

# Proper Orthogonal Decomposition for Linear-Quadratic Optimal Control

Martin Gubisch and Stefan Volkwein

**Mathematics Subject Classification (2010).** 35K20, 35K90, 49K20, 65K05, 65Nxx.

**Keywords.** Proper orthogonal decomposition, model-order reduction, a-priori and a-posteriori error analysis, linear-quadratic optimal control, primal-dual active set strategy.

## 1. Introduction

Optimal control problems for partial differential equation are often hard to tackle numerically because their discretization leads to very large scale optimization problems. Therefore, different techniques of model reduction were developed to approximate these problems by smaller ones that are tractable with less effort.

*Balanced truncation* [2, 66, 81] is one well studied model reduction technique for state-space systems. This method utilizes the solutions to two Lyapunov equations, the so-called controllability and observability Gramians. The balanced truncation method is based on transforming the state-space system into a balanced form so that its controllability and observability Gramians become diagonal and equal. Moreover, the states that are difficult to reach or to observe, are truncated. The advantage of this method is that it preserves the asymptotic stability in the reduced-order system. Furthermore, a-priori error bounds are available. Recently, the theory of balanced truncation model reduction was extended to descriptor systems; see, e.g., [50] and [21].

Recently the application of *reduced-order models* to linear time varying and nonlinear systems, in particular to nonlinear control systems, has received

---

The authors gratefully acknowledges support by the German Science Fund DFG grant VO 1658/2-1 *A-posteriori-POD Error Estimators for Nonlinear Optimal Control Problems governed by Partial Differential Equations*. The first author is further supported by the Landesgraduiertenförderung of Baden-Württemberg.

an increasing amount of attention. The reduced-order approach is based on projecting the dynamical system onto subspaces consisting of basis elements that contain characteristics of the expected solution. This is in contrast to, e.g., finite element techniques (see, e.g., [7], where the basis elements of the subspaces do not relate to the physical properties of the system that they approximate. The *reduced basis* (RB) method, as developed in [20, 56] and [32], is one such reduced-order method, where the basis elements correspond to the dynamics of expected control regimes. Let us refer to the [14, 23, 51, 55] for the successful use of reduced basis method in PDE constrained optimization problems. Currently, *Proper orthogonal decomposition* (POD) is probably the mostly used and most successful model reduction technique for nonlinear optimal control problems, where the basis functions contain information from the solutions of the dynamical system at pre-specified time-instances, so-called snapshots; see, e.g., [8, 31, 69, 77]. Due to a possible linear dependence or almost linear dependence the snapshots themselves are not appropriate as a basis. Hence a singular value decomposition is carried out and the leading generalized eigenfunctions are chosen as a basis, referred to as the POD basis. POD is successfully used in a variety of fields including fluid dynamics, coherent structures [1, 3] and inverse problems [6]. Moreover in [5] POD is successfully applied to compute reduced-order controllers. The relationship between POD and balancing was considered in [46, 63, 79]. An error analysis for nonlinear dynamical systems in finite dimensions were carried out in [60] and a missing point estimation in models described by POD was studied in [4].

Reduced order models are used in PDE-constrained optimization in various ways; see, e.g., [28, 65] for a survey. In optimal control problems it is sometimes necessary to compute a feedback control law instead of a fixed optimal control. In the implementation of these feedback laws models of reduced-order can play an important and very useful role, see [5, 45, 48, 61]. Another useful application is the use in optimization problems, where a PDE solver is part of the function evaluation. Obviously, thinking of a gradient evaluation or even a step-size rule in the optimization algorithm, an expensive function evaluation leads to an enormous amount of computing time. Here, the reduced-order model can replace the system given by a PDE in the objective function. It is quite common that a PDE can be replaced by a five- or ten-dimensional system of ordinary differential equations. This results computationally in a very fast method for optimization compared to the effort for the computation of a single solution of a PDE. There is a large amount of literature in engineering applications in this regard, we mention only the papers [49, 52]. Recent applications can also be found in finance using the RB model [58] and the POD model [64, 67] in the context of calibration for models in option pricing.

In the present work we use POD for deriving low order models of dynamical systems. These low order models then serve as surrogates for the dynamical system in the optimization process. We consider a linear-quadratic

optimal control problem in an abstract setting and prove error estimates for the POD Galerkin approximations of the optimal control. This is achieved by combining techniques from [11, 12, 25] and [40, 41]. For nonlinear problems we refer the reader to [28, 57, 65]. However, unless the snapshots are generating a sufficiently rich state space or are computed from the exact (unknown) optimal controls, it is not a-priorly clear how far the optimal solution of the POD problem is from the exact one. On the other hand, the POD method is a universal tool that is applicable also to problems with time-dependent coefficients or to nonlinear equations. Moreover, by generating snapshots from the real (large) model, a space is constructed that inherits the main and relevant physical properties of the state system. This, and its ease of use makes POD very competitive in practical use, despite of a certain heuristic flavor. In this work, we review results for a POD a-posteriori analysis; see, e.g., [73] and [18, 35, 36, 70, 71, 76, 78]. We use a fairly standard perturbation method to deduce how far the suboptimal control, computed on the basis of the POD model, is from the (unknown) exact one. This idea turned out to be very efficient in our examples. It is able to compensate for the lack of a priori analysis for POD methods. Let us also refer to the papers [13, 19, 51], where a-posteriori error bounds are computed for linear-quadratic optimal control problems approximated by the reduced basis method.

The manuscript is organised in the following manner: In Section 2 we introduce the method of POD in real, separable Hilbert spaces and discuss its relationship to the singular value decomposition. We distinguish between two versions of the POD method: the discrete and the continuous one. Reduced-order modelling with POD is carried out in Section 3. The error between the exact solution and its POD approximation is investigated by an a-priori error analysis. In Section 4 we study quadratic optimal control problems governed by linear evolution problems and bilateral inequality constraints. These problems are infinite dimensional, convex optimization problems. Their optimal solutions are characterised by first-order optimality conditions. POD Galerkin discretizations of the optimality conditions are introduced and analysed. By an a-priori error analysis the error of the exact optimal control and its POD suboptimal approximation is estimated. For the error control in the numerical realisations we make use of an a-posteriori error analysis, which turns out to be very efficient in our numerical examples, which are presented in Section 5.

## 2. The POD method

Throughout we suppose that  $X$  is a real Hilbert space endowed with the inner product  $\langle \cdot, \cdot \rangle_X$  and the associated induced norm  $\| \cdot \|_X = \langle \cdot, \cdot \rangle_X^{1/2}$ . Furthermore, we assume that  $X$  is *separable*, i.e.,  $X$  has a countable dense subset. This implies that  $X$  possesses a countable orthonormal basis; see, e.g., [62, p. 47]. For the POD method in complex Hilbert spaces we refer to [75], for instance.

## 2.1. The discrete variant of the POD method

For fixed  $n, \wp \in \mathbb{N}$  let the so-called *snapshots*  $y_1^k, \dots, y_n^k \in X$  be given for  $1 \leq k \leq \wp$ . To avoid a trivial case we suppose that at least one of the  $y_j^k$ 's is nonzero. Then, we introduce the finite dimensional, linear subspace

$$\mathcal{V}^n = \text{span} \left\{ y_j^k \mid 1 \leq j \leq n \text{ and } 1 \leq k \leq \wp \right\} \subset X \quad (2.1)$$

with dimension  $d^n \in \{1, \dots, n\wp\} < \infty$ . We call the set  $\mathcal{V}^n$  *snapshot subspace*. In Section 2.3 we consider the case, where the number  $n$  is varied. Therefore, we emphasize this dependence by using the super index  $n$ . We distinguish two cases:

- 1) The separable Hilbert space  $X$  has finite dimension  $m$ . Then,  $X$  is isomorphic to  $\mathbb{R}^m$ ; see, e.g., [62, p. 47]. We set  $J = \{1, \dots, m\}$ . Clearly, we have  $d^n \leq \min(n\wp, m)$ .
- 2) Since  $X$  is separable, each orthonormal basis of  $X$  has countably many elements. In this case  $X$  is isomorphic to the set  $\ell_2$  of sequences  $\{x_i\}_{i \in \mathbb{N}}$  of real numbers which satisfy  $\sum_{i=1}^{\infty} |x_i|^2 < \infty$ ; see [62, p. 47], for instance. Then, we define  $J = \mathbb{N}$ .

The method of POD consists in choosing a complete orthonormal basis  $\{\psi_i\}_{i \in J}$  in  $X$  such that for every  $\ell \in \{1, \dots, d^n\}$  the mean square error between the  $n\wp$  elements  $y_j^k$  and their corresponding  $\ell$ -th partial Fourier sum is minimized on average:

$$\begin{cases} \min \sum_{k=1}^{\wp} \sum_{j=1}^n \alpha_j^n \left\| y_j^k - \sum_{i=1}^{\ell} \langle y_j^k, \psi_i \rangle_X \psi_i \right\|_X^2 \\ \text{s.t. } \{\psi_i\}_{i=1}^{\ell} \subset X \text{ and } \langle \psi_i, \psi_j \rangle_X = \delta_{ij}, \quad 1 \leq i, j \leq \ell, \end{cases} \quad (\mathbf{P}_n^{\ell})$$

where the  $\alpha_j^n$ 's denote positive weighting parameters. Here, the symbol  $\delta_{ij}$  denotes the Kronecker symbol satisfying  $\delta_{ii} = 1$  and  $\delta_{ij} = 0$  for  $i \neq j$ . An optimal solution  $\{\bar{\psi}_i^n\}_{i=1}^{\ell}$  to  $(\mathbf{P}_n^{\ell})$  is called a *POD basis of rank  $\ell$* . Notice that

$$\begin{aligned} & \left\| y_j^k - \sum_{i=1}^{\ell} \langle y_j^k, \psi_i \rangle_X \psi_i \right\|_X^2 \\ &= \left\langle y_j^k - \sum_{i=1}^{\ell} \langle y_j^k, \psi_i \rangle_X \psi_i, y_j^k - \sum_{l=1}^{\ell} \langle y_j^k, \psi_l \rangle_X \psi_l \right\rangle_X \\ &= \|y_j^k\|_X^2 - 2 \sum_{i=1}^{\ell} \langle y_j^k, \psi_i \rangle_X^2 + \sum_{i=1}^{\ell} \sum_{l=1}^{\ell} \langle y_j^k, \psi_i \rangle_X \langle y_j^k, \psi_l \rangle_X \langle \psi_i, \psi_l \rangle_X \\ &= \|y_j^k\|_X^2 - \sum_{i=1}^{\ell} \langle y_j^k, \psi_i \rangle_X^2 \end{aligned} \quad (2.2)$$

holds for any set  $\{\psi_i\}_{i=1}^\ell \subset X$  satisfying  $\langle \psi_i, \psi_j \rangle_X = \delta_{ij}$ . Thus,  $(\mathbf{P}_n^\ell)$  is equivalent with the maximization problem

$$\begin{cases} \max \sum_{k=1}^{\wp} \sum_{j=1}^n \alpha_j^n \sum_{i=1}^{\ell} \langle y_j^k, \psi_i \rangle_X^2 \\ \text{s.t. } \{\psi_i\}_{i=1}^\ell \subset X \text{ and } \langle \psi_i, \psi_j \rangle_X = \delta_{ij}, 1 \leq i, j \leq \ell. \end{cases} \quad (\hat{\mathbf{P}}_n^\ell)$$

Suppose that  $\{\psi_i\}_{i \in \mathcal{J}}$  is a complete orthonormal basis in  $X$ . Since  $X$  is separable, any  $y_j^k \in X$ ,  $1 \leq j \leq n$  and  $1 \leq k \leq \wp$ , can be written as

$$y_j^k = \sum_{i \in \mathcal{J}} \langle y_j^k, \psi_i \rangle_X \psi_i \quad (2.3)$$

and the (probably infinite) sum converges for all snapshots (even for all elements in  $X$ ). Thus, the POD basis  $\{\psi_i^n\}_{i=1}^\ell$  of rank  $\ell$  maximizes the absolute values of the first  $\ell$  Fourier coefficients  $\langle y_j^k, \psi_i \rangle_X$  for all  $n\wp$  snapshots  $y_j^k$  in an average sense. Let us recall the following definition for linear operators in Banach spaces.

**Definition 2.1.** Let  $\mathcal{B}_1, \mathcal{B}_2$  be two real Banach spaces. The operator  $\mathcal{T} : \mathcal{B}_1 \rightarrow \mathcal{B}_2$  is called a linear, bounded operator if these conditions are satisfied:

- 1)  $\mathcal{T}(\alpha u + \beta v) = \alpha \mathcal{T}u + \beta \mathcal{T}v$  for all  $\alpha, \beta \in \mathbb{R}$  and  $u, v \in \mathcal{B}_1$ .
- 2) There exists a constant  $c > 0$  such that  $\|\mathcal{T}u\|_{\mathcal{B}_2} \leq c \|u\|_{\mathcal{B}_1}$  for all  $u \in \mathcal{B}_1$ .

The set of all linear, bounded operators from  $\mathcal{B}_1$  to  $\mathcal{B}_2$  is denoted by  $\mathcal{L}(\mathcal{B}_1, \mathcal{B}_2)$  which is a Banach space equipped with the operator norm [62, pp. 69-70]

$$\|\mathcal{T}\|_{\mathcal{L}(\mathcal{B}_1, \mathcal{B}_2)} = \sup_{\|u\|_{\mathcal{B}_1}=1} \|\mathcal{T}u\|_{\mathcal{B}_2} \quad \text{for } \mathcal{T} \in \mathcal{L}(\mathcal{B}_1, \mathcal{B}_2).$$

If  $\mathcal{B}_1 = \mathcal{B}_2$  holds, we briefly write  $\mathcal{L}(\mathcal{B}_1)$  instead of  $\mathcal{L}(\mathcal{B}_1, \mathcal{B}_2)$ . The dual mapping  $\mathcal{T}' : \mathcal{B}'_2 \rightarrow \mathcal{B}'_1$  of an operator  $\mathcal{T} \in \mathcal{L}(\mathcal{B}_1, \mathcal{B}_2)$  is defined as

$$\langle \mathcal{T}'f, u \rangle_{\mathcal{B}'_1, \mathcal{B}_1} = \langle f, \mathcal{T}u \rangle_{\mathcal{B}'_2, \mathcal{B}_2} \quad \text{for all } (u, f) \in \mathcal{B}_1 \times \mathcal{B}'_2,$$

where, for instance,  $\langle \cdot, \cdot \rangle_{\mathcal{B}'_1, \mathcal{B}_1}$  denotes the dual pairing of the space  $\mathcal{B}_1$  with its dual space  $\mathcal{B}'_1 = \mathcal{L}(\mathcal{B}_1, \mathbb{R})$ .

Let  $\mathcal{H}_1$  and  $\mathcal{H}_2$  denote two real Hilbert spaces. For a given  $\mathcal{T} \in \mathcal{L}(\mathcal{H}_1, \mathcal{H}_2)$  the adjoint operator  $\mathcal{T}^* : \mathcal{H}_2 \rightarrow \mathcal{H}_1$  is uniquely defined by

$$\langle \mathcal{T}^*v, u \rangle_{\mathcal{H}_1} = \langle v, \mathcal{T}u \rangle_{\mathcal{H}_2} = \langle \mathcal{T}u, v \rangle_{\mathcal{H}_2} \quad \text{for all } (u, v) \in \mathcal{H}_1 \times \mathcal{H}_2.$$

Let  $\mathcal{J}_i : \mathcal{H}_i \rightarrow \mathcal{H}'_i$ ,  $i = 1, 2$ , denote the Riesz isomorphisms satisfying

$$\langle u, v \rangle_{\mathcal{H}_i} = \langle \mathcal{J}_i u, v \rangle_{\mathcal{H}'_i, \mathcal{H}_i} \quad \text{for all } v \in \mathcal{H}_i.$$

Then, we have the representation  $\mathcal{T}^* = \mathcal{J}_1^{-1} \mathcal{T}' \mathcal{J}_2$ ; see [72, p. 186]. Moreover,  $(\mathcal{T}^*)^* = \mathcal{T}$  for every  $\mathcal{T} \in \mathcal{L}(\mathcal{H}_1, \mathcal{H}_2)$ . If  $\mathcal{T} = \mathcal{T}^*$  holds,  $\mathcal{T}$  is said to be *selfadjoint*. The operator  $\mathcal{T} \in \mathcal{L}(\mathcal{H}_1, \mathcal{H}_2)$  is called *nonnegative* if  $\langle \mathcal{T}u, u \rangle_{\mathcal{H}_2} \geq 0$  for all  $u \in \mathcal{H}_1$ . Finally,  $\mathcal{T} \in \mathcal{L}(\mathcal{H}_1, \mathcal{H}_2)$  is called *compact* if for every bounded sequence  $\{u_n\}_{n \in \mathbb{N}} \subset \mathcal{H}_1$  the sequence  $\{\mathcal{T}u_n\}_{n \in \mathbb{N}} \subset \mathcal{H}_2$  contains a convergent subsequence.

Now we turn to  $(\mathbf{P}_n^\ell)$  and  $(\hat{\mathbf{P}}_n^\ell)$ . We make use of the following lemma.

**Lemma 2.2.** *Let  $X$  be a (separable) real Hilbert space and  $y_1^k, \dots, y_n^k \in X$  are given snapshots for  $1 \leq k \leq \wp$ . Define the linear operator  $\mathcal{R}^n : X \rightarrow X$  as follows:*

$$\mathcal{R}^n \psi = \sum_{k=1}^{\wp} \sum_{j=1}^n \alpha_j^n \langle \psi, y_j^k \rangle_X y_j^k \quad \text{for } \psi \in X \quad (2.4)$$

with positive weights  $\alpha_1^n, \dots, \alpha_n^n$ . Then,  $\mathcal{R}^n$  is a compact, nonnegative and selfadjoint operator.

*Proof.* It is clear that  $\mathcal{R}^n$  is a linear operator. From

$$\|\mathcal{R}^n \psi\|_X \leq \sum_{k=1}^{\wp} \sum_{j=1}^n \alpha_j^n |\langle \psi, y_j^k \rangle_X| \|y_j^k\|_X \quad \text{for } \psi \in X$$

and the Cauchy-Schwarz inequality [62, p. 38]

$$|\langle \varphi, \phi \rangle_X| \leq \|\varphi\|_X \|\phi\|_X \quad \text{for } \varphi, \phi \in X$$

we conclude that  $\mathcal{R}^n$  is bounded. Since  $\mathcal{R}^n \psi \in \mathcal{V}^n$  holds for all  $\psi \in X$ , the range of  $\mathcal{R}^n$  is finite dimensional. Thus,  $\mathcal{R}^n$  is a finite rank operator which is compact; see [62, p. 199]. Next we show that  $\mathcal{R}^n$  is nonnegative. For that purpose we choose an arbitrary element  $\psi \in X$  and consider

$$\langle \mathcal{R}^n \psi, \psi \rangle_X = \sum_{k=1}^{\wp} \sum_{j=1}^n \alpha_j^n \langle \psi, y_j^k \rangle_X \langle y_j^k, \psi \rangle_X = \sum_{k=1}^{\wp} \sum_{j=1}^n \alpha_j^n \langle \psi, y_j^k \rangle_X^2 \geq 0.$$

Thus,  $\mathcal{R}^n$  is nonnegative. For any  $\psi, \tilde{\psi} \in X$  we derive

$$\begin{aligned} \langle \mathcal{R}^n \psi, \tilde{\psi} \rangle_X &= \sum_{k=1}^{\wp} \sum_{j=1}^n \alpha_j^n \langle \psi, y_j^k \rangle_X \langle y_j^k, \tilde{\psi} \rangle_X = \sum_{k=1}^{\wp} \sum_{j=1}^n \alpha_j^n \langle \tilde{\psi}, y_j^k \rangle_X \langle y_j^k, \psi \rangle_X \\ &= \langle \mathcal{R}^n \tilde{\psi}, \psi \rangle_X = \langle \psi, \mathcal{R}^n \tilde{\psi} \rangle_X. \end{aligned}$$

Thus,  $\mathcal{R}^n$  is selfadjoint. □

Next we recall some important results from the spectral theory of operators (on infinite dimensional spaces). We begin with the following definition; see [62, Section VI.3].

**Definition 2.3.** *Let  $\mathcal{H}$  be a real Hilbert space and  $\mathcal{T} \in \mathcal{L}(\mathcal{H})$ .*

- 1) *A complex number  $\lambda$  belongs to the resolvent set  $\rho(\mathcal{T})$  if  $\lambda \mathcal{I} - \mathcal{T}$  is a bijection with a bounded inverse. Here,  $\mathcal{I} \in \mathcal{L}(\mathcal{H})$  stands for the identity operator. If  $\lambda \notin \rho(\mathcal{T})$ , then  $\lambda$  is an element of the spectrum  $\sigma(\mathcal{T})$  of  $\mathcal{T}$ .*
- 2) *Let  $u \neq 0$  be a vector with  $\mathcal{T}u = \lambda u$  for some  $\lambda \in \mathbb{C}$ . Then,  $u$  is said to be an eigenvector of  $\mathcal{T}$ . We call  $\lambda$  the corresponding eigenvalue. If  $\lambda$  is an eigenvalue, then  $\lambda \mathcal{I} - \mathcal{T}$  is not injective. This implies  $\lambda \in \sigma(\mathcal{T})$ . The set of all eigenvalues is called the point spectrum of  $\mathcal{T}$ .*

We will make use of the next two essential theorems for compact operators; see [62, p. 203].

**Theorem 2.4 (Riesz-Schauder).** *Let  $\mathcal{H}$  be a real Hilbert space and  $\mathcal{T} : \mathcal{H} \rightarrow \mathcal{H}$  a linear, compact operator. Then the spectrum  $\sigma(\mathcal{T})$  is a discrete set having no limit points except perhaps 0. Furthermore, the space of eigenvectors corresponding to each nonzero  $\lambda \in \sigma(\mathcal{T})$  is finite dimensional.*

**Theorem 2.5 (Hilbert-Schmidt).** *Let  $\mathcal{H}$  be a real separable Hilbert space and  $\mathcal{T} : \mathcal{H} \rightarrow \mathcal{H}$  a linear, compact, selfadjoint operator. Then, there is a sequence of eigenvalues  $\{\lambda_i\}_{i \in \mathcal{J}}$  and of an associated complete orthonormal basis  $\{\psi_i\}_{i \in \mathcal{J}} \subset X$  satisfying*

$$\mathcal{T}\psi_i = \lambda_i\psi_i \quad \text{and} \quad \lambda_i \rightarrow 0 \text{ as } i \rightarrow \infty.$$

Since  $X$  is a separable real Hilbert space and  $\mathcal{R}^n : X \rightarrow X$  is a linear, compact, nonnegative, selfadjoint operator (see Lemma 2.2), we can utilize Theorems 2.4 and 2.5: there exist a complete countable orthonormal basis  $\{\bar{\psi}_i^n\}_{i \in \mathcal{J}}$  and a corresponding sequence of real eigenvalues  $\{\bar{\lambda}_i^n\}_{i \in \mathcal{J}}$  satisfying

$$\mathcal{R}^n \bar{\psi}_i^n = \bar{\lambda}_i^n \bar{\psi}_i^n, \quad \bar{\lambda}_1^n \geq \dots \geq \bar{\lambda}_{d^n} > \bar{\lambda}_{d^n+1} = \dots = 0. \quad (2.5)$$

The spectrum of  $\mathcal{R}^n$  is a pure point spectrum except for possibly 0. Each nonzero eigenvalue of  $\mathcal{R}^n$  has finite multiplicity and 0 is the only possible accumulation point of the spectrum of  $\mathcal{R}^n$ .

**Remark 2.6.** From (2.4), (2.5) and  $\|\psi\|_X = 1$  we infer that

$$\begin{aligned} \sum_{k=1}^{\wp} \sum_{j=1}^n \alpha_j^n \langle y_j^k, \bar{\psi}_i^n \rangle_X^2 &= \left\langle \sum_{k=1}^{\wp} \sum_{j=1}^n \alpha_j^n \langle y_j^k, \bar{\psi}_i^n \rangle_X y_j^k, \bar{\psi}_i^n \right\rangle_X \\ &= \langle \mathcal{R}^n \bar{\psi}_i^n, \bar{\psi}_i^n \rangle_X = \bar{\lambda}_i^n \quad \text{for any } i \in \mathcal{J}. \end{aligned} \quad (2.6)$$

In particular, it follows that

$$\sum_{k=1}^{\wp} \sum_{j=1}^n \alpha_j^n \langle y_j^k, \bar{\psi}_i^n \rangle_X^2 = 0 \quad \text{for all } i > d^n. \quad (2.7)$$

Since  $\{\bar{\psi}_i^n\}_{i \in \mathcal{J}}$  is a complete orthonormal basis and  $\|y_j^k\|_X < \infty$  holds for  $1 \leq k \leq \wp$ ,  $1 \leq j \leq n$ , we derive from (2.6) and (2.7) that

$$\begin{aligned} \sum_{k=1}^{\wp} \sum_{j=1}^n \alpha_j^n \|y_j^k\|_X^2 &= \sum_{k=1}^{\wp} \sum_{j=1}^n \alpha_j^n \sum_{\nu \in \mathcal{J}} \langle y_j^k, \bar{\psi}_\nu^n \rangle_X^2 \\ &= \sum_{\nu \in \mathcal{J}} \sum_{k=1}^{\wp} \sum_{j=1}^n \alpha_j^n \langle y_j^k, \bar{\psi}_\nu^n \rangle_X^2 = \sum_{i \in \mathcal{J}} \bar{\lambda}_i^n = \sum_{i=1}^{d^n} \bar{\lambda}_i^n. \end{aligned} \quad (2.8)$$

By (2.8) the (probably infinite) sum  $\sum_{i \in \mathcal{J}} \bar{\lambda}_i^n$  is bounded. It follows from (2.2) that the objective of  $(\mathbf{P}_n^\ell)$  can be written as

$$\sum_{k=1}^{\wp} \sum_{j=1}^n \alpha_j^n \left\| y_j^k - \sum_{i=1}^{\ell} \langle y_j^k, \psi_i \rangle_X \psi_i \right\|_X^2 = \sum_{i=1}^{d^n} \bar{\lambda}_i^n - \sum_{k=1}^{\wp} \sum_{j=1}^n \alpha_j^n \sum_{i=1}^{\ell} \langle y_j^k, \psi_i \rangle_X^2 \quad (2.9)$$

which we will use in the proof of Theorem 2.7.  $\diamond$

Now we can formulate the main result for  $(\mathbf{P}_n^\ell)$  and  $(\hat{\mathbf{P}}_n^\ell)$ .

**Theorem 2.7.** *Let  $X$  be a separable real Hilbert space,  $y_1^k, \dots, y_n^k \in X$  for  $1 \leq k \leq \wp$  and  $\mathcal{R}^n : X \rightarrow X$  be defined by (2.4). Suppose that  $\{\bar{\lambda}_i^n\}_{i \in \mathcal{J}}$  and  $\{\bar{\psi}_i^n\}_{i \in \mathcal{J}}$  denote the nonnegative eigenvalues and associated orthonormal eigenfunctions of  $\mathcal{R}^n$  satisfying (2.5). Then, for every  $\ell \in \{1, \dots, d^n\}$  the first  $\ell$  eigenfunctions  $\{\bar{\psi}_i^n\}_{i=1}^\ell$  solve  $(\mathbf{P}_n^\ell)$  and  $(\hat{\mathbf{P}}_n^\ell)$ . Moreover, the value of the cost evaluated at the optimal solution  $\{\bar{\psi}_i^n\}_{i=1}^\ell$  satisfies*

$$\sum_{k=1}^{\wp} \sum_{j=1}^n \alpha_j^n \left\| y_j^k - \sum_{i=1}^{\ell} \langle y_j^k, \bar{\psi}_i^n \rangle_X \bar{\psi}_i^n \right\|_X^2 = \sum_{i=\ell+1}^{d^n} \bar{\lambda}_i^n \quad (2.10)$$

and

$$\sum_{k=1}^{\wp} \sum_{j=1}^n \alpha_j^n \sum_{i=1}^{\ell} \langle y_j^k, \bar{\psi}_i^n \rangle_X^2 = \sum_{i=1}^{\ell} \bar{\lambda}_i^n. \quad (2.11)$$

*Proof.* We prove the claim for  $(\hat{\mathbf{P}}_n^\ell)$  by finite induction over  $\ell \in \{1, \dots, d^n\}$ .

- 1) The base case: Let  $\ell = 1$  and  $\psi \in X$  with  $\|\psi\|_X = 1$ . Since  $\{\bar{\psi}_\nu^n\}_{\nu \in \mathcal{J}}$  is a complete orthonormal basis in  $X$ , we have the representation

$$\psi = \sum_{\nu \in \mathcal{J}} \langle \psi, \bar{\psi}_\nu^n \rangle_X \bar{\psi}_\nu^n. \quad (2.12)$$

Inserting this expression for  $\psi$  in the objective of  $(\hat{\mathbf{P}}_n^\ell)$  we find that

$$\begin{aligned} \sum_{k=1}^{\wp} \sum_{j=1}^n \alpha_j^n \langle y_j^k, \psi \rangle_X^2 &= \sum_{k=1}^{\wp} \sum_{j=1}^n \alpha_j^n \left\langle y_j^k, \sum_{\nu \in \mathcal{J}} \langle \psi, \bar{\psi}_\nu^n \rangle_X \bar{\psi}_\nu^n \right\rangle_X^2 \\ &= \sum_{k=1}^{\wp} \sum_{j=1}^n \alpha_j^n \sum_{\nu \in \mathcal{J}} \sum_{\mu \in \mathcal{J}} \left( \langle y_j^k, \langle \psi, \bar{\psi}_\nu^n \rangle_X \bar{\psi}_\nu^n \rangle_X \langle y_j^k, \langle \psi, \bar{\psi}_\mu^n \rangle_X \bar{\psi}_\mu^n \rangle_X \right) \\ &= \sum_{k=1}^{\wp} \sum_{j=1}^n \alpha_j^n \sum_{\nu \in \mathcal{J}} \sum_{\mu \in \mathcal{J}} \left( \langle y_j^k, \bar{\psi}_\nu^n \rangle_X \langle y_j^k, \bar{\psi}_\mu^n \rangle_X \langle \psi, \bar{\psi}_\nu^n \rangle_X \langle \psi, \bar{\psi}_\mu^n \rangle_X \right) \\ &= \sum_{\nu \in \mathcal{J}} \sum_{\mu \in \mathcal{J}} \left( \left\langle \sum_{k=1}^{\wp} \sum_{j=1}^n \alpha_j^n \langle y_j^k, \bar{\psi}_\nu^n \rangle_X y_j^k, \bar{\psi}_\mu^n \right\rangle_X \langle \psi, \bar{\psi}_\nu^n \rangle_X \langle \psi, \bar{\psi}_\mu^n \rangle_X \right). \end{aligned}$$

Utilizing (2.4), (2.5) and  $\|\bar{\psi}_\nu^n\|_X = 1$  we find that

$$\begin{aligned} \sum_{k=1}^{\wp} \sum_{j=1}^n \alpha_j^n \langle y_j^k, \psi \rangle_X^2 &= \sum_{\nu \in \mathcal{J}} \sum_{\mu \in \mathcal{J}} \left( \langle \bar{\lambda}_\nu^n \bar{\psi}_\nu^n, \bar{\psi}_\mu^n \rangle_X \langle \psi, \bar{\psi}_\nu^n \rangle_X \langle \psi, \bar{\psi}_\mu^n \rangle_X \right) \\ &= \sum_{\nu \in \mathcal{J}} \bar{\lambda}_\nu^n \langle \psi, \bar{\psi}_\nu^n \rangle_X^2. \end{aligned}$$



From  $\bar{\lambda}_1^n \geq \bar{\lambda}_\nu^n$  for all  $\nu \in \mathcal{J}$  and (2.6) we infer that

$$\begin{aligned} \sum_{\nu \in \mathcal{J}} \bar{\lambda}_\nu^n \langle \psi, \bar{\psi}_\nu^n \rangle_X^2 &\leq \bar{\lambda}_1^n \sum_{\nu \in \mathcal{J}} \langle \psi, \bar{\psi}_\nu^n \rangle_X^2 = \bar{\lambda}_1^n \|\psi\|_X^2 = \bar{\lambda}_1^n \\ &= \sum_{k=1}^{\wp} \sum_{j=1}^n \alpha_j^n \langle y_j^k, \bar{\psi}_1^n \rangle_X^2, \end{aligned}$$

i.e.,  $\bar{\psi}_1^n$  solves  $(\hat{\mathbf{P}}_n^\ell)$  for  $\ell = 1$  and (2.11) holds. This gives the base case. Notice that (2.9) and (2.11) imply (2.10).

