot 10^{-3}$	$1.59 \cdot 10^{-3}$
$\ \bar{u}^\ell - \bar{u}^{FE}\ _{\mathcal{U}}$	$2.06 \cdot 10^{-2}$	$1.19 \cdot 10^{-2}$	$2.07 \cdot 10^{-3}$	$1.59 \cdot 10^{-3}$
Different u_a	96	67	15	16
Different u_b	63	38	6	4

Table 3 Run 1: number of active nodes

	u^1	u^2	\bar{u}^{FE}
$u^k = u_a$	1321 (1321)	1814 (1812)	2233
$u^k = u_b$	986 (986)	3632 (3627)	3891

In parenthesis the number of nodes, where $u^k = \bar{u}^{FE} = u_a$ or $u_k = \bar{u}^{FE} = u_b$, respectively

step, the tolerance ε_{apo} is not reached with the available $\ell_{\text{max}} = 40$ basis functions. Though we make an effort in direction of the optimal control, the algorithm seems to perform even worse than with the basis corresponding to the uncontrolled state. This can be seen in the higher control errors that cause the algorithm to run up to $\ell_{\text{max}} = 40$ ansatz functions. We can see, however, that the errors in the suboptimal state are one order smaller than without gradient steps, so the POD basis did improve after all. After $k = 2$ gradient steps, the performance is considerably better: The algorithm already terminates with a ROM rank of $\ell = 13$ like in the optimal case. In Table 1 we provide the required CPU times and final errors. Additionally regard Table 2 where we compare the errors for the POD suboptimal solutions for fixed rank $\ell = 15$. Here, we also provide the number of nodes that are restricted by the box constraints either in the suboptimal control \bar{u}^{15} or in the FE optimal control \bar{u}^{FE} , but not in both. It tells us, how many of the restricted nodes are mistaken. This number decreases to 21 by the gradient steps. Next we are interested in the approximation of the active sets. The computations are done with 68 triangulation nodes at the boundary and 251 time steps; that is a total amount of $68 \cdot 251 = 17068$ boundary nodes in the time interval $[0, T]$. The FE optimal control is restricted by u_a at 2233 and by u_b at 3891 nodes. In Table 3 we present the number of nodes where u_k is restricted to the lower or upper bound and, in parenthesis, how many of these nodes are actually restricted correctly, i.e. equal to \bar{u}^{FE} , what amounts to more than 99%. Finally, let us illustrate the changes achieved in the POD basis by the OS-POD steps. The left plot of Fig. 3 shows how the decay of normalized eigenvalues differs depending on the used control for snapshot generation. The

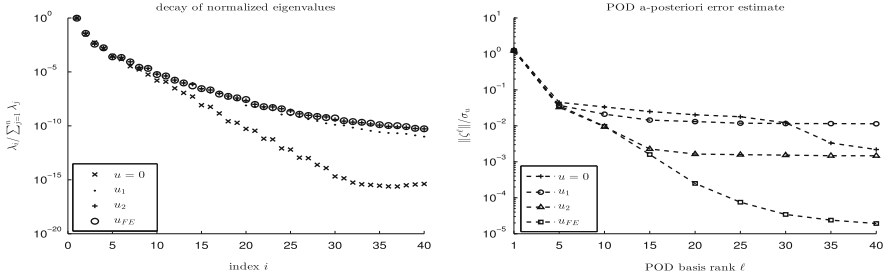


Fig. 3 Run 1: comparison of eigenvalue decay for POD basis generated with u_k after k gradient steps or with \bar{u}^{FE} (left) and a-posteriori error for increasing ℓ (right)

