

MODEL REDUCTION USING PROPER ORTHOGONAL DECOMPOSITION

S. VOLKWEIN

ABSTRACT. In this lecture notes an introduction to model reduction utilizing proper orthogonal decomposition (POD) is given. The close connection between POD and singular value decomposition (SVD) of rectangular matrices is emphasized. As an application POD is used to derive a reduced-order model for non-linear initial value problems. The strategy is extended to linear-quadratic optimal control problems governed by ordinary differential equations. The relationship to classical model reduction techniques like balanced truncation is studied.

1. The POD method in \mathbb{R}^m

In this section we introduce the POD method in the Euclidean space \mathbb{R}^m and study the close connection to the SVD of rectangular matrices; see [6]. We also refer to the monograph [3].

1.1. POD and SVD. Let $Y = [y_1, \dots, y_n]$ be a real-valued $m \times n$ matrix of rank $d \leq \min\{m, n\}$ with columns $y_j \in \mathbb{R}^m$, $1 \leq j \leq n$. Consequently,

$$(1.1) \quad \bar{y} = \frac{1}{n} \sum_{j=1}^n y_j$$

can be viewed as the column-averaged mean of the matrix Y .

SVD [10] guarantees the existence of real numbers $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_d > 0$ and orthogonal matrices $U \in \mathbb{R}^{m \times m}$ with columns $\{u_i\}_{i=1}^m$ and $V \in \mathbb{R}^{n \times n}$ with columns $\{v_i\}_{i=1}^n$ such that

$$(1.2) \quad U^T Y V = \begin{pmatrix} D & 0 \\ 0 & 0 \end{pmatrix} =: \Sigma \in \mathbb{R}^{m \times n},$$

where $D = \text{diag}(\sigma_1, \dots, \sigma_d) \in \mathbb{R}^{d \times d}$ and the zeros in (1.2) denote matrices of appropriate dimensions. Moreover the vectors $\{u_i\}_{i=1}^d$ and $\{v_i\}_{i=1}^d$ satisfy

$$(1.3) \quad Y v_i = \sigma_i u_i \quad \text{and} \quad Y^T u_i = \sigma_i v_i \quad \text{for } i = 1, \dots, d.$$

They are eigenvectors of $Y Y^T$ and $Y^T Y$, respectively, with eigenvalues $\lambda_i = \sigma_i^2 > 0$, $i = 1, \dots, d$. The vectors $\{u_i\}_{i=d+1}^m$ and $\{v_i\}_{i=d+1}^n$ (if $d < m$ respectively $d < n$) are eigenvectors of $Y Y^T$ and $Y^T Y$ with eigenvalue 0.

From (1.2) we deduce that

$$Y = U \Sigma V^T.$$

Date: December 7, 2011.

Key words and phrases. Proper orthogonal decomposition, singular value decomposition, model reduction, error estimates, linear-quadratic regulator problems, balanced truncation.

It follows that Y can also be expressed as

$$(1.4) \quad Y = U^d D (V^d)^T,$$

where $U^d \in \mathbb{R}^{m \times d}$ and $V^d \in \mathbb{R}^{n \times d}$ are given by

$$\begin{aligned} U_{ij}^d &= U_{ij} & \text{for } 1 \leq i \leq m, 1 \leq j \leq d, \\ V_{ij}^d &= V_{ij} & \text{for } 1 \leq i \leq n, 1 \leq j \leq d. \end{aligned}$$

Setting $B^d = D(V^d)^T \in \mathbb{R}^{d \times n}$ we can write (1.4) in the form

$$Y = U^d B^d \quad \text{with } B^d = D(V^d)^T \in \mathbb{R}^{d \times n}.$$

Thus, the column space of Y can be represented in terms of the d linearly independent columns of U^d . The coefficients in the expansion for the columns y_j , $j = 1, \dots, n$, in the basis $\{u_i\}_{i=1}^d$ are given by the j th-column of B^d . Since U is orthogonal, we find that

$$\begin{aligned} y_j &= \sum_{i=1}^d B_{ij}^d U_{\cdot,i}^d = \sum_{i=1}^d (D(V^d)^T)_{ij} u_i = \sum_{i=1}^d \underbrace{((U^d)^T U^d)}_{=I^d \in \mathbb{R}^{d \times d}} D(V^d)^T_{ij} u_i \\ &\stackrel{(1.4)}{=} \sum_{i=1}^d ((U^d)^T Y)_{ij} u_i = \sum_{i=1}^d \underbrace{\left(\sum_{k=1}^m U_{ki}^d Y_{kj} \right)}_{=u_i^T y_j} u_i = \sum_{i=1}^d \langle u_i, y_j \rangle_{\mathbb{R}^m} u_i, \end{aligned}$$

where $\langle \cdot, \cdot \rangle_{\mathbb{R}^m}$ denotes the canonical inner product in \mathbb{R}^m . Thus,

$$(1.5) \quad y_j = \sum_{i=1}^d \langle y_j, u_i \rangle_{\mathbb{R}^m} u_i \quad \text{for } j = 1, \dots, n$$

Let us now interpret SVD in terms of POD. One of the central issues of POD is the reduction of data expressing their *essential information* by means of a few basis vectors. The problem of approximating all spatial coordinate vectors y_j of Y simultaneously by a single, normalized vector as well as possible can be expressed as

$$(\mathbf{P}^1) \quad \max_{u \in \mathbb{R}^m} \sum_{j=1}^n |\langle y_j, u \rangle_{\mathbb{R}^m}|^2 \quad \text{subject to (s.t.) } \|u\|_{\mathbb{R}^m}^2 = 1,$$

where $\|u\|_{\mathbb{R}^m} = \sqrt{\langle u, u \rangle_{\mathbb{R}^m}}$ for $u \in \mathbb{R}^m$.

Note that (\mathbf{P}^1) is a constrained optimization problem that can be solved by considering first-order necessary optimality conditions. We introduce the function $e : \mathbb{R}^m \rightarrow \mathbb{R}$ by $e(u) = 1 - \|u\|_{\mathbb{R}^m}^2$ for $u \in \mathbb{R}^m$. Then, the equality constraint in (\mathbf{P}^1) can be expressed as $e(u) = 0$. Notice that $\nabla e(u) = 2u^T$ is linear independent if $u \neq 0$ holds. In particular, a solution to (\mathbf{P}^1) satisfies $u \neq 0$. Thus, any solution to (\mathbf{P}^1) is a *regular point*. Let $\mathcal{L} : \mathbb{R}^m \times \mathbb{R} \rightarrow \mathbb{R}$ be the Lagrange functional associated with (\mathbf{P}^1) , i.e.,

$$\mathcal{L}(u, \lambda) = \sum_{j=1}^n |\langle y_j, u \rangle_{\mathbb{R}^m}|^2 + \lambda(1 - \|u\|_{\mathbb{R}^m}^2) \quad \text{for } (u, \lambda) \in \mathbb{R}^m \times \mathbb{R}.$$

Suppose that $u \in \mathbb{R}^m$ is a solution to (\mathbf{P}^1) . Since u is regular, there exists a Lagrange multiplier satisfying the first-order necessary optimality condition

$$\nabla \mathcal{L}(u, \lambda) \stackrel{!}{=} 0 \quad \text{in } \mathbb{R}^m \times \mathbb{R}.$$

We compute the gradient of \mathcal{L} with respect to u :

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial u_i}(u, \lambda) &= \frac{\partial}{\partial u_i} \left(\sum_{j=1}^n \left| \sum_{k=1}^m Y_{kj} u_k \right|^2 + \lambda \left(1 - \sum_{k=1}^m u_k^2 \right) \right) \\ &= 2 \sum_{j=1}^n \left(\sum_{k=1}^m Y_{kj} u_k \right) Y_{ij} - 2\lambda u_i \\ &= 2 \sum_{k=1}^m \left(\underbrace{\sum_{j=1}^n Y_{ij} Y_{jk}^T}_{=(YY^T)_{ik}} u_k \right) - 2\lambda u_i. \end{aligned}$$

Thus,

$$(1.6) \quad \nabla_u \mathcal{L}(u, \lambda) = 2(YY^T u - \lambda u) \stackrel{!}{=} 0 \quad \text{in } \mathbb{R}^m.$$

Equation (1.6) yields the eigenvalue problem

$$(1.7a) \quad YY^T u = \lambda u \quad \text{in } \mathbb{R}^m.$$

Notice that $YY^T \in \mathbb{R}^{m \times m}$ is a symmetric matrix satisfying

$$u^T (YY^T) u = (Y^T u)^T Y^T u = \|Y^T u\|_{\mathbb{R}^n}^2 \geq 0 \quad \text{for all } u \in \mathbb{R}^m.$$

Thus, YY^T is positive semi-definite. It follows that YY^T possesses m non-negative eigenvalues $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_m \geq 0$ and the corresponding eigenvectors can be chosen such that they are pairwise orthonormal.

From $\frac{\partial \mathcal{L}}{\partial \lambda}(u, \lambda) \stackrel{!}{=} 0$ in \mathbb{R} we infer the constraint

$$(1.7b) \quad \|u\|_{\mathbb{R}^m} = 1.$$

Due to SVD the vector u_1 solves (1.7) and

$$\begin{aligned} \sum_{j=1}^n |\langle y_j, u_1 \rangle_{\mathbb{R}^m}|^2 &= \sum_{j=1}^n \langle y_j, u_1 \rangle_{\mathbb{R}^m} \langle y_j, u_1 \rangle_{\mathbb{R}^m} = \sum_{j=1}^n \langle \langle y_j, u_1 \rangle_{\mathbb{R}^m} y_j, u_1 \rangle_{\mathbb{R}^m} \\ &= \left\langle \sum_{j=1}^n \langle y_j, u_1 \rangle_{\mathbb{R}^m} y_j, u_1 \right\rangle_{\mathbb{R}^m} = \left\langle \sum_{j=1}^n \left(\sum_{k=1}^m Y_{kj} (u_1)_k \right) y_j, u_1 \right\rangle_{\mathbb{R}^m} \\ &= \left\langle \sum_{k=1}^m \left(\sum_{j=1}^n Y_{\cdot,j} Y_{jk}^T (u_1)_k \right), u_1 \right\rangle_{\mathbb{R}^m} = \langle YY^T u_1, u_1 \rangle_{\mathbb{R}^m} \\ &= \lambda_1 \langle u_1, u_1 \rangle_{\mathbb{R}^m} = \lambda_1 \|u_1\|_{\mathbb{R}^m}^2 = \lambda_1. \end{aligned}$$

We next prove that u_1 solves (\mathbf{P}^1) . Suppose that $\tilde{u} \in \mathbb{R}^m$ is an arbitrary vector with $\|\tilde{u}\|_{\mathbb{R}^m} = 1$. Since $\{u_i\}_{i=1}^m$ is an orthonormal basis in \mathbb{R}^m , we have

$$\tilde{u} = \sum_{i=1}^m \langle \tilde{u}, u_i \rangle_{\mathbb{R}^m} u_i.$$

Thus,

$$\begin{aligned}
\sum_{j=1}^n |\langle y_j, \tilde{u} \rangle_{\mathbb{R}^m}|^2 &= \sum_{j=1}^n \left| \left\langle y_j, \sum_{i=1}^m \langle \tilde{u}, u_i \rangle_{\mathbb{R}^m} u_i \right\rangle_{\mathbb{R}^m} \right|^2 \\
&= \sum_{j=1}^n \sum_{i=1}^m \sum_{k=1}^m (\langle y_j, \langle \tilde{u}, u_i \rangle_{\mathbb{R}^m} u_i \rangle_{\mathbb{R}^m} \langle y_j, \langle \tilde{u}, u_k \rangle_{\mathbb{R}^m} u_k \rangle_{\mathbb{R}^m}) \\
&= \sum_{j=1}^n \sum_{i=1}^m \sum_{k=1}^m (\langle y_j, u_i \rangle_{\mathbb{R}^m} \langle y_j, u_k \rangle_{\mathbb{R}^m} \langle \tilde{u}, u_i \rangle_{\mathbb{R}^m} \langle \tilde{u}, u_k \rangle_{\mathbb{R}^m}) \\
&= \sum_{i=1}^m \sum_{k=1}^m \left(\underbrace{\left\langle \sum_{j=1}^n \langle y_j, u_i \rangle_{\mathbb{R}^m} y_j, u_k \right\rangle_{\mathbb{R}^m}}_{=\lambda_i u_i} \langle \tilde{u}, u_i \rangle_{\mathbb{R}^m} \langle \tilde{u}, u_k \rangle_{\mathbb{R}^m} \right) \\
&= \sum_{i=1}^m \sum_{k=1}^m \left(\underbrace{\langle \lambda_i u_i, u_k \rangle_{\mathbb{R}^m}}_{=\lambda_i \delta_{ik}} \langle \tilde{u}, u_i \rangle_{\mathbb{R}^m} \langle \tilde{u}, u_k \rangle_{\mathbb{R}^m} \right) \\
&= \sum_{i=1}^m \lambda_i |\langle \tilde{u}, u_i \rangle_{\mathbb{R}^m}|^2 \leq \lambda_1 \sum_{i=1}^m |\langle \tilde{u}, u_i \rangle_{\mathbb{R}^m}|^2 = \lambda_1 \|\tilde{u}\|_{\mathbb{R}^m}^2 = \lambda_1 \\
&= \sum_{j=1}^n |\langle y_j, u_1 \rangle_{\mathbb{R}^m}|^2.
\end{aligned}$$

Consequently, u_1 solves (\mathbf{P}^1) and $\operatorname{argmax}(\mathbf{P}^1) = \sigma_1^2 = \lambda_1$.