2) The induction hypothesis: Now we suppose that

$$\left\{ \begin{array}{l} \text{for any } \ell \in \{1, \dots, d^n - 1\} \text{ the set } \{\bar{\psi}_i^n\}_{i=1}^\ell \subset X \text{ solve } (\hat{\mathbf{P}}_n^\ell) \\ \text{and } \sum_{k=1}^{\wp} \sum_{j=1}^n \alpha_j^n \sum_{i=1}^\ell \langle y_j^k, \bar{\psi}_i^n \rangle_X^2 = \sum_{i=1}^\ell \bar{\lambda}_i^n. \end{array} \right. \quad (2.13)$$

3) The induction step: We consider

$$\left\{ \begin{array}{l} \max \sum_{k=1}^{\wp} \sum_{j=1}^n \alpha_j^n \sum_{i=1}^{\ell+1} \langle y_j^k, \psi_i \rangle_X^2 \\ \text{s.t. } \{\psi_i\}_{i=1}^{\ell+1} \subset X \text{ and } \langle \psi_i, \psi_j \rangle_X = \delta_{ij}, 1 \leq i, j \leq \ell + 1. \end{array} \right. \quad (\hat{\mathbf{P}}_n^{\ell+1})$$

By (2.13) the elements  $\{\bar{\psi}_i^n\}_{i=1}^\ell$  maximize the term

$$\sum_{k=1}^{\wp} \sum_{j=1}^n \alpha_j^n \sum_{i=1}^\ell \langle y_j^k, \bar{\psi}_i^n \rangle_X^2.$$

Thus,  $(\hat{\mathbf{P}}_n^{\ell+1})$  is equivalent with

$$\left\{ \begin{array}{l} \max \sum_{k=1}^{\wp} \sum_{j=1}^n \alpha_j^n \langle y_j^k, \psi \rangle_X^2 \\ \text{s.t. } \psi \in X \text{ and } \|\psi\|_X = 1, \langle \psi, \bar{\psi}_i^n \rangle_X = 0, 1 \leq i \leq \ell. \end{array} \right. \quad (2.14)$$

Let  $\psi \in X$  be given satisfying  $\|\psi\|_X = 1$  and  $\langle \psi, \bar{\psi}_i^n \rangle_X = 0$  for  $i = 1 \dots, \ell$ . Then, using the representation (2.12) and  $\langle \psi, \bar{\psi}_i^n \rangle_X = 0$  for  $i = 1 \dots, \ell$ , we derive as above

$$\sum_{k=1}^{\wp} \sum_{j=1}^n \alpha_j^n \langle y_j^k, \psi \rangle_X^2 = \sum_{\nu \in \mathcal{J}} \bar{\lambda}_\nu^n \langle \psi, \bar{\psi}_\nu^n \rangle_X^2 = \sum_{\nu > \ell} \bar{\lambda}_\nu^n \langle \psi, \bar{\psi}_\nu^n \rangle_X^2.$$

From  $\bar{\lambda}_{\ell+1}^n \geq \bar{\lambda}_\nu^n$  for all  $\nu \geq \ell + 1$  and (2.6) we conclude that

$$\begin{aligned} \sum_{k=1}^{\wp} \sum_{j=1}^n \alpha_j^n \langle y_j^k, \psi \rangle_X^2 &\leq \bar{\lambda}_{\ell+1}^n \sum_{\nu > \ell} \langle \psi, \bar{\psi}_\nu^n \rangle_X^2 \leq \bar{\lambda}_{\ell+1}^n \sum_{\nu \in \mathcal{J}} \langle \psi, \bar{\psi}_\nu^n \rangle_X^2 \\ &= \bar{\lambda}_{\ell+1}^n \|\psi\|_X^2 = \bar{\lambda}_{\ell+1}^n = \sum_{k=1}^{\wp} \sum_{j=1}^n \alpha_j^n \langle y_j^k, \bar{\psi}_{\ell+1}^n \rangle_X^2. \end{aligned}$$

Thus,  $\bar{\psi}_{\ell+1}^n$  solves (2.14), which implies that  $\{\bar{\psi}_i^n\}_{i=1}^{\ell+1}$  is a solution to  $(\hat{\mathbf{P}}_n^{\ell+1})$  and

$$\sum_{k=1}^{\wp} \sum_{j=1}^n \alpha_j^n \sum_{i=1}^{\ell+1} \langle y_j^k, \bar{\psi}_i^n \rangle_X^2 = \sum_{i=1}^{\ell+1} \bar{\lambda}_i^n.$$

Again, (2.9) and (2.11) imply (2.10).  $\square$

**Remark 2.8.** Theorem 2.7 can also be proved by using the theory of nonlinear programming; see [31, 75], for instance. In this case  $(\hat{\mathbf{P}}_n^\ell)$  is considered as an equality constrained optimization problem. Applying a Lagrangian framework it turns out that (2.5) are first-order necessary optimality conditions for  $(\hat{\mathbf{P}}_n^\ell)$ .  $\diamond$

For the application of POD to concrete problems the choice of  $\ell$  is certainly of central importance for applying POD. It appears that no general a-priori rules are available. Rather the choice of  $\ell$  is based on heuristic considerations combined with observing the ratio of the modeled to the “total energy” contained in the snapshots  $y_1^k, \dots, y_n^k$ ,  $1 \leq k \leq \wp$ , which is expressed by

$$\mathcal{E}(\ell) = \frac{\sum_{i=1}^{\ell} \bar{\lambda}_i^n}{\sum_{i=1}^{d^n} \bar{\lambda}_i^n} \in [0, 1].$$

Utilizing (2.8) we have

$$\mathcal{E}(\ell) = \frac{\sum_{i=1}^{\ell} \bar{\lambda}_i^n}{\sum_{k=1}^{\wp} \sum_{j=1}^n \alpha_j^n \|y_j^k\|_X^2},$$

i.e., the computation of the eigenvalues  $\{\bar{\lambda}_i\}_{i=\ell+1}^{d_{\ell+1}}$  is not necessary. This is utilized in numerical implementations when iterative eigenvalue solver are applied like, e.g., the Lanczos method; see [2, Chapter 10], for instance.

In the following we will discuss three examples which illustrate that POD is strongly related to the singular value decomposition of matrices.

**Remark 2.9 (POD in Euclidean space  $\mathbb{R}^m$ ; see [39]).** Suppose that  $X = \mathbb{R}^m$  with  $m \in \mathbb{N}$  and  $\wp = 1$  hold. Then we have  $n$  snapshot vectors  $y_1, \dots, y_n$  and introduce the rectangular matrix  $Y = [y_1 | \dots | y_n] \in \mathbb{R}^{m \times n}$  with rank  $d^n \leq \min(m, n)$ . Choosing  $\alpha_j^n = 1$  for  $1 \leq j \leq n$  problem  $(\mathbf{P}_n^\ell)$  has the form

$$\begin{cases} \min \sum_{j=1}^n \left\| y_j - \sum_{i=1}^{\ell} (y_j^\top \psi_i) \psi_i \right\|_{\mathbb{R}^m}^2 \\ \text{s.t. } \{\psi_i\}_{i=1}^{\ell} \subset \mathbb{R}^m \text{ and } \psi_i^\top \psi_j = \delta_{ij}, \quad 1 \leq i, j \leq \ell, \end{cases} \quad (2.15)$$

where  $\|\cdot\|_{\mathbb{R}^m}$  stands for the Euclidean norm in  $\mathbb{R}^m$  and “ $\top$ ” denotes the transpose of a given vector (or matrix). From

$$(\mathcal{R}^n \psi)_i = \left( \sum_{j=1}^n (y_j^\top \psi) y_j \right)_i = \sum_{j=1}^n \sum_{l=1}^m Y_{lj} \psi_l Y_{ij} = (Y Y^\top \psi)_i, \quad \psi \in \mathbb{R}^m,$$

for each component  $1 \leq i \leq m$  we infer that (2.5) leads to the symmetric  $m \times m$  eigenvalue problem

$$YY^\top \bar{\psi}_i^n = \bar{\lambda}_i^n \bar{\psi}_i^n, \quad \bar{\lambda}_1^n \geq \dots \geq \bar{\lambda}_{d^n}^n > \bar{\lambda}_{d^n+1}^n = \dots = \bar{\lambda}_m^n = 0. \quad (2.16)$$

Recall that (2.16) can be solved by utilizing the singular value decomposition (SVD) [53]: There exist real numbers  $\bar{\sigma}_1^n \geq \bar{\sigma}_2^n \geq \dots \geq \bar{\sigma}_{d^n}^n > 0$  and orthogonal matrices  $\Psi \in \mathbb{R}^{m \times m}$  with column vectors  $\{\bar{\psi}_i^n\}_{i=1}^m$  and  $\Phi \in \mathbb{R}^{n \times n}$  with column vectors  $\{\bar{\phi}_i^n\}_{i=1}^n$  such that

$$\Psi^\top Y \Phi = \begin{pmatrix} D & 0 \\ 0 & 0 \end{pmatrix} =: \Sigma \in \mathbb{R}^{m \times n}, \quad (2.17)$$

where  $D = \text{diag}(\bar{\sigma}_1^n, \dots, \bar{\sigma}_{d^n}^n) \in \mathbb{R}^{d \times d}$  and the zeros in (2.17) denote matrices of appropriate dimensions. Moreover the vectors  $\{\bar{\psi}_i^n\}_{i=1}^d$  and  $\{\bar{\phi}_i^n\}_{i=1}^d$  satisfy

$$Y \bar{\phi}_i^n = \bar{\sigma}_i^n \bar{\psi}_i^n \quad \text{and} \quad Y^\top \bar{\psi}_i^n = \bar{\sigma}_i^n \bar{\phi}_i^n \quad \text{for } i = 1, \dots, d^n. \quad (2.18)$$

They are eigenvectors of  $YY^\top$  and  $Y^\top Y$ , respectively, with eigenvalues  $\bar{\lambda}_i^n = (\bar{\sigma}_i^n)^2 > 0$ ,  $i = 1, \dots, d^n$ . The vectors  $\{\bar{\psi}_i^n\}_{i=d^n+1}^m$  and  $\{\bar{\phi}_i^n\}_{i=d^n+1}^n$  (if  $d^n < m$  respectively  $d^n < n$ ) are eigenvectors of  $YY^\top$  and  $Y^\top Y$  with eigenvalue 0. Consequently, in the case  $n < m$  one can determine the POD basis of rank  $\ell$  as follows: Compute the eigenvectors  $\bar{\phi}_1^n, \dots, \bar{\phi}_\ell^n \in \mathbb{R}^n$  by solving the symmetric  $n \times n$  eigenvalue problem

$$Y^\top Y \bar{\phi}_i^n = \bar{\lambda}_i^n \bar{\phi}_i^n \quad \text{for } i = 1, \dots, \ell$$

and set, by (2.18),

$$\bar{\psi}_i^n = \frac{1}{(\bar{\lambda}_i^n)^{1/2}} Y \bar{\phi}_i^n \quad \text{for } i = 1, \dots, \ell.$$

For historical reasons this method of determining the POD-basis is sometimes called the *method of snapshots*; see [69]. On the other hand, if  $m < n$  holds, we can obtain the POD basis by solving the  $m \times m$  eigenvalue problem (2.16). If the matrix  $Y$  is badly scaled, we should avoid to build the matrix product  $YY^\top$  (or  $Y^\top Y$ ). In this case the SVD turns out to be more stable for the numerical computation of the POD basis of rank  $\ell$ .  $\diamond$

**Remark 2.10 (POD in  $\mathbb{R}^m$  with weighted inner product).** As in Remark 2.9 we choose  $X = \mathbb{R}^m$  with  $m \in \mathbb{R}^m$  and  $\wp = 1$ . Let  $W \in \mathbb{R}^{m \times m}$  be a given symmetric, positive definite matrix. We supply  $\mathbb{R}^m$  with the weighted inner product

$$\langle \psi, \tilde{\psi} \rangle_W = \psi^\top W \tilde{\psi} = \langle \psi, W \tilde{\psi} \rangle_{\mathbb{R}^m} = \langle W \psi, \tilde{\psi} \rangle_{\mathbb{R}^m} \quad \text{for } \psi, \tilde{\psi} \in \mathbb{R}^m.$$

Then, problem  $(\mathbf{P}_n^\ell)$  has the form

$$\begin{cases} \min \sum_{j=1}^n \alpha_j^n \left\| y_j - \sum_{i=1}^{\ell} \langle y_j, \psi_i \rangle_W \psi_i \right\|_W^2 \\ \text{s.t. } \{\psi_i\}_{i=1}^{\ell} \subset \mathbb{R}^m \text{ and } \langle \psi_i, \psi_j \rangle_W = \delta_{ij}, \quad 1 \leq i, j \leq \ell. \end{cases}$$

As in Remark 2.9 we introduce the matrix  $Y = [y_1 | \dots | y_n] \in \mathbb{R}^{m \times n}$  with rank  $d^n \leq \min(m, n)$ . Moreover, we define the diagonal matrix  $D = \text{diag}(\alpha_1^n, \dots, \alpha_n^n) \in \mathbb{R}^{n \times n}$ . We find that

$$\begin{aligned} (\mathcal{R}^n \psi)_i &= \left( \sum_{j=1}^n \alpha_j^n \langle y_j, \psi \rangle_W y_j \right)_i = \sum_{j=1}^n \sum_{l=1}^m \sum_{\nu=1}^m \alpha_j^n Y_{lj} W_{l\nu} \psi_\nu Y_{ij} \\ &= (YDY^\top W\psi)_i \quad \text{for } \psi \in \mathbb{R}^m, \end{aligned}$$

for each component  $1 \leq i \leq m$ . Consequently, (2.5) leads to the eigenvalue problem

$$YDY^\top W\bar{\psi}_i^n = \bar{\lambda}_i^n \bar{\psi}_i^n, \quad \bar{\lambda}_1^n \geq \dots \geq \bar{\lambda}_{d^n}^n > \bar{\lambda}_{d^n+1}^n = \dots = \bar{\lambda}_m^n = 0. \quad (2.19)$$

Since  $W$  is symmetric and positive definite,  $W$  possesses an eigenvalue decomposition of the form  $W = QBQ^\top$ , where  $B = \text{diag}(\beta_1, \dots, \beta_m)$  contains the eigenvalues  $\beta_1 \geq \dots \geq \beta_m > 0$  of  $W$  and  $Q \in \mathbb{R}^{m \times m}$  is an orthogonal matrix. We define

$$W^r = Q \text{diag}(\beta_1^r, \dots, \beta_m^r) Q^\top \quad \text{for } r \in \mathbb{R}.$$

Note that  $(W^r)^{-1} = W^{-r}$  and  $W^{r+s} = W^r W^s$  for  $r, s \in \mathbb{R}$ . Moreover, we have

$$\langle \psi, \tilde{\psi} \rangle_W = \langle W^{1/2} \psi, W^{1/2} \tilde{\psi} \rangle_{\mathbb{R}^m} \quad \text{for } \psi, \tilde{\psi} \in \mathbb{R}^m$$

and  $\|\psi\|_W = \|W^{1/2} \psi\|_{\mathbb{R}^m}$  for  $\psi \in \mathbb{R}^m$ . Analogously, the matrix  $D^{1/2}$  is defined. Inserting  $\psi_i^n = W^{1/2} \bar{\psi}_i^n$  in (2.19), multiplying (2.19) by  $W^{1/2}$  from the left and setting  $\hat{Y} = W^{1/2} Y D^{1/2}$  yield the symmetric  $m \times m$  eigenvalue problem

$$\hat{Y} \hat{Y}^\top \psi_i^n = \bar{\lambda}_i^n \psi_i^n, \quad 1 \leq i \leq \ell.$$

Note that

$$\hat{Y}^\top \hat{Y} = D^{1/2} Y^\top W Y D^{1/2} \in \mathbb{R}^{n \times n}. \quad (2.20)$$

Thus, the POD basis  $\{\bar{\psi}_i^n\}_{i=1}^\ell$  of rank  $\ell$  can also be computed by the methods of snapshots as follows: First solve the symmetric  $n \times n$  eigenvalue problem

$$\hat{Y}^\top \hat{Y} \phi_i^n = \bar{\lambda}_i^n \phi_i^n, \quad 1 \leq i \leq \ell \quad \text{and} \quad \langle \phi_i^n, \phi_j^n \rangle_{\mathbb{R}^n} = \delta_{ij}, \quad 1 \leq i, j \leq \ell.$$

Then we set (by using the SVD of  $\hat{Y}$ )

$$\bar{\psi}_i^n = W^{-1/2} \psi_i^n = \frac{1}{\bar{\sigma}_i^n} W^{-1/2} \hat{Y} \phi_i^n = \frac{1}{\bar{\sigma}_i^n} Y D^{1/2} \phi_i^n, \quad 1 \leq i \leq \ell. \quad (2.21)$$

Note that

$$\langle \bar{\psi}_i^n, \bar{\psi}_j^n \rangle_W = (\bar{\psi}_i^n)^\top W \bar{\psi}_j^n = \frac{1}{\bar{\sigma}_i^n \bar{\sigma}_j^n} (\phi_i^n)^\top \underbrace{D^{1/2} Y^\top W Y D^{1/2}}_{=\hat{Y}^\top \hat{Y}} \phi_j^n = \delta_{ij}$$

for  $1 \leq i, j \leq \ell$ . Thus, the POD basis  $\{\bar{\psi}_i^n\}_{i=1}^\ell$  of rank  $\ell$  is orthonormal in  $\mathbb{R}^m$  with respect to the inner product  $\langle \cdot, \cdot \rangle_W$ . We observe from (2.20) and (2.21) that the computation of  $W^{1/2}$  and  $W^{-1/2}$  is not required. For applications, where  $W$  is not just a diagonal matrix, the method of snapshots turns out to be more attractive with respect to the computational costs even if  $m > n$  holds.  $\diamond$

**Remark 2.11 (POD in  $\mathbb{R}^m$  with multiple snapshots).** Let us discuss the more general case  $\varphi = 2$  in the setting of Remark 2.10. The extension for  $\varphi > 2$  is straightforward. We introduce the matrix  $Y = [y_1^1 | \dots | y_n^1 | y_1^2 | \dots | y_n^2] \in \mathbb{R}^{m \times (n\varphi)}$  with rank  $d^n \leq \min(m, n\varphi)$ . Then we find

$$\begin{aligned} \mathcal{R}^n \psi &= \sum_{j=1}^n \left( \alpha_j^n \langle y_j^1, \psi \rangle_W y_j^1 + \alpha_j^n \langle y_j^2, \psi \rangle_W y_j^2 \right) \\ &= Y \underbrace{\begin{pmatrix} D & 0 \\ 0 & D \end{pmatrix}}_{=: \tilde{D} \in \mathbb{R}^{(n\varphi) \times (n\varphi)}} Y^\top W \psi = Y \tilde{D} Y^\top W \psi \quad \text{for } \psi \in \mathbb{R}^m. \end{aligned}$$

Hence, (2.5) corresponds to the eigenvalue problem

$$Y \tilde{D} Y^\top W \bar{\psi}_i^n = \bar{\lambda}_i^n \bar{\psi}_i^n, \quad \bar{\lambda}_1^n \geq \dots \geq \bar{\lambda}_{d^n}^n > \bar{\lambda}_{d^n+1}^n = \dots = \bar{\lambda}_m^n = 0. \quad (2.22)$$

Setting  $\psi_i^n = W^{1/2} \bar{\psi}_i^n$  in (2.22) and multiplying by  $W^{1/2}$  from the left yield

$$W^{1/2} Y \tilde{D} Y^\top W^{1/2} \psi_i^n = \bar{\lambda}_i^n \psi_i^n. \quad (2.23)$$

Let  $\hat{Y} = W^{1/2} Y \tilde{D}^{1/2} \in \mathbb{R}^{m \times (n\varphi)}$ . Using  $W^\top = W$  as well as  $\tilde{D}^\top = \tilde{D}$  we infer from (2.23) that the POD basis  $\{\bar{\psi}_i^n\}_{i=1}^\ell$  of rank  $\ell$  is given by the symmetric  $m \times m$  eigenvalue problem

$$\hat{Y} \hat{Y}^\top \psi_i^n = \bar{\lambda}_i^n \psi_i^n, \quad 1 \leq i \leq \ell, \quad \text{and} \quad \langle \psi_i^n, \psi_j^n \rangle_{\mathbb{R}^m} = \delta_{ij}, \quad 1 \leq i, j \leq \ell$$

and  $\bar{\psi}_i^n = W^{-1/2} \psi_i^n$ . Note that

$$\hat{Y}^\top \hat{Y} = \tilde{D}^{1/2} Y^\top W Y \tilde{D}^{1/2} \in \mathbb{R}^{(n\varphi) \times (n\varphi)}.$$

Thus, the POD basis of rank  $\ell$  can also be computed by the methods of snapshots as follows: First solve the symmetric  $(n\varphi) \times (n\varphi)$  eigenvalue problem

$$\hat{Y}^\top \hat{Y} \phi_i^n = \bar{\lambda}_i^n \phi_i^n, \quad 1 \leq i \leq \ell \quad \text{and} \quad \langle \phi_i^n, \phi_j^n \rangle_{\mathbb{R}^{n\varphi}} = \delta_{ij}, \quad 1 \leq i, j \leq \ell.$$

Then we set (by SVD)

$$\bar{\psi}_i^n = W^{-1/2} \psi_i^n = \frac{1}{\bar{\sigma}_i^n} W^{-1/2} \hat{Y} \phi_i^n = \frac{1}{\bar{\sigma}_i^n} Y \tilde{D}^{1/2} \phi_i^n$$

for  $1 \leq i \leq \ell$ . ◇

## 2.2. The continuous variant of the POD method

Let  $0 \leq t_1 < t_2 < \dots < t_n \leq T$  be a given time grid in the interval  $[0, T]$ . To simplify of the presentation, the time grid is assumed to be equidistant with step-size  $\Delta t = T/(n-1)$ , i.e.,  $t_j = (j-1)\Delta t$ . For nonequidistant grids we refer the reader to [41, 42, ]. Suppose that we have trajectories  $y^k \in C([0, T]; X)$ ,  $1 \leq k \leq \varphi$ . Here, the Banach space  $C([0, T]; X)$  contains all functions  $\varphi : [0, T] \rightarrow X$ , which are continuous on  $[0, T]$ ; see, e.g., [72, p. 142]. Let the snapshots be given as  $y_j^k = y^k(t_j) \in X$  or  $y_j^k \approx y^k(t_j) \in X$ . Then, the snapshot subspace  $\mathcal{V}^n$  introduced in (2.1) depends on the chosen time instances  $\{t_j\}_{j=1}^n$ . Consequently, the POD basis  $\{\bar{\psi}_i^n\}_{i=1}^\ell$  of rank  $\ell$  as well as the corresponding eigenvalues  $\{\bar{\lambda}_i^n\}_{i=1}^\ell$  depend also on the time instances (which has already been indicated by the superindex  $n$ ). Moreover, we have

not discussed so far what is the motivation to introduce the positive weights  $\{\alpha_j^n\}_{j=1}^n$  in  $(\mathbf{P}_n^\ell)$ . For this reason we proceed by investigating the following two questions:

- How to choose good time instances for the snapshots?
- What are appropriate positive weights  $\{\alpha_j^n\}_{j=1}^n$ ?

To address these two questions we will introduce a *continuous version* of POD. In Section 2.1 we have introduced the operator  $\mathcal{R}^n$  in (2.4). By  $\{\bar{\psi}_i^n\}_{i \in \mathcal{J}}$  and  $\{\bar{\lambda}_i^n\}_{i \in \mathcal{J}}$  we have denoted the eigenfunctions and eigenvalues for  $\mathcal{R}^n$  satisfying (2.5). Moreover, we have set  $d^n = \dim \mathcal{V}^n$  for the dimension of the snapshot set. Let us now introduce the snapshot set by

$$\mathcal{V} = \text{span} \left\{ y^k(t) \mid t \in [0, T] \text{ and } 1 \leq k \leq \wp \right\} \subset X$$

with dimension  $d \leq \infty$ . For any  $\ell \leq d$  we are interested in determining a POD basis of rank  $\ell$  which minimizes the mean square error between the trajectories  $y^k$  and the corresponding  $\ell$ -th partial Fourier sums on average in the time interval  $[0, T]$ :

$$\left\{ \begin{array}{l} \min \sum_{k=1}^{\wp} \int_0^T \left\| y^k(t) - \sum_{i=1}^{\ell} \langle y^k(t), \psi_i \rangle_X \psi_i \right\|_X^2 dt \\ \text{s.t. } \{\psi_i\}_{i=1}^{\ell} \subset X \text{ and } \langle \psi_i, \psi_j \rangle_X = \delta_{ij}, \quad 1 \leq i, j \leq \ell. \end{array} \right. \quad (\mathbf{P}^\ell)$$

An optimal solution  $\{\bar{\psi}_i\}_{i=1}^{\ell}$  to  $(\mathbf{P}^\ell)$  is called a *POD basis of rank  $\ell$* . Analogous to  $(\hat{\mathbf{P}}_n^\ell)$  we can – instead of  $(\mathbf{P}^\ell)$  – consider the problem

$$\left\{ \begin{array}{l} \max \sum_{k=1}^{\wp} \int_0^T \sum_{i=1}^{\ell} \langle y^k(t), \psi_i \rangle_X^2 dt \\ \text{s.t. } \{\psi_i\}_{i=1}^{\ell} \subset X \text{ and } \langle \psi_i, \psi_j \rangle_X = \delta_{ij}, \quad 1 \leq i, j \leq \ell. \end{array} \right. \quad (\hat{\mathbf{P}}^\ell)$$

A solution to  $(\mathbf{P}^\ell)$  and to  $(\hat{\mathbf{P}}^\ell)$  can be characterized by an eigenvalue problem for the linear integral operator  $\mathcal{R} : X \rightarrow X$  given as

$$\mathcal{R}\psi = \sum_{k=1}^{\wp} \int_0^T \langle y^k(t), \psi \rangle_X y^k(t) dt \quad \text{for } \psi \in X. \quad (2.24)$$

For the given real Hilbert space  $X$  we denote by  $L^2(0, T; X)$  the Hilbert space of square integrable functions  $t \mapsto \varphi(t) \in X$  so that [72, p. 143]

- the mapping  $t \mapsto \varphi(t)$  is measurable for  $t \in [0, T]$  and
- $\|\varphi\|_{L^2(0, T; X)} = \left( \int_0^T \|\varphi(t)\|_X^2 dt \right)^{1/2} < \infty$ .

Recall that  $\varphi : [0, T] \rightarrow X$  is called *measurable* if there exists a sequence  $\{\varphi_n\}_{n \in \mathbb{N}}$  of step functions  $\varphi_n : [0, T] \rightarrow X$  satisfying  $\varphi(t) = \lim_{n \rightarrow \infty} \varphi_n(t)$  for almost all  $t \in [0, T]$ . The standard inner product on  $L^2(0, T; X)$  is given by

$$\langle \varphi, \psi \rangle_{L^2(0, T; X)} = \int_0^T \langle \varphi(t), \psi(t) \rangle_X dt \quad \text{for } \varphi, \psi \in L^2(0, T; X).$$

**Lemma 2.12.** *Let  $X$  be a (separable) real Hilbert space and  $y^k \in L^2(0, T; X)$ ,  $1 \leq k \leq \wp$ , be given snapshot trajectories. Then, the operator  $\mathcal{R}$  introduced in (2.24) is compact, nonnegative and selfadjoint.*

*Proof.* First we write  $\mathcal{R}$  as a product of an operator and its Hilbert space adjoint. For that purpose let us define the linear operator  $\mathcal{Y} : L^2(0, T; \mathbb{R}^\wp) \rightarrow X$  by

$$\mathcal{Y}\phi = \sum_{k=1}^{\wp} \int_0^T \phi^k(t) y^k(t) dt \quad \text{for } \phi = (\phi^1, \dots, \phi^\wp) \in L^2(0, T; \mathbb{R}^\wp). \quad (2.25)$$

Utilizing the Cauchy-Schwarz inequality [62, p. 38] and  $y^k \in L^2(0, T; X)$  for  $1 \leq k \leq \wp$  we infer that

$$\begin{aligned} \|\mathcal{Y}\phi\|_X &\leq \sum_{k=1}^{\wp} \int_0^T |\phi^k(t)| \|y^k(t)\|_X dt \leq \sum_{k=1}^{\wp} \|\phi^k\|_{L^2(0, T)} \|y^k\|_{L^2(0, T; X)} \\ &\leq \left( \sum_{k=1}^{\wp} \|\phi^k\|_{L^2(0, T)}^2 \right)^{1/2} \left( \sum_{k=1}^{\wp} \|y^k(t)\|_X^2 \right)^{1/2} \\ &= C_{\mathcal{Y}} \|\phi\|_{L^2(0, T; \mathbb{R}^\wp)} \quad \text{for any } \phi \in L^2(0, T; \mathbb{R}^\wp), \end{aligned}$$

where we set  $C_{\mathcal{Y}} = (\sum_{k=1}^{\wp} \|y^k(t)\|_X^2)^{1/2} < \infty$ . Hence, the operator  $\mathcal{Y}$  is bounded. Its Hilbert space adjoint  $\mathcal{Y}^* : X \rightarrow L^2(0, T; \mathbb{R}^\wp)$  satisfies

$$\langle \mathcal{Y}^* \psi, \phi \rangle_{L^2(0, T; \mathbb{R}^\wp)} = \langle \psi, \mathcal{Y}\phi \rangle_X \quad \text{for } \psi \in X \text{ and } \phi \in L^2(0, T; \mathbb{R}^\wp).$$

Since we derive

$$\begin{aligned} \langle \mathcal{Y}^* \psi, \phi \rangle_{L^2(0, T; \mathbb{R}^\wp)} &= \langle \psi, \mathcal{Y}\phi \rangle_X = \left\langle \psi, \sum_{k=1}^{\wp} \int_0^T \phi^k(t) y^k(t) dt \right\rangle_X \\ &= \sum_{k=1}^{\wp} \int_0^T \langle \psi, y^k(t) \rangle_X \phi^k(t) dt = \left\langle \left( \langle \psi, y^k(\cdot) \rangle_X \right)_{1 \leq k \leq \wp}, \phi \right\rangle_{L^2(0, T; \mathbb{R}^\wp)} \end{aligned}$$

for  $\psi \in X$  and  $\phi \in L^2(0, T; \mathbb{R}^\wp)$ , the adjoint operator is given by

$$(\mathcal{Y}^* \psi)(t) = \begin{pmatrix} \langle \psi, y^1(t) \rangle_X \\ \vdots \\ \langle \psi, y^\wp(t) \rangle_X \end{pmatrix} \quad \text{for } \psi \in X \text{ and } t \in [0, T] \text{ a.e.,}$$

where ‘a.e.’ stands for ‘almost everywhere’. From (2.4) and

$$(\mathcal{Y}\mathcal{Y}^*)\psi = \mathcal{Y} \begin{pmatrix} \langle \psi, y^1(\cdot) \rangle_X \\ \vdots \\ \langle \psi, y^\wp(\cdot) \rangle_X \end{pmatrix} = \sum_{k=1}^{\wp} \int_0^T \langle \psi, y^k(t) \rangle_X y^k(t) dt \quad \text{for } \psi \in X$$

we infer that  $\mathcal{R} = \mathcal{Y}\mathcal{Y}^*$  holds. Moreover, let  $\mathcal{K} = \mathcal{Y}^*\mathcal{Y} : L^2(0, T; \mathbb{R}^\varphi) \rightarrow L^2(0, T; \mathbb{R}^\varphi)$ . We find that

$$(\mathcal{K}\phi)(t) = \begin{pmatrix} \sum_{k=1}^{\varphi} \int_0^T \langle y^k(s), y^1(t) \rangle_X \phi^k(s) ds \\ \vdots \\ \sum_{k=1}^{\varphi} \int_0^T \langle y^k(s), y^\varphi(t) \rangle_X \phi^k(s) ds \end{pmatrix}, \quad \phi \in L^2(0, T; \mathbb{R}^\varphi).$$

Since the operator  $\mathcal{Y}$  is bounded, its adjoint and therefore  $\mathcal{R} = \mathcal{Y}\mathcal{Y}^*$  are bounded operators. Notice that the kernel function

$$r_{ik}(s, t) = \langle y^k(s), y^i(t) \rangle_X, \quad (s, t) \in [0, T] \times [0, T] \text{ and } 1 \leq i, k \leq \varphi,$$

belongs to  $L^2(0, T) \times L^2(0, T)$ . Here, we shortly write  $L^2(0, T)$  for  $L^2(0, T; \mathbb{R})$ . Then, it follows from [80, pp. 197 and 277] that the linear integral operator  $\mathcal{K}_{ik} : L^2(0, T) \rightarrow L^2(0, T)$  defined by

$$\mathcal{K}_{ik}(t) = \int_0^T r_{ik}(s, t) \phi(s) ds, \quad \phi \in L^2(0, T),$$

is a compact operator. This implies, that the operator  $\sum_{k=1}^{\varphi} \mathcal{K}_{ik}$  is compact for  $1 \leq i \leq \varphi$  as well. Consequently,  $\mathcal{K}$  and therefore  $\mathcal{R} = \mathcal{K}^*$  are compact operators. From

$$\begin{aligned} \langle \mathcal{R}\psi, \psi \rangle_X &= \left\langle \sum_{k=1}^{\varphi} \int_0^T \langle \psi, y^k(t) \rangle_X y^k(t) dt, \psi \right\rangle_X \\ &= \sum_{k=1}^{\varphi} \int_0^T |\langle \psi, y^k(t) \rangle_X|^2 dt \geq 0 \quad \text{for all } \psi \in X \end{aligned}$$

we infer that  $\mathcal{R}$  is nonnegative. Finally, we have  $\mathcal{R}^* = (\mathcal{Y}\mathcal{Y}^*)^* = \mathcal{R}$ , i.e. the operator  $\mathcal{R}$  is selfadjoint.  $\square$

In the next theorem we formulate how the solution to  $(\mathbf{P}^\ell)$  and  $(\hat{\mathbf{P}}^\ell)$  can be found.