eigenvalues corresponding to the uncontrolled state decay faster and further than those corresponding to the more or less optimally controlled state; increasing the utilized rank further than $\ell = 35$ yields no more improvement. The difference caused by one gradient step is significant. A lot more basis functions contain still relevant information for the reduced order models. After the second gradient step the course is equal to the optimal situation, at least for the considered rank $\ell \leq 40$. The right plot of Fig. 3 shows the a-posteriori error for the suboptimal control. By one gradient step the control error first decreases, but then stagnates at this level. Though without any gradient step, the error is higher at the beginning, between 30 and 35 basis functions it jumps down once more and therefore the algorithm can reach the tolerance. However, the right plot shows that the absolute error in state stays far above the OS-POD results. In Fig. 4 we compare the first four POD basis functions obtained either with $u^0 = 0$, u_2 or \bar{u}^{FE} . In the first POD basis function associated with the uncontrolled equation ($u = 0$) we recognize the initial condition; see left plot of Fig. 1. The optimal state is richer in dynamics what is reflected by a different shape of the POD basis functions. After two OS-POD steps the basis has changed significantly and at least the first four modes can hardly be distinguished from the optimal ones. \diamond

Run 2 (Example 2) As a second test, we study a distributed control problem with control and state constraints. In Example 2 we choose $d = 1$, $\nu = 1$, $\beta = -5$, $N_t = 400$ time points in the time interval $[0, 1]$, $N_x = 600$ grid points in the domain $\Omega = [0, 3]$, $m = 50$ control components and $n = 800$, i.e. pointwise state constraints. For the data, we choose $f = 0$, $y_0 = \frac{1}{2}\chi_{[1.2, 1.8]}$ and $y_Q(t, x) = \frac{1}{9}(6x + 6tx - 2x^2)$ for $t < 1 - \frac{1}{3}x$, $y_Q(t, x) = 0$ elsewhere and $\sigma_Q = 1$, $\sigma_\Omega = 0$, $\sigma_u = \sigma_w = \varepsilon = 7.5e-02$. The control and state bounds are $u_a = -1$, $u_b = 4$ and $y_a = 0.05$, $y_b = 0.5$. Compared to the situation in Run 1 additional challenges arise here:

1. If the convection parameter β which resembles the dispersal speed of the initial profile is dominant, a rapid decay of the singular values of the POD operator \mathcal{R} is prevented. This results in a slower decay of the POD error $\ell \mapsto \|\bar{v} - \bar{v}^\ell\|_{\mathcal{V}}$, so larger POD basis ranks are required to ensure a good approximation.

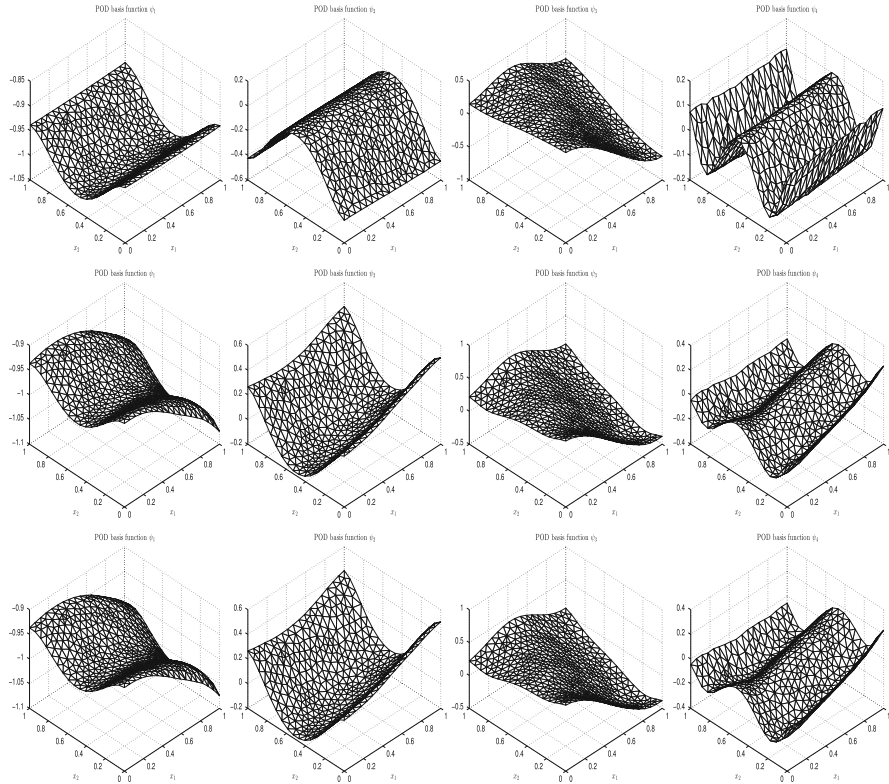


Fig. 4 Run 1: first four POD basis functions associated with the initial control $u^0 = 0$ (top), with the control gained after $k = 2$ OS-POD steps (middle) and with the optimal FE control \bar{u}^{FE} (bottom)