If we look for a second vector, orthogonal to u_1 that again describes the data set $\{y_i\}_{i=1}^n$ as well as possible then we need to solve

$$(\mathbf{P}^2) \quad \max_{u \in \mathbb{R}^m} \sum_{j=1}^n |\langle y_j, u \rangle_{\mathbb{R}^m}|^2 \quad \text{s.t.} \quad \|u\|_{\mathbb{R}^m} = 1 \text{ and } \langle u, u_1 \rangle_{\mathbb{R}^m} = 0.$$

SVD implies that u_2 is a solution to (\mathbf{P}^2) and $\operatorname{argmax}(\mathbf{P}^2) = \sigma_2^2 = \lambda_2$. In fact, u_2 solves the first-order necessary optimality conditions (1.7) and for

$$\tilde{u} = \sum_{i=2}^m \langle \tilde{u}, u_i \rangle_{\mathbb{R}^m} u_i \in \operatorname{span}\{u_1\}^\perp$$

we have

$$\sum_{j=1}^n |\langle y_j, \tilde{u} \rangle_{\mathbb{R}^m}|^2 \leq \lambda_2 = \sum_{j=1}^n |\langle y_j, u_2 \rangle_{\mathbb{R}^m}|^2.$$

Clearly this procedure can be continued by finite induction. We summarize our results in the following theorem.

Theorem 1.1. *Let $Y = [y_1, \dots, y_n] \in \mathbb{R}^{m \times n}$ be a given matrix with rank $d \leq \min\{m, n\}$. Further, let $Y = U\Sigma V^T$ be the singular value decomposition of Y , where $U = [u_1, \dots, u_m] \in \mathbb{R}^{m \times m}$, $V = [v_1, \dots, v_n] \in \mathbb{R}^{n \times n}$ are orthogonal matrices and the matrix $\Sigma \in \mathbb{R}^{m \times n}$ has the form as (1.2). Then, for any $\ell \in \{1, \dots, d\}$ the solution to*

$$(\mathbf{P}^\ell) \quad \max_{\tilde{u}_1, \dots, \tilde{u}_\ell \in \mathbb{R}^m} \sum_{i=1}^{\ell} \sum_{j=1}^n |\langle y_j, \tilde{u}_i \rangle_{\mathbb{R}^m}|^2 \quad \text{s.t.} \quad \langle \tilde{u}_i, \tilde{u}_j \rangle_{\mathbb{R}^m} = \delta_{ij} \text{ for } 1 \leq i, j \leq \ell$$

is given by the singular vectors $\{u_i\}_{i=1}^\ell$, i.e., by the first ℓ columns of U . Moreover,

$$(1.8) \quad \operatorname{argmax}(\mathbf{P}^\ell) = \sum_{i=1}^{\ell} \sigma_i^2 = \sum_{i=1}^{\ell} \lambda_i.$$

Proof. Since (\mathbf{P}^ℓ) is an equality constrained optimization problem, we introduce the Lagrangian

$$\mathcal{L} : \underbrace{\mathbb{R}^m \times \dots \times \mathbb{R}^m}_{\ell\text{-times}} \times \mathbb{R}^{\ell \times \ell}$$

by

$$\mathcal{L}(\psi_1, \dots, \psi_\ell, \Lambda) = \sum_{i=1}^{\ell} \sum_{j=1}^n |\langle y_j, \psi_i \rangle_{\mathbb{R}^m}|^2 + \sum_{i,j=1}^{\ell} \lambda_{ij} (\delta_{ij} - \langle \psi_i, \psi_j \rangle_{\mathbb{R}^m})$$

for $\psi_1, \dots, \psi_\ell \in \mathbb{R}^m$ and $\Lambda = ((\lambda_{ij})) \in \mathbb{R}^{\ell \times \ell}$. First-order necessary optimality conditions for (\mathbf{P}^ℓ) are given by

$$(1.9) \quad \frac{\partial \mathcal{L}}{\partial \psi_k}(\psi_1, \dots, \psi_\ell, \Lambda) \delta \psi_k = 0 \quad \text{for all } \delta \psi_k \in \mathbb{R}^m \text{ and } k \in \{1, \dots, \ell\}.$$

From

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial \psi_k}(\psi_1, \dots, \psi_\ell, \Lambda) \delta \psi_k &= 2 \sum_{i=1}^{\ell} \sum_{j=1}^n \langle y_j, \psi_i \rangle_{\mathbb{R}^m} \langle y_j, \delta \psi_k \rangle_{\mathbb{R}^m} \delta_{ik} \\ &\quad - \sum_{i,j=1}^{\ell} \lambda_{ij} \langle \psi_i, \delta \psi_k \rangle_{\mathbb{R}^m} \delta_{jk} - \sum_{i,j=1}^{\ell} \lambda_{ij} \langle \delta \psi_k, \psi_j \rangle_{\mathbb{R}^m} \delta_{ki} \\ &= 2 \sum_{j=1}^n \langle y_j, \psi_k \rangle_{\mathbb{R}^m} \langle y_j, \delta \psi_k \rangle_{\mathbb{R}^m} - \sum_{i=1}^{\ell} (\lambda_{ik} + \lambda_{ki}) \langle \psi_i, \delta \psi_k \rangle_{\mathbb{R}^m} \\ &= \left\langle 2 \sum_{j=1}^n \langle y_j, \psi_k \rangle_{\mathbb{R}^m} y_j - \sum_{i=1}^{\ell} (\lambda_{ik} + \lambda_{ki}) \psi_i, \delta \psi_k \right\rangle_{\mathbb{R}^m} \end{aligned}$$

and (1.9) we infer that

$$(1.10) \quad \sum_{j=1}^n \langle y_j, \psi_k \rangle_{\mathbb{R}^m} y_j = \frac{1}{2} \sum_{i=1}^{\ell} (\lambda_{ik} + \lambda_{ki}) \psi_i \quad \text{in } \mathbb{R}^m \text{ and for all } k \in \{1, \dots, \ell\}.$$

Note that

$$YY^T \psi = \sum_{j=1}^n \langle y_j, \psi \rangle_{\mathbb{R}^m} y_j \quad \text{for } \psi \in \mathbb{R}^m.$$

Thus, condition (1.10) can be expressed as

$$(1.11) \quad YY^T \psi_k = \frac{1}{2} \sum_{i=1}^{\ell} (\lambda_{ik} + \lambda_{ki}) \psi_i \quad \text{in } \mathbb{R}^m \text{ and for all } k \in \{1, \dots, \ell\}.$$

Now we proceed by induction. For $\ell = 1$ we have $k = 1$. It follows from (1.11) that

$$(1.12) \quad YY^T \psi_1 = \lambda_1 \psi_1 \quad \text{in } \mathbb{R}^m$$

with $\lambda_1 = \lambda_{11}$. Next we suppose that for $\ell \geq 1$ the first-order optimality conditions are given by

$$(1.13) \quad YY^T \psi_k = \lambda_k \psi_k \quad \text{in } \mathbb{R}^m \text{ and for all } k \in \{1, \dots, \ell\}.$$

We want to show that the first-order necessary optimality conditions for a POD basis $\{\psi_i\}_{i=1}^{\ell+1}$ of rank $\ell + 1$ are given by

$$(1.14) \quad YY^T \psi_k = \lambda_k \psi_k \quad \text{in } \mathbb{R}^m \text{ and for all } k \in \{1, \dots, \ell + 1\}.$$

By assumption we have (1.13). Thus, we only have to prove that

$$(1.15) \quad YY^T \psi_{\ell+1} = \lambda_{\ell+1} \psi_{\ell+1} \quad \text{in } \mathbb{R}^m.$$

Due to (1.11) we have

$$(1.16) \quad YY^T \psi_{\ell+1} = \frac{1}{2} \sum_{i=1}^{\ell+1} (\lambda_{i,\ell+1} + \lambda_{\ell+1,i}) \psi_i \quad \text{in } \mathbb{R}^m.$$

Since $\{\psi_i\}_{i=1}^{\ell+1}$ is a POD basis we have $\langle \psi_{\ell+1}, \psi_j \rangle_{\mathbb{R}^m} = 0$ for $1 \leq j \leq \ell$. Using (1.13) and the symmetry of YY^T we have for any $j \in \{1, \dots, \ell\}$

$$\begin{aligned} 0 &= \lambda_j \langle \psi_{\ell+1}, \psi_j \rangle_{\mathbb{R}^m} = \langle \psi_{\ell+1}, YY^T \psi_j \rangle_{\mathbb{R}^m} = \langle YY^T \psi_{\ell+1}, \psi_j \rangle_{\mathbb{R}^m} \\ &= \frac{1}{2} \sum_{i=1}^{\ell+1} (\lambda_{i,\ell+1} + \lambda_{\ell+1,i}) \langle \psi_i, \psi_j \rangle_{\mathbb{R}^m} = (\lambda_{j,\ell+1} + \lambda_{\ell+1,j}). \end{aligned}$$

This gives

$$(1.17) \quad \lambda_{\ell+1,i} = -\lambda_{i,\ell+1} \quad \text{for any } i \in \{1, \dots, \ell\}.$$

Inserting (1.17) into (1.16) we obtain

$$\begin{aligned} YY^T \psi_{\ell+1} &= \frac{1}{2} \sum_{i=1}^{\ell} (\lambda_{i,\ell+1} + \lambda_{\ell+1,i}) \psi_i + \lambda_{\ell+1,\ell+1} \psi_{\ell+1} \\ &= \frac{1}{2} \sum_{i=1}^{\ell} (\lambda_{i,\ell+1} - \lambda_{i,\ell+1}) \psi_i + \lambda_{\ell+1,\ell+1} \psi_{\ell+1} = \lambda_{\ell+1,\ell+1} \psi_{\ell+1}. \end{aligned}$$

Setting $\lambda_{\ell+1} = \lambda_{\ell+1,\ell+1}$ we obtain (1.15).

Summarizing, the necessary optimality conditions for (\mathbf{P}^ℓ) are given by the symmetric $m \times m$ eigenvalue problem

$$(1.18) \quad YY^T u_i = \lambda_i u_i \quad \text{for } i = 1, \dots, \ell.$$

It follows from SVD that $\{u_i\}_{i=1}^\ell$ solves (1.18). The proof that $\{u_i\}_{i=1}^\ell$ is a solution to (\mathbf{P}^ℓ) and that $\operatorname{argmax}(\mathbf{P}^\ell) = \sum_{i=1}^\ell \sigma_i^2$ holds is analogous to the proof for (\mathbf{P}^1) ; see Exercise 1.2). \square

Motivated by the previous theorem we give the next definition.

Definition 1.2. For $\ell \in \{1, \dots, d\}$ the vectors $\{u_i\}_{i=1}^\ell$ are called POD basis of rank ℓ .

The following result states that for every $\ell \leq d$ the approximation of the columns of Y by the first ℓ singular vectors $\{u_i\}_{i=1}^\ell$ is optimal in the mean among all rank ℓ approximations to the columns of Y .