**Theorem 2.13.** *Let  $X$  be a separable real Hilbert space and  $y^k \in L^2(0, T; X)$  are given trajectories for  $1 \leq k \leq \varphi$ . Suppose that the linear operator  $\mathcal{R}$  is defined by (2.24). Then, there exist nonnegative eigenvalues  $\{\bar{\lambda}_i\}_{i \in \mathcal{J}}$  and associated orthonormal eigenfunctions  $\{\bar{\psi}_i\}_{i \in \mathcal{J}}$  satisfying*

$$\mathcal{R}\bar{\psi}_i = \bar{\lambda}_i \bar{\psi}_i, \quad \bar{\lambda}_1 \geq \dots \geq \bar{\lambda}_d > \bar{\lambda}_{d+1} = \dots = 0. \quad (2.26)$$

For every  $\ell \in \{1, \dots, d\}$  the first  $\ell$  eigenfunctions  $\{\bar{\psi}_i\}_{i=1}^\ell$  solve  $(\mathbf{P}^\ell)$  and  $(\hat{\mathbf{P}}^\ell)$ . Moreover, the value of the objectives evaluated at the optimal solution  $\{\bar{\psi}_i\}_{i=1}^\ell$  satisfies

$$\sum_{k=1}^{\varphi} \int_0^T \left\| y^k(t) - \sum_{i=1}^{\ell} \langle y^k(t), \bar{\psi}_i \rangle_X \bar{\psi}_i \right\|_X^2 dt = \sum_{i=\ell+1}^d \bar{\lambda}_i \quad (2.27)$$



and

$$\sum_{k=1}^{\wp} \int_0^T \sum_{i=1}^{\ell} \langle y^k(t), \bar{\psi}_i \rangle_X^2 dt = \sum_{i=1}^{\ell} \bar{\lambda}_i, \quad (2.28)$$

respectively.

*Proof.* The existence of sequences  $\{\bar{\lambda}_i\}_{i \in \mathcal{J}}$  of eigenvalues and  $\{\bar{\psi}_i\}_{i \in \mathcal{J}}$  of associated eigenfunctions satisfying (2.26) follows from Lemma 2.12, Theorem 2.4 and Theorem 2.5. Analogous to the proof of Theorem 2.7 in Section 2.1 one can show that  $\{\bar{\psi}_i\}_{i=1}^{\ell}$  solves  $(\mathbf{P}^{\ell})$  as well as  $(\hat{\mathbf{P}}^{\ell})$  and that (2.27) respectively (2.28) are valid.  $\square$

**Remark 2.14.** Similar to (2.6) we have

$$\sum_{k=1}^{\wp} \int_0^T \|y^k(t)\|_X^2 dt = \sum_{i=1}^d \bar{\lambda}_i. \quad (2.29)$$

In fact,

$$\mathcal{R}\bar{\psi}_i = \sum_{k=1}^{\wp} \int_0^T \langle y^k(t), \bar{\psi}_i \rangle_X y^k(t) dt \quad \text{for every } i \in \mathcal{J}.$$

Taking the inner product with  $\bar{\psi}_i$ , using (2.26) and summing over  $i$  we get

$$\sum_{i=1}^d \sum_{k=1}^{\wp} \int_0^T \langle y^k(t), \bar{\psi}_i \rangle_X^2 dt = \sum_{i=1}^d \langle \mathcal{R}\bar{\psi}_i, \bar{\psi}_i \rangle_X = \sum_{i=1}^d \bar{\lambda}_i.$$

Expanding each  $y^k(t) \in X$  in terms of  $\{\bar{\psi}_i\}_{i \in \mathcal{J}}$  for each  $1 \leq k \leq \wp$  we have

$$y^k(t) = \sum_{i=1}^d \langle y^k(t), \bar{\psi}_i \rangle_X \bar{\psi}_i$$

and hence

$$\sum_{k=1}^{\wp} \int_0^T \|y^k(t)\|_X^2 dt = \sum_{k=1}^{\wp} \sum_{i=1}^d \int_0^T \langle y^k(t), \bar{\psi}_i \rangle_X^2 dt = \sum_{i=1}^d \bar{\lambda}_i,$$

which is (2.29).  $\diamond$

**Remark 2.15 (Singular value decomposition).** Suppose that  $y^k \in L^2(0, T; X)$  holds. By Theorem 2.13 there exist nonnegative eigenvalues  $\{\bar{\lambda}_i\}_{i \in \mathcal{J}}$  and associated orthonormal eigenfunctions  $\{\bar{\psi}_i\}_{i \in \mathcal{J}}$  satisfying (2.26). From  $\mathcal{K} = \mathcal{R}^*$  it follows that there is a sequence  $\{\bar{\phi}_i\}_{i \in \mathcal{J}}$  such that

$$\mathcal{K}\bar{\phi}_i = \bar{\lambda}_i \bar{\phi}_i, \quad 1 \dots, \ell.$$

We set  $\mathbb{R}_0^+ = \{s \in \mathbb{R} \mid s \geq 0\}$  and  $\bar{\sigma}_i = \bar{\lambda}_i^{1/2}$ . The sequence  $\{\bar{\sigma}_i, \bar{\phi}_i, \bar{\psi}_i\}_{i \in \mathcal{J}}$  in  $\mathbb{R}_0^+ \times L^2(0, T; \mathbb{R}^{\wp}) \times X$  can be interpreted as a singular value decomposition of the mapping  $\mathcal{Y} : L^2(0, T; \mathbb{R}^{\wp}) \rightarrow X$  introduced in (2.25). In fact, we have

$$\mathcal{Y}\bar{\phi}_i = \bar{\sigma}_i \bar{\psi}_i, \quad \mathcal{Y}^* \bar{\psi}_i = \bar{\sigma}_i \bar{\phi}_i, \quad i \in \mathcal{J}.$$

Since  $\bar{\sigma}_i > 0$  holds for  $1 = 1 \dots, d$ , we have  $\bar{\phi}_i = \bar{\lambda}_i / \bar{\sigma}_i$  for  $i = 1, \dots, d$ .  $\diamond$

### 2.3. Perturbation analysis for the POD basis

The eigenvalues  $\{\bar{\lambda}_i^n\}_{i \in \mathcal{J}}$  satisfying (2.5) depend on the time grid  $\{t_j\}_{j=1}^n$ . In this section we investigate the sum  $\sum_{i=\ell+1}^{d^n} \bar{\lambda}_i^n$ , the value of the cost in  $(\mathbf{P}_n^\ell)$  evaluated at the solution  $\{\bar{\psi}_i^n\}_{i=1}^\ell$  for  $n \rightarrow \infty$ . Clearly,  $n \rightarrow \infty$  is equivalent with  $\Delta t = T/(n-1) \rightarrow 0$ .

In general the spectrum  $\sigma(\mathcal{T})$  of an operator  $\mathcal{T} \in \mathcal{L}(X)$  does not depend continuously on  $\mathcal{T}$ . This is an essential difference to the finite dimensional case. For the compact and selfadjoint operator  $\mathcal{R}$  we have  $\sigma(\mathcal{R}) = \{\bar{\lambda}_i\}_{i \in \mathcal{J}}$ . Suppose that for  $\ell \in \mathbb{N}$  we have  $\bar{\lambda}_\ell > \bar{\lambda}_{\ell+1}$  so that we can separate the spectrum as follows:  $\sigma(\mathcal{R}) = \mathcal{S}_\ell \cup \mathcal{S}'_\ell$  with  $\mathcal{S}_\ell = \{\bar{\lambda}_1, \dots, \bar{\lambda}_\ell\}$  and  $\mathcal{S}'_\ell = \sigma(\mathcal{R}) \setminus \mathcal{S}_\ell$ . Then,  $\mathcal{S}_\ell \cap \mathcal{S}'_\ell = \emptyset$ . Moreover, setting  $V^\ell = \text{span}\{\bar{\psi}_1, \dots, \bar{\psi}_\ell\}$  we have  $X = V^\ell \oplus (V^\ell)^\perp$ , where the linear space  $(V^\ell)^\perp$  stands for the  $X$ -orthogonal complement of  $V^\ell$ . Let us assume that

$$\lim_{n \rightarrow \infty} \|\mathcal{R}^n - \mathcal{R}\|_{\mathcal{L}(X)} = 0 \quad (2.30)$$

holds. Then it follows from the perturbation theory of the spectrum of linear operators [37, pp. 212-214] that the space  $V_n^\ell = \text{span}\{\bar{\psi}_1^n, \dots, \bar{\psi}_\ell^n\}$  is isomorphic to  $V^\ell$  if  $n$  is sufficiently large. Furthermore, the change of a finite set of eigenvalues of  $\mathcal{R}$  is small provided  $\|\mathcal{R}^n - \mathcal{R}\|_{\mathcal{L}(X)}$  is sufficiently small. Summarizing, the behavior of the spectrum is much the same as in the finite dimensional case if we can ensure (2.30). Therefore, we start this section by investigating the convergence of  $\mathcal{R}^n - \mathcal{R}$  in the operator norm.

Recall that the Sobolev space  $H^1(0, T; X)$  is given by

$$H^1(0, T; X) = \{\varphi \in L^2(0, T; X) \mid \varphi_t \in L^2(0, T; X)\},$$

where  $\varphi_t$  denotes the weak derivative of  $\varphi$ . The space  $H^1(0, T; X)$  is a Hilbert space with the inner product

$$\langle \varphi, \phi \rangle_{H^1(0, T; X)} = \int_0^T \langle \varphi(t), \phi(t) \rangle_X + \langle \varphi_t(t), \phi_t(t) \rangle_X dt \text{ for } \varphi, \phi \in H^1(0, T; X)$$

and the induced norm  $\|\varphi\|_{H^1(0, T; X)} = \langle \varphi, \varphi \rangle_{H^1(0, T; X)}^{1/2}$ .

Let us choose the trapezoidal weights

$$\alpha_1^n = \frac{T}{2(n-1)}, \quad \alpha_j^n = \frac{T}{n-1} \text{ for } 2 \leq j \leq n-1, \quad \alpha_n^n = \frac{T}{2(n-1)}. \quad (2.31)$$

For this choice we observe that for every  $\psi \in X$  the element  $\mathcal{R}^n \psi$  is a trapezoidal approximation for  $\mathcal{R} \psi$ . We will make use of the following lemma.

**Lemma 2.16.** *Suppose that  $X$  is a (separable) real Hilbert space and that the snapshot trajectories  $y^k$  belong to  $H^1(0, T; X)$  for  $1 \leq k \leq \wp$ . Then, (2.30) holds true.*

*Proof.* For an arbitrary  $\psi \in X$  with  $\|\psi\|_X = 1$  we define  $F : [0, T] \rightarrow X$  by

$$F(t) = \sum_{k=1}^{\wp} \langle y^k(t), \psi \rangle_X y^k(t) \text{ for } t \in [0, T].$$

It follows that

$$\begin{aligned}\mathcal{R}\psi &= \int_0^T F(t) dt = \sum_{j=1}^{n-1} \int_{t_j}^{t_{j+1}} F(t) dt, \\ \mathcal{R}^n \psi &= \sum_{j=1}^n \alpha_j F(t_j) = \frac{\Delta t}{2} \sum_{j=1}^{n-1} (F(t_j) + F(t_{j+1})).\end{aligned}\tag{2.32}$$

Then, we infer from  $\|\psi\|_X = 1$  that

$$\|F(t)\|_X^2 \leq \left( \sum_{k=1}^{\wp} \|y^k(t)\|_X^2 \right)^2.\tag{2.33}$$

Now we show that  $F$  belongs to  $H^1(0, T; X)$  and its norm is bounded independently of  $\psi$ . Recall that  $y^k \in H^1(0, T; X)$  imply that  $y^k \in C([0, T]; X)$  holds for  $1 \leq k \leq \wp$ . Using (2.33) we have

$$\|F\|_{L^2(0, T; X)}^2 \leq \int_0^T \left( \sum_{k=1}^{\wp} \|y^k\|_{C([0, T]; X)}^2 \right)^2 dt \leq C_1$$

with  $C_1 = T(\sum_{k=1}^{\wp} \|y^k\|_{C([0, T]; X)}^2)^2$ . Moreover,  $F \in H^1(0, T; X)$  with

$$F_t(t) = \sum_{k=1}^{\wp} \langle y_t^k(t), \psi \rangle_X y^k(t) + \langle y^k(t), \psi \rangle_X y_t^k(t) \quad \text{f.a.a. } t \in [0, T],$$

where ‘f.a.a.’ stands for ‘for almost all’. Thus, we derive

$$\|F_t\|_{L^2(0, T; X)}^2 \leq 4 \int_0^T \left( \sum_{k=1}^{\wp} \|y^k(t)\|_X \|y_t^k(t)\|_X \right)^2 dt \leq C_2$$

with  $C_2 = 4 \sum_{k=1}^{\wp} \|y^k\|_{C([0, T]; X)}^2 \sum_{l=1}^{\wp} \|y_l^l\|_{L^2(0, T; X)}^2 < \infty$ . Consequently,

$$\|F\|_{H^1(0, T; X)} = \left( \int_0^T \|F(t)\|_X^2 + \|F_t(t)\|_X^2 dt \right)^{1/2} \leq C_3\tag{2.34}$$

with the constant  $C_3 = (C_1 + C_2)^{1/2}$ , which is independent of  $\psi$ . To evaluate  $\mathcal{R}^n \psi$  we notice that

$$\begin{aligned}\int_{t_j}^{t_{j+1}} F(t) dt &= \frac{1}{2} \int_{t_j}^{t_{j+1}} \left( F(t_j) + \int_{t_j}^t F_t(s) ds \right) dt \\ &\quad + \frac{1}{2} \int_{t_j}^{t_{j+1}} \left( F(t_{j+1}) + \int_{t_{j+1}}^t F_t(s) ds \right) dt \\ &= \frac{\Delta t}{2} (F(t_j) + F(t_{j+1})) \\ &\quad + \frac{1}{2} \int_{t_j}^{t_{j+1}} \left( \int_{t_{j+1}}^t F_t(s) ds + \int_{t_j}^t F_t(s) ds \right) dt.\end{aligned}\tag{2.35}$$

Utilizing (2.32) and (2.35) we obtain

$$\begin{aligned} \|\mathcal{R}^n\psi - \mathcal{R}\psi\|_X &= \left\| \sum_{j=1}^{n-1} \left( \frac{\Delta t}{2} (F(t_j) + F(t_{j+1})) - \int_{t_j}^{t_{j+1}} F(t) dt \right) \right\|_X \\ &\leq \frac{1}{2} \sum_{j=1}^{n-1} \left\| \int_{t_j}^{t_{j+1}} \int_{t_j}^t F_t(s) ds dt \right\|_X + \frac{1}{2} \sum_{j=1}^{n-1} \left\| \int_{t_j}^{t_{j+1}} \int_{t_{j+1}}^t F_t(s) ds dt \right\|_X. \end{aligned}$$

From the Cauchy-Schwarz inequality [62, p. 38] we deduce that

$$\begin{aligned} \sum_{j=1}^{n-1} \left\| \int_{t_j}^{t_{j+1}} \int_{t_j}^t F_t(s) ds dt \right\|_X &\leq \sum_{j=1}^{n-1} \int_{t_j}^{t_{j+1}} \left\| \int_{t_j}^t F_t(s) ds \right\|_X dt \\ &\leq \sqrt{\Delta t} \sum_{j=1}^{n-1} \left( \int_{t_j}^{t_{j+1}} \left\| \int_{t_j}^t F_t(s) ds \right\|_X^2 dt \right)^{1/2} \\ &\leq \sqrt{\Delta t} \sum_{j=1}^{n-1} \left( \int_{t_j}^{t_{j+1}} \left( \int_{t_j}^t \|F_t(s)\|_X ds \right)^2 dt \right)^{1/2} \\ &\leq \Delta t \sum_{j=1}^{n-1} \left( \int_{t_j}^{t_{j+1}} \int_{t_j}^t \|F_t(s)\|_X^2 ds dt \right)^{1/2} \leq T\sqrt{\Delta t} \|F\|_{H^1(0,T;X)}. \end{aligned} \tag{2.36}$$

Analogously, we derive

$$\sum_{j=1}^{n-1} \left\| \int_{t_j}^{t_{j+1}} \int_{t_{j+1}}^t F_t(s) ds dt \right\|_X \leq T\sqrt{\Delta t} \|F\|_{H^1(0,T;X)}. \tag{2.37}$$

From (2.34), (2.36) and (2.37) it follows that

$$\|\mathcal{R}^n\psi - \mathcal{R}\psi\|_X \leq \frac{C_4}{\sqrt{n}},$$

where  $C_4 = C_3 T^{3/2}$  is independent of  $n$  and  $\psi$ . Consequently,

$$\|\mathcal{R}^n - \mathcal{R}\|_{\mathcal{L}(X)} = \sup_{\|\psi\|_X=1} \|\mathcal{R}^n\psi - \mathcal{R}\psi\|_X \xrightarrow{n \rightarrow \infty} 0$$

which gives the claim.  $\square$

Now we follow [41, Section 3.2]. We suppose that  $y^k \in H^1(0, T; X)$  for  $1 \leq k \leq \wp$ . Thus  $y^k \in C([0, T]; X)$  holds, which implies that

$$\sum_{k=1}^{\wp} \sum_{j=1}^n \alpha_j^n \|y^k(t_j)\|_X^2 \rightarrow \sum_{k=1}^{\wp} \int_0^T \|y^k(t)\|_X^2 dt \quad \text{as } n \rightarrow \infty. \tag{2.38}$$

Combining (2.38) with (2.8) and (2.29) we find

$$\sum_{i=1}^{d^n} \bar{\lambda}_i^n \rightarrow \sum_{i=1}^d \bar{\lambda}_i \quad \text{as } n \rightarrow \infty. \tag{2.39}$$

Now choose and fix

$$\ell \quad \text{such that} \quad \bar{\lambda}_\ell \neq \bar{\lambda}_{\ell+1}. \tag{2.40}$$

Then, by spectral analysis of compact operators and Lemma 2.16 it follows that

$$\bar{\lambda}_i^n \rightarrow \bar{\lambda}_i \quad \text{for } 1 \leq i \leq \ell \text{ as } n \rightarrow \infty. \quad (2.41)$$

Combining (2.39) and (2.41) we derive

$$\sum_{i=\ell+1}^{d^n} \bar{\lambda}_i^n \rightarrow \sum_{i=\ell+1}^d \bar{\lambda}_i \quad \text{as } n \rightarrow \infty.$$

As a consequence of (2.40) and Lemma 2.16 we have  $\lim_{n \rightarrow \infty} \|\bar{\psi}_i^n - \bar{\psi}_i\|_X = 0$  for  $i = 1, \dots, \ell$ . Summarizing the following theorem has been shown.

**Theorem 2.17.** *Let  $X$  be a separable real Hilbert space, the weighting parameters  $\{\alpha_j^n\}_{j=1}^n$  be given by (2.31) and  $y^k \in H^1(0, T; X)$  for  $1 \leq k \leq \wp$ . Let  $\{(\bar{\psi}_i^n, \bar{\lambda}_i^n)\}_{i \in \mathcal{J}}$  and  $\{(\bar{\psi}_i, \bar{\lambda}_i)\}_{i \in \mathcal{J}}$  be eigenvector-eigenvalue pairs satisfying (2.5) and (2.26), respectively. Suppose that  $\ell \in \mathbb{N}$  is fixed such that (2.40) holds. Then we have*

$$\lim_{n \rightarrow \infty} |\bar{\lambda}_i^n - \bar{\lambda}_i| = \lim_{n \rightarrow \infty} \|\bar{\psi}_i^n - \bar{\psi}_i\|_X = 0 \quad \text{for } 1 \leq i \leq \ell,$$

and

$$\lim_{n \rightarrow \infty} \sum_{i=\ell+1}^{d^n} \bar{\lambda}_i^n = \sum_{i=\ell+1}^d \bar{\lambda}_i.$$

**Remark 2.18.** Theorem 2.17 gives an answer to the two questions posed at the beginning of Section 2.2: The time instances  $\{t_j\}_{j=1}^n$  and the associated positive weights  $\{\alpha_j^n\}_{j=1}^n$  should be chosen such that  $\mathcal{R}^n$  is a quadrature approximation of  $\mathcal{R}$  and  $\|\mathcal{R}^n - \mathcal{R}\|_{\mathcal{L}(X)}$  is small (for reasonable  $n$ ). A different strategy is applied in [44], where the time instances  $\{t_j\}_{j=1}^n$  are chosen by an optimization approach. Clearly, other choices for the weights  $\{\alpha_j^n\}_{j=1}^n$  are also possible provided (2.30) is guaranteed. For instance, we can choose the Simpson weights.  $\diamond$

### 3. Reduced-order modelling for evolution problems

In this section error estimates for POD Galerkin schemes for linear evolution problems are presented. The resulting error bounds depend on the number of POD basis functions. Let us refer, e.g., to [18, 22, 30, 40, 41, 42, 64] and [34], where POD Galerkin schemes for parabolic equations and elliptic equations are studied. Moreover, we would like to mention the recent papers [9] and [68], where improved rates of convergence results are derived.

#### 3.1. The abstract evolution problem

Let  $V$  and  $H$  be real, separable Hilbert spaces and suppose that  $V$  is dense in  $H$  with compact embedding. By  $\langle \cdot, \cdot \rangle_H$  and  $\langle \cdot, \cdot \rangle_V$  we denote the inner products in  $H$  and  $V$ , respectively. Let  $T > 0$  the final time. For  $t \in [0, T]$

we define a time-dependent symmetric bilinear form  $a(t; \cdot, \cdot) : V \times V \rightarrow \mathbb{R}$  satisfying

$$|a(t; \varphi, \psi)| \leq \gamma \|\varphi\|_V \|\psi\|_V \quad \forall \varphi \in V \text{ a.e. in } [0, T], \quad (3.1a)$$

$$a(t; \varphi, \varphi) \geq \gamma_1 \|\varphi\|_V^2 - \gamma_2 \|\varphi\|_H^2 \quad \forall \varphi \in V \text{ a.e. in } [0, T] \quad (3.1b)$$

for constants  $\gamma, \gamma_1 > 0$  and  $\gamma_2 \geq 0$  which do not depend on  $t$ . In (3.1), the abbreviation ‘‘a.e.’’ stands for ‘‘almost everywhere’’. By identifying  $H$  with its dual  $H'$  it follows that  $V \hookrightarrow H = H' \hookrightarrow V'$  each embedding being continuous and dense. Recall that the function space (see [10, pp. 472-479] and [72, pp. 146-148], for instance)

$$W(0, T) = \{\varphi \in L^2(0, T; V) \mid \varphi_t \in L^2(0, T; V')\}$$

is a Hilbert space endowed with the inner product

$$\langle \varphi, \phi \rangle_{W(0, T)} = \int_0^T \langle \varphi(t), \phi_t(t) \rangle_V + \langle \varphi_t(t), \phi_t(t) \rangle_{V'} dt \text{ for } \varphi, \phi \in W(0, T)$$

and the induced norm  $\|\varphi\|_{W(0, T)} = \langle \varphi, \varphi \rangle_{W(0, T)}^{1/2}$ . Furthermore,  $W(0, T)$  is continuously embedded into the space  $C([0, T]; H)$ . Hence,  $\varphi(0)$  and  $\varphi(T)$  are meaningful in  $H$  for an element  $\varphi \in W(0, T)$ . The integration by parts formula reads

$$\begin{aligned} \int_0^T \langle \varphi_t(t), \phi(t) \rangle_{V', V} dt + \int_0^T \langle \phi_t(t), \varphi(t) \rangle_{V', V} dt &= \frac{d}{dt} \int_0^T \langle \varphi(t), \psi(t) \rangle_H dt \\ &= \varphi(T)\phi(T) - \varphi(0)\phi(0) \end{aligned}$$

for  $\varphi, \phi \in W(0, T)$ . Moreover, we have the formula

$$\langle \varphi_t(t), \phi \rangle_{V', V} = \frac{d}{dt} \langle \varphi(t), \phi \rangle_H \quad \text{for } (\varphi, \phi) \in W(0, T) \times V \text{ and f.a.a. } t \in [0, T].$$

We suppose that for  $N_u \in \mathbb{N}$  the input space  $U = L^2(0, T; \mathbb{R}^{N_u})$  is chosen. In particular, we identify  $U$  with its dual space  $U'$ . For  $u \in U$ ,  $y_0 \in H$  and  $f \in L^2(0, T; V')$  we consider the linear evolution problem

$$\begin{aligned} \frac{d}{dt} \langle y(t), \varphi \rangle_H + a(t; y(t), \varphi) &= \langle (f + \mathcal{B}u)(t), \varphi \rangle_{V', V} \\ &\quad \forall \varphi \in V \text{ a.e. in } (0, T], \quad (3.2) \\ \langle y(0), \varphi \rangle_H &= \langle y_0, \varphi \rangle_H \quad \forall \varphi \in H, \end{aligned}$$

where  $\langle \cdot, \cdot \rangle_{V', V}$  stands for the dual pairing between  $V$  and its dual space  $V'$  and  $\mathcal{B} : U \rightarrow L^2(0, T; V')$  is a continuous, linear operator.

**Remark 3.1.** Notice that the techniques presented in this work can be adapted for problems, where the input space  $U$  is given by  $L^2(0, T; L^2(\mathcal{D}))$  for some open and bounded domain  $\mathcal{D} \subset \mathbb{R}^{\tilde{N}_u}$  for an  $\tilde{N}_u \in \mathbb{N}$ .  $\diamond$

**Theorem 3.2.** For  $t \in [0, T]$  let  $a(t; \cdot, \cdot) : V \times V \rightarrow \mathbb{R}$  be a time-dependent symmetric bilinear form satisfying (3.1). Then, for every  $u \in U$ ,  $f \in L^2(0, T; V')$  and  $y_\circ \in H$  there is a unique weak solution  $y \in W(0, T)$  satisfying (3.2) and

$$\|y\|_{W(0, T)} \leq C \left( \|y_\circ\|_H + \|f\|_{L^2(0, T; V')} + \|u\|_U \right) \quad (3.3)$$

for a constant  $C > 0$  which is independent of  $u$ ,  $y_\circ$  and  $f$ . If  $f \in L^2(0, T; H)$ ,  $a(t; \cdot, \cdot) = a(\cdot, \cdot)$  (independent of  $t$ ) and  $y_\circ \in V$  hold, we even have  $y \in L^\infty(0, T; V) \cap H^1(0, T; H)$ . Here,  $L^\infty(0, T; V)$  stands for the Banach space of all measurable functions  $\varphi : [0, T] \rightarrow V$  with  $\text{esssup}_{t \in [0, T]} \|\varphi(t)\|_V < \infty$  (see [72, p. 143], for instance).