2. The transport term βy_x requires further considerations for the full-order solution techniques. For instance, central differences lead to a stable discretization if $\nu \Delta t \leq \Delta x^2/2$ holds true, but nevertheless, strong oscillations of the discrete solution may occur if the condition $|\beta \Delta x/\nu| < 2$ is violated; see, e.g., [21]. An upwind scheme for βy_x which combines forward and backward differences prevents oscillations, but is only of convergence order one.
3. By evaluation of the a-posteriori error estimator, the active set equations $\bar{u}^\ell = u_a$ and $\bar{u}^\ell = u_b$ defining the control perturbation ζ_u^ℓ are fulfilled exactly by construction since $\bar{v}_1^\ell = \bar{u}^\ell$ holds. This is not the case for the state perturbation ζ_w^ℓ : Here, a high-order solution operation is required to calculate $\bar{v}_2^\ell = \varepsilon \bar{w}^\ell + \mathcal{I} \mathcal{I} \mathcal{B} \bar{u}^\ell$ and to determine the active sets $\bar{v}_2^\ell = \hat{y}_a$ and $\bar{v}_2^\ell = \hat{y}_b$, respectively. We propose to replace the active set equalities by $\|\bar{v}_2^\ell - \hat{y}_{a,b}\|_W < \varepsilon_{acc}$, where ε_{acc} is the accuracy of the full-order model.
4. If the penalized state constraint shall resemble a pointwise pure state constraint, one may choose a fine partition $(\Omega_j)_{1 \leq j \leq n_y} \subseteq \Omega$ of Ω and $\pi_j(x) = |\Omega_j|^{-1}$ for $x \in$

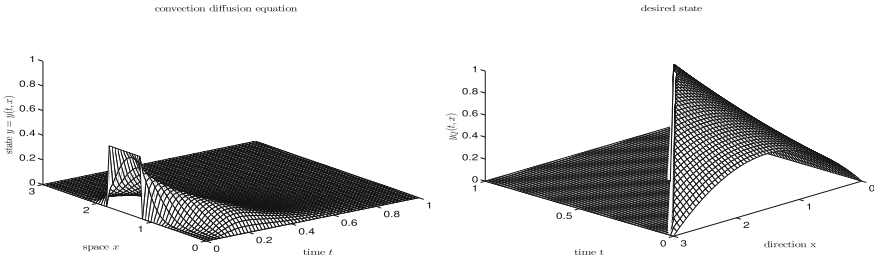


Fig. 5 Run 2: the uncontrolled state (left) and the desired state y_Q (right)

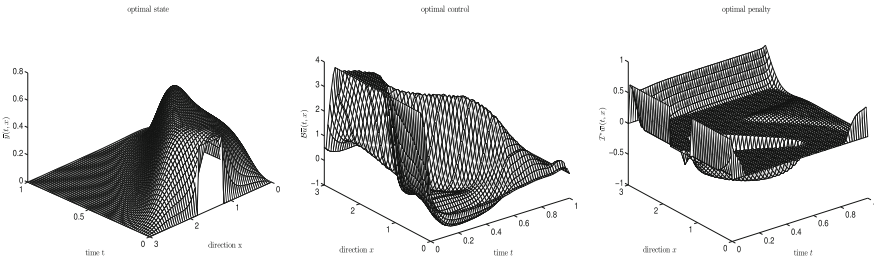


Fig. 6 Run 2: the optimal FE state \bar{y} (left), the optimal FE control \bar{u} and the optimal FE penalty \bar{w} (right) to (\mathbf{P}^ε)