Corollary 1.3 (Optimality of the POD basis). *Let all hypotheses of Theorem 1.1 be satisfied. Suppose that $\hat{U}^d \in \mathbb{R}^{m \times d}$ denotes a matrix with pairwise orthonormal vectors \hat{u}_i and that the expansion of the columns of Y in the basis $\{\hat{u}_i\}_{i=1}^d$ be given by*

$$Y = \hat{U}^d C^d, \quad \text{where } C_{ij}^d = \langle \hat{u}_i, y_j \rangle_{\mathbb{R}^m} \text{ for } 1 \leq i \leq d, 1 \leq j \leq n.$$

Then for every $\ell \in \{1, \dots, d\}$ we have

$$(1.19) \quad \|Y - U^\ell B^\ell\|_F \leq \|Y - \hat{U}^\ell C^\ell\|_F.$$

In (1.19), $\|\cdot\|_F$ denotes the Frobenius norm given by

$$\|A\|_F = \sqrt{\sum_{i=1}^m \sum_{j=1}^n |A_{ij}|^2} = \sqrt{\text{trace}(A^T A)} \quad \text{for } A \in \mathbb{R}^{m \times n},$$

the matrix U^ℓ denotes the first ℓ columns of U , B^ℓ the first ℓ rows of B and similarly for \hat{U}^ℓ and C^ℓ .

Remark 1.4. Notice that

$$\begin{aligned} \|Y - \hat{U}^\ell C^\ell\|_F^2 &= \sum_{i=1}^m \sum_{j=1}^n \left| Y_{ij} - \sum_{k=1}^{\ell} \hat{U}_{ik}^\ell C_{kj} \right|^2 = \sum_{j=1}^n \sum_{i=1}^m \left| Y_{ij} - \sum_{k=1}^{\ell} \langle \hat{u}_k, y_j \rangle_{\mathbb{R}^m} \hat{U}_{ik}^\ell \right|^2 \\ &= \sum_{j=1}^n \left\| y_j - \sum_{k=1}^{\ell} \langle y_j, \hat{u}_k \rangle_{\mathbb{R}^m} \hat{u}_k \right\|_{\mathbb{R}^m}^2. \end{aligned}$$

Analogously,

$$\|Y - U^\ell B^\ell\|_F^2 = \sum_{j=1}^n \left\| y_j - \sum_{k=1}^{\ell} \langle y_j, u_k \rangle_{\mathbb{R}^m} u_k \right\|_{\mathbb{R}^m}^2.$$

Thus, (1.19) implies that

$$\sum_{j=1}^n \left\| y_j - \sum_{k=1}^{\ell} \langle y_j, u_k \rangle_{\mathbb{R}^m} u_k \right\|_{\mathbb{R}^m}^2 \leq \sum_{j=1}^n \left\| y_j - \sum_{k=1}^{\ell} \langle y_j, \hat{u}_k \rangle_{\mathbb{R}^m} \hat{u}_k \right\|_{\mathbb{R}^m}^2$$

for any other set $\{\hat{u}_i\}_{i=1}^{\ell}$ of ℓ pairwise orthonormal vectors. Hence, the POD basis of rank ℓ can also be determined by solving

$$(1.20) \quad \min_{\tilde{u}_1, \dots, \tilde{u}_\ell \in \mathbb{R}^m} \sum_{j=1}^n \left\| y_j - \sum_{i=1}^{\ell} \langle y_j, \tilde{u}_i \rangle_{\mathbb{R}^m} \tilde{u}_i \right\|_{\mathbb{R}^m}^2 \quad \text{s.t. } \langle \tilde{u}_i, \tilde{u}_j \rangle_{\mathbb{R}^m} = \delta_{ij}, \quad 1 \leq i, j \leq \ell.$$

◇

Proof of Corollary 1.3. Note that (see Exercise 1.3))

$$\|Y - \hat{U}^\ell C^\ell\|_F^2 = \|\hat{U}^d (C^d - C_0^\ell)\|_F^2 = \|C^d - C_0^\ell\|_F^2 = \sum_{i=\ell+1}^d \sum_{j=1}^n |C_{ij}^d|^2,$$

where $C_0^\ell \in \mathbb{R}^{d \times n}$ results from $C \in \mathbb{R}^{d \times n}$ by replacing the last $d - \ell$ rows by 0. Similarly,

$$\begin{aligned}
\|Y - U^\ell B^\ell\|_F^2 &= \|U^k(B^d - B_0^\ell)\|_F^2 = \|B^d - B_0^\ell\|_F^2 = \sum_{i=\ell+1}^d \sum_{j=1}^n |B_{ij}^d|^2 \\
(1.21) \quad &= \sum_{i=\ell+1}^d \sum_{j=1}^n |\langle y_j, u_i \rangle_{\mathbb{R}^m}|^2 \\
&= \sum_{i=\ell+1}^d \sum_{j=1}^n \langle \langle y_j, u_i \rangle_{\mathbb{R}^m} y_j, u_i \rangle_{\mathbb{R}^m} = \sum_{i=\ell+1}^d \langle Y Y^T u_i, u_i \rangle_{\mathbb{R}^m} \\
&= \sum_{i=\ell+1}^d \sigma_i^2,
\end{aligned}$$

By Theorem 1.1 the vectors u_1, \dots, u_ℓ solve (\mathbf{P}^ℓ) . From (1.21),

$$\|Y\|_F^2 = \|\hat{U}^d C^d\|_F^2 = \|C^d\|_F^2 = \sum_{i=1}^d \sum_{j=1}^n |C_{ij}^d|^2$$

and

$$\|Y\|_F^2 = \|U^d B^d\|_F^2 = \|B^d\|_F^2 = \sum_{i=1}^d \sum_{j=1}^n |B_{ij}^d|^2 = \sum_{i=1}^d \sigma_i^2$$

we infer that

$$\begin{aligned}
\|Y - U^\ell B^\ell\|_F^2 &= \sum_{i=\ell+1}^d \sigma_i^2 = \sum_{i=1}^d \sigma_i^2 - \sum_{i=1}^{\ell} \sigma_i^2 = \|Y\|_F^2 - \sum_{i=1}^{\ell} \sum_{j=1}^n |\langle y_j, u_i \rangle_{\mathbb{R}^m}|^2 \\
&\leq \|Y\|_F^2 - \sum_{i=1}^{\ell} \sum_{j=1}^n |\langle y_j, \hat{u}_i \rangle_{\mathbb{R}^m}|^2 = \sum_{i=1}^d \sum_{j=1}^n |C_{ij}^d|^2 - \sum_{i=1}^{\ell} \sum_{j=1}^n |C_{ij}^d|^2 \\
&= \sum_{i=\ell+1}^d \sum_{j=1}^n |C_{ij}^d|^2 = \|Y - \hat{U}^\ell C^\ell\|_F^2,
\end{aligned}$$

which gives (1.19). \square

Remark 1.5. It follows from Corollary 1.3 that the POD basis of rank ℓ is optimal in the sense of representing in the mean the columns $\{y_j\}_{j=1}^n$ of Y as a linear combination by an orthonormal basis of rank ℓ :

$$\sum_{i=1}^{\ell} \sum_{j=1}^n |\langle y_j, u_i \rangle_{\mathbb{R}^m}|^2 = \sum_{i=1}^{\ell} \sigma_i^2 = \sum_{i=1}^{\ell} \lambda_i \geq \sum_{i=1}^{\ell} \sum_{j=1}^n |\langle y_j, \hat{u}_i \rangle_{\mathbb{R}^m}|^2$$

for any other set of orthonormal vectors $\{\hat{u}_i\}_{i=1}^{\ell}$. \diamond

The next corollary states that the POD coefficients are uncorrelated.

Corollary 1.6 (Uncorrelated POD coefficients). *Let all hypotheses of Theorem 1.1 hold. Then.*

$$\sum_{j=1}^n \langle y_j, u_i \rangle_{\mathbb{R}^m} \langle y_j, u_k \rangle_{\mathbb{R}^m} = \sum_{j=1}^n B_{ij}^\ell B_{kj}^\ell = \sigma_i^2 \delta_{ik} \quad \text{for } 1 \leq i, k \leq \ell.$$

Proof. The claim follows from (1.18) and $\langle u_i, u_k \rangle_{\mathbb{R}^m} = \delta_{ik}$ for $1 \leq i, k \leq \ell$:

$$\sum_{j=1}^n \langle y_j, u_i \rangle_{\mathbb{R}^m} \langle y_j, u_k \rangle_{\mathbb{R}^m} = \left\langle \underbrace{\sum_{j=1}^n \langle y_j, u_i \rangle_{\mathbb{R}^m} y_j}_{=Y Y^T u_i}, u_k \right\rangle_{\mathbb{R}^m} = \langle \sigma_i^2 u_i, u_k \rangle_{\mathbb{R}^m} = \sigma_i^2 \delta_{ik}.$$

□

Next we turn to the practical computation of a POD-basis of rank ℓ . If $n < m$ then one can determine the POD basis of rank ℓ as follows: Compute the eigenvectors $v_1, \dots, v_\ell \in \mathbb{R}^n$ by solving the symmetric $n \times n$ eigenvalue problem

$$(1.22) \quad Y^T Y v_i = \lambda_i v_i \quad \text{for } i = 1, \dots, \ell$$

and set, by (1.3),

$$u_i = \frac{1}{\sqrt{\lambda_i}} Y v_i \quad \text{for } i = 1, \dots, \ell.$$

For historical reasons [13] this method of determining the POD-basis is sometimes called the *method of snapshots*. On the other hand, if $m < n$ holds, we can obtain the POD basis by solving the $m \times m$ eigenvalue problem (1.18).

For the application of POD to concrete problems the choice of ℓ is certainly of central importance for applying POD. It appears that no general a-priori rules are available. Rather the choice of ℓ is based on heuristic considerations combined with observing the ratio of the modeled to the total energy contained in the system Y , which is expressed by

$$\mathcal{E}(\ell) = \frac{\sum_{i=1}^{\ell} \lambda_i}{\sum_{i=1}^d \lambda_i}.$$

Let us mention that POD is also called *Principal Component Analysis* (PCA) and *Karhunen-Loève Decomposition*.

1.2. The POD method with a weighted inner product. Let us endow the Euclidean space \mathbb{R}^m with the weighted inner product

$$(1.23) \quad \langle u, \tilde{u} \rangle_W = u^T W \tilde{u} = \langle u, W \tilde{u} \rangle_{\mathbb{R}^m} = \langle W u, \tilde{u} \rangle_{\mathbb{R}^m} \quad \text{for } u, \tilde{u} \in \mathbb{R}^m,$$

where $W \in \mathbb{R}^{m \times m}$ is a symmetric, positive-definite matrix. Furthermore, let $\|u\|_W = \sqrt{\langle u, u \rangle_W}$ for $u \in \mathbb{R}^m$ be the associated induced norm. For the choice $W = I$, the inner product (1.23) coincides the Euclidean inner product.