*Proof.* For a proof of the existence of a unique solution we refer to [10, pp. 512-520]. The a-priori error estimate follows from standard variational techniques and energy estimates. The regularity result follows from [10, pp. 532-533] and [17, pp. 360-364].  $\square$

**Remark 3.3.** We split the solution to (3.2) in one part, which depends on the fixed initial condition  $y_\circ$  and right-hand  $f$ , and another part depending linearly on the input variable  $u$ . Let  $\hat{y} \in W(0, T)$  be the unique solution to

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

We define the subspace

$$W_0(0, T) = \{ \varphi \in W(0, T) \mid \varphi(0) = 0 \text{ in } H \}$$

endowed with the topology of  $W(0, T)$ . Let us now introduce the linear solution operator  $\mathcal{S} : U \rightarrow W_0(0, T)$ : for  $u \in U$  the function  $y = \mathcal{S}u \in W_0(0, T)$  is the unique solution to

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

From  $y \in W_0(0, T)$  we infer  $y(t_b) = 0$  in  $H$ . The boundedness of  $\mathcal{S}$  follows from (3.3). Now, the solution to (3.2) can be expressed as  $y = \hat{y} + \mathcal{S}u$ .  $\diamond$

### 3.2. The POD method for the evolution problem

Let  $u \in U$ ,  $f \in L^2(0, T; V')$  and  $y_\circ \in H$  be given and  $y = \hat{y} + \mathcal{S}u$ . To keep the notation simple we apply only a spatial discretization with POD basis functions, but no time integration by, e.g., the implicit Euler method. Therefore, we utilize the continuous version of the POD method introduced in Section 2.2. In this section we distinguish two choices for  $X$ :  $X = H$  and  $X = V$ . We suppose that the snapshots  $y^k$ ,  $k = 1, \dots, \varphi$ , belong to  $L^2(0, T; V)$

and introduce the following notations:

$$\begin{aligned}\mathcal{R}_V \psi &= \sum_{k=1}^{\wp} \int_0^T \langle \psi, y^k(t) \rangle_V y^k(t) dt && \text{for } \psi \in V, \\ \mathcal{R}_H \psi &= \sum_{k=1}^{\wp} \int_0^T \langle \psi, y^k(t) \rangle_H y^k(t) dt && \text{for } \psi \in H.\end{aligned}\quad (3.4)$$

Moreover, we set  $\mathcal{K}_V = \mathcal{R}_V^*$  and  $\mathcal{K}_H = \mathcal{R}_H^*$ . In Remark 2.15 we have introduced the singular value decomposition of the operator  $\mathcal{Y}$  defined by (2.25). To distinguish the two choices for the Hilbert space  $X$  we denote by the sequence  $\{(\sigma_i^V, \psi_i^V, \phi_i^V)\}_{i \in \mathcal{J}}^\ell \subset \mathbb{R}_0^+ \times V \times L^2(0, T; \mathbb{R}^\wp)$  of triples the singular value decomposition for  $X = V$ , i.e., we have that

$$\mathcal{R}_V \psi_i^V = \lambda_i^V \psi_i^V, \quad \mathcal{K}_V \phi_i^V = \lambda_i^V \phi_i^V, \quad \sigma_i^V = \sqrt{\lambda_i^V}, \quad i \in \mathcal{J}.$$

Furthermore, let the sequence  $\{(\sigma_i^H, \psi_i^H, \phi_i^H)\}_{i \in \mathcal{J}}^\ell \subset \mathbb{R}_0^+ \times H \times L^2(0, T; \mathbb{R}^\wp)$  in satisfy

$$\mathcal{R}_H \psi_i^H = \lambda_i^H \psi_i^H, \quad \mathcal{K}_H \phi_i^H = \lambda_i^H \phi_i^H, \quad \sigma_i^H = \sqrt{\lambda_i^H}, \quad i \in \mathcal{J}.\quad (3.5)$$

The relationship between the singular values  $\sigma_i^H$  and  $\sigma_i^V$  is investigated in the next lemma, which is taken from [68].

**Lemma 3.4.** *Suppose that the snapshots  $y^k \in L^2(0, T; V)$ ,  $k = 1, \dots, \wp$ . Then we have:*

- 1) For all  $i \in \mathcal{J}$  with  $\sigma_i^H > 0$  we have  $\psi_i^H \in V$ .
- 2)  $\sigma_i^V = 0$  for all  $i > d$  with some  $d \in \mathbb{N}$  if and only if  $\sigma_i^H = 0$  for all  $i > d$ , i.e., we have  $d_H = d_V$  if the rank of  $\mathcal{R}_V$  is finite.
- 3)  $\sigma_i^V > 0$  for all  $i \in \mathcal{J}$  if and only if  $\sigma_i^H > 0$  for all  $i \in \mathcal{J}$ .

*Proof.* We argue similarly as in the proof of Lemma 3.1 in [68].

- 1) Let  $\sigma_i^H > 0$  hold. Then, it follows that  $\lambda_i^H > 0$ . We infer from  $y^k \in L^2(0, T; V)$  that  $\mathcal{R}_H \psi \in V$  for any  $\psi \in H$ . Hence, we infer from (3.5) and that  $\psi_i^H = \mathcal{R}_H \psi_i^H / \lambda_i^H \in V$ .
- 2) Assume that  $\sigma_i^V = 0$  for all  $i > d$  with some  $d \in \mathbb{N}$ . Then, we deduce from (2.27) that

$$y^k(t) = \sum_{i=1}^d \langle y^k(t), \psi_i^V \rangle_V \psi_i^V \quad \text{for every } k = 1, \dots, \wp.\quad (3.6)$$

From

$$\begin{aligned}\mathcal{R}_H \psi_j^H &= \sum_{k=1}^{\wp} \int_0^T \langle \psi_j^H, y^k(t) \rangle_H y^k(t) dt \\ &= \sum_{i=1}^d \left( \sum_{k=1}^{\wp} \int_0^T \langle \psi_j^H, y^k(t) \rangle_H \langle y^k(t), \psi_i^V \rangle_V dt \right) \psi_i^V, \quad j \in \mathcal{J},\end{aligned}$$



we conclude that the range of  $\mathcal{R}_H$  is at most  $d$ , which implies that  $\lambda_i^H = 0$  for all  $i > d$ . Analogously, we deduce from  $\sigma_i^H = 0$  for all  $i > d$  that the range of  $\mathcal{R}_V$  is at most  $d$ .

3) The claim follows directly from part 2).  $\square$

Next we recall an *inverse inequality* from [40, Lemma 2].

**Lemma 3.5.** *For all  $v \in \mathcal{V} = \text{span} \{y^k(t) \mid t \in [0, T] \text{ and } 1 \leq k \leq \wp\}$  we*

$$\|v\|_V \leq \sqrt{\|(M^\ell)^{-1}\|_2 \|S^\ell\|_2} \|v\|_H, \quad (3.7)$$

where

$$M^\ell = ((\langle \psi_j, \psi_i \rangle_H)) \in \mathbb{R}^{d \times d} \quad \text{and} \quad S^\ell = ((\langle \psi_j, \psi_i \rangle_V)) \in \mathbb{R}^{d \times d}$$

denote the mass and stiffness matrix, respectively, with  $\psi_i = \psi_i^V$  for  $X = V$  and  $\psi_i = \psi_i^H$  for  $X = H$ . Moreover,  $\|\cdot\|_2$  denotes the spectral norm for symmetric matrices.

*Proof.* Let  $v \in \mathcal{V} \subset V$  be chosen arbitrarily. Then,

$$v = \sum_{i=1}^d \langle v, \psi_i \rangle_X \psi_i$$

with  $\psi_i = \psi_i^V$  for  $X = V$  and  $\psi_i = \psi_i^H$  for  $X = H$ . Defining the vector  $\mathbf{v} = (\langle v, \psi_1 \rangle_X, \dots, \langle v, \psi_d \rangle_X) \in \mathbb{R}^d$  we get

$$\begin{aligned} \|v\|_V^2 &= \mathbf{v}^\top S^\ell \mathbf{v} \leq \|S^\ell\|_2 \mathbf{v}^\top \mathbf{v} \\ &\leq \|S^\ell\|_2 \|(M^\ell)^{-1}\|_2 \mathbf{v}^\top M^\ell \mathbf{v} = \|S^\ell\|_2 \|(M^\ell)^{-1}\|_2 \|v\|_H^2 \end{aligned}$$

which gives (3.7).  $\square$

**Remark 3.6.** In the case  $X = H$  the mass matrix  $M^\ell$  is the identity, whereas  $S^\ell$  is the identity for the choice  $X = V$ . Thus, we have

$$\|v\|_V \leq \sqrt{\|S^\ell\|_2} \|v\|_H \quad \text{and} \quad \|v\|_V \leq \sqrt{\|(M^\ell)^{-1}\|_2} \|v\|_H$$

for  $X = H$  and  $X = V$ , respectively.  $\diamond$

Let us define the two POD subspaces

$$V^\ell = \text{span} \{\psi_1^V, \dots, \psi_\ell^V\} \subset V, \quad H^\ell = \text{span} \{\psi_1^H, \dots, \psi_\ell^H\} \subset V \subset H,$$

where  $H^\ell \subset V$  follows from part 1) of Lemma 3.4. Moreover, we introduce the orthogonal projection operators  $\mathcal{P}_H^\ell : V \rightarrow H^\ell \subset V$  and  $\mathcal{P}^\ell : V \rightarrow V^\ell \subset V$  as follows:

$$\begin{aligned} v^\ell &= \mathcal{P}_H^\ell \varphi \text{ for any } \varphi \in V \quad \text{iff } v^\ell \text{ solves } \min_{w^\ell \in H^\ell} \|\varphi - w^\ell\|_V, \\ v^\ell &= \mathcal{P}_V^\ell \varphi \text{ for any } \varphi \in V \quad \text{iff } v^\ell \text{ solves } \min_{w^\ell \in V^\ell} \|\varphi - w^\ell\|_V. \end{aligned} \quad (3.8)$$

It follows from the first-order optimality conditions that  $v^\ell = \mathcal{P}_H^\ell \varphi$  satisfies

$$\langle v^\ell, \psi_i^H \rangle_V = \langle \varphi, \psi_i^H \rangle_V, \quad 1 \leq i \leq \ell. \quad (3.9)$$

Writing  $v^\ell \in H^\ell$  in the form  $v^\ell = \sum_{j=1}^{\ell} v_j^\ell \psi_j^H$  we derive from (3.9) that the vector  $v^\ell = (v_1^\ell, \dots, v_\ell^\ell)^\top \in \mathbb{R}^\ell$  satisfies the linear system

$$\sum_{j=1}^{\ell} \langle \psi_j^H, \psi_i^H \rangle_V v_j^\ell = \langle \varphi, \psi_i^H \rangle_V, \quad 1 \leq i \leq \ell.$$

For the operator  $\mathcal{P}_V^\ell$  we have the explicit representation

$$\mathcal{P}_V^\ell \varphi = \sum_{i=1}^{\ell} \langle \varphi, \psi_i^V \rangle_V \psi_i^V \text{ for } \varphi \in V.$$

Since the linear operators  $\mathcal{P}_V^\ell$  and  $\mathcal{P}_H^\ell$  are orthogonal projections, we have  $\|\mathcal{P}_V^\ell\|_{\mathcal{L}(V)} = \|\mathcal{P}_H^\ell\|_{\mathcal{L}(V)} = 1$ . As  $\{\psi_i^V\}_{i \in \mathcal{J}}$  is a complete orthonormal basis in  $V$ , we have

$$\lim_{\ell \rightarrow \infty} \int_0^T \|w(t) - \mathcal{P}_V^\ell w(t)\|_V^2 dt = 0 \quad \text{for all } w \in L^2(0, T; V). \quad (3.10)$$

Next we review an essential result from [68, Theorem 6.2], which we will use in our a-priori error analysis for the choice  $X = H$ . Recall that  $\psi_i^H \in V$  holds for  $1 \leq i \leq d$  and the image of  $\mathcal{P}_H^\ell$  belongs to  $V$ . Consequently,  $\|\psi_i^H - \mathcal{P}_H^\ell \psi_i^H\|_V$  is well-defined for  $1 \leq i \leq d$ .

**Theorem 3.7.** *Suppose that  $y^k \in L^2(0, T; V)$  for  $1 \leq k \leq \wp$ . Then,*

$$\sum_{k=1}^{\wp} \int_0^T \|y^k(t) - \mathcal{P}_H^\ell y^k(t)\|_V^2 dt = \sum_{i=\ell+1}^{d_H} \lambda_i^H \|\psi_i^H - \mathcal{P}_H^\ell \psi_i^H\|_V^2.$$

Here,  $d_H$  is the rank of the operator  $\mathcal{R}_H$ , which may be infinite. Moreover,  $\mathcal{P}_H^\ell y^k$  converges to  $y^k$  in  $L^2(0, T; V)$  as  $\ell$  tends to  $\infty$  for each  $k \in \{1, \dots, \wp\}$ .

*Proof.* We sketch the proof. For more details we refer the reader to [68]. Suppose that  $1 \leq \ell \leq d_H$  and  $1 \leq \ell_o < \infty$  hold. Then,  $\lambda_i^H > 0$  for  $1 \leq i \leq \ell$ . Let  $\mathcal{I}_V : V \rightarrow V$  denote the identity operator. As  $\mathcal{I}_V - \mathcal{P}_H^\ell$  is an orthonormal projection on  $V$ , we conclude  $\|\mathcal{I}_V - \mathcal{P}_H^\ell\|_{\mathcal{L}(V)} = 1$ . Furthermore,  $y^k \in L^2(0, T; V)$  holds for each  $k \in \{1, \dots, \wp\}$ . Thus, (3.10) implies that  $\mathcal{P}_V^{\ell_o} y^k \rightarrow y^k$  in  $L^2(0, T; V)$  as  $\ell_o \rightarrow \infty$  for each  $k$ . Hence, we obtain

$$\begin{aligned} & \sum_{k=1}^{\wp} \int_0^T \|(\mathcal{I}_V - \mathcal{P}_H^\ell)(y^k(t) - \mathcal{P}_V^{\ell_o} y^k(t))\|_V^2 dt \\ & \leq \sum_{k=1}^{\wp} \int_0^T \|y^k(t) - \mathcal{P}_V^{\ell_o} y^k(t)\|_V^2 dt = \sum_{i=\ell_o+1}^{d_V} \lambda_i^V \rightarrow 0 \text{ as } \ell_o \rightarrow \infty, \end{aligned}$$

where,  $d_V$  is the rank of the operator  $\mathcal{R}_V$ , which may be infinite. This implies, that  $(\mathcal{I}_V - \mathcal{P}_H^\ell) \mathcal{P}_V^{\ell_o} y^k$  converges to  $(\mathcal{I}_V - \mathcal{P}_H^\ell) y^k$  in  $L^2(0, T; V)$  as  $\ell_o \rightarrow \infty$

for each  $k$ . Hence,

$$\begin{aligned} & \sum_{k=1}^{\wp} \int_0^T \|y^k(t) - \mathcal{P}_H^\ell y^k(t)\|_V^2 dt \\ &= \lim_{\ell_o \rightarrow \infty} \sum_{k=1}^{\wp} \int_0^T \|(\mathcal{I}_V - \mathcal{P}_H^\ell) \mathcal{P}_V^{\ell_o} y^k(t)\|_V^2 dt. \end{aligned} \quad (3.11)$$

Now, we apply the following result [68, Lemma 6.1]:

$$\sum_{k=1}^{\wp} \int_0^T \|(\mathcal{I} - \mathcal{P}_H^\ell) \mathcal{P}_V^{\ell_o} y^k(t)\|_V^2 dt = \sum_{i=1}^{\ell_o} \lambda_i^V \|\psi_i^V - \mathcal{P}_H^\ell \psi_i^V\|_V^2. \quad (3.12)$$

Combining (3.11) and (3.12) we get the error formula:

$$\begin{aligned} & \sum_{k=1}^{\wp} \int_0^T \|y^k(t) - \mathcal{P}_H^\ell y^k(t)\|_V^2 dt \\ &= \lim_{\ell_o \rightarrow \infty} \sum_{i=1}^{\ell_o} \lambda_i^V \|\psi_i^V - \mathcal{P}_H^\ell \psi_i^V\|_V^2 = \sum_{i \in \mathcal{J}} \lambda_i^V \|\psi_i^V - \mathcal{P}_H^\ell \psi_i^V\|_V^2. \end{aligned} \quad (3.13)$$

From  $\|\mathcal{I}_V - \mathcal{P}_H^\ell\|_{\mathcal{L}(V)} = 1$ ,  $\|\psi_i^V\|_V = 1$  for all  $i \in \mathcal{J}$  and  $\sum_{i \in \mathcal{J}} \lambda_i < \infty$  we infer that sum on the right-hand side in (3.13) is finite. Now, the claim follows by arguments from the Hilbert-Schmidt theory.  $\square$

We will also need the following result, which follows from the continuous embedding  $V \hookrightarrow H$ . For a proof we refer to [68, Proposition 6.5].

**Lemma 3.8.** *Let  $y^k \in L^2(0, T; V)$  for each  $k \in \{1, \dots, \wp\}$  and  $\lambda_i^H > 0$  for all  $i \in \mathcal{J}$ . Then,*

$$\lim_{\ell \rightarrow \infty} \|\varphi - \mathcal{P}_H^\ell \varphi\|_V = 0 \quad \text{for all } \varphi \in V.$$

### 3.3. The POD Galerkin approximation

In the context of Section 2.2 we choose  $\wp = 1$ ,  $y^1 = \mathcal{S}u$  and compute a POD basis  $\{\psi_i\}_{i=1}^\ell$  of rank  $\ell$  by solving  $(\mathbf{P}^\ell)$  with  $\psi_i = \psi_i^V$  for  $X = V$  and  $\psi_i = \psi_i^H$  for  $X = H$ . Then, we define the subspace  $X^\ell = \text{span}\{\psi_1, \dots, \psi_\ell\}$ , i.e.,  $X^\ell = V^\ell$  for  $X = V$  and  $X^\ell = H^\ell$  for  $X = H$ . Now we approximate the state variable  $y$  by the Galerkin expansion

$$y^\ell(t) = \hat{y}(t) + \sum_{i=1}^{\ell} y_i^\ell(t) \psi_i \in V \quad \text{a.e. in } [0, T] \quad (3.14)$$

with coefficient functions  $y_i^\ell : [0, T] \rightarrow \mathbb{R}$ . We introduce the vector-valued coefficient function

$$y^\ell = (y_1^\ell, \dots, y_\ell^\ell) : [0, T] \rightarrow \mathbb{R}^\ell.$$

Since  $\hat{y}(0) = y_o$  holds, we suppose that  $y^\ell(0) = 0$ . Then,  $y^\ell(0) = y_o$  is valid, i.e., the POD state matches exactly the initial condition. Inserting (3.14) into

(3.2) and using the test space in  $V^\ell$  for  $1 \leq i \leq \ell$  we obtain the following POD Galerkin scheme for (3.2):  $y^\ell \in W(0, T)$  solves

$$\begin{aligned} \frac{d}{dt} \langle y^\ell(t), \psi \rangle_H + a(t; y^\ell(t), \psi) &= \langle (f + \mathcal{B}u)(t), \psi \rangle_{V', V} \quad \forall \psi \in X^\ell \text{ a.e.}, \\ y^\ell(0) &= 0. \end{aligned} \quad (3.15)$$

We call (3.15) a *low dimensional* or *reduced-order model* for (3.2).

**Proposition 3.9.** *Let all assumptions of Theorem 3.2 be satisfied and the POD basis of rank  $\ell$  be computed as described at the beginning of Section 3.1. Then, there exists a unique solution  $y^\ell \in H^1(0, T; \mathbb{R}^\ell) \hookrightarrow W(0, T)$  solving (3.15).*

*Proof.* Choosing  $\psi = \psi_i$ ,  $1 \leq i \leq \ell$ , and applying (3.14) we infer from (3.15) that the coefficient vector  $y^\ell$  satisfies

$$M^\ell \dot{y}^\ell(t) + A^\ell(t)y(t) = \hat{F}^\ell(t) \text{ a.e. in } [0, T], \quad y^\ell(0) = 0, \quad (3.16)$$

where we have set

$$\begin{aligned} M^\ell &= ((\langle \psi_i, \psi_j \rangle_X)) \in \mathbb{R}^{\ell \times \ell}, \quad A^\ell(t) = ((a(t; \psi_i, \psi_j))) \in \mathbb{R}^{\ell \times \ell}, \\ \hat{F}^\ell(t) &= (\langle (f + \mathcal{B}u)(t) - \hat{y}_t(t), \psi_i \rangle_{V', V} - a(t; \hat{y}_t(t), \psi_i)) \in \mathbb{R}^\ell \end{aligned} \quad (3.17)$$

with  $\psi_i = \psi_i^V$  for  $X = V$  and  $\psi_i = \psi_i^H$  for  $X = H$ . Since (3.16) is a linear ordinary differential equation system the existence of a unique  $y^\ell \in H^1(0, T; \mathbb{R}^\ell)$  follows by standard arguments.  $\square$

**Remark 3.10.** 1) In contrast to [29, 73], for instance, the POD approximation does not belong to  $X^\ell$ , but to the affine space  $\hat{y} + X^\ell$  provided  $\hat{y} \neq 0$ . The benefit of this approach is that  $y^\ell(0) = y_0$  – and not  $y^\ell(0) = \mathcal{P}_H^\ell y_0$  or  $y^\ell(0) = \mathcal{P}_V^\ell y_0$ . This improves the approximation quality of the POD basis which is illustrated in our numerical tests.

2) We proceed analogously to Remark 3.3 and introduce the linear and bounded solution operator  $\mathcal{S}^\ell : U \rightarrow W_0(0, T)$ : for  $u \in U$  the function  $w^\ell = \mathcal{S}^\ell u \in W(0, T)$  satisfies  $w^\ell(0) = 0$  and

$$\frac{d}{dt} \langle w^\ell(t), \psi \rangle_H + a(t; w^\ell(t), \psi) = \langle (\mathcal{B}u)(t), \psi \rangle_{V', V} \quad \forall \psi \in X^\ell \text{ a.e.}$$

Then, the solution to (3.15) is given by  $y^\ell = \hat{y} + \mathcal{S}^\ell u$ . Analogous to the proof of (3.3) we derive that there exists a positive constant  $C_2$  which does not depend on  $\ell$  or  $u$  so that

$$\|\mathcal{S}^\ell u\|_{W(0, T)} \leq C \|u\|_U.$$

Thus,  $\mathcal{S}^\ell$  is bounded uniformly with respect to  $\ell$ .  $\diamond$

To investigate the convergence of the error  $y - y^\ell$  we make use of the following two inequalities:

1) *Gronwall's inequality:* For  $T > 0$  let  $v : [0, T] \rightarrow \mathbb{R}$  be a nonnegative, differentiable function satisfying

$$v'(t) \leq \varphi(t)v(t) + \chi(t) \quad \text{for all } t \in [0, T],$$

where  $\varphi$  and  $\chi$  are real-valued, nonnegative, integrable functions on  $[0, T]$ . Then

$$v(t) \leq \exp\left(\int_0^t \varphi(s) ds\right) \left(v(0) + \int_0^t \chi(s) ds\right) \quad \text{for all } t \in [0, T]. \quad (3.18)$$

In particular, if

$$v' \leq \varphi v \text{ in } [0, T] \quad \text{and} \quad v(0) = 0$$

hold, then  $v = 0$  in  $[0, T]$ .

2) *Young's inequality*: For every  $a, b \in \mathbb{R}$  and for every  $\varepsilon > 0$  we have

$$ab \leq \frac{\varepsilon a^2}{2} + \frac{b^2}{2\varepsilon}.$$

**Theorem 3.11.** *Let  $u \in U$  be chosen arbitrarily so that  $Su \neq 0$ .*

1) *To compute a POD basis  $\{\psi_i\}_{i=1}^\ell$  of rank  $\ell$  we choose  $\varphi = 1$  and  $y^1 = Su$ . Then,  $y = \hat{y} + Su$  and  $y^\ell = \hat{y} + \mathcal{S}^\ell u$  satisfies the a-priori error estimate*

$$\|y^\ell - y\|_{W(0,T)}^2 \leq \begin{cases} 2 \sum_{i=\ell+1}^{d_V} \lambda_i^V + C_1 \|y_t^1 - \mathcal{P}_V^\ell y_t^1\|_{L^2(0,T;V')}^2 & \text{if } X = V, \\ 2 \sum_{i=\ell+1}^{d_H} \lambda_i^H \|\psi_i^H - \mathcal{P}_H^\ell \psi_i^H\|_V^2 \\ \quad + C_1 \|y_t^1 - \mathcal{P}_H^\ell y_t^1\|_{L^2(0,T;V')}^2 & \text{if } X = H, \end{cases} \quad (3.19)$$

where the constant  $C_1$  depends on the terminal time  $T$  and the constants  $\gamma, \gamma_1, \gamma_2$  introduced in (3.1).

2) *Suppose that  $Su \in H^1(0, T; V)$  holds true. If we set  $\varphi = 2$  and compute a POD basis of rank  $\ell$  using the trajectories  $y^1 = Su$  and  $y^2 = (Su)_t$ , it follows that*

$$\|y^\ell - y\|_{W(0,T)}^2 \leq \begin{cases} C_2 \sum_{i=\ell+1}^{d_V} \lambda_i^V & \text{for } X = V, \\ C_2 \sum_{i=\ell+1}^{d_H} \lambda_i^H \|\psi_i^V - \mathcal{P}_H^\ell \psi_i^H\|_V^2 & \text{for } X = H, \end{cases} \quad (3.20)$$

for a constant  $C_2$  which depends on  $\gamma, \gamma_1, \gamma_2$ , and  $T$ .

3) *If  $S\tilde{u}$  belongs to  $H^1(0, T; V)$  for every  $\tilde{u} \in U$  and if  $\lambda_i^H > 0$  for all  $i \in \mathcal{J}$ , then we have*

$$\lim_{\ell \rightarrow \infty} \|\mathcal{S} - \mathcal{S}^\ell\|_{\mathcal{L}(U, W(0,T))} = 0. \quad (3.21)$$

*Proof.* 1) For almost all  $t \in [0, T]$  we make use of the decomposition

$$\begin{aligned} y^\ell(t) - y(t) &= \hat{y}(t) + (\mathcal{S}^\ell u)(t) - \hat{y}(t) - (Su)(t) \\ &= (\mathcal{S}^\ell u)(t) - \mathcal{P}^\ell((\mathcal{S}^\ell u)(t)) + \mathcal{P}^\ell((\mathcal{S}^\ell u)(t)) - (Su)(t) \\ &= \vartheta^\ell(t) - \varrho^\ell(t), \end{aligned} \quad (3.22)$$

where  $\mathcal{P}^\ell = \mathcal{P}_V^\ell$  for  $X = V$ ,  $\mathcal{P}^\ell = \mathcal{P}_H^\ell$  for  $X = H$ ,  $\vartheta^\ell(t) = (\mathcal{S}^\ell u)(t) - \mathcal{P}^\ell((\mathcal{S}^\ell u)(t)) \in X^\ell$  and  $\varrho^\ell(t) = \mathcal{P}^\ell((\mathcal{S}^\ell u)(t)) - (\mathcal{S}u)(t)$ . From  $y^1(t) = (\mathcal{S}u)(t)$  and (2.27) we infer that

$$\begin{aligned} \|\varrho^\ell\|_{W(0,T)}^2 &= \|y^1 - \mathcal{P}_V^\ell y^1(t)\|_{L^2(0,T;V)}^2 + \|y_t^1 - \mathcal{P}_V^\ell y_t^1(t)\|_{L^2(0,T;V')}^2 \\ &= \sum_{i=\ell+1}^{d_V} \lambda_i + \|y_t^1 - \mathcal{P}_V^\ell y_t^1(t)\|_{L^2(0,T;V')}^2 \end{aligned} \quad (3.23)$$

in case of  $X = V$ , where  $d_V$  stands for rank of  $\mathcal{R}_V$ . For the choice  $X = H$  we derive from Theorem 3.7 that

$$\|\varrho^\ell\|_{W(0,T)}^2 = \sum_{i=\ell+1}^{d_H} \lambda_i^H \|\psi_i^H - \mathcal{P}_H^\ell \psi_i^H\|_V^2 + \|y_t^1 - \mathcal{P}_V^\ell y_t^1(t)\|_{L^2(0,T;V')}^2. \quad (3.24)$$

Here,  $d_H$  denotes for rank of  $\mathcal{R}_H$ . Using (3.2), (3.15), and

$$\vartheta^\ell(t) = y^\ell - \hat{y}(t) - \mathcal{P}^\ell((\mathcal{S}u)(t)) = y^\ell(t) - y(t) + (\mathcal{S}u)(t) - \mathcal{P}^\ell((\mathcal{S}u)(t))$$

we derive that

$$\frac{d}{dt} \langle \vartheta^\ell(t), \psi \rangle_H + a(t; \vartheta^\ell(t), \psi) = \langle y_t^1(t) - \mathcal{P}^\ell y_t^1(t), \psi \rangle_{V',V} \quad (3.25)$$

for all  $\psi \in X^\ell$  and for almost all  $t \in [0, T]$ . From choosing  $\psi = \vartheta^\ell(t)$ , (3.1b) and (3.21) we find

$$\frac{d}{dt} \|\vartheta^\ell(t)\|_H^2 + \gamma_1 \|\vartheta^\ell(t)\|_V^2 - 2\gamma_2 \|\vartheta^\ell(t)\|_H^2 \leq \frac{1}{\gamma_1} \|y_t^1(t) - \mathcal{P}^\ell y_t^1(t)\|_{V'}^2.$$

From (3.18) – setting  $v(t) = \|\vartheta^\ell(t)\|_H^2 \geq 0$ ,  $\varphi(t) = \gamma_2 > 0$ ,  $\chi(t) = \|y_t^1(t) - \mathcal{P}^\ell y_t^1(t)\|_{L^2(0,T;V')}^2 \geq 0$  – and  $\vartheta^\ell(0) = 0$  it follows that

$$\|\vartheta^\ell(t)\|_H^2 \leq c_1 \|y_t^1 - \mathcal{P}^\ell y_t^1\|_{L^2(0,T;V')}^2 \quad \text{for almost all } t \in [0, T]$$

with  $c_1 = \exp(\gamma_2 T)$ , so that

$$\begin{aligned} \|\vartheta^\ell\|_{L^2(0,T;V)}^2 &\leq \frac{2\gamma_2}{\gamma_1} \|\vartheta^\ell\|_{L^2(0,T;H)}^2 + \frac{1}{\gamma_1^2} \|y_t^1 - \mathcal{P}^\ell y_t^1\|_{L^2(0,T;V')}^2 \\ &\leq c_2 \|y_t^1 - \mathcal{P}^\ell y_t^1\|_{L^2(0,T;V')}^2 \end{aligned} \quad (3.26)$$

with  $c_2 = \max(2\gamma_2 T c_1, 1/\gamma_1)/\gamma_1$ . Moreover, we conclude from (3.1a), (3.19) and (3.26) that

$$\begin{aligned} \|\vartheta_t^\ell\|_{L^2(0,T;V')}^2 &\leq \frac{\gamma}{2} \|\vartheta^\ell\|_{L^2(0,T;V)}^2 + \frac{1}{2} \|y_t^1 - \mathcal{P}^\ell y_t^1\|_{L^2(0,T;V')}^2 \\ &\leq c_3 \|y_t^1 - \mathcal{P}^\ell y_t^1\|_{L^2(0,T;V')}^2 \end{aligned} \quad (3.27)$$

with  $c_3 = \max(\gamma c_2, 1)/2$ . Combining (3.22), (3.23), (3.26) and (3.27) we obtain (3.22) with  $C_1 = 2 \max(1, c_2, c_3)$ .