Ω_j as well as 0 otherwise. In this case, we have $(\mathcal{J}_j y)(t) = |\Omega_j|^{-1} \int_{\Omega_j} y(t, \mathbf{x}) \, d\mathbf{x} \approx y(t, \mathbf{x}_j)$. Now, choosing $\varepsilon \ll 1$ and $\sigma_w \gg 1$ ensures $\varepsilon w + \mathcal{J}y \approx y$: The penalty w cannot compensate strong violations of the state constraint any more. A small ε leads to bad condition numbers of the optimality system matrices already for the full-order model which causes not only stability problems, but also less regular state solutions. Since the convergence of POD solutions to the full-order ones require an additional regularity of the snapshot ensemble, a good accuracy of the POD model can be expected only if additional effort is conducted for finding appropriate snapshots.

The uncontrolled FE state is plotted in the left plot of Fig. 5. The discontinuous desired state y_Q is presented in the right plot of Fig. 5. The optimal FE solution to (\mathbf{P}^ε) is shown in Fig. 6. The primal-dual active set strategy (PDASS) required a rather large number of iterations to converge. The complex structure of the active and inactive sets is given in Fig. 7. In this example, 39 updates of the active sets are conducted until the iteration stops after 1217 s. Due to $\sigma_u \gg 0$ as well as the control constraints which prevent that \bar{u} develops singularities and \bar{y} loses regularity, the state solution is smooth. However, $\varepsilon \ll 1$ causes a plateau, where the upper state constraint $y_b = 0.5$ is active. This dynamics which do not occur in the

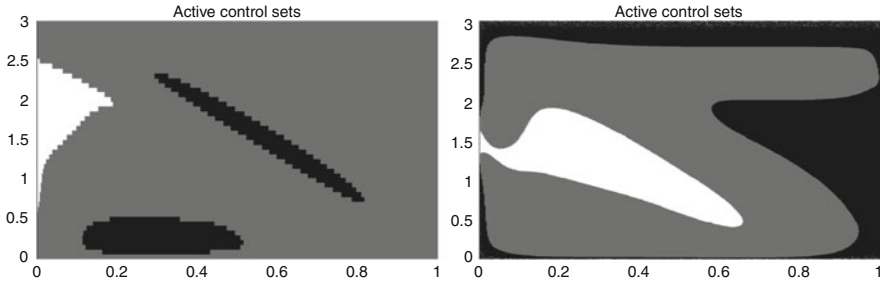


Fig. 7 Run 2: the active sets of the upper bounds (*white*) and the lower bounds (*black*) as well as the inactive regions (*grey*) for the control constraints (*left*) and the mixed penalty-state constraints (*right*)

Table 4 Run 2: error analysis for different numbers of initializing gradient steps

GradSteps	CPU time	Err(v^∇)	$\hat{J}^\ell(v^\nabla)$	$\ell(v^\nabla)$	Iter(ψ^ℓ)	$\ \bar{v}^\ell - \bar{v}^{\text{FE}}\ _\mathcal{V}$	$\ \mathcal{F}^* \zeta^\ell\ _\mathcal{V} / \sigma$
0	10.24 s	–	–	16	12	3.5 e+00	4.4 e+00
1	22.08 s	4.36	37.18	26	18	1.7e–01	2.1e–01
2	29.32 s	3.07	36.52	32	19	1.6e–01	1.9e–01
3	28.64 s	2.65	31.55	32	18	1.6e–01	1.9e–01
4	13.42 s	2.51	31.02	32	8	1.6e–01	1.8e–01
5	15.53 s	2.42	29.59	33	8	1.6e–01	1.9e–01
10	24.55 s	2.17	28.53	33	8	1.5e–01	1.8e–01
25	53.46 s	2.15	28.44	33	8	1.6e–01	1.8e–01