Example 1.7. Let us motivate the weighted inner product by an example. Suppose that $\Omega = (0, 1) \subset \mathbb{R}$ holds. We consider the space $L^2(\Omega)$ of square integrable functions on Ω :

$$L^2(\Omega) = \left\{ \varphi : \Omega \rightarrow \mathbb{R} \mid \int_{\Omega} |\varphi|^2 dx < \infty \right\}.$$

Recall that $L^2(\Omega)$ is a Hilbert space endowed with the inner product

$$\langle \varphi, \tilde{\varphi} \rangle_{L^2(\Omega)} = \int_{\Omega} \varphi \tilde{\varphi} dx \quad \text{for } \varphi, \tilde{\varphi} \in L^2(\Omega)$$

and the induced norm $\|\varphi\|_{L^2(\Omega)} = \sqrt{\langle \varphi, \varphi \rangle_{L^2(\Omega)}}$ for $\varphi \in L^2(\Omega)$. For the step size $h = 1/(m-1)$ let us introduce a spatial grid in Ω by

$$x_i = (i-1)h \quad \text{for } i = 1, \dots, m.$$

For any $\varphi, \tilde{\varphi} \in L^2(\Omega)$ we introduce a discrete inner product by trapezoidal approximation:

$$(1.24) \quad \langle \varphi, \tilde{\varphi} \rangle_{L_h^2(\Omega)} = h \left(\frac{\varphi_1^h \tilde{\varphi}_1^h}{2} + \sum_{i=2}^{m-1} (\varphi_i^h \tilde{\varphi}_i^h) + \frac{\varphi_m^h \tilde{\varphi}_m^h}{2} \right),$$

where

$$\varphi_i^h = \begin{cases} \frac{2}{h} \int_0^{h/2} \varphi(x) \, dx & \text{for } i = 1, \\ \frac{1}{h} \int_{x_i-h/2}^{x_i+h/2} \varphi(x) \, dx & \text{for } i = 2, \dots, m-1, \\ \frac{2}{h} \int_{1-h/2}^1 \varphi(x) \, dx & \text{for } i = m \end{cases}$$

and the $\tilde{\varphi}_i^h$'s are defined analogously. Setting $W = \text{diag}(h/2, h, \dots, h, h/2) \in \mathbb{R}^{m \times m}$, $\varphi^h = (\varphi_1^h, \dots, \varphi_m^h)^T \in \mathbb{R}^m$ and $\tilde{\varphi}^h = (\tilde{\varphi}_1^h, \dots, \tilde{\varphi}_m^h)^T \in \mathbb{R}^m$ we find

$$\langle \varphi, \tilde{\varphi} \rangle_{L_h^2(\Omega)} = \langle \varphi^h, \tilde{\varphi}^h \rangle_W = (\varphi^h)^T W \tilde{\varphi}^h.$$

Thus, the discrete L^2 -inner product can be written as a weighted inner product of the form (1.23). \diamond

Now we replace (\mathbf{P}^1) by

$$(\mathbf{P}_W^1) \quad \max_{u \in \mathbb{R}^m} \sum_{j=1}^n |\langle y_j, u \rangle_W|^2 \quad \text{s.t.} \quad \|u\|_W = 1.$$

Analogously to Section 1.1 we treat (\mathbf{P}_W^1) as an equality constrained optimization problem. The Lagrangian $\mathcal{L} : \mathbb{R}^m \times \mathbb{R} \rightarrow \mathbb{R}$ for (\mathbf{P}_W^1) is given by

$$\mathcal{L}(u, \lambda) = \sum_{j=1}^n |\langle y_j, u \rangle_W|^2 + \lambda(1 - \|u\|_W^2) \quad \text{for } (u, \lambda) \in \mathbb{R}^m \times \mathbb{R}.$$

Suppose that $u \in \mathbb{R}^m$ is a solution to (\mathbf{P}_W^1) . Then, a first-order necessary optimality condition is given by

$$\nabla \mathcal{L}(u, \lambda) \stackrel{!}{=} 0 \quad \text{in } \mathbb{R}^m \times \mathbb{R}.$$

We compute the gradient of \mathcal{L} with respect to u : Since W is symmetric, we derive

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial u_i}(u, \lambda) &= \frac{\partial}{\partial u_i} \left(\sum_{j=1}^n \left| \sum_{k=1}^m \sum_{\nu=1}^m Y_{j\nu}^T W_{\nu k} u_k \right|^2 + \lambda \left(1 - \sum_{k=1}^m \sum_{\nu=1}^m u_\nu W_{\nu k} u_k \right) \right) \\ &= 2 \sum_{j=1}^n \left(\sum_{k=1}^m \sum_{\nu=1}^m Y_{j\nu}^T W_{\nu k} u_k \right) \left(\sum_{\mu=1}^m Y_{j\mu}^T W_{\mu i} \right) \\ &\quad - \lambda \left(\sum_{\nu=1}^m u_\nu W_{\nu i} + \sum_{k=1}^m W_{ik} u_k \right) \\ &= 2 \sum_{k=1}^m \sum_{\nu=1}^m \sum_{\mu=1}^m W_{i\mu} \sum_{j=1}^n Y_{\mu j} Y_{j\nu}^T W_{\nu k} u_k - 2\lambda \left(\sum_{k=1}^m W_{ik} u_k \right) \\ &= 2 \left(W Y Y^T W u - \lambda W u \right)_i. \end{aligned}$$

Thus,

$$(1.25) \quad \nabla_u \mathcal{L}(u, \lambda) = 2(WY Y^T W u - \lambda W u) \stackrel{!}{=} 0 \quad \text{in } \mathbb{R}^m.$$

Equation (1.25) yields the generalized eigenvalue problem

$$(1.26) \quad (WY)(WY)^T u = \lambda W u.$$

Since W is symmetric and positive definite, W possesses an eigenvalue decomposition of the form $W = Q D Q^T$, where $D = \text{diag}(\eta_1, \dots, \eta_m)$ contains the eigenvalues $\eta_1 \geq \dots \geq \eta_m > 0$ of W and $Q \in \mathbb{R}^{m \times m}$ is an orthogonal matrix. We define

$$W^\alpha = Q \text{diag}(\eta_1^\alpha, \dots, \eta_m^\alpha) Q^T \quad \text{for } \alpha \in \mathbb{R}.$$

Note that $(W^\alpha)^{-1} = W^{-\alpha}$ and $W^{\alpha+\beta} = W^\alpha W^\beta$ for $\alpha, \beta \in \mathbb{R}$; see Exercise 1.4). Moreover, we have

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

and $\|u\|_W = \|W^{1/2} u\|_{\mathbb{R}^m}$ for $u \in \mathbb{R}^m$.

Setting $\bar{u} = W^{1/2} u \in \mathbb{R}^m$ and $\bar{Y} = W^{1/2} Y \in \mathbb{R}^{m \times n}$ and multiplying (1.26) by $W^{-1/2}$ from the left we deduce the symmetric, $m \times m$ eigenvalue problem

$$(1.27a) \quad \bar{Y} \bar{Y}^T \bar{u} = \lambda \bar{u} \quad \text{in } \mathbb{R}^m.$$

From $\frac{\partial \mathcal{L}}{\partial \lambda}(u, \lambda) \stackrel{!}{=} 0$ in \mathbb{R} we infer the constraint $\|u\|_W = 1$ that can be expressed as

$$(1.27b) \quad \|\bar{u}\|_{\mathbb{R}^m} = 1.$$

Thus, the first-order optimality conditions (1.27) for (\mathbf{P}_W^1) are — as for (\mathbf{P}^1) (compare (1.7)) — an $m \times m$ eigenvalue problem, but the matrix Y as well as the vector u have to be weighted by the matrix $W^{1/2}$.

It can be shown (see Exercise 1.4.1)) that

$$u_1 = W^{-1/2} \bar{u}_1$$

solves (\mathbf{P}_W^1) , where \bar{u}_1 is an eigenvector of $\bar{Y} \bar{Y}^T$ corresponding to the largest eigenvalue λ_1 with $\|\bar{u}_1\|_{\mathbb{R}^m} = 1$. Due to SVD the vector u_1 can be also determined by solving the symmetric $n \times n$ eigenvalue problem

$$\bar{Y}^T \bar{Y} \bar{v}_1 = \lambda_1 \bar{v}_1$$

where $\bar{Y}^T \bar{Y} = Y^T W Y$, and setting

$$(1.28) \quad u_1 = W^{-1/2} \bar{u}_1 = \frac{1}{\sqrt{\lambda_1}} W^{-1/2} \bar{Y} \bar{v}_1 = \frac{1}{\sqrt{\lambda_1}} Y \bar{v}_1.$$

As in Section 1.1 we can continue by looking at a second vector $u \in \mathbb{R}^m$ with $\langle u, u_1 \rangle_W = 0$ that maximizes $\sum_{j=1}^n |\langle y_j, u \rangle_W|^2$. Let us generalize Theorem 1.1 as follows.

Theorem 1.8. *Let $Y \in \mathbb{R}^{m \times n}$ be a given matrix with rank $d \leq \min\{m, n\}$, W a symmetric, positive definite matrix, $\bar{Y} = W^{1/2} Y$ and $\ell \in \{1, \dots, d\}$. Further, let $\bar{Y} = \bar{U} \Sigma \bar{V}^T$ be the singular value decomposition of \bar{Y} , where $\bar{U} = [\bar{u}_1, \dots, \bar{u}_m] \in \mathbb{R}^{m \times m}$, $\bar{V} = [\bar{v}_1, \dots, \bar{v}_n] \in \mathbb{R}^{n \times n}$ are orthogonal matrices and the matrix Σ has the form*

$$\bar{U}^T \bar{Y} \bar{V} = \begin{pmatrix} D & 0 \\ 0 & 0 \end{pmatrix} = \Sigma \in \mathbb{R}^{m \times n}.$$

Then the solution to

$$(\mathbf{P}_W^\ell) \quad \max_{\tilde{u}_1, \dots, \tilde{u}_\ell \in \mathbb{R}^m} \sum_{i=1}^{\ell} \sum_{j=1}^n | \langle y_j, \tilde{u}_i \rangle_W |^2 \quad \text{s.t.} \quad \langle \tilde{u}_i, \tilde{u}_j \rangle_W = \delta_{ij} \text{ for } 1 \leq i, j \leq \ell$$

is given by the vectors $u_i = W^{-1/2} \tilde{u}_i$, $i = 1, \dots, \ell$. Moreover,

$$(1.29) \quad \operatorname{argmax}(\mathbf{P}_W^\ell) = \sum_{i=1}^{\ell} \sigma_i^2 = \sum_{i=1}^{\ell} \lambda_i.$$

Proof. Using similar arguments as in the proof of Theorem 1.1 one can prove that $\{u_i\}_{i=1}^{\ell}$ solves (\mathbf{P}_W^ℓ) ; see Exercice 1.4. \square

Remark 1.9. Due to SVD and $\bar{Y}^T \bar{Y} = Y^T W Y$ the POD basis $\{u_i\}_{i=1}^{\ell}$ of rank ℓ can be determined by the method of snapshots as follows: Solve the symmetric $n \times n$ eigenvalue problem

$$Y^T W Y \bar{v}_i = \lambda_i \bar{v}_i \quad \text{for } i = 1, \dots, \ell,$$

and set

$$u_i = W^{-1/2} \tilde{u}_i = \frac{1}{\sqrt{\lambda_i}} W^{-1/2} (\bar{Y} \bar{v}_i) = \frac{1}{\sqrt{\lambda_i}} W^{-1/2} W^{1/2} Y \bar{v}_i = \frac{1}{\sqrt{\lambda_i}} Y \bar{v}_i$$

for $i = 1, \dots, \ell$. Notice that

$$\langle u_i, u_j \rangle_W = u_i^T W u_j = \frac{\delta_{ij} \lambda_j}{\sqrt{\lambda_i \lambda_j}} \quad \text{for } 1 \leq i, j \leq \ell.$$

For $m \gg n$ the method of snapshots turns out to be faster than computing the POD basis via (1.27). Notice that the matrix $W^{1/2}$ is also not required for the method of snapshots. \diamond

1.3. Application to time-dependent systems. For $T > 0$ we consider the semi-linear initial value problem

$$(1.30a) \quad \dot{y}(t) = Ay(t) + f(t, y(t)) \quad \text{for } t \in (0, T],$$

$$(1.30b) \quad y(0) = y_0,$$

where $y_0 \in \mathbb{R}^m$ is a chosen initial condition, $A \in \mathbb{R}^{m \times m}$ is a given matrix, $f : [0, T] \times \mathbb{R}^m \rightarrow \mathbb{R}^m$ is continuous in both arguments and locally Lipschitz-continuous with respect to the second argument. It is well known that (1.30) has a unique (classical) solution $y \in C^1(0, T; \mathbb{R}^m) \cap C([0, T]; \mathbb{R}^m)$ given by the implicit integral representation

$$y(t) = e^{tA} y_0 + \int_0^t e^{(t-s)A} f(s, y(s)) ds$$

with $e^{tA} = \sum_{i=0}^{\infty} t^i A^i / (i!)$. Let $0 \leq t_1 < t_2 < \dots < t_n \leq T$ be a given time grid in the interval $[0, T]$. For simplicity of the presentation, the time grid is assumed to be equidistant with step-size $\Delta t = T/(n-1)$, i.e., $t_j = (j-1)\Delta t$. We suppose that we know the solution to (1.30) at the given time instances t_j , $j \in \{1, \dots, n\}$. Our goal is to determine a POD basis of rank $\ell \leq n$ that describes the ensemble

$$y_j = y(t_j) = e^{t_j A} y_0 + \int_0^{t_j} e^{(t_j-s)A} f(s, y(s)) ds, \quad j = 1, \dots, n,$$

as well as possible with respect to the weighted inner product:

$$(\hat{\mathbf{P}}_W^{n,\ell}) \quad \min_{\tilde{u}_1, \dots, \tilde{u}_\ell \in \mathbb{R}^m} \sum_{j=1}^n \alpha_j \left\| y_j - \sum_{i=1}^{\ell} \langle y_j, \tilde{u}_i \rangle_W \tilde{u}_i \right\|_W^2 \quad \text{s.t.} \quad \langle \tilde{u}_i, \tilde{u}_j \rangle_W = \delta_{ij} \text{ for } 1 \leq i, j \leq \ell,$$

where the α_j 's denote non-negative weights which will be specified later on. Note that for $\alpha_j = 1$ for $j = 1, \dots, n$ and $W = I$ problem $(\hat{\mathbf{P}}_W^{n,\ell})$ coincides with (1.20).