2) The claim follows directly from

$$\begin{aligned} \|(\mathcal{S}u)_t - \mathcal{P}^\ell(\mathcal{S}u)_t\|_{L^2(0,T;V)}^2 &= \|y^2 - \mathcal{P}^\ell y^2\|_{L^2(0,T;V)}^2 \\ &\leq \begin{cases} \sum_{i=\ell+1}^{d_V} \lambda_i^V & \text{if } X = V, \\ \sum_{i=\ell+1}^{d_H} \lambda_i^H \|\psi_i^H - \mathcal{P}_H^\ell \psi_i^H\|_V^2 & \text{if } X = H. \end{cases} \end{aligned}$$

3) Using  $\mathcal{S}\tilde{u} \in H^1(0, T; V)$  for any  $\tilde{u} \in U$ , Remark 2.8 and applying the proof of Proposition 4.4 in [73] we infer that there exists a constant  $C_3$  which is independent of  $\ell$  satisfying

$$\begin{aligned} \|\mathcal{S} - \mathcal{S}^\ell\|_{\mathcal{L}(U, W(0, T))} &= \sup_{\|\tilde{u}\|_U=1} \|(\mathcal{S} - \mathcal{S}^\ell)\tilde{u}\|_{W(0, T)} \\ &\leq c_3 \sup_{\|\tilde{u}\|_U=1} \int_0^T \|\tilde{y}(t) - \mathcal{P}^\ell \tilde{y}(t)\|_V^2 + \|\tilde{y}_t(t) - \mathcal{P}^\ell \tilde{y}_t(t)\|_V^2 dt \xrightarrow{\ell \rightarrow \infty} 0 \end{aligned}$$

with  $\tilde{y} = \mathcal{S}\tilde{u}$ . By assumption, the elements  $\tilde{y}(t)$  and  $\tilde{y}_t(t)$  belong to  $L^2(0, T; V)$ . Now the claim follows for  $X = V$  from (3.10) and for  $X = H$  from Lemma 3.8.  $\square$

**Remark 3.12.** 1) Note that the a-priori error estimates (3.19) and (3.20) depend on the arbitrarily chosen, but fixed control  $u \in U$ , which is also utilized to compute the POD basis. Moreover, these a-priori estimates do not involve errors by the POD discretization of the initial condition  $y_\circ$  – in contrast to the error analysis presented in [29, 40, 41, 64, 73], for instance.

2) From (3.21) we infer

$$\|\hat{y} + \mathcal{S}^\ell \tilde{u} - \hat{y} - \mathcal{S}\tilde{u}\|_{W(0, T)} \leq \|\mathcal{S} - \mathcal{S}^\ell\|_{\mathcal{L}(U, W(0, T))} \|\tilde{u}\|_U \xrightarrow{\ell \rightarrow \infty} 0$$

for any  $\tilde{u} \in U$ .

3) For the numerical realization we have to utilize also a time integration method like, e.g., the implicit Euler or the Crank-Nicolson method. We refer the reader to [40, 41, 42], where different time discretization schemes are considered. Moreover, in [47, 64] also a finite element discretization of the ansatz space  $V$  is incorporated in the a-priori error analysis.  $\diamond$

**Example 3.13.** Accurate approximation results are achieved if the subspace spanned by the snapshots is (approximately) of low dimension. Let  $T > 0$ ,  $\Omega = (0, 2) \subset \mathbb{R}$  and  $Q = (0, T) \times \Omega$ . We set  $f(t, \mathbf{x}) = e^{-t}(\pi^2 - 1) \sin(\pi \mathbf{x})$  for  $(t, \mathbf{x}) \in Q$  and  $y_\circ(\mathbf{x}) = \sin(\pi \mathbf{x})$  for  $\mathbf{x} \in \Omega$ . Let  $H = L^2(\Omega)$ ,  $V = H_0^1(\Omega)$  and

$$a(t; \varphi, \phi) = \int_\Omega \varphi'(\mathbf{x}) \phi'(\mathbf{x}) d\mathbf{x} \quad \text{for } \varphi, \phi \in V,$$

i.e., the bilinear form  $a$  is independent of  $t$ . Finally, we choose  $u = 0$ . Then, the exact solution to (3.2) is given by  $y(t, \mathbf{x}) = e^{-t} \sin(\pi \mathbf{x})$  spans the one-dimensional space  $\{\alpha \psi \mid \alpha \in \mathbb{R}\}$  with  $\psi(\mathbf{x}) = \sin(\pi \mathbf{x})$ . Choosing the space  $X = H$ , this implies that all eigenvalues of the operator  $\mathcal{R}_H$  introduced in (3.4) except of the first one are zero and  $\psi_1 = \psi \in V$  is the single POD element corresponding to a nontrivial eigenvalue of  $\mathcal{R}_H$ . Further, the reduced order model of the rank-1 POD-Galerkin ansatz

$$\begin{aligned} \dot{y}^1(t) + \|\psi_1'\|_H^2 y^1(t) &= \langle f(t), \psi_1 \rangle_H \quad \text{for } t \in (0, T], \\ y^1(0) &= \langle y_0, \psi_1 \rangle_H \end{aligned}$$

has the solution  $y^1(t) = e^{-t}$ , so both the projection

$$(\mathcal{P}^1 y)(t, \mathbf{x}) = \langle y(t), \psi_1 \rangle_X \psi_1(\mathbf{x}), \quad (t, \mathbf{x}) \in \bar{Q},$$

of the state  $y$  on the POD-Galerkin space and the reduced-order solution  $y^1(t) = y^1(t)\psi_1$  coincide with the exact solution  $y$ . In the latter case, this is due to the fact that the data functions  $f$  and  $y_0$  as well as all time derivative snapshots  $\dot{y}(t)$  are already elements of  $\text{span}(\psi_1)$ , so no projection error occurs here, cp. the a priori error bounds given in (3.20). In the case  $X = V$ , we get the same results with  $\psi_1(\mathbf{x}) = \sin(\pi \mathbf{x})/2$  and  $y^1(t) = 2e^{-t}$ .  $\diamond$

Utilizing the techniques as in the proof of Theorem 7.5 in [68] one can derive an a-priori error bound without including the time derivatives into the snapshot subspace. In the next proposition we formulate the a-priori error estimate.

**Proposition 3.14.** *Let  $y_0 \in V$  and  $u \in U$  be chosen arbitrarily so that  $Su \neq 0$ . To compute a POD basis  $\{\psi_i\}_{i=1}^\ell$  of rank  $\ell$  we choose  $\varphi = 1$  and  $y^1 = Su$ . Then,  $y = \hat{y} + Su$  and  $y^\ell = \hat{y} + \mathcal{S}^\ell u$  satisfies the a-priori error estimate*

$$\|y^\ell - y\|_{W(0,T)}^2 \leq \begin{cases} C \sum_{i=\ell+1}^{d_V} \lambda_i^V \|\psi_i^V - \mathcal{P}^\ell \psi_i^V\|_V^2 & \text{if } X = V, \\ C \sum_{i=\ell+1}^{d_H} \lambda_i^H \|\psi_i^H\|_V^2 & \text{if } X = H, \end{cases} \quad (3.28)$$

where the constant  $C$  depends on the terminal time  $T$  and the constants  $\gamma, \gamma_1, \gamma_2$  introduced in (3.1). Moreover,  $\mathcal{P}^\ell : V \rightarrow V^\ell$  is the orthogonal projection given as follows:

$$v^\ell = \mathcal{P}_V^\ell \varphi \text{ for any } \varphi \in V \text{ iff } v^\ell \text{ solves } \min_{w^\ell \in V^\ell} \|\varphi - w^\ell\|_H.$$

In particular, we have  $y^\ell \rightarrow y$  in  $W(0, T)$  as  $\ell \rightarrow \infty$ .

## 4. The linear-quadratic optimal control problem

In this section we apply a POD Galerkin approximation to linear-quadratic optimal control problems. Linear-quadratic problems are interesting in several respects: in particular, they occur in each level of a sequential quadratic



programming (SQP) methods; see, e.g., [54]. In contrast to methods of balanced truncation type, the POD method is somehow lacking a reliable a-priori error analysis. Unless its snapshots are generating a sufficiently rich state space, it is not a-priorily clear how far the optimal solution of the POD problem is from the exact one. On the other hand, the POD method is a universal tool that is applicable also to problems with time-dependent coefficients or to nonlinear equations. By generating snapshots from the real (large) model, a space is constructed that inherits the main and relevant physical properties of the state system. This, and its ease of use makes POD very competitive in practical use, despite of certain heuristic.

Here we prove convergence and derive a-priori error estimates for the optimal control problem. The error estimates rely on the (unrealistic) assumption that the POD basis is computed from the (exact) optimal solution. However, these estimates are utilized to develop an a-posteriori error analysis for the POD Galerkin approximation of the optimal control problem. Using a perturbation method [16] we deduce how far the suboptimal control, computed by the POD Galerkin approximation, is from the (unknown) exact one. This idea turned out to be very efficient in our numerical examples. Thus, we are able to compensate for the lack of an a-priori error analysis for the POD method.

#### 4.1. Problem formulation

In this section we introduce our optimal control problem, which is a constrained optimization problem in a Hilbert space. The objective is a quadratic function. The evolution problem (3.2) serves as an equality constraint. Moreover, bilateral control bounds lead to inequality constraints in the minimization. For the readers convenience we recall (3.2) here. Let  $U = L^2(0, T; \mathbb{R}^{N_u})$  denote the control space with  $N_u \in \mathbb{N}$ . For  $u \in U$ ,  $y_0 \in H$  and  $f \in L^2(0, T; V')$  we consider the state equation

$$\begin{aligned} \frac{d}{dt} \langle y(t), \varphi \rangle_H + a(t; y(t), \varphi) &= \langle (f + \mathcal{B}u)(t), \varphi \rangle_{V', V} \\ &\quad \forall \varphi \in V \text{ a.e. in } (0, T], \quad (4.1) \\ \langle y(0), \varphi \rangle_H &= \langle y_0, \varphi \rangle_H \quad \forall \varphi \in H, \end{aligned}$$

where  $\mathcal{B} : U \rightarrow L^2(0, T; V')$  is a continuous, linear operator. Due to Theorem 3.2 there exists a unique solution  $y \in W(0, T)$  to (4.1).

We introduce the Hilbert space

$$X = W(0, T) \times U$$

endowed with the natural product topology, i.e., with the inner product

$$\langle x, \tilde{x} \rangle_X = \langle y, \tilde{y} \rangle_{W(0, T)} + \langle u, \tilde{u} \rangle_U \quad \text{for } x = (y, u), \tilde{x} = (\tilde{y}, \tilde{u}) \in X$$

and the norm  $\|x\|_X = (\|y\|_{W(0, T)}^2 + \|u\|_U^2)^{1/2}$  for  $x = (y, u) \in X$ .

**Assumption 1.** For  $t \in [0, T]$  let  $a(t; \cdot, \cdot) : V \times V \rightarrow \mathbb{R}$  be a time-dependent symmetric bilinear form satisfying (3.1). Moreover,  $f \in L^2(0, T; V')$ ,  $y_0 \in H$  and  $\mathcal{B} \in \mathcal{L}(U, L^2(0, T; V'))$  holds.

In Remark 3.3 we have introduced the particular solution  $\hat{y} \in W(0, T)$  as well as the linear, bounded solution operator  $\mathcal{S}$ . Then, the solution to (4.1) can be expressed as  $y = \hat{y} + \mathcal{S}u$ . By  $X_{\text{ad}}$  we denote the closed, convex and bounded set of admissible solutions for the optimization problem as

$$X_{\text{ad}} = \{(\hat{y} + \mathcal{S}u, u) \in X \mid u_a \leq u \leq u_b \text{ in } \mathbb{R}^m \text{ a.e. in } [0, T]\},$$

where  $u_a = (u_{a,1}, \dots, u_{a,N_u})$ ,  $u_b = (u_{b,1}, \dots, u_{b,N_u}) \in U$  satisfy  $u_{a,i} \leq u_{b,i}$  for  $1 \leq i \leq N_u$  a.e. in  $[0, T]$ . Since  $u_{a,i} \leq u_{b,i}$  holds for  $1 \leq i \leq N_u$ , we infer from Theorem 3.2 that the set  $X_{\text{ad}}$  is nonempty.

The quadratic objective  $J : X \rightarrow \mathbb{R}$  is given by

$$J(x) = \frac{\sigma_Q}{2} \int_0^T \|y(t) - y_Q(t)\|_H^2 dt + \frac{\sigma_\Omega}{2} \|y(T) - y_\Omega\|_H^2 + \frac{\sigma}{2} \|u\|_U^2 \quad (4.2)$$

for  $x = (y, u) \in X$ , where  $(y_Q, y_\Omega) \in L^2(0, T; H) \times H$  are given desired states. Furthermore,  $\sigma_Q, \sigma_\Omega \geq 0$  and  $\sigma > 0$ . Of course, more general cost functionals can be treated analogously.

Now the quadratic programming problem is given by

$$\min J(x) \quad \text{subject to (s.t.) } x \in X_{\text{ad}}. \quad (\mathbf{P})$$

From  $x = (y, u) \in X_{\text{ad}}$  we infer that  $y = \hat{y} + \mathcal{S}u$  holds. Hence,  $y$  is a dependent variable. We call  $u$  the *control* and  $y$  the *state*. In this way,  $(\mathbf{P})$  becomes an *optimal control problem*. Utilizing the relationship  $y = \hat{y} + \mathcal{S}u$  we define a so-called *reduced cost functional*  $\hat{J} : U \rightarrow \mathbb{R}$  by

$$\hat{J}(u) = J(\hat{y} + \mathcal{S}u, u) \quad \text{for } u \in U.$$

Moreover, the set of admissible controls is given as

$$U_{\text{ad}} = \{u \in U \mid u_a \leq u \leq u_b \text{ in } \mathbb{R}^m \text{ a.e. in } [0, T]\},$$

which is convex, closed and bounded in  $U$ . Then, we consider the reduced optimal control problem:

$$\min \hat{J}(u) \quad \text{s.t. } u \in U_{\text{ad}}. \quad (\hat{\mathbf{P}})$$

Clearly, if  $\bar{u}$  is the optimal solution to  $(\hat{\mathbf{P}})$ , then  $\bar{x} = (\hat{y} + \mathcal{S}\bar{u}, \bar{u})$  is the optimal solution to  $(\mathbf{P})$ . On the other hand, if  $\bar{x} = (\bar{y}, \bar{u})$  is the solution to  $(\mathbf{P})$ , then  $\bar{u}$  solves  $(\hat{\mathbf{P}})$ .

**Example 4.1.** We introduce an example for  $(\mathbf{P})$  and discuss the presented theory for this application. Let  $\Omega \subset \mathbb{R}^d$ ,  $d \in \{1, 2, 3\}$ , be an open and bounded domain with Lipschitz-continuous boundary  $\Gamma = \partial\Omega$ . For  $T > 0$  we set  $Q = (0, T) \times \Omega$  and  $\Sigma = (0, T) \times \Gamma$ . We choose  $H = L^2(\Omega)$  and  $V = H_0^1(\Omega)$  endowed with the usual inner products

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

and their induced norms, respectively. Let  $\chi_i \in H$ ,  $1 \leq i \leq m$ , denote given control shape functions. Then, for given control  $u \in U$ , initial condition

$y_\circ \in H$  and inhomogeneity  $f \in L^2(0, T; H)$  we consider the linear heat equation

$$\begin{aligned} y_t(t, \mathbf{x}) - \Delta y(t, \mathbf{x}) &= f(t, \mathbf{x}) + \sum_{i=1}^m u_i(t) \chi_i(\mathbf{x}), & \text{a.e. in } Q, \\ y(t, \mathbf{x}) &= 0, & \text{a.e. in } \Sigma, \\ y(0, \mathbf{x}) &= y_\circ(\mathbf{x}), & \text{a.e. in } \Omega. \end{aligned} \tag{4.3}$$

We introduce the time-independent, symmetric bilinear form

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

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

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

Hence, we have  $\gamma = \gamma_1 = \gamma_2 = 1$  in (3.1). It follows that the weak formulation of (4.3) can be expressed in the form (3.2). Moreover, the unique weak solution to (4.3) belongs to the space  $L^\infty(0, T; V)$  provided  $y_\circ \in V$  holds.  $\diamond$

#### 4.2. Existence of a unique optimal solution

We suppose the following hypothesis for the objective.

**Assumption 2.** *In (4.2) the desired states  $(y_Q, y_\Omega)$  belong to  $L^2(0, T; H) \times H$ . Furthermore,  $\sigma_Q, \sigma_\Omega \geq 0$  and  $\sigma > 0$  are satisfied.*

Let us review the following result for quadratic optimization problems in Hilbert spaces; see [72, pp. 50-51].

**Theorem 4.2.** *Suppose that  $\mathcal{U}$  and  $\mathcal{H}$  are given Hilbert spaces with norms  $\|\cdot\|_{\mathcal{U}}$  and  $\|\cdot\|_{\mathcal{H}}$ , respectively. Furthermore, let  $\mathcal{U}_{\text{ad}} \subset \mathcal{U}$  be non-empty, bounded, closed, convex and  $z_d \in \mathcal{H}$ ,  $\kappa \geq 0$ . The mapping  $\mathcal{G} : \mathcal{U} \rightarrow \mathcal{H}$  is assumed to be a linear and continuous operator. Then there exists an optimal control  $\bar{u}$  solving*

$$\min_{u \in \mathcal{U}_{\text{ad}}} \mathcal{J}(u) := \frac{1}{2} \|\mathcal{G}u - z_d\|_{\mathcal{H}}^2 + \frac{\kappa}{2} \|u\|_{\mathcal{U}}^2. \tag{4.4}$$

*If  $\kappa > 0$  holds or if  $\mathcal{G}$  is injective, then  $\bar{u}$  is uniquely determined.*

**Remark 4.3.** In the proof of Theorem 4.2 it is only used that  $\mathcal{J}$  is continuous and convex. Therefore, the existence of an optimal control follows for general convex, continuous cost functionals  $\mathcal{J} : \mathcal{U} \rightarrow \mathbb{R}$  with a Hilbert space  $\mathcal{U}$ .  $\diamond$

Next we can use Theorem 4.2 to obtain an existence result for the optimal control problem  $(\hat{\mathbf{P}})$ , which imply the existence of an optimal solution to  $(\mathbf{P})$ .

**Theorem 4.4.** *Let Assumptions 1 and 2 be valid. Moreover, let the bilateral control constraints  $u_a, u_b \in U$  satisfy  $u_a \leq u_b$  componentwise in  $\mathbb{R}^m$  a.e. in  $[0, T]$ . Then,  $(\hat{\mathbf{P}})$  has a unique optimal solution  $\bar{u}$ .*

*Proof.* Let us choose the Hilbert spaces  $\mathcal{H} = L^2(0, T; H) \times H$  and  $\mathcal{U} = U$ . Moreover,  $\mathcal{E} : W(0, T) \rightarrow L^2(0, T; H)$  is the canonical embedding operator, which is linear and bounded. We define the operator  $\mathcal{E}_2 : W(0, T) \rightarrow H$  by  $\mathcal{E}_2\varphi = \varphi(T)$  for  $\varphi \in W(0, T)$ . Since  $W(0, T)$  is continuously embedded into  $C([0, T]; H)$ , the linear operator  $\mathcal{E}_2$  is continuous. Finally, we set

$$\mathcal{G} = \begin{pmatrix} \sqrt{\sigma_Q} \mathcal{E}_1 \mathcal{S} \\ \sqrt{\sigma_\Omega} \mathcal{E}_2 \mathcal{S} \end{pmatrix} \in \mathcal{L}(\mathcal{U}, \mathcal{H}), \quad z_d = \begin{pmatrix} \sqrt{\sigma_Q} (y_Q - \hat{y}) \\ \sqrt{\sigma_\Omega} (y_\Omega - \hat{y}(T)) \end{pmatrix} \in \mathcal{H} \quad (4.5)$$

and  $\mathcal{U}_{\text{ad}} = U_{\text{ad}}$ . Then,  $(\hat{\mathbf{P}})$  and (4.4) coincide. Consequently, the claim follows from Theorem 4.2 and  $\sigma > 0$ .  $\square$

Next we consider the case that  $u_a = -\infty$  or/and  $u_b = +\infty$ . In this case  $U_{\text{ad}}$  is not bounded. However, we have the following result [72, p. 52].

**Theorem 4.5.** *Let Assumptions 1 and 2 be satisfied. If  $u_a = -\infty$  or/and  $u_b = +\infty$ , problem  $(\hat{\mathbf{P}})$  admits a unique solution.*

*Proof.* We utilize the setting of the proof of Theorem 4.4. By assumption there exists an element  $u_0 \in U_{\text{ad}}$ . For  $u \in U$  with  $\|u\|_U^2 > 2\hat{J}(u_0)/\sigma$  we have

$$\hat{J}(u) = \mathcal{J}(u) = \frac{1}{2} \|\mathcal{G}u - z_d\|_{\mathcal{H}}^2 + \frac{\sigma}{2} \|u\|_U^2 \geq \frac{\sigma}{2} \|u\|_U^2 > \hat{J}(u_0).$$

Thus, the minimization of  $\hat{J}$  over  $U_{\text{ad}}$  is equivalent with the minimization of  $\hat{J}$  over the bounded, convex and closed set

$$U_{\text{ad}} \cap \left\{ u \in U \mid \|u\|_U^2 \leq \frac{2\hat{J}(u_0)}{\sigma} \right\}.$$

Now the claim follows from Theorem 4.2.  $\square$

### 4.3. First-order necessary optimality conditions

In (4.4) we have introduced the quadratic programming problem

$$\min_{u \in U_{\text{ad}}} \mathcal{J}(u) = \frac{1}{2} \|\mathcal{G}u - z_d\|_{\mathcal{H}}^2 + \frac{\sigma}{2} \|u\|_U^2. \quad (4.6)$$

Existence of a unique solution has been investigated in Section 4.2. In this section we characterize the solution to (4.6) by first-order optimality conditions, which are essential to prove convergence and rate of convergence results for the POD approximations in Section 4.4. To derive first-order conditions we require the notion of derivatives in function spaces. Therefore, we recall the following definition [72, pp. 56-57].

**Definition 4.6.** *Suppose that  $\mathcal{B}_1$  and  $\mathcal{B}_2$  are real Banach spaces,  $\mathcal{U} \subset \mathcal{B}_1$  be an open subset and  $\mathcal{F} : \mathcal{U} \subset \mathcal{B}_1 \rightarrow \mathcal{B}_2$  a given mapping. The directional derivative of  $\mathcal{F}$  at a point  $u \in \mathcal{U}$  in the direction  $h \in \mathcal{B}_2$  is defined by*

$$D\mathcal{F}(u; h) := \lim_{\varepsilon \searrow 0} \frac{1}{\varepsilon} (\mathcal{F}(u + \varepsilon h) - \mathcal{F}(u))$$

provided the limit exists in  $\mathcal{B}_2$ . Suppose that the directional derivative exists for all  $h \in \mathcal{B}_1$  and there is a linear, continuous operator  $\mathcal{T} : \mathcal{U} \rightarrow \mathcal{B}_2$  satisfying

$$D\mathcal{F}(u; h) = \mathcal{T}h \quad \text{for all } h \in \mathcal{U}.$$

Then,  $\mathcal{F}$  is said to be Gâteaux-differentiable at  $u$  and  $\mathcal{T}$  is the Gâteaux derivative of  $\mathcal{F}$  at  $u$ . We write  $\mathcal{T} = \mathcal{F}'(u)$ .

**Remark 4.7.** Let  $\mathcal{H}$  be a real Hilbert space and  $\mathcal{F} : \mathcal{H} \rightarrow \mathbb{R}$  be Gâteaux-differentiable at  $u \in \mathcal{H}$ . Then, its Gâteaux derivative  $\mathcal{F}'(u)$  at  $u$  belongs to  $\mathcal{H}' = \mathcal{L}(\mathcal{H}, \mathbb{R})$ . Due to Riesz theorem there exists a unique element  $\nabla \mathcal{F}(u) \in \mathcal{H}$  satisfying

$$\langle \nabla \mathcal{F}(u), v \rangle_{\mathcal{H}} = \langle \mathcal{F}'(u), v \rangle_{\mathcal{H}', \mathcal{H}} \quad \text{for all } v \in \mathcal{H}.$$

We call  $\nabla \mathcal{F}(u)$  the (Gâteaux) gradient of  $\mathcal{F}$  at  $u$ . ◇

**Theorem 4.8.** Let  $\mathcal{U}$  be a real Hilbert space and  $\mathcal{U}_{\text{ad}}$  be convex subset. Suppose that  $\bar{u} \in \mathcal{U}_{\text{ad}}$  is a solution to (4.6)

$$\min_{u \in \mathcal{U}_{\text{ad}}} \mathcal{J}(u).$$

Then the following variational inequality holds

$$\langle \nabla \mathcal{J}(\bar{u}), u - \bar{u} \rangle_{\mathcal{U}} \geq 0 \quad \text{for all } u \in \mathcal{U}_{\text{ad}}, \quad (4.7)$$

where the gradient of  $\mathcal{J}$  is given by

$$\nabla \mathcal{J}(\bar{u}) = \mathcal{G}^*(\mathcal{G}\bar{u} - z_d) + \sigma \bar{u} \quad \text{for } u \in \mathcal{U}.$$

If  $\bar{u} \in \mathcal{U}_{\text{ad}}$  solves (4.7), then  $\bar{u}$  is a solution to (4.6).

*Proof.* Since  $\mathcal{J}$  is Gâteaux-differentiable and convex in  $\mathcal{U}$ , the result follows directly from [72, pp. 63-63]. □

Inequality (4.7) is a first-order necessary and sufficient condition for (4.6), which can be expressed as

$$\langle \mathcal{G}\bar{u} - z_d, \mathcal{G}\bar{u} - \mathcal{G}\bar{u} \rangle_{\mathcal{H}} + \langle \sigma \bar{u}, u - \bar{u} \rangle_{\mathcal{U}} \geq 0 \quad \text{for all } u \in \mathcal{U}_{\text{ad}}. \quad (4.8)$$

Next we study (4.8) for  $(\hat{\mathbf{P}})$ . Utilizing the setting from (4.5) we obtain

$$\begin{aligned} & \langle \mathcal{G}\bar{u} - z_d, \mathcal{G}\bar{v} \rangle_{\mathcal{H}} \\ &= \sigma_Q \langle \mathcal{S}\bar{u} - (y_Q - \hat{y}), \mathcal{S}(u - \bar{u}) \rangle_{L^2(0, T; H)} \\ & \quad + \sigma_\Omega \langle (\mathcal{S}\bar{u})(T) - (y_\Omega - \hat{y}(T)), (\mathcal{S}(u - \bar{u}))(T) \rangle_H \\ &= \sigma_Q \langle \mathcal{S}\bar{u}, \mathcal{S}(u - \bar{u}) \rangle_{L^2(0, T; H)} + \sigma_\Omega \langle (\mathcal{S}\bar{u})(T), (\mathcal{S}(u - \bar{u}))(T) \rangle_H \\ & \quad - \sigma_Q \langle y_Q - \hat{y}, \mathcal{S}(u - \bar{u}) \rangle_{L^2(0, T; H)} - \sigma_\Omega \langle y_\Omega - \hat{y}(T), (\mathcal{S}(u - \bar{u}))(T) \rangle_H. \end{aligned}$$

Let us define the two linear, bounded operators  $\Theta : W_0(0, T) \rightarrow W_0(0, T)'$  and  $\Xi : L^2(0, T; H) \times H \rightarrow W_0(0, T)'$  by

$$\begin{aligned} \langle \Theta \varphi, \phi \rangle_{W_0(0, T)', W_0(0, T)} &= \int_0^T \langle \sigma_Q \varphi(t), \phi(t) \rangle_H dt + \langle \sigma_\Omega \varphi(T), \phi(T) \rangle_H, \\ \langle \Xi z, \phi \rangle_{W_0(0, T)', W_0(0, T)} &= \int_0^T \langle \sigma_Q z_Q(t), \phi(t) \rangle_H dt + \langle \sigma_\Omega z_\Omega, \phi(T) \rangle_H \end{aligned} \quad (4.9)$$

for  $\varphi, \phi \in W_0(0, T)$  and  $z = (z_Q, z_\Omega) \in L^2(0, T; H) \times H$ . Then, we find

$$\begin{aligned} & \langle \mathcal{G}\bar{u} - z_d, \mathcal{G}\bar{v} \rangle_{\mathcal{J}\mathcal{L}} \\ &= \langle \Theta(\mathcal{S}\bar{u}) - \Xi(y_Q - \hat{y}, y_\Omega - \hat{y}(T)), \mathcal{S}(u - \bar{u}) \rangle_{W_0(0, T)', W_0(0, T)} \\ &= \langle \mathcal{S}'\Theta\mathcal{S}\bar{u}, u - \bar{u} \rangle_U - \langle \mathcal{S}'\Xi(y_Q - \hat{y}, y_\Omega - \hat{y}(T)), u - \bar{u} \rangle_U. \end{aligned} \quad (4.10)$$

Let us define the linear  $\mathcal{A} : U \rightarrow W(0, T)$  as follows: for given  $u \in U$  the function  $p = \mathcal{A}u \in W(0, T)$  is the unique solution to

$$\begin{aligned} -\frac{d}{dt} \langle p(t), \varphi \rangle_H + a(t; p(t), \varphi) &= -\sigma_Q \langle (\mathcal{S}u)(t), \varphi \rangle_H \quad \forall \varphi \in V \text{ a.e.}, \\ p(T) &= -\sigma_\Omega (\mathcal{S}u)(T) \quad \text{in } H. \end{aligned} \quad (4.11)$$

It follows from (3.1) and  $\mathcal{S}u \in W(0, T)$  that the operator  $\mathcal{A}$  is well-defined and bounded.

**Lemma 4.9.** *Let Assumption 1 be satisfied and  $u, v \in U$ . We set  $y = \mathcal{S}u \in W_0(0, T)$ ,  $w = \mathcal{S}v \in W_0(0, T)$ , and  $p = \mathcal{A}v \in W(0, T)$ . Then,*

$$\int_0^T \langle (\mathcal{B}u)(t), p(t) \rangle_{V', V} dt = - \int_0^T \sigma_Q \langle w(t), y(t) \rangle_H dt - \sigma_\Omega \langle w(T), y(T) \rangle_H.$$

*Proof.* We derive from  $y = \mathcal{S}u$ ,  $p = \mathcal{A}u$ ,  $y \in W_0(0, T)$  and integration by parts

$$\begin{aligned} & \int_0^T \langle (\mathcal{B}u)(t), p(t) \rangle_{V', V} dt = \int_0^T \langle y_t(t), p(t) \rangle_{V', V} + a(t; y(t), p(t)) dt \\ &= \int_0^T -\langle p_t(t), y(t) \rangle_{V', V} + a(t; p(t), y(t)) dt + \langle p(T), y(T) \rangle_H \\ &= - \int_0^T \sigma_Q \langle w(t), y(t) \rangle_H dt - \sigma_\Omega \langle w(T), y(T) \rangle_H \end{aligned}$$

which is the claim.  $\square$

We define  $\hat{p} \in W(0, T)$  as the unique solution to

$$\begin{aligned} -\frac{d}{dt} \langle \hat{p}(t), \varphi \rangle_H + a(t; \hat{p}(t), \varphi) &= \sigma_Q \langle y_Q(t) - \hat{y}(t), \varphi \rangle_H \quad \forall \varphi \in V \text{ a.e.}, \\ \hat{p}(T) &= \sigma_\Omega (y_\Omega - \hat{y}(T)) \quad \text{in } H. \end{aligned} \quad (4.12)$$

Then, for every  $u \in U$  the function  $p = \hat{p} + \mathcal{A}u$  is the unique solution to

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

with  $y = \hat{y} + \mathcal{S}u$ . Moreover, we have the following result.