uncontrolled state \hat{y} have to be included in appropriate snapshots to generate an accurate POD basis. Due to the strong convection, projections even on the optimal POD space spanned by the POD elements of the optimal snapshots \bar{y} and \bar{p} cause significant approximation errors if the POD basis rank ℓ is not chosen sufficiently large. Table 4 shows that this procedure does not lead to an adequate model error if state constraints are taken into account. The first row presents the gradient-based error indicator $\text{Err}(v) = \|v - \mathbb{P}_v(v + d_v)\|_\mathcal{V}$ which is our termination criterion for the projected gradient method [15]; the value almost stagnates after circa eight iterations such as the corresponding objective value $\hat{J}^\ell(v)$. The third line presents the POD basis ranks used for the active set strategy. We choose $\ell = \min\{\max\{v \mid \lambda_v > \lambda_{\min}\}, \ell_{\max}\}$ where we set $\lambda_{\min} = 10^{-4}$ and ℓ_{\max} ensures that the model reduction effect does not vanish by using too many POD elements. We see that at least two gradient steps are required to get a sufficiently rich snapshot sample. The next row shows the number of active set updates in the reduced model. Four initializing gradient steps lead to a fast termination of this routine. However, the corresponding errors do not decay below the value 0.15 independent of the number of gradient steps or the chosen basis rank: Here, the gradient steps do not lead to a control u^∇ which is close enough to \bar{u} to guarantee good snapshots for the POD basis. The a-posteriori error bounds $\|\mathcal{F}^* \zeta^\ell\|_\mathcal{V} / \sigma$ turn out to be of the same order as

Table 5 Run 2: error analysis for different numbers of gradient step/active set interactions

PDASS steps	1	2	3	4	(\bar{u}, \bar{w})
CPU time	15.63 s	32.19 s	51.61 s	102.30 s	49.32 s
$\ \bar{v}^\ell - \bar{v}^{\text{FE}}\ _{\mathcal{V}}$	$1.61 \cdot 10^{-1}$	$3.22 \cdot 10^{-2}$	$6.11 \cdot 10^{-3}$	$8.49 \cdot 10^{-4}$	$5.20 \cdot 10^{-4}$
$\ \mathcal{S}^* \zeta^\ell\ _{\mathcal{V}}/\sigma$	$1.90 \cdot 10^{-1}$	$3.53 \cdot 10^{-2}$	$8.24 \cdot 10^{-3}$	$8.61 \cdot 10^{-4}$	$5.26 \cdot 10^{-4}$

the errors themselves. Finally, the calculation times show that the model reduction would be very efficient if the quality of the snapshots could be improved. Table 5 shows that the additional effort leads both to sufficiently small reduction errors and still very efficient calculation times: With three steps, the a-posteriori error estimator guarantees that the reduced order model error is below the discretization error of the full order model. Solving the reduced-order problem lasts 51.61 s with this strategy which is just 4.24 % of the full-order calculation time. \diamond

9 Conclusions

We have presented a combination of adaptive OS-POD basis computation and a-posteriori error estimation for solving linear-quadratic optimal control problems with bilaterally control and state constraints. The considerations started from a basic POD Galerkin approach, where the quality of the reduced order model is controlled by an a-posteriori error estimate. In the context of optimal control it turned out to be important that the POD basis is not computed from arbitrary control and state data, but models more or less their optimal course. We succeed in providing convincing numerical tests for the combination of OSPOD and a-posteriori error analysis.

Acknowledgement This work was supported by the DFG project *A-Posteriori-POD Error Estimators for Nonlinear Optimal Control Problems governed by Partial Differential Equations*, grant VO 1658/2-1.

References

1. Arian, E., Fahl, M., Sachs, E.W.: Trust-region proper orthogonal decomposition for flow control. Technical Report 2000-25, ICASE (2000)
2. Benner, P., Mehrmann, V., Sorensen, D.C.: Dimension Reduction of Large-Scale Systems. Lecture Notes in Computational Science and Engineering, vol. 45. Springer, Berlin (2005)
3. Chapelle, D., Gariah, A., Saint-Marie, J.: Galerkin approximation with proper orthogonal decomposition: new error estimates and illustrative examples. ESAIM: Math. Model. Numer. Anal. **46**, 731–757 (2012)
4. Dautray, R., Lions, J.-L.: Mathematical Analysis and Numerical Methods for Science and Technology. Volume 5: Evolution Problems I. Springer, Berlin (2000)
5. Dontchev, A.L., Hager, W.W., Poore, A.B., Yang, B.: Optimality, stability, and convergence in nonlinear control. Appl. Math. Optim. **31**, 297–326 (1995)