Example 1.10. Let us consider the following one-dimensional heat equation:

$$(1.31a) \quad \theta_t(t, x) = \theta_{xx}(t, x) \quad \text{for all } (t, x) \in Q = (0, T) \times \Omega,$$

$$(1.31b) \quad \theta_x(t, 0) = \theta_x(t, 1) = 0 \quad \text{for all } t \in (0, T),$$

$$(1.31c) \quad \theta(0, x) = \theta_0(x) \quad \text{for all } x \in \Omega = (0, 1) \subseteq \mathbb{R},$$

where $\theta_0 \in C(\bar{\Omega})$ is a given initial condition. To solve (1.31) numerically we apply a classical finite difference approximation for the spatial variable x . In Example 1.7 we have introduced the spatial grid $\{x_i\}_{i=1}^m$ in the interval $[0, 1]$. Let us denote by $y_i : [0, T] \rightarrow \mathbb{R}$ the numerical approximation for $\theta(\cdot, x_i)$ for $i = 1, \dots, m$. The second partial derivative θ_{xx} in (1.31a) and the boundary conditions (1.31b) are discretized by centered difference quotients of second-order so that we obtain the following ordinary differential equations for the time-dependent functions y_i :

$$(1.32a) \quad \begin{cases} \dot{y}_1(t) = \frac{-2y_1(t) + 2y_2(t)}{h^2}, \\ \dot{y}_i(t) = \frac{y_{i-1}(t) - 2y_i(t) + y_{i+1}(t)}{h^2}, \quad i = 2, \dots, m-1, \\ \dot{y}_m(t) = \frac{-2y_m(t) + 2y_{m-1}(t)}{h^2} \end{cases}$$

for $t \in (0, T]$. From (1.31c) we infer the initial condition

$$(1.32b) \quad y_i(0) = \theta_0(x_i), \quad i = 1, \dots, m.$$

Introducing the matrix

$$A = \frac{1}{h^2} \begin{pmatrix} -2 & 2 & & & 0 \\ 1 & -2 & 1 & & \\ & \ddots & \ddots & \ddots & \\ & & & 1 & -2 & 1 \\ 0 & & & & 2 & -2 \end{pmatrix} \in \mathbb{R}^{m \times m}$$

and the vectors

$$y(t) = \begin{pmatrix} y_1(t) \\ \vdots \\ y_m(t) \end{pmatrix} \text{ for } t \in [0, T], \quad y_0 = \begin{pmatrix} \theta_0(x_1) \\ \vdots \\ \theta_0(x_m) \end{pmatrix} \in \mathbb{R}^m$$

we can express (1.32) in the form

$$(1.33) \quad \begin{aligned} \dot{y}(t) &= Ay(t) \quad \text{for } t \in (0, T], \\ y(0) &= y_0 \end{aligned}$$

Setting $f \equiv 0$ the linear initial-value problem coincides with (1.30). Note that now the vector $y(t)$, $t \in [0, T]$, represents a function in Ω evaluated at m grid points. Therefore, we should supply \mathbb{R}^m by a weighted inner product representing

a discretized inner product in an appropriate function space. Here we choose the inner product introduced in (1.24); see Example 1.7. Next we choose a time grid $\{t_j\}_{j=1}^n$ in the interval $[0, T]$ and define $y_j = y(t_j)$ for $j = 1, \dots, n$. If we are interested in finding a POD basis of rank $\ell \leq n$ that describes the ensemble $\{y_j\}_{j=1}^n$ as well as possible, we end up with $(\hat{\mathbf{P}}_W^{n, \ell})$. \diamond

To solve $(\hat{\mathbf{P}}_W^{n, \ell})$ we apply the techniques used in Sections 1.1 and 1.2, i.e., we use the Lagrangian framework. Thus, we introduce the Lagrange functional

$$\mathcal{L} : \underbrace{\mathbb{R}^m \times \dots \times \mathbb{R}^m}_{\ell\text{-times}} \times \mathbb{R}^{\ell \times \ell} \rightarrow \mathbb{R}$$

by

$$\mathcal{L}(u_1, \dots, u_\ell, \Lambda) = \sum_{j=1}^n \alpha_j \left\| y_j - \sum_{i=1}^{\ell} \langle y_j, u_i \rangle_W u_i \right\|_W^2 + \sum_{i=1}^{\ell} \sum_{j=1}^{\ell} \Lambda_{ij} (1 - \langle u_i, u_j \rangle_W)$$

for $u_1, \dots, u_\ell \in \mathbb{R}^m$ and $\Lambda \in \mathbb{R}^{\ell \times \ell}$ with elements Λ_{ij} , $1 \leq i, j \leq \ell$. It turns out that the solution to $(\hat{\mathbf{P}}_W^{n, \ell})$ is given by the first-order necessary optimality conditions

$$(1.34a) \quad \nabla_{u_i} \mathcal{L}(u_1, \dots, u_\ell, \Lambda) \stackrel{!}{=} 0 \quad \text{in } \mathbb{R}^m, \quad 1 \leq i \leq \ell,$$

and

$$(1.34b) \quad \langle u_i, u_j \rangle_W \stackrel{!}{=} \delta_{ij}, \quad 1 \leq i, j \leq \ell.$$

From (1.34a) we derive

$$(1.35) \quad YDY^T W u_i = \lambda_i u_i \quad \text{for } i = 1, \dots, \ell,$$

where $D = \text{diag}(\alpha_1, \dots, \alpha_n) \in \mathbb{R}^{n \times n}$. Inserting $u_i = W^{-1/2} \bar{u}_i$ in (1.35) and multiplying (1.35) by $W^{1/2}$ from the left yield

$$(1.36a) \quad W^{1/2} Y D Y^T W^{1/2} \bar{u}_i = \lambda_i \bar{u}_i.$$

From (1.34b) we find

$$(1.36b) \quad \langle \bar{u}_i, \bar{u}_j \rangle_{\mathbb{R}^m} = \bar{u}_i^T \bar{u}_j = u_i^T W u_j = \langle u_i, u_j \rangle_W = \delta_{ij}, \quad 1 \leq i, j \leq \ell.$$

Setting $\bar{Y} = W^{1/2} Y D^{1/2} \in \mathbb{R}^{m \times n}$ and using $W^T = W$ as well as $D^T = D$ we infer from (1.36) that the solution $\{u_i\}_{i=1}^{\ell}$ to $(\hat{\mathbf{P}}_W^{n, \ell})$ is given by the symmetric $m \times m$ eigenvalue problem

$$\bar{Y} \bar{Y}^T \bar{u}_i = \lambda_i \bar{u}_i, \quad 1 \leq i \leq \ell \quad \text{and} \quad \langle \bar{u}_i, \bar{u}_j \rangle_{\mathbb{R}^m} = \delta_{ij}, \quad 1 \leq i, j \leq \ell.$$

Note that

$$\bar{Y}^T \bar{Y} = D^{1/2} Y^T W Y D^{1/2} \in \mathbb{R}^{n \times n}.$$

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

$$\bar{Y}^T \bar{Y} \bar{v}_i = \lambda_i \bar{v}_i, \quad 1 \leq i \leq \ell \quad \text{and} \quad \langle \bar{v}_i, \bar{v}_j \rangle_{\mathbb{R}^n} = \delta_{ij}, \quad 1 \leq i, j \leq \ell.$$

Then we set (by SVD)

$$u_i = W^{-1/2} \bar{u}_i = \frac{1}{\sqrt{\lambda_i}} W^{-1/2} \bar{Y} \bar{v}_i = \frac{1}{\sqrt{\lambda_i}} Y D^{1/2} \bar{v}_i, \quad 1 \leq i \leq \ell;$$

compare (1.28).

Note that

$$\langle u_i, u_j \rangle_W = u_i^T W u_j = \frac{1}{\sqrt{\lambda_i \lambda_j}} \bar{v}_i^T \underbrace{D^{1/2} Y^T W Y D^{1/2}}_{= \bar{Y}^T \bar{Y}} \bar{v}_j = \frac{\lambda_i}{\sqrt{\lambda_i \lambda_j}} \bar{v}_i^T \bar{v}_j = \frac{\lambda_i \delta_{ij}}{\sqrt{\lambda_i \lambda_j}}$$

for $1 \leq i, j \leq \ell$, i.e., the POD basis vectors u_1, \dots, u_ℓ are orthonormal in \mathbb{R}^m with respect to the inner product $\langle \cdot, \cdot \rangle_W$.

Of course, the snapshot ensemble $\{y_j\}_{j=1}^n$ for $(\hat{\mathbf{P}}_W^{n, \ell})$ and therefore the snapshot set $\text{span}\{y_1, \dots, y_n\}$ depend on the chosen time instances $\{t_j\}_{j=1}^n$. Consequently, the POD basis vectors $\{u_i\}_{i=1}^\ell$ and the corresponding eigenvalues $\{\lambda_i\}_{i=1}^\ell$ depend also on the time instances, i.e.,

$$u_i = u_i^n \quad \text{and} \quad \lambda_i = \lambda_i^n, \quad 1 \leq i \leq \ell.$$

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

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

To address these two questions we will introduce a *continuous version* of POD. Let $y : [0, T] \rightarrow \mathbb{R}^m$ be the unique solution to (1.30). If we are interested to find a POD basis of rank ℓ that describes the whole trajectory $\{y(t) \mid t \in [0, T]\} \subset \mathbb{R}^m$ as good as possible we have to consider the following minimization problem

$$\begin{aligned} (\hat{\mathbf{P}}_W^\ell) \quad & \min_{\tilde{u}_1, \dots, \tilde{u}_\ell \in \mathbb{R}^m} \int_0^T \left\| y(t) - \sum_{i=1}^{\ell} \langle y(t), \tilde{u}_i \rangle_W \tilde{u}_i \right\|_W^2 dt \\ & \text{s.t. } \langle \tilde{u}_i, \tilde{u}_j \rangle_W = \delta_{ij}, \quad 1 \leq i, j \leq \ell, \end{aligned}$$

To solve $(\hat{\mathbf{P}}_W^\ell)$ we use similar arguments as in Sections 1.1 and 1.2. For $\ell = 1$ we obtain instead of $(\hat{\mathbf{P}}_W^\ell)$ the minimization problem

$$(1.37) \quad \min_{\tilde{u} \in \mathbb{R}^m} \int_0^T \left\| y(t) - \langle y(t), \tilde{u} \rangle_W \tilde{u} \right\|_W^2 dt \quad \text{s.t.} \quad \|\tilde{u}\|_W^2 = 1,$$

Suppose that $\{\tilde{u}_i\}_{i=2}^m$ are chosen in such a way that $\{\tilde{u}, \tilde{u}_2, \dots, \tilde{u}_m\}$ is an orthonormal basis in \mathbb{R}^m with respect to the inner product $\langle \cdot, \cdot \rangle_W$. Then we have

$$y(t) = \langle y(t), \tilde{u} \rangle_W \tilde{u} + \sum_{i=2}^m \langle y(t), \tilde{u}_i \rangle_W \tilde{u}_i \quad \text{for all } t \in [0, T].$$

Thus,

$$\begin{aligned} \int_0^T \left\| y(t) - \langle y(t), \tilde{u} \rangle_W \tilde{u} \right\|_W^2 dt &= \int_0^T \left\| \sum_{i=2}^m \langle y(t), \tilde{u}_i \rangle_W \tilde{u}_i \right\|_W^2 dt \\ &= \sum_{i=2}^m \int_0^T |\langle y(t), \tilde{u}_i \rangle_W|^2 dt \end{aligned}$$

we conclude that (1.37) is equivalent with the following maximization problem

$$(1.38) \quad \max_{\tilde{u} \in \mathbb{R}^m} \int_0^T |\langle y(t), \tilde{u} \rangle_W|^2 dt \quad \text{s.t.} \quad \|\tilde{u}\|_W^2 = 1.$$