**Lemma 4.10.** *Let Assumption 1 be satisfied. Then,  $\mathcal{B}'\mathcal{A} = -\mathcal{S}'\Theta\mathcal{S} \in \mathcal{L}(U)$ , where linear and bounded operator  $\Theta$  has been defined in (4.9). Moreover,  $\mathcal{B}'\hat{p} = \mathcal{S}'\Xi(y_Q - \hat{y}, y_\Omega - \hat{y}(T))$ , where  $\hat{p}$  is the solution to (4.12).*

*Proof.* Let  $u, v \in \mathcal{U}$  be chosen arbitrarily. We set  $y = \mathcal{S}u \in W_0(0, T)$  and  $w = \mathcal{S}v \in W_0(0, T)$ . Recall that we identify  $U$  with its dual space  $U'$ . From the integration by parts formula and Lemma 4.9 we infer that

$$\begin{aligned} \langle \mathcal{S}'\Theta\mathcal{S}v, u \rangle_U &= \langle \Theta\mathcal{S}v, \mathcal{S}u \rangle_{W_0(0, T)', W_0(0, T)} = \langle \Theta w, y \rangle_{W_0(0, T)', W_0(0, T)} \\ &= \int_0^T \sigma_Q \langle w(t), y(t) \rangle_H dt + \sigma_\Omega \langle w(T), y(T) \rangle_H \\ &= -\langle \mathcal{B}u, p \rangle_{L^2(0, T; V'), L^2(0, T; V)} = -\langle u, \mathcal{B}'p \rangle_U = -\langle \mathcal{B}'\mathcal{A}v, u \rangle_U. \end{aligned}$$

Since  $u, v \in U$  are chosen arbitrarily, we have  $\mathcal{B}'\mathcal{A} = \mathcal{S}'\Theta\mathcal{S}$ . Further, we find

$$\begin{aligned} \langle \mathcal{S}'\Xi(y_Q - \hat{y}, y_\Omega - \hat{y}(T)), u \rangle_U &= \langle \Xi(y_Q - \hat{y}), y_\Omega - \hat{y}(T) \rangle_{W_0(0, T)', W_0(0, T)} \\ &= \int_0^T \sigma_Q \langle y_Q - \hat{y}(t), y(t) \rangle_H dt + \sigma_\Omega \langle y_\Omega - \hat{y}(T), y(T) \rangle_H \\ &= \int_0^T -\langle \hat{p}_t(t), y(t) \rangle_H + a(t; \hat{p}(t), y(t)) dt + \langle \hat{p}(T), y(T) \rangle_H \\ &= \int_0^T \langle y_t(t), \hat{p}(t) \rangle_H + a(t; y(t), \hat{p}(t)) dt = \int_0^T \langle (\mathcal{B}u)(t), \hat{p}(t) \rangle_{V', V} dt \\ &= \langle \mathcal{B}'\hat{p}, u \rangle_U. \end{aligned}$$

which gives the claim.  $\square$

We infer from (4.10) and Lemma 4.10 that

$$\langle \mathcal{G}\bar{u} - z_d, \mathcal{G}\bar{v} \rangle_{\mathcal{H}} = -\langle \mathcal{B}'(\hat{p} + \mathcal{A}\bar{u}), u - \bar{u} \rangle_U.$$

This implies the following variational inequality for  $(\hat{\mathbf{P}})$

$$\begin{aligned} \langle \mathcal{G}\bar{u} - z_d, \mathcal{G}u - \mathcal{G}\bar{u} \rangle_{\mathcal{H}} + \sigma \langle \bar{u}, u - \bar{u} \rangle_{\mathcal{U}} \\ = \langle \sigma\bar{u} - \mathcal{B}'(\hat{p} + \mathcal{A}\bar{u}), u - \bar{u} \rangle_U \geq 0 \quad \text{for all } u \in U_{\text{ad}}. \end{aligned}$$

Summarizing we have proved the following result.

**Theorem 4.11.** *Suppose that Assumptions 1 and 2 hold. Then,  $(\bar{y}, \bar{u})$  is a solution to  $(\mathbf{P})$  if and only if  $(\bar{y}, \bar{u})$  satisfy together with the adjoint variable  $\bar{p}$  the first-order optimality system*

$$\bar{y} = \hat{y} + \mathcal{S}\bar{u}, \quad \bar{p} = \hat{p} + \mathcal{A}\bar{u}, \quad u_a \leq \bar{u} \leq u_b \quad (4.13a)$$

$$\langle \sigma\bar{u} - \mathcal{B}'\bar{p}, u - \bar{u} \rangle_U \geq 0 \quad \text{for all } u \in U_{\text{ad}}. \quad (4.13b)$$

**Remark 4.12.** By using a Lagrangian framework it follows from Theorem 4.11 and [72] that the variational inequality (4.13b) is equivalent to the existence of two functions  $\bar{\mu}_a, \bar{\mu}_b \in U$  satisfying  $\bar{\mu}_a, \bar{\mu}_b \geq 0$ ,

$$\sigma\bar{u} - \mathcal{B}'\bar{p} + \bar{\mu}_b - \bar{\mu}_a = 0$$

and the complementarity condition

$$\bar{\mu}_a(t)^\top (u_a(t) - \bar{u}(t)) = \bar{\mu}_b(t)^\top (\bar{u}(t) - u_b(t)) = 0 \quad \text{f.a.a. } t \in [0, T].$$

Thus, (4.13) is equivalent to the system

$$\begin{aligned} \bar{y} &= \hat{y} + \mathcal{S}\bar{u}, & \bar{p} &= \hat{p} + \mathcal{A}\bar{u}, & \sigma\bar{u} - \mathcal{B}'\bar{p} + \bar{\mu}_b - \bar{\mu}_a &= 0, \\ u_a &\leq \bar{u} \leq u_b, & 0 &\leq \bar{\mu}_a, & 0 &\leq \bar{\mu}_b, \\ \bar{\mu}_a(t)^\top (u_a(t) - \bar{u}(t)) &= \bar{\mu}_b(t)^\top (\bar{u}(t) - u_b(t)) &= 0 &\text{ a.e. in } [0, T]. \end{aligned} \quad (4.14)$$

Utilizing a complementarity function it can be shown that (4.14) is equivalent with

$$\begin{aligned} \bar{y} &= \hat{y} + \mathcal{S}\bar{u}, & \bar{p} &= \hat{p} + \mathcal{A}\bar{u}, & \sigma\bar{u} - \mathcal{B}'\bar{p} + \bar{\mu}_b - \bar{\mu}_a &= 0, & u_a &\leq \bar{u} \leq u_b, \\ \bar{\mu}_a &= \max(0, \bar{\mu}_a + \eta(\bar{u} - u_a)), & \bar{\mu}_b &= \max(0, \bar{\mu}_b + \eta(\bar{u} - u_b)), \end{aligned} \quad (4.15)$$

where  $\eta > 0$  is an arbitrary real number. The max-and min-operations are interpreted componentwise in the pointwise everywhere sense.  $\diamond$

The gradient  $\nabla \hat{J} : U \rightarrow U$  of the reduced cost functional  $\hat{J}$  is given by

$$\nabla J(u) = \sigma u - \mathcal{B}^* p, \quad u \in U,$$

where  $p = \hat{p} + \mathcal{A}u$  holds true; see, e.g., [26]. Thus, a first-order sufficient optimality condition for  $(\hat{\mathbf{P}})$  is given by the variational inequality

$$\langle \sigma\bar{u} - \mathcal{B}'\bar{p}, u - \bar{u} \rangle_U \geq 0 \quad \text{for all } u \in U_{\text{ad}}, \quad (4.16)$$

with  $\bar{p} = \hat{p} + \mathcal{A}\bar{u}$ .

Problem  $(\hat{\mathbf{P}})$  can be solved numerically by a primal-dual active set strategy with the choice  $\eta = \sigma$ . In this case the method is equivalent to a locally superlinearly convergent semi-smooth Newton algorithm applied to (4.15); see [24, 26, 74]. In Algorithm 4.1 we formulate the method in the context of our application. In Section 5 we compare Algorithm 4.1 with the Banach fixed point iteration as well as with the projected gradient method [38, 54].

---

**Algorithm 4.1** (Primal-dual active set strategy)

---

**Require:** Starting value  $(u^0, \lambda^0)$  and maximal iteration number  $k_{\text{max}}$ .

1: Set  $k = 0$ . For  $i = 1, \dots, m$  determine the active sets

$$\begin{aligned} \mathcal{A}_{ai}^k &= \{t \in [0, T] \mid \sigma u_i^k + \lambda_i^k < u_{ai} \text{ a.e.}\}, \\ \mathcal{A}_{bi}^k &= \{t \in [0, T] \mid \sigma u_i^k + \lambda_i^k > u_{bi} \text{ a.e.}\} \end{aligned}$$

and the inactive set  $\mathcal{J}_i^k = [0, T] \setminus \mathcal{A}_i^k$  with  $\mathcal{A}_i^k = \mathcal{A}_{ai}^k \cup \mathcal{A}_{bi}^k$ .

2: **repeat**

3: Compute the solution  $(y, p, u)$  to the optimality system

$$y = \hat{y} + \mathcal{S}u, \quad p = \hat{p} + \mathcal{A}u, \quad u_i = \begin{cases} u_{ai} & \text{on } \mathcal{A}_{ia}^k, \\ u_{bi} & \text{on } \mathcal{A}_{ib}^k, \\ (\mathcal{B}'p)_i / \sigma & \text{on } \mathcal{J}_i^k, \end{cases} \quad (1 \leq i \leq m)$$

4: Set  $(y^{k+1}, u^{k+1}, p^{k+1}) = (y, u, p)$ ,  $\lambda^{k+1} = \mathcal{B}'p^{k+1} - \sigma u^{k+1}$  and  $k = k+1$ .

5: Compute the active and inactive sets according to step 1.

6: **until**  $(\mathcal{A}_{ai}^k = \mathcal{A}_{ai}^{k-1} \text{ and } \mathcal{A}_{bi}^k = \mathcal{A}_{bi}^{k-1})$  **or**  $k = k_{\text{max}}$ .

---



#### 4.4. The POD Galerkin approximation for $(\hat{\mathbf{P}})$

In this subsection we introduce the POD Galerkin schemes for the variational inequality (4.16) using a POD Galerkin approximation for the state and dual variables. Moreover, we study the convergence of the POD discretizations, where we make use of the analysis in [29, 40, 41, 42, 68, 73]. For a general introduction we also refer the reader to the survey paper [28].

In Section 3.3 we have introduced a POD Galerkin scheme for the state equation (4.1). Suppose that  $\{\psi_i\}_{i=1}^\ell$  be a POD basis of rank  $\ell$  computed from  $(\mathbf{P}^\ell)$  with  $\psi_i = \psi_i^V$  in case of  $X = V$  and  $\psi_i = \psi_i^H$  in case of  $X = H$ . We set  $X^\ell = \text{span}\{\psi_1, \dots, \psi_\ell\} \subset V$ . Let the linear and bounded projection operator  $\mathcal{P}^\ell$  denote  $\mathcal{P}_V^\ell$  for  $X = V$  and  $\mathcal{P}_H^\ell$  for  $X = H$ ; see (3.8).

Recall the POD Galerkin ansatz (3.14) for the state variable. Analogously, we approximate the adjoint variable  $p = \hat{p} + \mathcal{A}u$  by the Galerkin expansion

$$p^\ell(t) = \hat{p}(t) + \sum_{i=1}^{\ell} p_i^\ell(t) \psi_i \in V \quad \text{for } t \in [0, T] \quad (4.17)$$

with coefficient functions  $p_i^\ell : [0, T] \rightarrow \mathbb{R}$  and with  $\hat{p}$  from (4.12). Let the vector-valued coefficient function given by

$$p^\ell = (p_1^\ell, \dots, p_\ell^\ell) : [0, T] \rightarrow \mathbb{R}^\ell$$

If we assume that  $p^\ell(T) = -\sigma_\Omega y^\ell(T)$  holds, then we infer from  $\hat{p}(T) = \sigma_\Omega(y_\Omega - \hat{y}(T))$  and (4.17) that

$$p^\ell(T) = \hat{p}(T) - \sigma_\Omega \sum_{i=1}^{\ell} y_i^\ell(T) \psi_i = \sigma_\Omega(y_\Omega - y^\ell(T)).$$

This motivates the following POD scheme for the approximation of  $p = \hat{p} + \mathcal{A}u$  is given as follows:  $p^\ell \in W(0, T)$  satisfies

$$\begin{aligned} -\frac{d}{dt} \langle p^\ell(t), \psi \rangle_H + a(t; p^\ell(t), \psi) &= \sigma_Q \langle (y_Q - y^\ell)(t), \psi \rangle_H \quad \forall \psi \in X^\ell \text{ a.e.}, \\ p^\ell(T) &= -\sigma_\Omega y^\ell(T). \end{aligned} \quad (4.18)$$

It follows by similar arguments as for (3.15) that there is a unique solution  $p^\ell \in W(0, T)$ .

**Remark 4.13.** Recall that we have introduced the linear and bounded solution operator  $\mathcal{S}^\ell : U \rightarrow W(0, T)$  as an approximation for the state solution operator  $\mathcal{S}$ ; see Remark 3.10-2). Analogously, we define an approximation of the adjoint solution operator  $\mathcal{A}$  as follows: Let  $\mathcal{A}^\ell : U \rightarrow W(0, T)$  denote the solution operator to

$$\begin{aligned} -\frac{d}{dt} \langle w^\ell(t), \psi \rangle_H + a(t; w^\ell(t), \psi) &= -\sigma_1 \langle (\mathcal{S}^\ell u)(t), \psi \rangle_H \quad \forall \psi \in X^\ell \text{ a.e.}, \\ w^\ell(T) &= -\sigma_2 (\mathcal{S}^\ell u)(T). \end{aligned}$$

Then  $p^\ell = \hat{p} + \mathcal{A}^\ell u$  is the unique solution to (4.18).  $\diamond$

**Lemma 4.14.** *Let Assumption 1 on page 33 be satisfied and  $u, v \in U$ . We set  $y^\ell = \mathcal{S}^\ell u \in W_0(0, T)$ ,  $w^\ell = \mathcal{S}^\ell v \in W_0(0, T)$ , and  $p^\ell = \mathcal{A}^\ell v \in W(0, T)$ . Then,*

$$\int_0^T \langle (\mathcal{B}u)(t), p^\ell(t) \rangle_{V', V} dt = - \int_0^T \sigma_Q \langle w^\ell(t), y^\ell(t) \rangle_H dt - \sigma_\Omega \langle w^\ell(T), y^\ell(T) \rangle_H.$$

Moreover,  $\mathcal{B}'\mathcal{A}^\ell = -(\mathcal{S}^\ell)' \Theta \mathcal{S}^\ell \in \mathcal{L}(U)$ , where linear and bounded operator  $\Theta$  has been defined in (4.9).

*Proof.* Since the POD basis for the state and adjoint coincide, the claim follows by the same arguments used to prove Lemmas 4.9 and 4.10.  $\square$

**Theorem 4.15.** *Suppose that Assumptions 1 and 2 hold. Let  $X = V$  and  $u \in U$  be arbitrarily given so that  $\mathcal{S}u, \mathcal{A}u \in H^1(0, T; V) \setminus \{0\}$ .*

- 1) *To compute a POD basis  $\{\psi_i\}_{i=1}^\ell$  of rank  $\ell$  we choose  $\wp = 4$ ,  $y^1 = \mathcal{S}u$ ,  $y^2 = (\mathcal{S}u)_t$ ,  $y^3 = \mathcal{A}u$  and  $y^4 = (\mathcal{A}u)_t$ . Then,  $p = \hat{p} + \mathcal{A}u$  and  $p^\ell = \hat{p} + \mathcal{A}^\ell u$  satisfies thea-priori error estimate*

$$\|p^\ell - p\|_{W(0, T)}^2 \leq \begin{cases} C \sum_{i=\ell+1}^{d_V} \lambda_i^V & \text{if } X = V, \\ C \sum_{i=\ell+1}^{d_H} \lambda_i^H \|\psi_i^H - \mathcal{P}_H^\ell \psi_i^H\|_V^2 & \text{if } X = H \end{cases} \quad (4.19)$$

for a constant  $C$  which depends on  $\gamma, \gamma_1, \gamma_2, T, \sigma_\Omega$  and  $\sigma_Q$ .

- 2) *If  $\mathcal{S}\tilde{u}$  and  $\mathcal{A}\tilde{u}$  belong to  $H^1(0, T; V)$  for every  $\tilde{u} \in U$  and if  $\lambda_i^H > 0$  for all  $i \in \mathcal{J}$ , then we have*

$$\lim_{\ell \rightarrow \infty} \|\mathcal{A} - \mathcal{A}^\ell\|_{\mathcal{L}(U, W(0, T))} = 0. \quad (4.20)$$

*Proof.* Analogous to (3.22) we have  $p^\ell(t) - p(t) = \theta^\ell(t) + \rho^\ell(t)$  for almost all  $t \in [0, T]$  with  $\theta^\ell(t) = (\mathcal{A}^\ell u)(t) - \mathcal{P}^\ell(\mathcal{A}u)(t)$  and  $\rho^\ell(t) = \mathcal{P}^\ell(\mathcal{A}u)(t) - (\mathcal{A}u)(t)$ . Here,  $\mathcal{P}^\ell = \mathcal{P}_V^\ell$  for  $X = V$  and  $\mathcal{P}^\ell = \mathcal{P}_H^\ell$  for  $X = H$ . Now, the proof of the claims follows by similar arguments as the proofs of Theorem 3.11, Proposition 4.7 in [29], Proposition 4.6 in [73] and Theorem 7.3 in [68]. To estimate the terminal term  $\theta^\ell(T)$  we use observe that

$$\begin{aligned} \|\theta^\ell(T)\|_H &= \|\mathcal{P}^\ell((\mathcal{A}u)(T)) - (\mathcal{A}^\ell u)(T)\|_H \\ &\leq \sigma_\Omega \left( \|\mathcal{P}^\ell((\mathcal{S}u)(T)) - (\mathcal{S}u)(T)\|_H + \|(\mathcal{S}u)(T) - (\mathcal{S}^\ell u)(T)\|_H \right) \\ &\leq \sigma_\Omega \left( \|\mathcal{P}^\ell((\mathcal{S}u)(T)) - (\mathcal{S}u)(T)\|_H + \|y(T) - y^\ell(T)\|_H \right) \\ &\leq \sigma_\Omega \left( \|\mathcal{P}^\ell(\mathcal{S}u) - (\mathcal{S}u)\|_{C([0, T]; H)} + \|y - y^\ell\|_{C([0, T]; H)} \right) \\ &\leq \sigma_\Omega C_E \left( \|\mathcal{P}^\ell(\mathcal{S}u) - (\mathcal{S}u)\|_{H^1(0, T; V)} + \|y - y^\ell\|_{W(0, T)} \right) \end{aligned}$$

with an embedding constant  $C_E$ . The first term on the right-hand can be handled by (2.27), the second term is estimated in Theorem 3.11. Finally, (4.20) follows from (3.21) and the fact that the operator  $\mathcal{S}^\ell$  is bounded uniformly with respect to  $\ell$ .  $\square$

**Remark 4.16.** 1) The inclusion of adjoint information into the snapshot ensemble improves the approximation quality also for nonlinear problems; see [15].

2) Analogous to Remark 3.12-2) the a-priori estimate (4.19) holds for an arbitrarily chosen, but fixed control  $u \in U$ . Furthermore, (4.20) implies that

$$\lim_{\ell \rightarrow \infty} \|\hat{p} + \mathcal{A}^\ell \tilde{u} - \hat{p} - \mathcal{A}\tilde{u}\|_{W(0,T)} = 0$$

for any  $\tilde{u} \in U$ .

3) We can also extend the results in Proposition 3.14 for the adjoint equation and get an a-priori error estimate choosing  $\wp = 2$ ,  $y^1 = \mathcal{S}u$  and  $y^2 = \mathcal{A}u$ .  $\diamond$

The POD Galerkin approximation for  $(\hat{\mathbf{P}})$  is as follows:

$$\min \hat{J}^\ell(u) \quad \text{s.t.} \quad u \in U_{\text{ad}}, \quad (\hat{\mathbf{P}}^\ell)$$

where the cost is defined by  $\hat{J}^\ell(u) = J(\hat{y} + \mathcal{S}^\ell u, u)$  for  $u \in U$ . Let  $\bar{u}^\ell$  be the solution to  $(\hat{\mathbf{P}}^\ell)$ . Then, a first-order sufficient optimality condition is given by the variational inequality

$$\langle \sigma \bar{u}^\ell - \mathcal{B}' \bar{p}^\ell, u - \bar{u}^\ell \rangle_U \geq 0 \quad \text{for all } u \in U_{\text{ad}}, \quad (4.21)$$

where  $\bar{p}^\ell = \hat{p}^\ell + \mathcal{A}^\ell \bar{u}^\ell$  holds.

**Theorem 4.17.** *Suppose that Assumptions 1 and 2 hold. Let  $u \in U$  be arbitrarily given so that  $\mathcal{S}u, \mathcal{A}u \in H^1(0, T; V) \setminus \{0\}$ .*

1) *To compute a POD basis  $\{\psi_i\}_{i=1}^\ell$  of rank  $\ell$  we choose  $\wp = 4$ ,  $y^1 = \mathcal{S}u$ ,  $y^2 = (\mathcal{S}u)_t$ ,  $y^3 = \mathcal{A}u$  and  $y^4 = (\mathcal{A}u)_t$ . Then, the optimal solution  $\bar{u}$  to  $(\hat{\mathbf{P}})$  and the associated POD suboptimal solution  $\bar{u}^\ell$  to  $(\hat{\mathbf{P}}^\ell)$  satisfy*

$$\lim_{\ell \rightarrow \infty} \|\bar{u}^\ell - \bar{u}\|_U = 0 \quad (4.22)$$

for  $X = V$  and  $X = H$ .

2) *If an optimal POD basis of rank  $\ell$  is computed by choosing  $\wp = 4$ ,  $y^1 = \mathcal{S}\bar{u}$ ,  $y^2 = (\mathcal{S}\bar{u})_t$ ,  $y^3 = \mathcal{A}\bar{u}$  and  $y^4 = (\mathcal{A}\bar{u})_t$ , then we have*

$$\|\bar{u}^\ell - \bar{u}\|_U \leq \begin{cases} \frac{C}{\sigma} \sum_{i=\ell+1}^{d_V} \lambda_i^V & \text{if } X = V, \\ \frac{C}{\sigma} \sum_{i=\ell+1}^{d_H} \lambda_i^H \|\psi_i^H - \mathcal{P}_H^\ell \psi_i^H\|_V^2 & \text{if } X = H, \end{cases} \quad (4.23)$$

where the constant  $C$  which depends on  $\gamma$ ,  $\gamma_1$ ,  $\gamma_2$ ,  $T$ ,  $\sigma_\Omega$ ,  $\sigma_Q$  and the norm  $\|\mathcal{B}'\|_{\mathcal{L}(L^2(0,T;V),U)}$ .

*Proof.* Choosing  $u = \bar{u}^\ell$  in (4.16) and  $u = \bar{u}$  in (4.21) we get the variational inequality

$$0 \leq \langle \sigma(\bar{u} - \bar{u}^\ell) - \mathcal{B}'(\bar{p} - \bar{p}^\ell), \bar{u}^\ell - \bar{u} \rangle_U. \quad (4.24)$$

Utilizing Lemma 4.14 and  $\langle \Theta\varphi, \varphi \rangle_{W_0(0,T)', W_0(0,T)} \geq 0$  for all  $\varphi \in W_0(0,T)$  we infer from (4.24) that

$$\begin{aligned} 0 &\leq \langle \mathcal{B}'\mathcal{A}^\ell \bar{u}^\ell - \mathcal{B}'\mathcal{A}\bar{u}, \bar{u}^\ell - \bar{u} \rangle_U - \sigma \|\bar{u} - \bar{u}^\ell\|_U^2 \\ &= \langle \mathcal{B}'\mathcal{A}^\ell(\bar{u}^\ell - \bar{u}) + \mathcal{B}'(\mathcal{A}^\ell - \mathcal{A})\bar{u}, \bar{u}^\ell - \bar{u} \rangle_U - \sigma \|\bar{u} - \bar{u}^\ell\|_U^2 \\ &\leq \langle \Theta\mathcal{S}^\ell(\bar{u} - \bar{u}^\ell), \mathcal{S}^\ell(\bar{u}^\ell - \bar{u}) \rangle_U + \|\mathcal{B}'(\mathcal{A}^\ell - \mathcal{A})\bar{u}\|_U \|\bar{u}^\ell - \bar{u}\|_U - \sigma \|\bar{u} - \bar{u}^\ell\|_U^2 \\ &\leq \|\mathcal{B}'(\mathcal{A}^\ell - \mathcal{A})\bar{u}\|_U \|\bar{u}^\ell - \bar{u}\|_U - \sigma \|\bar{u} - \bar{u}^\ell\|_U^2. \end{aligned}$$

Consequently,

$$\|\bar{u} - \bar{u}^\ell\|_U \leq \frac{1}{\sigma} \|\mathcal{B}'(\mathcal{A}^\ell - \mathcal{A})\bar{u}\|_U.$$

Now (4.22) and (4.23) follow from (4.20) and (4.19), respectively.  $\square$

**Remark 4.18.** It follows from Proposition 3.14 and Remark 4.16-3) that (4.22) holds true provided we choose  $\wp = 2$ ,  $y^1 = \mathcal{S}u$ ,  $y^2 = \mathcal{A}u$  or even  $\wp = 1$  and  $y^1 = \mathcal{S}u$ .  $\diamond$

In Algorithm 4.2 we formulate a discrete version of the primal-dual active set method (see Algorithm 4.1) which is utilized to solve  $(\hat{\mathbf{P}}^\ell)$  in Section 5.

---

**Algorithm 4.2** (POD discretized primal-dual active set strategy)

---

**Require:** POD basis  $\{\psi_i\}_{i=1}^\ell$ , starting value  $(u^{\ell 0}, \lambda^{\ell 0})$  and maximal iteration number  $k_{\max}$ .

- 1: Set  $k = 0$ , determine the active sets

$$\mathcal{A}_{ai}^{\ell k} = \{t \in [0, T] \mid \sigma u_i^{k\ell} + \lambda_i^{k\ell} < u_{ai} \text{ a.e.}\},$$

$$\mathcal{A}_{bi}^{\ell k} = \{t \in [0, T] \mid \sigma u_i^{k\ell} + \lambda_i^{k\ell}(t) > u_{bi}(t)\}$$

and the inactive sets  $\mathcal{J}_i^{\ell k} = [0, T] \setminus \mathcal{A}_i^{\ell k}$  with  $\mathcal{A}_i^{\ell k} = \mathcal{A}_{ai}^{\ell k} \cup \mathcal{A}_{bi}^{\ell k}$ .

- 2: **repeat**

- 3: Determine the solution  $(y^\ell, u^\ell, p^\ell)$  to the optimality system

$$y^\ell = \hat{y} + \mathcal{S}^\ell u^\ell, \quad p^\ell = \hat{p} + \mathcal{A}^\ell u^\ell, \quad u^\ell = \begin{cases} u_a & \text{on } \mathcal{A}_a^{k\ell}, \\ u_b & \text{on } \mathcal{A}_b^{k\ell}, \\ \mathcal{B}'p^\ell/\sigma & \text{on } \mathcal{J}^{k\ell}. \end{cases}$$

- 4: Set  $(y^{\ell, k+1}, u^{\ell, k+1}, p^{\ell, k+1}) = (y^\ell, u^\ell, p^\ell)$ ,  $\lambda^{\ell, k+1} = \mathcal{B}'p^{\ell, k+1} - \sigma u^{\ell, k+1}$  and  $k = k + 1$ .

- 5: Compute the active and inactive sets according to step 1.

- 6: **until**  $(\mathcal{A}_a^{\ell k} = \mathcal{A}_a^{\ell, k-1} \text{ and } \mathcal{A}_b^{\ell k} = \mathcal{A}_b^{\ell, k-1})$  **or**  $k = k_{\max}$ .
- 

#### 4.5. POD a-posteriori error analysis

In [73] a POD a-posteriori error estimates are presented which can be applied to our optimal control problem as well. Based on a perturbation method [16] it is deduced how far the suboptimal control  $\bar{u}^\ell$  is from the (unknown) exact

optimal control  $\bar{u}$ . Thus, our goal is to estimate the norm  $\|\bar{u} - \bar{u}^\ell\|_U$  without the knowledge of the optimal solution  $\bar{u}$ . In general,  $\bar{u}^\ell \neq \bar{u}$  holds, so that  $\bar{u}^\ell$  does not satisfy the variational inequality (4.16). However, there exists a function  $\zeta^\ell \in U$  such that

$$\langle \sigma \bar{u}^\ell - \mathcal{B}'\tilde{p}^\ell + \zeta^\ell, u - \bar{u}^\ell \rangle_U \geq 0 \quad \forall u \in U_{\text{ad}}, \quad (4.25)$$

with  $\tilde{p}^\ell = \hat{p} + \mathcal{A}\bar{u}^\ell$ . Therefore,  $\bar{u}^\ell$  satisfies the optimality condition of the perturbed parabolic optimal control problem

$$\min_{u \in U_{\text{ad}}} \tilde{J}(u) = J(\hat{y} + \mathcal{S}u, u) + \langle \zeta^\ell, u \rangle_U$$

with ‘‘perturbation’’  $\zeta^\ell$ . The smaller  $\zeta^\ell$  is, the closer  $\bar{u}^\ell$  is to  $\bar{u}$ . Next we estimate  $\|\bar{u} - \bar{u}^\ell\|_U$  in terms of  $\|\zeta^\ell\|_U$ . By Lemma 4.10 we have

$$\mathcal{B}'(\bar{p} - \tilde{p}^\ell) = \mathcal{B}'\mathcal{A}(\bar{u} - \bar{u}^\ell) = -\mathcal{S}'\Theta\mathcal{S}(\bar{u} - \bar{u}^\ell) = \mathcal{S}'\Theta(\tilde{y}^\ell - \bar{y}) \quad (4.26)$$

with  $\tilde{y}^\ell = \hat{y} + \mathcal{S}\bar{u}^\ell$ . Choosing  $u = \bar{u}^\ell$  in (4.16),  $u = \bar{u}$  in (4.25) and using (4.26) we obtain

$$\begin{aligned} 0 &\leq \langle -\sigma(\bar{u} - \bar{u}^\ell) + \mathcal{B}'(\bar{p} - \tilde{p}^\ell) + \zeta^\ell, \bar{u} - \bar{u}^\ell \rangle_U \\ &= -\sigma \|\bar{u} - \bar{u}^\ell\|_U^2 + \langle \mathcal{S}'\Theta(\tilde{y}^\ell - \bar{y}), \bar{u} - \bar{u}^\ell \rangle_U + \langle \zeta^\ell, \bar{u} - \bar{u}^\ell \rangle_U \\ &= -\sigma \|\bar{u} - u_p\|_U^2 - \langle \Theta(\bar{y} - \tilde{y}^\ell), \bar{y} - \tilde{y}^\ell \rangle_{W_0(0,T)', W_0(0,T)} + \langle \zeta^\ell, \bar{u} - \bar{u}^\ell \rangle_U \\ &= -\sigma \|\bar{u} - \bar{u}^\ell\|_U^2 + \langle \zeta^\ell, \bar{u}^\ell - \bar{u}^\ell \rangle_U \leq -\sigma \|\bar{u} - \bar{u}^\ell\|_U^2 + \|\zeta^\ell\|_U \|\bar{u} - \bar{u}^\ell\|_U. \end{aligned}$$

Hence, we get the a-posteriori error estimation

$$\|\bar{u} - \bar{u}^\ell\|_U \leq \frac{1}{\sigma} \|\zeta^\ell\|_U.$$

**Theorem 4.19.** *Suppose that Assumptions 1 and 2 hold. Let  $u \in U$  be arbitrarily given so that  $\mathcal{S}u, \mathcal{A}u \in H^1(0, T; V) \setminus \{0\}$ . To compute a POD basis  $\{\psi_i\}_{i=1}^\ell$  of rank  $\ell$  we choose  $\wp = 4$ ,  $y^1 = \mathcal{S}u$ ,  $y^2 = (\mathcal{S}u)_t$ ,  $y^3 = \mathcal{A}u$  and  $y^4 = (\mathcal{A}u)_t$ . Define the function  $\zeta^\ell \in U$  by*

$$\zeta_i^\ell(t) = \begin{cases} -\min(0, \xi_i^\ell(t)) & \text{a.e. in } \mathcal{A}_{ai}^\ell = \{t \in [0, T] \mid \bar{u}_i^\ell(t) = u_{ai}(t)\}, \\ \max(0, \xi_i^\ell(t)) & \text{a.e. in } \mathcal{A}_{bi}^\ell = \{t \in [0, T] \mid \bar{u}_i^\ell(t) = u_{bi}(t)\}, \\ -\xi_i^\ell(t) & \text{a.e. in } [0, T] \setminus (\mathcal{A}_{ai}^\ell \cup \mathcal{A}_{bi}^\ell), \end{cases}$$

where  $\xi^\ell = \sigma \bar{u}^\ell - \mathcal{B}'(\hat{p} + \mathcal{A}\bar{u}^\ell)$  in  $U$ . Then, the a-posteriori error estimate

$$\|\bar{u} - \bar{u}^\ell\|_U \leq \frac{1}{\sigma} \|\zeta^\ell\|_U. \quad (4.27)$$

In particular,  $\lim_{\ell \rightarrow \infty} \|\zeta^\ell\|_U = 0$ .