6. Grimm, E.: Optimality system POD and a-posteriori error analysis for linear-quadratic optimal control problems. Master Thesis, University of Konstanz (2013). <https://kops.uni-konstanz.de/handle/123456789/27761>
7. Gubisch, M., Volkwein, S.: Proper orthogonal decomposition for linear-quadratic optimal control. (2013, submitted). <http://kops.uni-konstanz.de/handle/123456789/25037>
8. Gubisch, M., Volkwein, S.: POD a-posteriori error analysis for optimal control problems with mixed control-state constraints. *Comput. Optim. Appl.* **58**, 619–644 (2014)
9. Hintermüller, M., Ito, K., Kunisch, K.: The primal-dual active set strategy as a semismooth Newton method. *SIAM J. Optim.* **13**, 865–888 (2003)
10. Hintermüller, M., Kopacka, I., Volkwein, S.: Mesh-independence and preconditioning for solving control problems with mixed control-state constraints. *ESAIM: Control Optim. Calc. Var.* **15**, 626–652 (2009)
11. Holmes, P., Lumley, J.L., Berkooz, G., Rowley, C.W.: *Turbulence, Coherent Structures, Dynamical Systems and Symmetry*. Cambridge Monographs on Mechanics, 2nd edn. Cambridge University Press, Cambridge (2012)
12. Hinze, M., Volkwein, S.: Proper orthogonal decomposition surrogate models for nonlinear dynamical systems: error estimates and suboptimal control. Chapter 10 of [2]
13. Hinze, M., Volkwein, S.: Error estimates for abstract linear-quadratic optimal control problems using proper orthogonal decomposition. *Comput. Optim. Appl.* **39**, 319–345 (2008)
14. Hinze, M., Pinnau, R., Ulbrich, M., Ulbrich, S.: *Optimization with PDE Constraints*. Springer, Berlin (2009)
15. Kelley, C.T.: *Iterative Methods for Optimization*. SIAM Frontiers in Applied Mathematics, Philadelphia (1999)
16. Kunisch, K., Volkwein, S.: Galerkin proper orthogonal decomposition methods for a general equation in fluid dynamics. *SIAM J. Numer. Anal.* **40**, 492–515 (2002)
17. Kunisch, K., Volkwein, S.: Proper orthogonal decomposition for optimality systems. *ESAIM: Math. Model. Numer. Anal.* **42**, 1–23 (2008)
18. Müller, M.: *Uniform convergence of the POD method and applications to optimal control*. Ph.D Thesis, University of Graz (2011)
19. Schilders, W.H.A., van der Vorst, H.A., Rommes, J.: *Model Order Reduction: Theory, Research Aspects and Applications*. Mathematics in Industry, vol. 13. Springer, Berlin (2008)
20. Singler, J.R.: New POD expressions, error bounds, and asymptotic results for reduced order models of parabolic PDEs. *SIAM J. Numer. Anal.* **52**, 852–876 (2014)
21. Strikwerda, J.: *Finite Difference Schemes and Partial Differential Equations*. SIAM, Philadelphia (2004)
22. Studinger, A., Volkwein, S.: Numerical analysis of POD a-posteriori error estimation for optimal control. *Int. Ser. Numer. Math.* **164**, 137–158 (2013)
23. Tröltzsch, F.: Regular Lagrange multipliers for control problems with mixed pointwise control-state constraints. *SIAM J. Optim.* **22**, 616–634 (2005)
24. Tröltzsch, F.: *Optimal Control of Partial Differential Equations. Theory, Methods and Applications*, vol. 112. American Mathematical Society, Providence (2010)
25. Tröltzsch, F., Volkwein, S.: POD a-posteriori error estimates for linear-quadratic optimal control problems. *Comput. Optim. Appl.* **44**, 83–115 (2009)
26. Volkwein, S.: Optimality system POD and a-posteriori error analysis for linear-quadratic problems. *Control. Cybern.* **40**, 1109–1125 (2011)