The Lagrange functional $\mathcal{L} : \mathbb{R}^m \times \mathbb{R} \rightarrow \mathbb{R}$ associated with (1.38) is given by

$$\mathcal{L}(u, \lambda) = \int_0^T |\langle y(t), u \rangle_W|^2 dt + \lambda(1 - \|u\|_W^2) \quad \text{for } (u, \lambda) \in \mathbb{R}^m \times \mathbb{R}.$$

First-order necessary optimality conditions are given by

$$\nabla \mathcal{L}(u, \lambda) \stackrel{!}{=} 0 \quad \text{in } \mathbb{R}^m \times \mathbb{R}.$$

Therefore, we compute the partial derivative of \mathcal{L} with respect to the i th component u_i of the vector u :

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial u_i}(u, \lambda) &= \frac{\partial}{\partial u_i} \left(\int_0^T \left| \sum_{k=1}^m \sum_{\nu=1}^m y_k(t) W_{k\nu} u_\nu \right|^2 dt + \lambda \left(1 - \sum_{k=1}^m \sum_{\nu=1}^m u_k W_{k\nu} u_\nu \right) \right) \\ &= 2 \int_0^T \left(\sum_{k=1}^m \sum_{\nu=1}^m y_k(t) W_{k\nu} u_\nu \right) \sum_{\mu=1}^m y_\mu(t) W_{\mu i} dt - 2\lambda \sum_{k=1}^m W_{ik} u_k \\ &= 2 \left(\int_0^T \langle y(t), u \rangle_W W y(t) dt - \lambda W u \right)_i \end{aligned}$$

for $i \in \{1, \dots, m\}$. Thus,

$$\nabla_u \mathcal{L}(u, \lambda) = 2 \left(\int_0^T \langle y(t), u \rangle_W W y(t) dt - \lambda W u \right) \stackrel{!}{=} 0 \quad \text{in } \mathbb{R}^m,$$

which gives

$$(1.39) \quad \int_0^T \langle y(t), u \rangle_W W y(t) dt = \lambda W u \quad \text{in } \mathbb{R}^m.$$

Multiplying (1.39) by W^{-1} from the left yields

$$(1.40) \quad \int_0^T \langle y(t), u \rangle_W y(t) dt = \lambda u \quad \text{in } \mathbb{R}^m.$$

We define the operator $\mathcal{R} : \mathbb{R}^m \rightarrow \mathbb{R}^m$ as

$$(1.41) \quad \mathcal{R}u = \int_0^T \langle y(t), u \rangle_W y(t) dt \quad \text{for } u \in \mathbb{R}^m.$$

Lemma 1.11. *The operator \mathcal{R} is linear and bounded (i.e., continuous). Moreover,*

1) \mathcal{R} is non-negative:

$$\langle \mathcal{R}u, u \rangle_W \geq 0 \quad \text{for all } u \in \mathbb{R}^m.$$

2) \mathcal{R} is self-adjoint (or symmetric):

$$\langle \mathcal{R}u, \tilde{u} \rangle_W = \langle u, \mathcal{R}\tilde{u} \rangle_W \quad \text{for all } u, \tilde{u} \in \mathbb{R}^m.$$

Proof. For arbitrary $u, \tilde{u} \in \mathbb{R}^m$ and $\alpha, \tilde{\alpha} \in \mathbb{R}$ we have

$$\begin{aligned} \mathcal{R}(\alpha u + \tilde{\alpha} \tilde{u}) &= \int_0^T \langle y(t), \alpha u + \tilde{\alpha} \tilde{u} \rangle_W y(t) dt \\ &= \int_0^T (\alpha \langle y(t), u \rangle_W + \tilde{\alpha} \langle y(t), \tilde{u} \rangle_W) y(t) dt \\ &= \alpha \int_0^T \langle y(t), u \rangle_W y(t) dt + \tilde{\alpha} \int_0^T \langle y(t), \tilde{u} \rangle_W y(t) dt = \alpha \mathcal{R}u + \tilde{\alpha} \mathcal{R}\tilde{u}, \end{aligned}$$

so that \mathcal{R} is linear. From the Cauchy-Schwarz inequality we derive

$$\begin{aligned} \|\mathcal{R}u\|_W &\leq \int_0^T \|\langle y(t), u \rangle_W y(t)\|_W dt = \int_0^T |\langle y(t), u \rangle_W| \|y(t)\|_W dt \\ &\leq \int_0^T \|y(t)\|_W^2 \|u\|_W dt = \left(\int_0^T \|y(t)\|_W^2 dt \right) \|u\|_W = \|y\|_{L^2(0,T;\mathbb{R}^m)}^2 \|u\|_W \end{aligned}$$

for an arbitrary $u \in \mathbb{R}^m$. Since $y \in C([0, T]; \mathbb{R}^m) \subset L^2(0, T; \mathbb{R}^m)$ holds, the norm $\|y\|_{L^2(0,T;\mathbb{R}^m)}$ is bounded. Therefore, \mathcal{R} is bounded. Since

$$\begin{aligned} \langle \mathcal{R}u, u \rangle_W &= \left(\int_0^T \langle y(t), u \rangle_W y(t) dt \right)^T W u = \int_0^T \langle y(t), u \rangle_W y(t)^T W u dt \\ &= \int_0^T |\langle y(t), u \rangle_W|^2 dt \geq 0 \end{aligned}$$

for all $u \in \mathbb{R}^m$ holds, \mathcal{R} is non-negative. Finally, we infer from

$$\begin{aligned} \langle \mathcal{R}u, \tilde{u} \rangle_W &= \int_0^T \langle y(t), u \rangle_W \langle y(t), \tilde{u} \rangle_W dt = \left\langle \int_0^T \langle y(t), \tilde{u} \rangle_W y(t) dt, u \right\rangle_W \\ &= \langle \mathcal{R}\tilde{u}, u \rangle_W = \langle u, \mathcal{R}\tilde{u} \rangle_W \end{aligned}$$

for all $u, \tilde{u} \in \mathbb{R}^m$ that \mathcal{R} is self-adjoint. \square

Utilizing the operator \mathcal{R} we can write (1.40) as the eigenvalue problem

$$\mathcal{R}u = \lambda u \quad \text{in } \mathbb{R}^m.$$

It follows from Lemma 1.11 that \mathcal{R} possesses eigenvectors $\{u_i\}_{i=1}^m$ and associated real eigenvalues $\{\lambda_i\}_{i=1}^m$ such that

$$(1.42) \quad \mathcal{R}u_i = \lambda_i u_i \quad \text{for } 1 \leq i \leq m \quad \text{and} \quad \lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_m \geq 0.$$

Note that

$$\int_0^T |\langle y(t), u_i \rangle_W|^2 dt = \int_0^T \langle \langle y(t), u_i \rangle_W y(t), u_i \rangle_W dt = \langle \mathcal{R}u_i, u_i \rangle_W = \lambda_i \|u_i\|_W^2 = \lambda_i$$

for $i \in \{1, \dots, m\}$ so that u_1 solves (1.37).

Proceeding as in Sections 1.1 and 1.2 we obtain the following result.

Theorem 1.12. *Let $y \in C([0, T]; \mathbb{R}^m)$ be the unique solution to (1.30). Then the POD basis of rank ℓ solving the minimization problem $(\hat{\mathbf{P}}_W^\ell)$ is given by the eigenvectors $\{u_i\}_{i=1}^\ell$ of \mathcal{R} corresponding to the ℓ largest eigenvalues $\lambda_1 \geq \dots \geq \lambda_\ell$.*

Remark 1.13 (Methods of snapshots). Let us introduce the linear and bounded operator $\mathcal{Y} : L^2(0, T) \rightarrow \mathbb{R}^m$ by

$$\mathcal{Y}v = \int_0^T v(t)y(t) dt \quad \text{for } v \in L^2(0, T).$$

The adjoint $\mathcal{Y}^* : \mathbb{R}^m \rightarrow L^2(0, T)$ satisfying

$$\langle \mathcal{Y}^*u, v \rangle_{L^2(0,T)} = \langle u, \mathcal{Y}v \rangle_W \quad \text{for all } (u, v) \in \mathbb{R}^m \times L^2(0, T)$$

is given as

$$(\mathcal{Y}^*u)(t) = \langle u, y(t) \rangle_W \quad \text{for } u \in \mathbb{R}^m \text{ and almost all } t \in [0, T].$$

Then we have

$$\mathcal{Y}\mathcal{Y}^*u = \int_0^T \langle u, y(t) \rangle_W y(t) dt = \int_0^T \langle y(t), u \rangle_W y(t) dt = \mathcal{R}u$$

for all $u \in \mathbb{R}^m$, i.e., $\mathcal{R} = \mathcal{Y}\mathcal{Y}^*$ holds. Furthermore,

$$(\mathcal{Y}^*\mathcal{Y}v)(t) = \left\langle \int_0^T v(s)y(s) ds, y(t) \right\rangle_W = \int_0^T \langle y(s), y(t) \rangle_W v(s) ds =: (\mathcal{K}v)(t)$$

for all $v \in L^2(0, T)$ and almost all $t \in [0, T]$. Thus, $\mathcal{K} = \mathcal{Y}^*\mathcal{Y}$. It can be shown that the operator \mathcal{K} is linear, bounded, non-negative and self-adjoint. Moreover, \mathcal{K} is compact. Therefore, the POD basis can also be computed as follows: Solve

$$(1.43) \quad \mathcal{K}v_i = \lambda_i v_i \text{ for } 1 \leq i \leq \ell, \quad \lambda_1 \geq \dots \geq \lambda_\ell > 0, \quad \int_0^T v_i(t)v_j(t) dt = \delta_{ij}$$

and set

$$u_i = \frac{1}{\sqrt{\lambda_i}} \mathcal{Y}v_i = \frac{1}{\sqrt{\lambda_i}} \int_0^T v_i(t)y(t) dt \quad \text{for } i = 1, \dots, \ell.$$

Note that (1.43) is a symmetric eigenvalue problem in the infinite-dimensional function space $L^2(0, T)$. For the functional analytic theory we refer, e.g., to [11]. \diamond

Let us turn back to the optimality conditions (1.35). For any $u \in \mathbb{R}^m$ and $i \in \{1, \dots, m\}$ we derive

$$\begin{aligned} (YDY^TWu)_i &= \sum_{\nu=1}^m \sum_{j=1}^m \sum_{k=1}^m \alpha_j Y_{ij} Y_{kj} W_{k\nu} u_\nu = \sum_{j=1}^n \alpha_j Y_{ij} \langle y_j, u \rangle_W \\ &= \sum_{j=1}^n \alpha_j \langle y_j, u \rangle_W (y_j)_i, \end{aligned}$$

where $(y_j)_i$ stands for the i th component of the vector $y_j \in \mathbb{R}^m$. Thus,

$$YDY^TWu = \sum_{j=1}^n \alpha_j \langle y_j, u \rangle_W y_j =: \mathcal{R}^n u.$$

Note that the operator $\mathcal{R}^n : \mathbb{R}^m \rightarrow \mathbb{R}^m$ is linear and bounded. Moreover,

$$\langle \mathcal{R}^n u, u \rangle_W = \left\langle \sum_{j=1}^n \alpha_j \langle y_j, u \rangle_W y_j, u \right\rangle_W = \sum_{j=1}^n \alpha_j |\langle y_j, u \rangle_W|^2 \geq 0$$

holds for all $u \in \mathbb{R}^m$ so that \mathcal{R}^n is non-negative. Further,

$$\begin{aligned} \langle \mathcal{R}^n u, \tilde{u} \rangle_W &= \left\langle \sum_{j=1}^n \alpha_j \langle y_j, u \rangle_W y_j, \tilde{u} \right\rangle_W = \sum_{j=1}^n \alpha_j \langle y_j, u \rangle_W \langle y_j, \tilde{u} \rangle_W \\ &= \left\langle \sum_{j=1}^n \alpha_j \langle y_j, \tilde{u} \rangle_W y_j, u \right\rangle_W = \langle \mathcal{R}^n \tilde{u}, u \rangle_W = \langle u, \mathcal{R}^n \tilde{u} \rangle_W \end{aligned}$$

for all $u, \tilde{u} \in \mathbb{R}^m$, i.e., \mathcal{R}^n is self-adjoint. Therefore, \mathcal{R}^n has the same properties as the operator \mathcal{R} . Summarizing, we have

$$(1.44a) \quad \mathcal{R}^n u_i^n = \lambda_i^n u_i^n, \quad \lambda_1^n \geq \dots \geq \lambda_\ell^n \geq \dots \geq \lambda_{d(n)}^n > \lambda_{d(n)+1}^n = \dots = \lambda_m^n = 0,$$

$$(1.44b) \quad \mathcal{R}u_i = \lambda_i u_i, \quad \lambda_1 \geq \dots \geq \lambda_\ell \geq \dots \geq \lambda_d > \lambda_{d+1} = \dots = \lambda_m = 0.$$