*Proof.* Estimate (4.27) has already be shown. We proceed by constructing the function  $\zeta^\ell$ . Here we adapt the lines of the proof of Proposition 3.2 in [73] to our optimal control problem. Suppose that we know  $\bar{u}^\ell$  and  $\tilde{p}^\ell = \hat{p} + \mathcal{A}\bar{u}^\ell$ . The goal is to determine  $\zeta^\ell \in U$  satisfying (4.25). We distinguish three different cases.

- Case  $\bar{u}_i^\ell(t) = u_{ai}(t)$  for fixed  $t \in [0, T]$  and  $i \in \{1, \dots, N_u\}$ : Then,  $u_i(t) - \bar{u}_i^\ell(t) = u_i(t) - u_{ai}(t) \geq 0$  for all  $u \in U_{ad}$ . Hence,  $\zeta_i^\ell(t)$  has to satisfy

$$(\sigma \bar{u}^\ell - \mathcal{B}' \bar{p}^\ell)_i(t) + \zeta_i^\ell(t) \geq 0. \quad (4.28)$$

Setting  $\zeta_i^\ell(t) = -\min(0, (\sigma \bar{u}^\ell - \mathcal{B}' \bar{p}^\ell)_i(t))$  the value  $\zeta_i^\ell(t)$  satisfies (4.28).

- Case  $\bar{u}_i^\ell(t) = u_{bi}(t)$  for fixed  $t \in [0, T]$  and  $i \in \{1, \dots, N_u\}$ : Now,  $u_i(t) - \bar{u}_i^\ell(t) = u_i(t) - u_{bi}(t) \leq 0$  for all  $u \in U_{ad}$ . Analogously to the first case we define  $\zeta_i^\ell(t) = \max(0, (\sigma \bar{u}^\ell - \mathcal{B}' \bar{p}^\ell)_i(t))$  to ensure (4.28).
- Case  $u_{ai}(t) < \bar{u}_i^\ell(t) < u_{bi}(t)$  for fixed  $t \in [0, T]$  and  $i \in \{1, \dots, N_u\}$ : Consequently,  $(\sigma \bar{u}^\ell - \mathcal{B}' \bar{p}^\ell)_i(t) + \zeta_i^\ell(t) = 0$  holds so that  $\zeta_i^\ell(t) = -(\sigma \bar{u}^\ell - \mathcal{B}' \bar{p}^\ell)_i(t)$  guarantees (4.28).

It remains to prove that  $\zeta^\ell$  tends to zero for  $\ell \rightarrow \infty$ . Here we follow adapt the proof of Theorem 4.11 in [73]. By Theorem 4.17-1), the sequence  $\{\bar{u}^\ell\}_{\ell \in \mathbb{N}}$  converges to  $\bar{u}$  in  $U$ . Since the linear operator  $\mathcal{B}'\mathcal{A}$  is bounded and  $\bar{p}^\ell = \hat{p} + \mathcal{A}\bar{u}^\ell$  holds,  $\{\mathcal{B}'\bar{p}^\ell\}_{\ell \in \mathbb{N}}$  tends to  $\mathcal{B}'\bar{p} = \mathcal{B}'\mathcal{A}\bar{u}$  as well. Hence, there exist subsequences  $\{\bar{u}^{\ell_k}\}_{k \in \mathbb{N}}$  and  $\{\mathcal{B}'\bar{p}^{\ell_k}\}_{k \in \mathbb{N}}$  satisfying

$$\lim_{k \rightarrow \infty} \bar{u}_i^{\ell_k}(t) = \bar{u}_i(t) \quad \text{and} \quad \lim_{k \rightarrow \infty} (\mathcal{B}'\bar{p}^{\ell_k})_i(t) = (\mathcal{B}'\bar{p})_i(t)$$

f.a.a.  $t \in [0, T]$  and for  $1 \leq i \leq N_u$ . Next we consider the active and inactive sets for  $\bar{u}$ .

- Let  $t \in \mathcal{J}_i = \{t \in [0, T] \mid u_{ai}(t) < \bar{u}_i(t) < u_{bi}(t)\}$  for  $i \in \{1, \dots, N_u\}$ . For  $k_\circ = k_\circ(t) \in \mathbb{N}$  sufficiently large,  $\bar{u}_i^{\ell_k}(t) \in (u_{ai}(t), u_{bi}(t))$  for all  $k \geq k_\circ$  and f.a.a.  $t \in \mathcal{J}_i$ . Thus,  $(\sigma \bar{u}^{\ell_k} - \mathcal{B}'\bar{p}^{\ell_k})_i(t) = 0$  for all  $k \geq k_\circ(t)$  in  $\mathcal{J}_i$  a.e. This implies

$$\zeta_i^{\ell_k}(t) = 0 \quad \forall k \geq k_\circ(t) \text{ and f.a.a. } t \in \mathcal{J}_i. \quad (4.29)$$

- Suppose that  $t \in \mathcal{A}_{ai} = \{t \in [0, T] \mid u_{ai}(t) = \bar{u}_i(t)\}$  for  $i \in \{1, \dots, N_u\}$ . From  $(\sigma \bar{u}_i - \mathcal{B}'\bar{p})_i(t) \geq 0$  in  $\mathcal{A}_{ai}$  a.e. we deduce

$$\lim_{k \rightarrow \infty} \zeta_i^{\ell_k}(t) = -\min(0, (\sigma \bar{u}^\ell - \mathcal{B}'\bar{p}^\ell)_i(t)) = 0 \quad \text{f.a.a. } t \in \mathcal{A}_{ai}.$$

- Suppose that  $t \in \mathcal{A}_{bi} = \{t \in [0, T] \mid u_{bi}(t) = \bar{u}_i(t)\}$ . Analogously to the second case we find

$$\lim_{k \rightarrow \infty} \zeta_i^{\ell_k}(t) = 0 \quad \text{f.a.a. } t \in \mathcal{A}_{bi}. \quad (4.30)$$

Combining (4.29)-(4.30) we conclude that  $\lim_{k \rightarrow \infty} \zeta_i^{\ell_k} = 0$  a.e. in  $[0, T]$  and for  $1 \leq i \leq N_u$ . Utilizing the dominated convergence theorem [62, p. 24] we have

$$\lim_{k \rightarrow \infty} \|\zeta^{\ell_k}\|_U = 0.$$

Since all subsequences contain a subsequence converging to zero, the claim follows from a standard argument.  $\square$

**Remark 4.20.** 1) Theorem 4.19 shows that  $\|\zeta^\ell\|_U$  can be expected smaller than any tolerance  $\epsilon > 0$  provided that  $\ell$  is taken sufficiently large. Motivated by this result we set up Algorithm 4.3. Note that the quality of the POD Galerkin scheme is improved by only increasing the number

of POD basis elements. Another approach is to update the POD basis in the optimization process; see, e.g., [1, 3, 43].

- 2) We infer from Remark 4.18 that Theorem 4.19 holds still true if we take  $\wp = 2$ ,  $y^1 = \mathcal{S}u$  and  $y^2 = \mathcal{A}u$ .
- 3) In [70] POD a-posteriori error estimates are tested numerically for a linear-quadratic optimal control problem. It turns out that in certain cases not only an increase of the number of POD ansatz functions decreases the error in the reduced-order solution. In this case a change of the POD basis is needed; see, [43, 76], for instance.
- 4) Let us refer to [35], where POD a-posteriori error estimates are combined with an sequential quadratic programming method in order to solve a nonlinear PDE constrained optimal control problem. Furthermore, the presented analysis for linear-quadratic problems can be extended to semilinear optimal control problems by a second-order analysis; see in [36].  $\diamond$

---

**Algorithm 4.3** (POD reduced-order method with a-posteriori estimator)

---

**Require:** Initial control  $u^{0\ell} \in U$ , initial number  $\ell$  for the POD ansatz functions, a maximal number  $\ell_{\max} > \ell$  of POD ansatz functions, and a stopping tolerance  $\epsilon > 0$ .

- 1: Determine  $\hat{y}$ ,  $\hat{p}$ ,  $y^1 = \mathcal{S}u^{0\ell}$ ,  $y^2 = \mathcal{A}u^{0\ell}$ .
  - 2: Compute a POD basis  $\{\psi_i\}_{i=1}^{\ell_{\max}}$  choosing  $y^1$  and  $y^2$ . Set  $\ell = 1$ .
  - 3: **repeat**
  - 4:   Establish the POD Galerkin discretization using  $\{\psi_i\}_{i=1}^{\ell}$ .
  - 5:   Call Algorithm 4.2 to compute suboptimal control  $\bar{u}^{\ell}$ .
  - 6:   Determine  $\zeta^{\ell}$  according to Theorem 4.15 and compute  $\epsilon_{\text{ape}} = \|\zeta^{\ell}\|_U/\sigma$ .
  - 7:   **if**  $\epsilon_{\text{ape}} < \epsilon$  **or**  $\ell = \ell_{\max}$  **then**
  - 8:     Return  $\ell$  and suboptimal control  $\bar{u}^{\ell}$  and STOP.
  - 9:   **end if**
  - 10:   Set  $\ell = \ell + 1$ .
  - 11: **until**  $\ell > \ell_{\max}$
- 

## 5. Numerical experiments

In this section we present numerical test examples to illustrate our theoretical findings. The programs are written in MATLAB utilizing the PARTIAL DIFFERENTIAL EQUATION TOOLBOX for the computation of the finite element (FE) discretization. For the temporal integration the implicit Euler method is applied based on the equidistant time grid  $t_j = (j-1)\Delta t$ ,  $j = 1, \dots, n$  and  $\Delta t = T/(n-1)$ .

**Run 5.1 (POD for the heat equation).** Let us apply the setting of Example 4.1. We choose the final time  $T = 3$ , the spatial domain  $\Omega = (0, 2) \subset \mathbb{R}$ , the Hilbert spaces  $H = L^2(\Omega)$ ,  $V = H_0^1(\Omega)$ , the source term  $f(t, \mathbf{x}) = t^3 - \mathbf{x}^2$

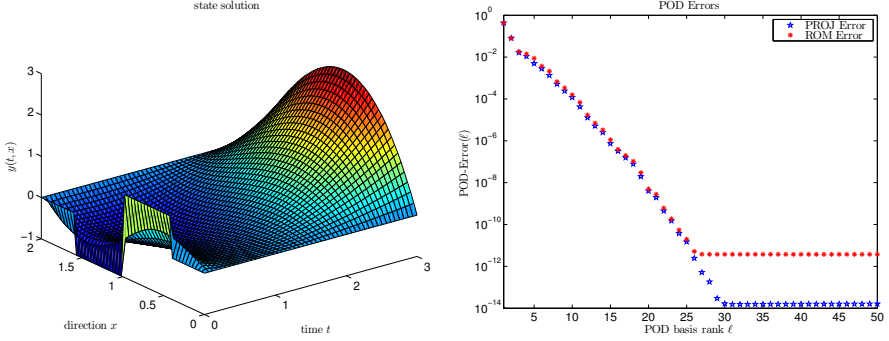


FIGURE 5.1. Run 5.1: The FE solution  $y^h$  (left) and the residuals corresponding to the POD basis rank  $\ell$  (right).

for  $(t, \mathbf{x}) \in Q$  and the discontinuous initial value  $y_\circ(\mathbf{x}) = \chi_{(0.5,1,0)} - \chi_{(1,1.5)}$  for  $\mathbf{x} \in \Omega$ , where, e.g.,  $\chi_{(0.5,1)}$  denotes the characteristic function on the subdomain  $(0.5, 1) \subset \Omega$ ,  $\chi_{(0.5,1)}(\mathbf{x}) = 1$  for  $\mathbf{x} \in (0.5, 1)$  and  $\chi_{(0.5,1)}(\mathbf{x}) = 0$  otherwise. To obtain an accurate approximation of the exact solution we choose  $n = 4000$  so that  $\Delta t \approx 7.5 \cdot 10^{-4}$  holds. For the FE discretization we choose  $m = 500$  spatial grid points and the equidistant mesh size  $h = 2/(m+1) \approx 4 \cdot 10^{-3}$ . Thus, the FE error – measured in the  $H$ -norm – is of the order  $10^{-4}$ . In the left graphic of Figure 5.1, the FE solution  $y^h$  to the state equation (4.3) is visualized. To compute a POD basis  $\{\psi_i\}_{i=1}^\ell$  of rank  $\ell$  we utilize the multiple discrete snapshots  $y_j^1 = y^h(t_j)$  for  $1 \leq j \leq n$  as well  $y_1^2 = 0$  and  $y_j^2 = (y^h(t_j) - y^h(t_{j-1}))/\Delta t$ ,  $j = 2, \dots, n$ , i.e., we include the temporal difference quotients. We choose  $X = H$  and utilize the (stable) singular value decomposition to determine the POD basis of rank  $\ell$ ; compare Remark 2.11. We address this issue in a more detail in Run 5.4. Since the snapshots are FE functions, the POD basis elements are also FE functions. In the right plot of Figure 5.1, the projection and reduced-order error given by

$$\text{PROJ Error}(\ell) = \left( \sum_{j=1}^n \alpha_j \left\| y^h(t_j) - \sum_{i=1}^{\ell} \langle y^h(t_j), \psi_i \rangle_H \psi_i \right\|_H^2 \right)^{1/2},$$

$$\text{ROM Error}(\ell) = \left( \sum_{j=1}^n \alpha_j \left\| y^h(t_j) - y^\ell(t_j) \right\|_H^2 \right)^{1/2}$$

are plotted for different POD basis ranks  $\ell$ . The chosen trapezoidal weights  $\alpha_j$  have been introduced in (2.31). We observe that both errors decay rapidly and coincide until the accuracy  $10^{-12}$ , which is already significant smaller than the FE discretisation error. This numerical results reflects the a-priori error estimates of Theorem 3.11.  $\diamond$

**Run 5.2 (POD for a convection dominated parabolic problem).** To present a more challenging situation, we study a convection-reaction-diffusion equation



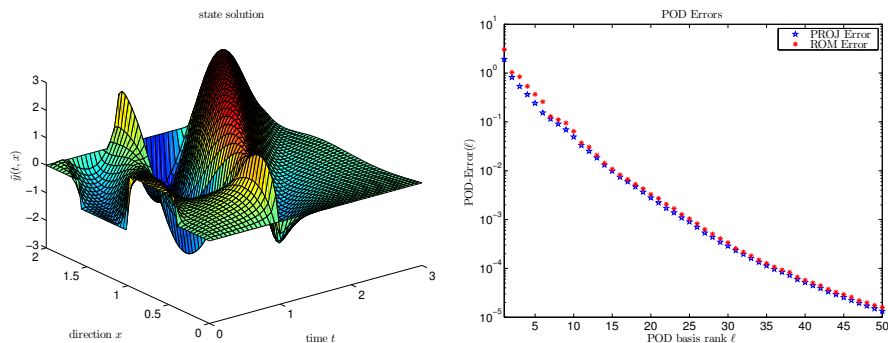


FIGURE 5.2. Run 5.2: The FE solution  $y^h$  (left) and the residuals corresponding the POD basis rank  $\ell$  (right).

with a source term which is close to be singular: Let  $T$ ,  $\Omega$ ,  $y_\circ$ ,  $H$  and  $V$  be given as in Run 5.1. The time-independent bilinear form  $a$  is given by

$$a(\phi, \varphi) = \eta_2 \langle \phi', \varphi' \rangle_H + \eta_1 \langle \phi', \varphi \rangle_H + \eta_0 \langle \phi, \varphi \rangle_H \quad \text{for } \varphi, \phi \in V.$$

We choose the diffusivity  $\sigma_2 = 0.025$ , the velocity  $\sigma_1 = 1.0$  that determines the speed in which the initial profile  $y_\circ$  is shifted to the boundary and the reaction rate  $\sigma_0 = -0.001$ . Finally,  $f(t, \mathbf{x}) = \mathbb{P}(\frac{1}{1-t}) \cos(\pi \mathbf{x})$  for  $(t, \mathbf{x}) \in Q$ , where  $(\mathbb{P}z)(t) = \min(+l, \max(-l, z(t)))$  restricts the image of  $z$  on a bounded interval. In this situation, the state solution  $y$  develops a jump at  $t = 1$  for  $l \rightarrow \infty$ ; see the left plot of Figure 5.2. The right plot of Figure 5.2 demonstrates that in this case, the decay of the reconstruction residuals and the decay of the errors are much slower. The manifold dynamics of the state solution require an inconvenient large number of POD basis elements. Since the supports of these ansatz functions in general cover the whole domain  $\Omega$ , the corresponding system matrices  $M^\ell$  and  $A^\ell$  of the reduced model (compare (3.17)) are not sparse in contrast to the matrices arising in the finite element Galerkin framework, so the model order reduction cannot be provided efficiently for this example if a good accuracy of the solution function  $y^\ell$  is required.  $\diamond$

**Run 5.3 (True and exact approximation error).** Let us consider the setting of Run 5.1 again. The exact solution to (4.3) does not possess a representation by elementary functions. Hence, the presented reconstruction and reduction errors actually are the residuals with respect to a high-order finite element solution  $y^h$ . To compute an approximation  $y$  of the exact solution  $y_{ex}$  we apply a Crank-Nicolson method (with Rannacher smoothing [59]) ensuring  $\|y - y_{ex}\|_{L^2(0,T;H)} = \mathcal{O}(\Delta t^2 + h^2) \approx 10^{-5}$ . In the context of model reduction, such a state is sometimes called the “true” solution. To compute the FE state  $y^h$  we apply the implicit Euler method. In the left plot of Figure 5.3 we compare the true solution with the associated POD approximation for different values  $n = Nt \in \{64, 128, 256, \dots, 8192\}$  of the time integration and

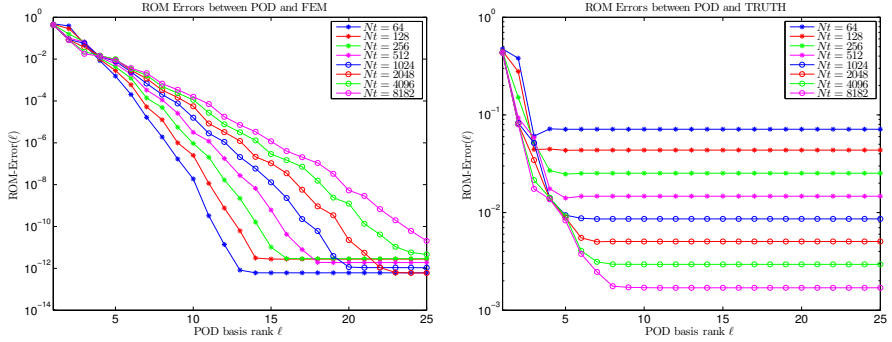


FIGURE 5.3. Run 5.3: The ROM errors with respect to the true solution (left) and the exact one (right).

for the spatial mesh size  $h = 4 \cdot 10^{-3}$ . For the norm we apply a discrete  $L^2(0, T; H)$ -norm as in Run 5.1. Let us mention that we compute for every  $n$  a corresponding FE solution  $y^h$ . We observe that the residuals ignore the errors arising by the application of time and space discretization schemes for the full-order model. The errors decay below the discretization error  $10^{-5}$ . If these discretization errors are taken into account, the residuals stagnate at the level of the full-order model accuracy instead of decaying to zero; see the right plot of Figure 5.3. Due to the implicit Euler method we have  $\|y^h - y_{ex}\|_{L^2(0, T; H)} = \mathcal{O}(\Delta t + h^2)$  with  $h = 4 \cdot 10^{-3}$ . In particular, from  $n \in \{64, 128, 256, \dots, 8192\}$  it follows that  $\Delta t > 3 \cdot 10^{-4} > h^2 = 1.6 \cdot 10^{-5}$ . Therefore, the spatial error is dominated by the time error for all values of  $n$ . We can observe that the exact residuals do not decay below a limit of the order  $\Delta t$ . One can observe that for fixed POD basis rank  $\ell$ , the residuals with respect to the true solution increase if the high-order accuracy is improved by enlarging  $n$ , since the reduced order model has to approximate a more complex system in this case, where the residuals with respect to the exact solution decrease due to the lower limit of stagnation  $\Delta t = 3/(n - 1)$ .  $\diamond$

**Run 5.4 (Different strategies for the POD basis computation).** Let  $Y \in \mathbb{R}^{m \times n}$  denote the matrix of snapshots in the discrete setting,  $W = (\langle \varphi_i, \varphi_j \rangle_X) \in \mathbb{R}^{m \times m}$  be the (sparse) spatial weight matrix arising from the finite element basis  $\{\varphi_i\}_{i=1}^m$  and  $D = \Delta t \text{diag}(\frac{1}{2}, 1, \dots, 1, \frac{1}{2}) \in \mathbb{R}^{n \times n}$  be the trapezoidal time integration matrix fitting to implicit Euler discretization. As it is stated in Remark 2.10, the POD basis  $\{\psi_i\}_{i=1}^\ell$  of rank  $\ell$  can be determined by providing an eigenvalue decomposition of the matrix  $\hat{Y} \hat{Y}^\top = W^{1/2} Y D Y^\top W^{1/2} \in \mathbb{R}^{m \times m}$ , one of  $\hat{Y}^\top \hat{Y} = D^{1/2} Y^\top W Y D^{1/2} \in \mathbb{R}^{n \times n}$ , or a singular value decomposition of  $\hat{Y} = W^{1/2} Y D^{1/2} \in \mathbb{R}^{m \times n}$ . Since  $n \gg m$  in Runs 5.1-5.3, the first variant is the cheapest one from a computational point of view. In case of multiple space dimensions or if a second-order time integration scheme such as some Crank-Nicolson technique is applied, the situation is converse. On the other hand, a singular value decomposition is more accurate than

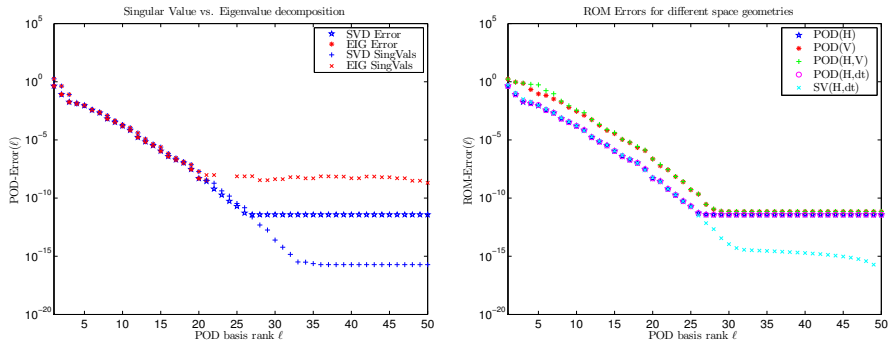


FIGURE 5.4. Run 5.4: Singular values  $\sigma_i$  using the SVD (SVD Vals) or the eigenvalue decomposition (EIG Vals) and the associated ROM errors (SVD error and EIG Error, respectively) (left); ROM errors for different the choices for  $X$ , the error norm and the snapshot ensembles (right).

an eigenvalue decomposition if the POD elements corresponding to eigenvalues/singular values which are close to zero are taken into account: Since  $\lambda_i = \sigma_i^2$  holds for all eigenvalues  $\lambda_i$  and singular values  $\sigma_i$ , the singular values are able to decay to machine precision, where the eigenvalues stagnate significantly above. This is illustrated in the left graphic of Figure 5.4. Indeed, for  $\ell > 20$  the EIG-ROM system matrices become singular due to the numerical errors in the eigenfunctions and the reduced order system is ill-posed in this case while the SVD-ROM model remains stable. In the right plot of Figure 5.4 POD elements are constructed with respect to different scalar products and the resulting ROM errors are compared:  $\|\cdot\|_H$ -residuals for  $X = H$  (denoted by POD(H)),  $\|\cdot\|_V$ -residuals for  $X = V$  (denoted by POD(V)),  $\|\cdot\|_V$ -residuals for  $X = H$  (denoted by POD(H,V)), which also works quite well, the consideration of time derivatives in the snapshot sample (denoted by POD(H,dt)) which allows to apply the a priori error estimate given in (3.20) and the corresponding sums of singular values (denoted by SV(H,dt)) corresponding to the unused eigenfunctions in the latter case which indeed nearly coincide with the ROM errors. In many applications, the quality of the reduced order model does not vary significantly if the weights matrix  $W$  refers to the space  $X = H$  or  $X = V$  and if time derivatives of the used snapshots are taken into account or not. Especially, the ROM residual decays with the same order as the sum over the remaining singular values,  $\|y - y^\ell\|_{W(0,T)} \sim \sum_{i=\ell+1}^{\infty} \sigma_i$  independent of the chosen geometrical framework.  $\diamond$

**Run 5.5 (Iterative methods for the optimal control problem).** In this numerical test we consider solution techniques for the linear-quadratic optimal control problem (P). We define the weights  $\sigma_Q = 1$ ,  $\sigma_\Omega = 0$ , the desired state  $y_Q(t, \mathbf{x}) = t(1 - (\mathbf{x} - 1)^2)$  for  $(t, \mathbf{x}) \in Q$ , the desired final state  $y_\Omega = 0$  (which is redundant due to  $\sigma_\Omega = 0$ , of course), the upper and lower bounds  $u_a = 0.25$ ,

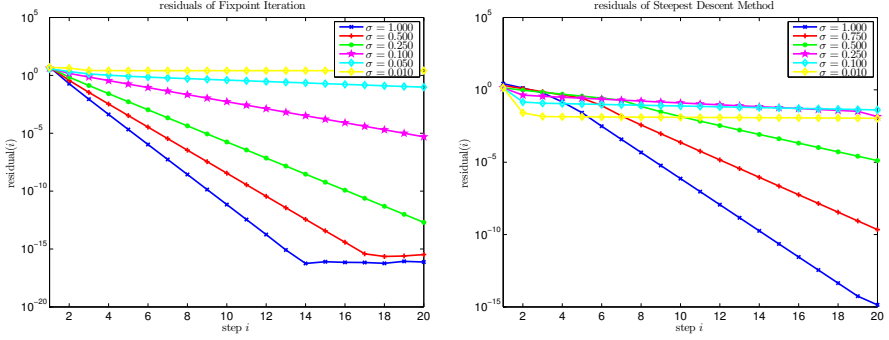


FIGURE 5.5. Residuals of the Banach fixpoint iteration (left) and the projected gradient method (right) for different regularization parameters  $\sigma$ .

$u_b = 0.75$ , the control operator  $(\mathcal{B}u)(t, \mathbf{x}) = u_1(t)\chi_{\Omega_1}(\mathbf{x}) + \dots + u_{10}(t)\chi_{\Omega_{10}}(\mathbf{x})$ , where  $\{\Omega_i \mid i = 1, \dots, 10\}$  is a uniform partition of  $\Omega$  (especially,  $(\mathcal{B}^*p)_i(t) = \int_{\Omega} \chi_i(\mathbf{x})p(t, \mathbf{x}) d\mathbf{x}$  holds) and initial control  $u_o(t) \equiv 1$ .

- 1) *Banach fixed point method*: The first-order necessary and sufficient optimality conditions (4.13) can be reformulated as the equivalent fixpoint problem

$$u = \mathbb{P}\left(\frac{1}{\sigma}\left(\mathcal{B}'\mathcal{A}u - \mathcal{B}'\hat{p}\right)\right) =: F(u),$$

where  $\mathbb{P}(u) = \min(\max(u, u_a), u_b)$  is the orthogonal projection on the set of admissible points  $U_{\text{ad}}$ . The optimal control  $\bar{u} \in U$  can therefore be determined by the Banach fixpoint iteration  $u_{k+1} = F(u_k)$  ( $k > 0$ ) with arbitrary initialization  $u_0 \in U_{\text{ad}}$  provided that  $F$  is a contraction. Since  $\mathbb{P}$  is Lipschitz-continuous with respect to the Lipschitz constant 1, we get

$$\|F(u) - F(v)\|_U \leq \frac{\|\mathcal{B}'\mathcal{A}\|_{\mathcal{L}(U)}}{\sigma} \|u - v\|_U \quad \text{for all } u, v \in U,$$

so the contraction of  $F$  is guaranteed if the regularization parameter  $\sigma$  is sufficiently large. Except of matrix multiplications, each iteration step requires the forward solving of the state equation for  $\tilde{y}(u) = \mathcal{S}u$  and the backwards solving of the adjoint equation for  $\tilde{p}(\tilde{y}) = \mathcal{A}u$ . As it can be observed in the left plot of Figure 5.5, the iteration indeed does not converge if  $\sigma$  is smaller than some critical value  $\sigma_o \approx 0.02$ . Furthermore, the convergence speed of the iteration loop tends to zero for  $\sigma \downarrow \sigma_o$ . We therefore can make use of this method if the control term  $\|u\|_U^2/2$  in the objective functional  $J$  models a control cost such as the required energy and hence shall be small. On the other hand, if we just penalize the objective functional to enforce the strict convexity property and are interested in the case  $\sigma \rightarrow 0$  (the resulting controls

usually are of “bang-bang”-type in this case, i.e.  $u(t) \in \{u_a, u_b\}$  for almost all  $t \in [0, T]$ ), we shall apply some other optimization technique.

- 2) *Projected gradient method*: A suitable steepest descent method for the control-constrained optimization problem is the projected gradient algorithm; see, [38], for instance. Here, the next iteration point is given by the formula  $u_{k+1} = \mathbb{P}(u_k + s_k d_k)$ , where  $d_k = -\nabla J(u_k) = \sigma u_k - \mathcal{B}'(\mathcal{A}u_k + \hat{p})$  is the direction of the steepest descent of  $J$  in the current iteration point  $u_k$  and  $s_k > 0$  is chosen by Algorithm 5.1. This

---

**Algorithm 5.1** (Backtracking strategy)

---

**Require:** Maximal number  $j_{\max}$  of iterations and parameter  $c \in (0, 1)$ .

- 1: Set  $s^{(0)} = 1$  and  $j = 1$ .
  - 2: **while**  $\hat{J}(u_k + s^{(j)} d_k) > \hat{J}(u_k) - cs^{(j)} \|d_k\|_U$  **and**  $j < j_{\max}$  **do**
  - 3:   Set  $s^{(j+1)} = s^{(j)}/2$  and  $j = j + 1$ .
  - 4: **end while**
  - 5: **return**  $s_k = s^{(j)}$
- 

procedure is globally convergent. However, as before, the convergence speed becomes extremely slow for  $\sigma \rightarrow 0$ . In addition, if the step size condition  $\hat{J}(u_k + s^{(j)} d_k) \leq \hat{J}(u_k) - cs^{(j)} \|d_k\|_U$  is just fulfilled for very small step sizes  $s^{(j)}$ , many evaluations of the reduced objective functional are required to test whether  $\hat{J}(u_k + s^{(j)} d_k) \leq \hat{J}(u_k) - cs^{(j)} \|d_k\|_U$  is satisfied. Here, each evaluation requires to solve the state equation. Therefore, also the single iteration steps may become quite expensive. The right plot of Figure 5.5 demonstrates that also the projected gradient method cannot deal with small regularizations. In contrast to the Banach iteration, the residuals decay for arbitrarily small values of  $\sigma$ , but the numerical effort explodes if  $\sigma$  tends to zero.