Let us note that

$$(1.45) \quad \int_0^T \|y(t)\|_W^2 dt = \sum_{i=1}^d \lambda_i = \sum_{i=1}^m \lambda_i.$$

In fact,

$$\mathcal{R}u_i = \int_0^T \langle y(t), u_i \rangle_W y(t) dt \quad \text{for every } i \in \{1, \dots, m\}.$$

Taking the inner product with u_i , using (1.44b) and summing over i we arrive at

$$\sum_{i=1}^d \int_0^T |\langle y(t), u_i \rangle_W|^2 dt = \sum_{i=1}^d \langle \mathcal{R}u_i, u_i \rangle_W = \sum_{i=1}^d \lambda_i = \sum_{i=1}^m \lambda_i.$$

Expanding $y(t) \in \mathbb{R}^m$ in terms of $\{u_i\}_{i=1}^m$ we have

$$y(t) = \sum_{i=1}^m \langle y(t), u_i \rangle_W u_i$$

and hence

$$\int_0^T \|y(t)\|_W^2 dt = \sum_{i=1}^m \int_0^T |\langle y(t), u_i \rangle_W|^2 dt = \sum_{i=1}^m \lambda_i,$$

which is (1.45). Analogously, we obtain

$$(1.46) \quad \sum_{j=1}^n \alpha_j \|y(t_j)\|_W^2 = \sum_{i=1}^{d(n)} \lambda_i^n = \sum_{i=1}^m \lambda_i^n \quad \text{for every } n \in \mathbb{N}.$$

For convenience we do not indicate the dependence of α_j on n . Let $y \in C([0, T]; \mathbb{R}^m)$ hold. To ensure

$$(1.47) \quad \sum_{j=1}^n \alpha_j \|y(t_j)\|_W^2 \rightarrow \int_0^T \|y(t)\|_W^2 dt \quad \text{as } \Delta t \rightarrow 0$$

we have to choose the α_j 's appropriately. Here we take the trapezoidal weights

$$(1.48) \quad \alpha_1 = \frac{\Delta t}{2}, \quad \alpha_j = \Delta t \text{ for } 2 \leq j \leq n-1, \quad \alpha_n = \frac{\Delta t}{2}.$$

Suppose that we have

$$(1.49) \quad \lim_{n \rightarrow \infty} \|\mathcal{R}^n - \mathcal{R}\|_{L(\mathbb{R}^m)} = \lim_{n \rightarrow \infty} \sup_{\|u\|_W=1} \|\mathcal{R}^n u - \mathcal{R}u\|_W = 0$$

provided $y \in C^1([0, T]; \mathbb{R}^m)$ is satisfied. In (1.49) $L(\mathbb{R}^m)$ denotes the Banach space of all linear and bounded operators mapping from \mathbb{R}^m into itself. Combining (1.47) with (1.45) and (1.46) we find

$$(1.50) \quad \sum_{i=1}^m \lambda_i^n \rightarrow \sum_{i=1}^m \lambda_i \quad \text{as } n \rightarrow \infty.$$

Now choose and fix

$$(1.51) \quad \ell \quad \text{such that} \quad \lambda_\ell \neq \lambda_{\ell+1}.$$

Then by spectral analysis of compact operators ([5, pp. 212–214]) and (1.49) it follows that

$$(1.52) \quad \lambda_i^n \rightarrow \lambda_i \quad \text{for } 1 \leq i \leq \ell \text{ as } n \rightarrow \infty.$$

Combining (1.50) and (1.52) there exists $\bar{n} \in \mathbb{N}$ such that

$$(1.53) \quad \sum_{i=\ell+1}^m \lambda_i^n \leq 2 \sum_{i=\ell+1}^m \lambda_i \quad \text{for all } n \geq \bar{n},$$

if $\sum_{i=\ell+1}^m \lambda_i \neq 0$. Moreover, for ℓ as above, \bar{n} can also be chosen such that

$$(1.54) \quad \sum_{i=\ell+1}^{d(n)} |\langle y_0, u_i^n \rangle_W|^2 \leq 2 \sum_{i=\ell+1}^m |\langle y_0, u_i \rangle_W|^2 \quad \text{for all } n \geq \bar{n},$$

provided that $\sum_{i=\ell+1}^m |\langle y_0, u_i \rangle_W|^2 \neq 0$ (1.49) hold. Recall that the vector $y_0 \in \mathbb{R}^m$ stands for the initial condition in (1.30b). Then we have

$$(1.55) \quad \|y_0\|_W^2 = \sum_{i=1}^m |\langle y_0, u_i \rangle_W|^2.$$

If $t_1 = 0$ holds, we have $y_0 \in \text{span}\{y_j\}_{j=1}^n$ for every n and

$$(1.56) \quad \|y_0\|_W^2 = \sum_{i=1}^{d(n)} |\langle y_0, u_i^n \rangle_W|^2.$$

Therefore, for $\ell < d(n)$ by (1.55) and (1.56)

$$\begin{aligned} \sum_{i=\ell+1}^{d(n)} |\langle y_0, u_i^n \rangle_W|^2 &= \sum_{i=1}^{d(n)} |\langle y_0, u_i^n \rangle_W|^2 - \sum_{i=1}^{\ell} |\langle y_0, u_i^n \rangle_W|^2 + \sum_{i=1}^{\ell} |\langle y_0, u_i \rangle_W|^2 \\ &\quad + \sum_{i=\ell+1}^m |\langle y_0, u_i \rangle_W|^2 - \sum_{i=1}^m |\langle y_0, u_i \rangle_W|^2 \\ &= \sum_{i=1}^{\ell} \left(|\langle y_0, u_i \rangle_W|^2 - |\langle y_0, u_i^n \rangle_W|^2 \right) + \sum_{i=\ell+1}^m |\langle y_0, u_i \rangle_W|^2. \end{aligned}$$

As a consequence of (1.49) and (1.51) we have $\lim_{n \rightarrow \infty} \|u_i^n - u_i\|_W = 0$ for $i = 1, \dots, \ell$ and hence (1.54) follows.

Summarizing we have the following theorem.

Theorem 1.14. *Assume that $y \in C^1([0, T]; \mathbb{R}^m)$ is the unique solution to (1.30). Let $\{(u_i^n, \lambda_i^n)\}_{i=1}^m$ and $\{(u_i, \lambda_i)\}_{i=1}^m$ be the eigenvector-eigenvalue pairs given by (1.44). Suppose that $\ell \in \{1, \dots, m\}$ is fixed such that (1.51) and*

$$\sum_{i=\ell+1}^m \lambda_i \neq 0, \quad \sum_{i=\ell+1}^m |\langle y_0, u_i \rangle_W|^2 \neq 0$$

hold. Then we have

$$(1.57) \quad \lim_{n \rightarrow \infty} \|\mathcal{R}^n - \mathcal{R}\|_{L(\mathbb{R}^m)} = 0.$$

This implies

$$\begin{aligned} \lim_{n \rightarrow \infty} |\lambda_i^n - \lambda_i| &= \lim_{n \rightarrow \infty} \|u_i^n - u_i\|_W = 0 \quad \text{for } 1 \leq i \leq \ell, \\ \lim_{n \rightarrow \infty} \sum_{i=\ell+1}^m (\lambda_i^n - \lambda_i) &= 0 \quad \text{and} \quad \lim_{n \rightarrow \infty} \sum_{i=\ell+1}^m |\langle y_0, u_i^n \rangle_W|^2 = \sum_{i=\ell+1}^m |\langle y_0, u_i \rangle_W|^2. \end{aligned}$$

Proof. We only have to verify (1.57). For that purpose we choose an arbitrary $u \in \mathbb{R}^m$ with $\|u\|_W = 1$ and introduce $f_u : [0, T] \rightarrow \mathbb{R}^m$ by

$$f_u(t) = \langle y(t), u \rangle_W y(t) \quad \text{for } t \in [0, T].$$

Then, we have $f_u \in C^1([0, T]; \mathbb{R}^m)$ with

$$\dot{f}_u(t) = \langle \dot{y}(t), u \rangle_W y(t) + \langle y(t), u \rangle_W \dot{y}(t) \quad \text{for } t \in [0, T]$$

By Taylor expansion there exist $\tau_{j1}(t), \tau_{j2}(t) \in [t_j, t_{j+1}]$ depending on t

$$\begin{aligned} \int_{t_j}^{t_{j+1}} f_u(t) dt &= \frac{1}{2} \int_{t_j}^{t_{j+1}} f_u(t_j) + \dot{f}_u(\tau_{j1}(t))(t - t_j) dt \\ &\quad + \frac{1}{2} \int_{t_j}^{t_{j+1}} f_u(t_{j+1}) + \dot{f}_u(\tau_{j2}(t))(t - t_{j+1}) dt \\ &= \frac{\Delta t}{2} (f_u(t_j) + f_u(t_{j+1})) + \frac{1}{2} \int_{t_j}^{t_{j+1}} \dot{f}_u(\tau_{j1}(t))(t - t_j) dt \\ &\quad + \frac{1}{2} \int_{t_j}^{t_{j+1}} \dot{f}_u(\tau_{j2}(t))(t - t_{j+1}) dt. \end{aligned}$$

Hence,

$$\begin{aligned} \|\mathcal{R}^n u - \mathcal{R}u\|_W &= \left\| \sum_{j=1}^n \alpha_j f_u(t_j) - \int_0^T f_u(t) dt \right\|_W \\ &= \left\| \sum_{j=1}^{n-1} \left(\frac{\Delta t}{2} (f_u(t_j) + f_u(t_{j+1})) - \int_{t_j}^{t_{j+1}} f_u(t) dt \right) \right\|_W \\ &\leq \frac{1}{2} \sum_{j=1}^{n-1} \int_{t_j}^{t_{j+1}} \|\dot{f}_u(\tau_{j1}(t))\|_W |t - t_j| + \|\dot{f}_u(\tau_{j2}(t))\|_W |t - t_{j+1}| dt \\ &\leq \frac{1}{2} \max_{t \in [0, T]} \|\dot{f}_u(t)\|_W \sum_{j=1}^{n-1} \left(\frac{(t - t_j)^2}{2} - \frac{(t_{j+1} - t)^2}{2} \Big|_{t=t_j}^{t=t_{j+1}} \right) \\ &= \frac{\Delta t}{2} \max_{t \in [0, T]} \|\dot{f}_u(t)\|_W \sum_{j=1}^{n-1} \Delta t = \frac{\Delta t T}{2} \max_{t \in [0, T]} \|\dot{f}_u(t)\|_W \\ &\leq \frac{\Delta t T}{2} \max_{t \in [0, T]} \|\dot{f}_u(t)\|_W \\ &= \frac{\Delta t T}{2} \max_{t \in [0, T]} \|\langle \dot{y}(t), u \rangle_W y(t) + \langle y(t), u \rangle_W \dot{y}(t)\|_W \\ &= \Delta t T \max_{t \in [0, T]} \|\dot{y}(t)\|_W \|y(t)\|_W \leq \Delta t T \|y\|_{C^1([0, T]; \mathbb{R}^m)}^2. \end{aligned}$$

Consequently,

$$\|\mathcal{R}^n - \mathcal{R}\|_{L(\mathbb{R}^m)} = \sup_{\|u\|_W=1} \|\mathcal{R}^n u - \mathcal{R}u\|_W \leq 2\Delta t \|y\|_{C^1([0, T]; \mathbb{R}^m)}^2 \xrightarrow{\Delta t \rightarrow 0} 0$$

which is (1.57). \square

1.4. Exercises.

- 1.1) Show that any optimal solution to (\mathbf{P}^ℓ) is a regular point.
 1.2) Verify the claim in Theorem 1.1 that $\operatorname{argmax}(\mathbf{P}^\ell) = \sum_{i=1}^\ell \sigma_i^2$ holds true.
 1.3) Show that the Frobenius norm is a matrix norm and that

$$\|AB\|_F \leq \|A\|_F \|B\|_F \quad \text{for any } A, B \in \mathbb{R}^{n \times n}$$

is valid. Suppose that $U^d \in \mathbb{R}^{m \times d}$ is a matrix with pairwise orthonormal vectors $u_i \in \mathbb{R}^m$, $1 \leq i \leq d$. Prove that

$$\|UA\|_F = \|A\|_F \quad \text{for any matrix } A \in \mathbb{R}^{d \times n}.$$

- 1.4) Suppose that $W \in \mathbb{R}^{m \times m}$ is symmetric and positive definite. Let $\eta_1 \geq \dots \geq \eta_m > 0$ denote the eigenvalues of W and $W^\alpha = Q \operatorname{diag}(\eta_1^\alpha, \dots, \eta_m^\alpha) Q^T$ be the eigenvalue decomposition of W . We define

$$W^\alpha = Q \operatorname{diag}(\eta_1^\alpha, \dots, \eta_m^\alpha) Q^T \quad \text{for } \alpha \in \mathbb{R}.$$

Show that $(W^\alpha)^{-1}$ exists and $(W^\alpha)^{-1} = W^{-\alpha}$. Prove that $W^{\alpha+\beta} = W^\alpha W^\beta$ holds for $\alpha, \beta \in \mathbb{R}$.