- 3) *Primal-dual active set strategy*: This method – see Algorithm 4.1 for the infinite-dimensional case and Algorithm 4.2 for the POD discretization – solves the state and the adjoint equation simultaneously within the implicit linear scheme

$$u_{k+1}(t) = \chi_{\mathcal{A}_a^k}(t)u_a + \chi_{\mathcal{A}_b^k}(t)u_b(t) + \chi_{\mathcal{I}^k}(t) \frac{1}{\sigma}(\mathcal{B}'\mathcal{A}u_{k+1})(t)$$

f.a.a.  $t \in [0, T]$ . Since this technique is equivalent to a semismooth Newton procedure [24] locally superlinear convergence rates are provided. Further, the algorithm is able to deal with smaller regularizations than the other two methods presented: Reasonable computation times are provided for all  $\sigma > \sigma_o \approx 0.0002$ , see Figure 5.6. For parameters below this critical value, the bang-bang-control  $u$  oscillates between  $u_a$  and  $u_b$  at the boundary grid points of the active sets. Notice that both the critical  $\sigma_o$  and the error between the exact solution and the suboptimal final iteration depend on the number of discretization points. The numerical effort of the simultaneous solving operations in each iteration

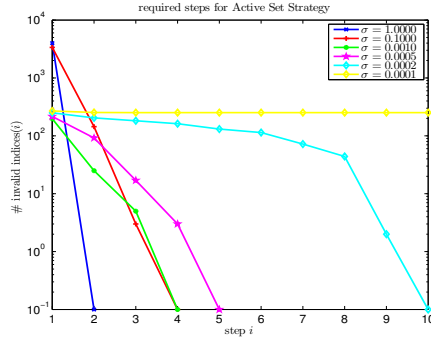


FIGURE 5.6. Run 5.5: The number of grid points, where the previous and the actual active sets differ for different regularization parameters  $\sigma$ .

step is significantly larger than the single solution since the initial condition for the state and the final condition for the adjoint state prevent to iteratively solve  $n$  times a system of dimension  $2m$ ; instead, all time and space values  $(y(t_i, \mathbf{x}_j), p(t_i, \mathbf{x}_j))$  are determined by solving a linear system of the dimension  $2nm$ . Here, the model order reduction techniques come into play which will lead to formidable calculation time reductions (or even make an execution of the primal-dual active set strategy just possible). In the following, we will make use of this optimization procedure.  $\diamond$

**Run 5.6 (Different Galerkin expansions).** In this run we compare the *modified* POD Galerkin expansions (3.14) for the state variable and (4.17) for the dual variable with the *standard* Galerkin approximations:

$$y^\ell(t) = \sum_{i=1}^{\ell} y_i^\ell(t) \psi_i, \quad p^\ell(t) = \sum_{i=1}^{\ell} p_i^\ell(t) \psi_i \quad \text{for } t \in [0, T]. \quad (5.1)$$

We choose the same setting as in Run 5.5. Let  $\sigma = 0.1$ . In Figures 5.7 and 5.8 we plot the optimal FE solution components  $(\bar{y}^h, \bar{u}^h, \bar{p}^h, \bar{\lambda}^h)$  obtained by using the primal-dual active set strategy. We observe that the support of the multiplier  $\mathcal{B}\bar{\lambda}^h$  coincides with the active set for the control variable  $\mathcal{B}\bar{u}^h$ . Further, the relation  $\bar{u}^h = \mathbb{P}(\mathcal{B}'\bar{p}^h/\sigma)$  can be observed. As it is stated in Remark 3.10-1) the advantage of the modified Galerkin ansatz is that the ROM errors do not include the projection of the initial value on the POD space. Figure 5.9 illustrates the impact of homogenization, where we not only plot the ROM errors, but also the a-posteriori error estimate for different  $\ell$ ; compare Section 4.5. First we see that the ROM errors and the a-posteriori error estimate nearly coincide in all scenarios. In the left plot of Figure 5.9 the POD basis is computed from snapshots of the state equation taking the control guess  $u_0 \equiv 1$ . One observes that the dynamics of the corresponding homogeneous snapshots in the modified ansatz are not sufficient to decrease

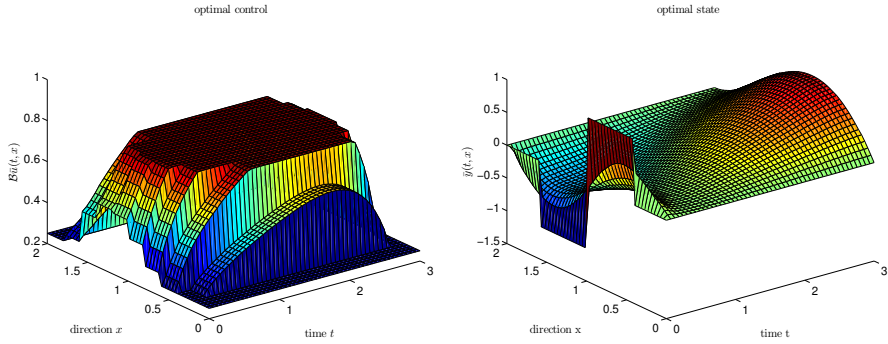


FIGURE 5.7. Run 5.6: The optimal FE control  $\mathcal{B}\bar{u}^h$  (left) and the optimal FE state  $\bar{y}^h$  (right).

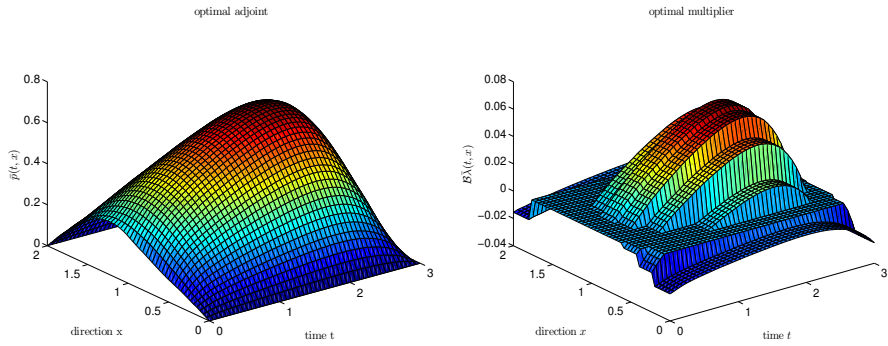


FIGURE 5.8. Run 5.6: The optimal FE adjoint state  $\bar{p}^h$  and the optimal FE Lagrange multiplier  $\mathcal{B}\bar{\lambda}^h$ .

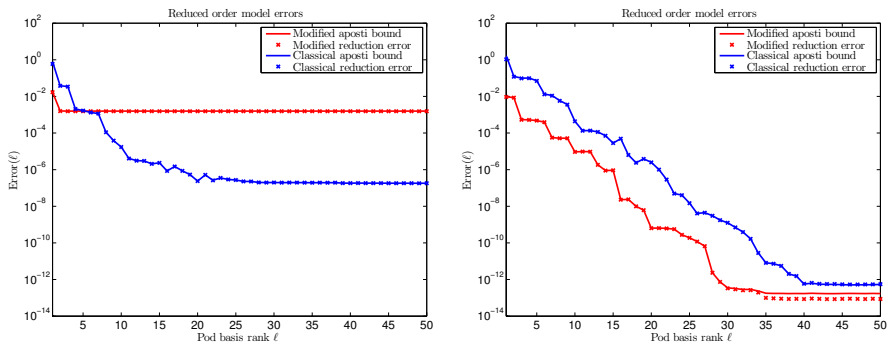


FIGURE 5.9. Run 5.6: The ROM errors for the standard and the modified POD ansatz for initial control guesses  $u_0 = 1$  (left) and  $u_0 = \bar{u}^h$  (right).

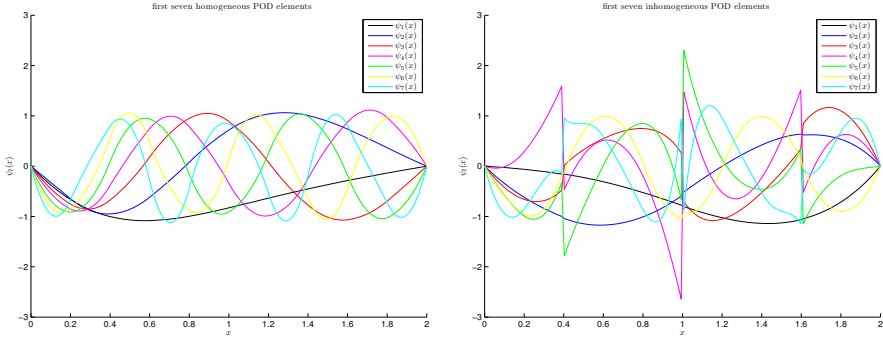


FIGURE 5.10. Run 5.6: The first POD basis elements for the modified (left) and the standard (right) Galerkin expansion.

the control error below a level of  $10^{-3}$  while the standard Galerkin ansatz, exploiting also the dynamics of the initial value and the inhomogeneity, induces a higher dimensional POD space and leads to an error order below  $10^{-6}$ . In the right plot of Figure 5.9 the optimal FE control  $\bar{u}^h$  creates the snapshots. Here, the modified Galerkin ansatz pays: The approximation error in the standard Galerkin ansatz is dominated by the projection error of the initial value  $y_o$  on the POD Galerkin ansatz space. This example also shows that good approximations of the reduced order model are just guaranteed in the case that the snapshots which build up the POD basis include the dynamics of the optimal state solution; otherwise, enlarging the POD basis rank does not necessarily improve the accuracy of the results. Algorithm 4.3 proposes a solution for this problem: Here, basis updates are provided if the a posteriori error estimator presented in Theorem 4.19 indicate that the control error does not decay in the current POD model. Figure 5.9 shows that these error bounds are sharp. Indeed, if the algorithm is initialized with the control guess  $u_o \equiv 1$  and a single basis update is provided, i.e., a new POD basis is calculated with respect to the achieved suboptimal POD control  $u_1^{\ell_{\max}}$ . This new POD basis coincides with the POD basis associated with the best (but usually unknown) control guess  $\bar{u}^h$ . Thus, the resulting error decay by enlarging  $\ell$  is the same one as in the right graphic of Figure 5.9. In Figure 5.10 the first POD basis functions are presented for the modified and standard Galerkin expansions. Consequently, the reconstruction of the initial condition  $y_o$  with the standard Galerkin ansatz works quite well as it is demonstrated in Figure 5.11 – especially, due to the shape of the POD basis functions, no oscillations at the jump points occur as can be observed by trigonometric Fourier approximations, for instance. For the modified POD Galerkin ansatz it is neither required nor possible to build up the initial value  $y_o$  accurately. But this is not needed, because the initial condition is explicitly included in the initial condition; see (3.14). If the model data is perturbed by noise, the improvement of homogenization is even significantly stronger. For the following simulation, we add random data onto the initial value  $y_o$ . The controls



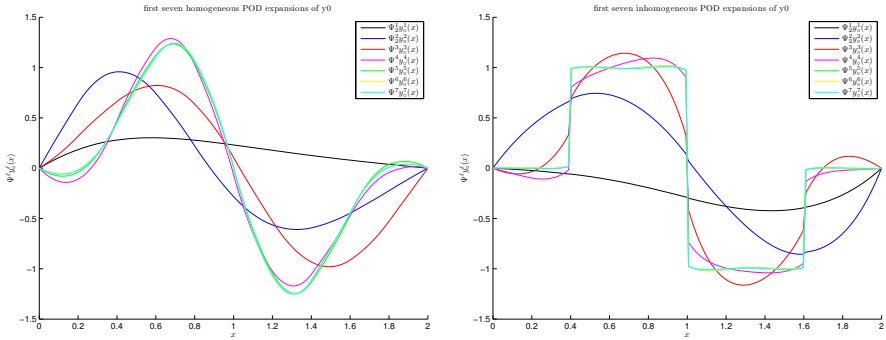


FIGURE 5.11. Run 5.6: The reconstruction error  $\Psi^\ell y_0 = \sum_{i=1}^\ell \langle y_0, \psi_i \rangle_H \psi_i$  for the initial condition  $y_0$  for the modified (left) and the standard (right) POD Galerkin expansions.

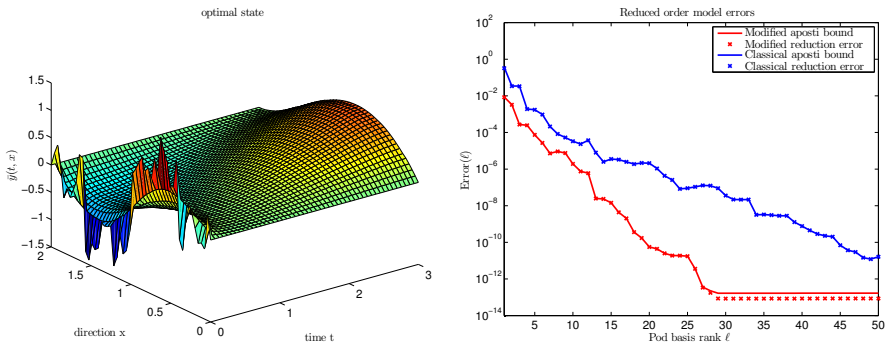


FIGURE 5.12. Run 5.6: The optimal state solution for perturbed initial data (left) and the ROM errors for the two POD ansatze (right).

gained by the modified model then reach the optimal precision  $10^{-13}$  with 29 POD basis functions, where even 50 basis elements are not sufficient in the standard ansatz to decrease the error below a level of  $10^{-11}$ , see Figure 5.12. We observe that the noise in the initial value is inherited to the POD basis elements of the modified Galerkin ansatz; despite of this perturbation, their shape does not differ much from those of the POD basis for the unperturbed initial conditions standard Galerkin ansatz. This is different for the standard POD Galerkin ansatz; compare Figures 5.10 and 5.13.  $\diamond$

## References

[1] K. Afanasiev and M. Hinze. Adaptive control of a wake flow using proper orthogonal decomposition. Lecture Notes in Pure and Applied Mathematics, 216:317-332, 2001.

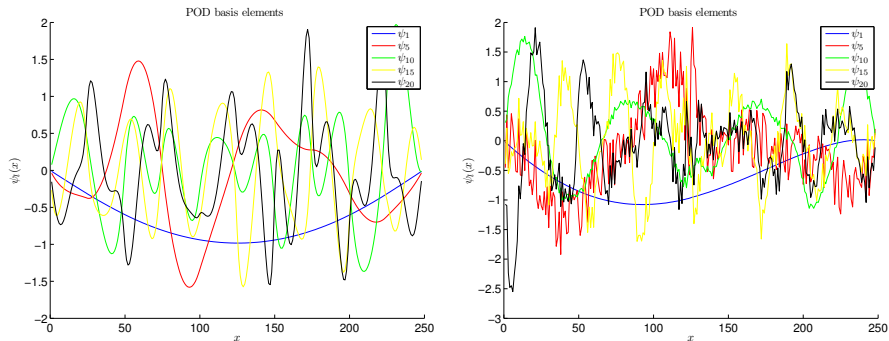


FIGURE 5.13. Run 5.6: The first POD basis elements for the modified (left) and the standard (right) Galerkin expansion in case of the perturbed initial condition.

- [2] A.C. Antoulas. *Approximation of Large-Scale Dynamical Systems*. Advances in Design and Control, SIAM, Philadelphia, 2005.
- [3] E. Arian, M. Fahl, and E.W. Sachs. Trust-region proper orthogonal decomposition for flow control. Technical Report 2000-25, ICASE, 2000.
- [4] P. Astrid, S. Weiland, K. Willcox, and T. Backx. Missing point estimation in models described by proper orthogonal decomposition. Proceedings 43rd IEEE Conference on Decision and Control, December, 2004.
- [5] J.A. Atwell, J.T. Borggaard, and B.B. King. Reduced-order controllers for Burgers' equation with a nonlinear observer. *Int. J. Appl. Math. Comp. Sci.*, 11:1311-1330, 2001.
- [6] H.T. Banks, M.L. Joyner, B. Winchesky, and W.P. Winfree. Nondestructive evaluation using a reduced-order computational methodology. *Inverse Problems*, 16:1-17, 2000.
- [7] S.C. Brenner and L.R. Scott. *The Mathematical Theory of Finite Element Methods*. Texts in Applied Mathematics, vol. 15, Springer, 2008.
- [8] A. Chatterjee. *An introduction to the proper orthogonal decomposition*. *Current Science*, 78:539-575, 2000.
- [9] D. Chapelle, A. Gariah, and J. Saint-Marie. Galerkin approximation with proper orthogonal decomposition: new error estimates and illustrative examples. *ESAIM: Math. Model. Numer. Anal.*, 46:731-757, 2012.
- [10] R. Dautray and J.-L. Lions. *Mathematical Analysis and Numerical Methods for Science and Technology. Volume 5: Evolution Problems I*. Springer, Berlin, 2000.
- [11] K. Deckelnick and M. Hinze. Error estimates in space and time for tracking-type control of the instationary Stokes system. *International Series of Numerical Mathematics*, 143:87-103, 2002).
- [12] K. Deckelnick and M. Hinze. Semidiscretization and error estimates for distributed control of the instationary Navier-Stokes equations. *Numerische Mathematik*, 97:297-320, 2004.

- [13] L. Dede. Reduced basis method and a posteriori error estimation for parametrized linear-quadratic optimal control problems. *SIAM Journal on Scientific Computing*, 32:997-1019, 2010.
- [14] M. Dihlmann and B. Haasdonk. Certified nonlinear parameter optimization with reduced basis surrogate models. Submitted, 2013.
- [15] F. Diwoky and S. Volkwein. Nonlinear boundary control for the heat equation utilizing proper orthogonal decomposition. In K.-H. Hoffmann, R. H. W. Hoppe, V. Schulz, editors, *Fast Solution of Discretized Optimization Problems*, International Series of Numerical Mathematics, 138:73-87, 2001.
- [16] A. L. Dontchev, W. W. Hager, A. B. Poore, and B. Yang. Optimality, stability, and convergence in nonlinear control. *Applied Math. and Optimization*, 31:297-326, 1995.
- [17] L.C. Evans. *Partial Differential Equations*. American Math. Society, Providence, Rhode Island, 2008.
- [18] M. Gubisch and S. Volkwein. POD a-posteriori error analysis for optimal control problems with mixed control-state constraints. Submitted, 2013.
- [19] M. Grepl and M. Kärcher. A posteriori error estimation for reduced order solutions of parametrized parabolic optimal control problems. Submitted, 2013.
- [20] M.A. Grepl, Y. Maday, N.C. Nguyen and A.T. Patera. Efficient reduced-basis treatment of nonaffine and nonlinear partial differential equations. *ESAIM: Math. Model. Numer. Anal.*, 41:575-605, 2007.
- [21] M. Heinkenschloss, D.C. Sorensen, and K. Sun. Balanced truncation model reduction for a class of descriptor systems with application to the Oseen equations. *SIAM J. Sci. Comput.*, 30:1038-1063, 2008.
- [22] S. Herkt, M. Hinze, and R. Pinnau. *Convergence analysis of Galerkin POD for linear second order evolution equations*. *Electronic Transactions on Numerical Analysis*, 40:321-337, 2013.
- [23] C. Himpe and M. Ohlberger. Cross-gramian based combined state and parameter reduction. Submitted, 2013.
- [24] M. Hintermüller, K. Ito, and K. Kunisch. The primal-dual active set strategy as a semismooth Newton method. *SIAM Journal on Optimization*, 13:865-888, 2003.
- [25] M. Hinze. A variational discretization concept in control constrained optimization: the linear-quadratic case. *Computational Optim. Appl.*, 30:45-61, 2005.
- [26] M. Hinze, R. Pinnau, M. Ulbrich, and S. Ulbrich. *Optimization with PDE Constraints*. Springer, 2009.
- [27] M. Hinze and F. Tröltzsch. Discrete concepts versus error analysis in pde constrained optimization. *GAMM-Mitteilungen*, 33:148-162, 2010.
- [28] M. Hinze and S. Volkwein. Proper orthogonal decomposition surrogate models for nonlinear dynamical systems: error estimates and suboptimal control. *Lecture Notes in Computational Science and Engineering*, 45:261-306, 2005.
- [29] M. Hinze and S. Volkwein. Error estimates for abstract linear-quadratic optimal control problems using proper orthogonal decomposition. *Computational Optimization and Applications*, 39:319-345, 2008.
- [30] D. Hömberg and S. Volkwein. Control of laser surface hardening by a reduced-order approach using proper orthogonal decomposition. *Mathematical and Computer Modelling*, 38:1003-1028, 2003.

- [31] P. Holmes, J.L. Lumley, G. Berkooz, and C.W. Rowley. *Turbulence, Coherent Structures, Dynamical Systems and Symmetry*. Cambridge Monographs on Mechanics, Cambridge University Press, second edition, 2012.
- [32] K. Ito and S.S. Ravindran. A reduced basis method for control problems governed by PDEs. In W. Desch, F. Kappel, and K. Kunisch, eds., *Control and Estimation of Distributed Parameter Systems*. Proceedings of the International Conference in Vorau, 1996, Birkhäuser-Verlag, Basel, 126:153-168, 1998.
- [33] M. Grepl and M. Kärcher. A posteriori error estimation for reduced order solutions of parametrized parabolic optimal control problems. Submitted, 2013.
- [34] M. Kahlbacher and S. Volkwein. Galerkin proper orthogonal decomposition methods for parameter dependent elliptic systems. *Discussiones Mathematicae: Differential Inclusions, Control and Optimization*, 27:95-117, 2007.
- [35] M. Kahlbacher and S. Volkwein. POD a-posteriori error based inexact SQP method for bilinear elliptic optimal control problems. *ESAIM: Mathematical Modelling and Numerical Analysis*, 46:491-511, 2012.
- [36] E. Kammann, F. Tröltzsch, and S. Volkwein. A method of a-posteriori error estimation with application to proper orthogonal decomposition. *ESAIM: Mathematical Modelling and Numerical Analysis*, 47:555-581, 2013.
- [37] T. Kato. *Perturbation Theory for Linear Operators*. Springer-Verlag, Berlin, 1980.
- [38] C.T. Kelley. *Iterative Methods for Optimization*. SIAM Frontiers in Applied Mathematics, Philadelphia, 1999.
- [39] K. Kunisch and S. Volkwein. Control of Burgers' equation by a reduced order approach using proper orthogonal decomposition. *Journal on Optimization Theory and Applications*, 102:345-371, 1999.
- [40] K. Kunisch and S. Volkwein. Galerkin proper orthogonal decomposition methods for parabolic problems. *Numerische Mathematik*, 90:117-148, 2001.
- [41] K. Kunisch and S. Volkwein. Galerkin proper orthogonal decomposition methods for a general equation in fluid dynamics. *SIAM Journal on Numerical Analysis*, 40:492-515, 2002.
- [42] K. Kunisch and S. Volkwein. Crank-Nicolson Galerkin proper orthogonal decomposition approximations for a general equation in fluid dynamics. Proceedings of the 18th GAMM Seminar on *Multigrid and related methods for optimization problems*, Leipzig, 97-114, 2002.
- [43] K. Kunisch and S. Volkwein. Proper orthogonal decomposition for optimality systems. *ESAIM: Mathematical Modelling and Numerical Analysis*, 42:1-23, 2008.
- [44] K. Kunisch and S. Volkwein. Optimal snapshot location for computing POD basis functions. *ESAIM: Mathematical Modelling and Numerical Analysis*, 44:509-529, 2010.
- [45] K. Kunisch and S. Volkwein, and L. Xie. HJB-POD based feedback design for the optimal control of evolution problems. *SIAM Journal on Applied Dynamical Systems*, 3:701-722, 2004.
- [46] S. Lall, J.E. Marsden, and S. Glavaski. A subspace approach to balanced truncation for model reduction of nonlinear control systems. *Int. J. Robust Nonlinear Control*, 12:519-535, 2002.

- [47] O. Lass and S. Volkwein. POD Galerkin schemes for nonlinear elliptic-parabolic systems. *SIAM Journal on Scientific Computing*, 35(3):A1271-A1298, 2013.
- [48] F. Leibfritz and S. Volkwein. Reduced order output feedback control design for PDE systems using proper orthogonal decomposition and nonlinear semidefinite programming. *Linear Algebra and Its Applications*, 415:542-575, 2006.
- [49] H.V. Ly and H.T. Tran. Modeling and control of physical processes using proper orthogonal decomposition. *Mathematical and Computer Modeling*, 33:223-236, 2001.
- [50] V. Mehrmann and T. Stykel. Balanced truncation model reduction for large-scale systems in descriptor form. In *Dimension Reduction of Large-Scale Systems, Lecture Notes Comput. Sci. Eng.*, 45,, 2005.
- [51] F. Negri, G. Rozza, A. Manzoni, and A. Quateroni. Reduced basis method for parametrized elliptic optimal control problems. *SIAM Journal on Scientific Computing*, to appear.
- [52] B. Noack, K. Afanasiev, M. Morzynski, G. Tadmor, and F. Thiele. A hierarchy of low-dimensional models for the transient and post-transient cylinder wake. *J. Fluid. Mech.*, 497:335-363, 2003.
- [53] B. Noble. *Applied Linear Algebra*. Englewood Cliffs, NJ : Prentice-Hall, 1969.
- [54] J. Nocedal and S.J. Wright. *Numerical Optimization*. Springer Series in Operation Research, second edition, 2006.
- [55] M. Ohlberger and M. Schäfer. Error control based model reduction for parameter optimization of elliptic homogenization problems. To appear in the *Proceedings of the 1st IFAC Workshop on Control of Systems Governed by Partial Differential Equations* , 2013
- [56] A.T. Patera and G. Rozza. Reduced Basis Approximation and A Posteriori Error Estimation for Parametrized Partial Differential Equations. MIT Papalardo Graduate Monographs in Mechanical Engineering, 2006.
- [57] R. Pinnau. Model reduction via proper orthogonal decomposition. In *Model Order Reduction: Theory, Research Aspects and Applications*, Mathematics in Industry, vol. 13, pp. 95-109, Springer, 2008.
- [58] O. Pironneau. Calibration of options on a reduced basis. *Journal of Computational and Applied Mathematics*, 232:139-147, 2009.
- [59] R. Rannacher. Finite element solution of diffusion problems with irregular data. *Numerische Mathematik*, 43: 309-327, 1984.
- [60] M. Rathinam and L. Petzold. Dynamic iteration using reduced order models: a method for simulation of large scale modular systems. *SIAM J. Numer. Anal.*, 40:1446-1474, 2002.
- [61] S.S. Ravindran. Reduced-order adaptive controllers for fluid flows using POD. *SIAM J. Sci. Comput.*, 15:457-478, 2000.
- [62] M. Reed and B. Simon. *Methods of Modern Mathematical Physics I: Functional Analysis*. Academic Press, New York, 1980.
- [63] C.W. Rowley. Model reduction for fluids, using balanced proper orthogonal decomposition. *Int. J. on Bifurcation and Chaos*, 15:997-1013, 2005.
- [64] E.W. Sachs and M. Schu. A-priori error estimates for reduced order models in finance. *ESAIM: Mathematical Modelling and Numerical Analysis*, 47:449-469, 2013.

- [65] E.W. Sachs and S. Volkwein. POD Galerkin approximations in PDE-constrained optimization. *Gamm-Mitteilungen*, 33:194-208, 2010.
- [66] W.H.A. Schilders, H.A. van der Vorst, and J. Rommes. *Model Order Reduction: Theory, Research Aspects and Applications*. Mathematics in Industry, vol. 13, Springer, 2008.
- [67] M. Schu. *Adaptive Trust-Region POD Methods and their Application in Finance*. Ph.D thesis, University of Trier, 2013.
- [68] J.R. Singler. New POD expressions, error bounds, and asymptotic results for reduced order models of parabolic PDEs. Submitted, 2013.
- [69] L. Sirovich. Turbulence and the dynamics of coherent structures. Parts I-II. *Quarterly of Applied Mathematics*, XVI:561-590, 1987.
- [70] A. Studinger and S. Volkwein. Numerical analysis of POD a-posteriori error estimation for optimal control. *International Series of Numerical Mathematics*, 164:137-158, 2013
- [71] T. Tonn, K. Urban, and S. Volkwein. Comparison of the reduced-basis and POD a-posteriori error estimators for an elliptic linear quadratic optimal control problem. *Mathematical and Computer Modelling of Dynamical Systems*, 17:355-369, 2011.
- [72] F. Tröltzsch. *Optimal Control of Partial Differential Equations. Theory, Methods and Applications*. American Math. Society, Providence, volume 112, 2010.
- [73] F. Tröltzsch and S. Volkwein. POD a-posteriori error estimates for linear-quadratic optimal control problems. *Computational Optimization and Applications*, 44:83-115, 2009.
- [74] M. Ulbrich. *Semismooth Newton Methods for Variational Inequalities and Constrained Optimization Problems in Function Spaces*. MOS-SIAM Series on Optimization, vol. 11, SIAM, 2011.
- [75] S. Volkwein. Optimal control of a phase-field model using proper orthogonal decomposition. *Zeitschrift für Angewandte Mathematik und Mechanik*, 81:83-97, 2001.
- [76] S. Volkwein. Optimality system POD and a-posteriori error analysis for linear-quadratic problems. *Control and Cybernetics*, 40:1109-1125, 2011.
- [77] S. Volkwein. *Proper Orthogonal Decomposition: Theory and Reduced-Order Modelling*. Lecture Notes, University of Konstanz, 2012.
- [78] G. Vossen and S. Volkwein. Model reduction techniques with a-posteriori error analysis for linear-quadratic optimal control problems. *Numerical Algebra, Control and Optimization*, 2:465-485, 2012.
- [79] K. Willcox and J. Peraire. Balanced model reduction via the proper orthogonal decomposition. *American Institute of Aeronautics and Astronautics (AIAA)*, 2323-2330, 2002.
- [80] K. Yosida. *Functional Analysis* Classics in Mathematics, Reprint of the 1980 edition, Springer-Verlag, Heidelberg, 1995.
- [81] K. Zhou, J.C. Doyle, and K. Glover. *Robust and Optimal Control*. Prentice-Hall, Upper Saddle River, NJ, 1996.

Martin Gubisch  
University of Konstanz  
Department of Mathematics and Statistics  
Universitätsstraße 10  
D-78457 Konstanz  
Germany  
e-mail: [Martin.Gubisch@uni-konstanz.de](mailto:Martin.Gubisch@uni-konstanz.de)

Stefan Volkwein  
University of Konstanz  
Department of Mathematics and Statistics  
Universitätsstraße 10  
D-78457 Konstanz  
Germany  
e-mail: [Stefan.Volkwein@uni-konstanz.de](mailto:Stefan.Volkwein@uni-konstanz.de)