- 1.5) Verify the claims of Theorem 1.8.
 1.5.1) Prove that $u_i = W^{-1/2} \bar{u}_i$, $1 \leq i \leq \ell$, solves (\mathbf{P}_W^ℓ) , where the matrix W and the vectors $\bar{u}_1, \dots, \bar{u}_m$ are introduced in Theorem 1.8.
 1.5.2) Show that (1.29) holds.
 1.6) Prove that u_1 given by (1.42) is a global solution to (1.37).
 1.7) Verify (1.46).

2. Reduced-order modeling (ROM)

In Section 1 we have introduced the POD basis of rank ℓ in \mathbb{R}^m and discussed its application to initial-value problems. If the POD basis is computed, it can be used to derive a so-called *low-dimensional approximation* or a *reduced-order model* for (1.30). This is the focus of this section.

2.1. ROM for time-dependent systems. Suppose that we have determined a POD basis $\{u_j\}_{j=1}^\ell$ of rank $\ell \in \{1, \dots, m\}$ in \mathbb{R}^m . Then we make the ansatz

$$(2.1) \quad y^\ell(t) = \sum_{j=1}^\ell \underbrace{\langle y^\ell(t), u_j \rangle_W}_{=: y_j^\ell(t)} u_j \quad \text{for all } t \in [0, T],$$

where the Fourier coefficients y_j^ℓ , $1 \leq j \leq \ell$, are functions mapping $[0, T]$ into \mathbb{R} . Since

$$y(t) = \sum_{j=1}^m \langle y(t), u_j \rangle_W u_j \quad \text{for all } t \in [0, T]$$

holds, $y^\ell(t)$ is an approximation for $y(t)$ provided $\ell < m$. Inserting (2.1) into (1.30) yields

$$(2.2a) \quad \sum_{j=1}^\ell \dot{y}_j^\ell(t) u_j = \sum_{j=1}^\ell y_j^\ell(t) A u_j + f(t, y^\ell(t)), \quad t \in (0, T],$$

$$(2.2b) \quad \sum_{j=1}^\ell y_j^\ell(0) u_j = y_0$$

Note that (2.2) is an initial-value problem in \mathbb{R}^m for $\ell \leq m$ coefficient functions $y_j^\ell(t)$, $1 \leq j \leq \ell$ and $t \in [0, T]$, so that the coefficients are overdetermined. Therefore, we assume that (2.2) holds after projection on the ℓ dimensional subspace $V^\ell = \text{span}\{u_j\}_{j=1}^\ell$. From (2.2a) and $\langle u_j, u_i \rangle_W = \delta_{ij}$ we infer that

$$(2.3) \quad \dot{y}_i^\ell(t) = \sum_{j=1}^{\ell} y_j^\ell(t) \langle Au_j, u_i \rangle_W + \langle f(t, y^\ell(t)), u_i \rangle_W$$

for $1 \leq i \leq \ell$ and $t \in (0, T]$. Let us introduce the matrix

$$A = ((a_{ij})) \in \mathbb{R}^{\ell \times \ell} \quad \text{with} \quad a_{ij} = \langle Au_j, u_i \rangle_W,$$

the vector-valued mapping

$$y^\ell = \begin{pmatrix} y_1^\ell \\ \vdots \\ y_\ell^\ell \end{pmatrix} : [0, T] \rightarrow \mathbb{R}^\ell$$

and the non-linearity $F = (F_1, \dots, F_\ell)^T : [0, T] \times \mathbb{R}^\ell \rightarrow \mathbb{R}^\ell$ by

$$F_i(t, y) = \left\langle f\left(t, \sum_{j=1}^{\ell} y_j u_j\right), u_i \right\rangle_W \quad \text{for } t \in [0, T] \text{ and } y = (y_1, \dots, y_\ell) \in \mathbb{R}^\ell.$$

Then, (2.3) can be expressed as

$$(2.4a) \quad \dot{y}^\ell(t) = Ay^\ell(t) + F(t, y^\ell(t)) \quad \text{for } t \in (0, T]$$

From (2.2b) we derive

$$(2.4b) \quad y^\ell(0) = y_0,$$

where

$$y_0 = \begin{pmatrix} \langle y_0, u_1 \rangle_W \\ \vdots \\ \langle y_0, u_\ell \rangle_W \end{pmatrix} \in \mathbb{R}^\ell$$

holds. System (2.4) is called the *POD-Galerkin projection* for (1.30). In case of $\ell \ll m$ the ℓ -dimensional system (2.4) is a low-dimensional approximation for (1.30). Therefore, (2.4) is a reduced-order model for (1.30).

2.2. Error analysis for the reduced-order model. In this section we focus on error analysis for POD Galerkin approximations. For a more detailed presentation we refer the reader to [7, 8, 9] and [4].

Let us suppose that $y \in C([0, T]; \mathbb{R}^m) \cap C^1(0, T; \mathbb{R}^m)$ is the unique solution to (1.30) and $\{u_i\}_{i=1}^\ell$ the POD basis of rank ℓ solving

$$(2.5) \quad \min \int_0^T \left\| y(t) - \sum_{i=1}^{\ell} \langle y(t), u_i \rangle_W u_i \right\|_W^2 dt \quad \text{s.t.} \quad \langle u_j, u_i \rangle_W = \delta_{ij}, \quad 1 \leq i, j \leq \ell.$$

The reduced-order model for (1.30) is given by (2.4). We are interested in estimating the error

$$\int_0^T \|y(t) - y^\ell(t)\|_W^2 dt.$$

Let us introduce the finite-dimensional space

$$V^\ell = \text{span}\{u_1, \dots, u_\ell\} \subset \mathbb{R}^m$$

and the projection $\mathcal{P}^\ell : \mathbb{R}^m \rightarrow V^\ell$ by

$$\mathcal{P}^\ell u = \sum_{i=1}^{\ell} \langle u, u_i \rangle_W u_i \quad \text{for } u \in \mathbb{R}^m.$$

Then,

$$\begin{aligned} \mathcal{P}^\ell(\alpha u + \tilde{\alpha} \tilde{u}) &= \sum_{i=1}^{\ell} \langle \alpha u + \tilde{\alpha} \tilde{u}, u_i \rangle_W u_i = \sum_{i=1}^{\ell} \left(\alpha \langle u, u_i \rangle_W + \tilde{\alpha} \langle \tilde{u}, u_i \rangle_W \right) u_i \\ &= \alpha \mathcal{P}^\ell u + \tilde{\alpha} \mathcal{P}^\ell \tilde{u} \end{aligned}$$

for all $\alpha, \tilde{\alpha} \in \mathbb{R}$ and $u, \tilde{u} \in \mathbb{R}^m$ so that \mathcal{P}^ℓ is linear. Further,

$$(2.6) \quad \begin{aligned} \|\mathcal{P}^\ell\|_{L(\mathbb{R}^m)}^2 &= \sup_{\|u\|_W=1} \|\mathcal{P}^\ell u\|_W^2 = \sup_{\|u\|_W=1} \sum_{i=1}^{\ell} |\langle u, u_i \rangle_W|^2 \\ &\leq \sup_{\|u\|_W=1} \sum_{i=1}^m |\langle u, u_i \rangle_W|^2 = \sup_{\|u\|_W=1} \|u\|_W^2 = 1, \end{aligned}$$

i.e., \mathcal{P}^ℓ is bounded and therefore continuous. In particular, (2.6) and $\|\mathcal{P}^\ell u\|_W = \|u\|_W$ for any $u \in V^\ell$ imply $\|\mathcal{P}^\ell\|_{L(\mathbb{R}^m)} = 1$.

Throughout we shall use the decomposition

$$(2.7) \quad y(t) - y^\ell(t) = y(t) - \mathcal{P}^\ell y(t) + \mathcal{P}^\ell y(t) - y^\ell(t) = \varrho^\ell(t) + \vartheta^\ell(t),$$

where $\varrho^\ell(t) = y(t) - \mathcal{P}^\ell y(t)$ and $\vartheta^\ell(t) = \mathcal{P}^\ell y(t) - y^\ell(t)$. Note that

$$\int_0^T \left\| y(t) - \sum_{i=1}^{\ell} \langle y(t), u_i \rangle_W u_i \right\|_W^2 dt = \int_0^T \|y(t) - \mathcal{P}^\ell y(t)\|_W^2 dt = \int_0^T \|\varrho^\ell(t)\|_W^2 dt.$$

Since $\{u_i\}_{i=1}^{\ell}$ is a POD basis of rank ℓ we have

$$(2.8) \quad \int_0^T \|\varrho^\ell(t)\|_W^2 dt = \sum_{i=\ell+1}^m \lambda_i.$$

Next we estimate the term $\vartheta^\ell(t)$. Utilizing (1.30a) and (2.4) we obtain for every $u^\ell \in V^\ell$ and $t \in (0, T]$

$$(2.9) \quad \begin{aligned} \langle \dot{\vartheta}^\ell(t), u^\ell \rangle_W &= \langle \mathcal{P}^\ell \dot{y}(t) - \dot{y}(t), u^\ell \rangle_W + \langle \dot{y}(t) - \dot{y}^\ell(t), u^\ell \rangle_W \\ &= \langle \mathcal{P}^\ell \dot{y}(t) - \dot{y}(t), u^\ell \rangle_W \\ &\quad + \langle A(y(t) - y^\ell(t)) + f(t, y(t)) - f(t, y^\ell(t)), u^\ell \rangle_W \end{aligned}$$

We choose $u^\ell = \vartheta^\ell(t) \in V^\ell$. Let

$$\|A\| = \max_{\|u\|_W=1} \|Au\|_W$$

the matrix norm induced by the vector norm $\|\cdot\|_W$. Further,

$$\frac{1}{2} \frac{d}{dt} \|\vartheta^\ell(t)\|_W^2 = \langle \dot{\vartheta}^\ell(t), \vartheta^\ell(t) \rangle_W \quad \text{for every } t \in (0, T].$$

holds. Then, we infer from (2.9)

$$(2.10) \quad \begin{aligned} \frac{1}{2} \frac{d}{dt} \|\vartheta^\ell(t)\|_W^2 &\leq \|A\| (\|\varrho^\ell(t)\|_W + \|\vartheta^\ell(t)\|_W) \|\vartheta^\ell(t)\|_W \\ &\quad + \|f(t, y(t)) - f(t, y^\ell(t))\|_W \|\vartheta^\ell(t)\|_W \\ &\quad + \|\mathcal{P}^\ell \dot{y}(t) - \dot{y}(t)\|_W \|\vartheta^\ell(t)\|_W. \end{aligned}$$

Suppose that f is Lipschitz-continuous with respect to the second argument, i.e., there exists a constant $L_f \geq 0$ satisfying

$$\|f(t, u) - f(t, \tilde{u})\|_W \leq L_f \|u - \tilde{u}\|_W \quad \text{for all } u, \tilde{u} \in \mathbb{R}^m \text{ and } t \in [0, T].$$

Moreover, we have

$$\|\mathcal{P}^\ell \dot{y}(t) - \dot{y}(t)\|_W^2 = \left\| \sum_{i=\ell+1}^m \langle \dot{y}(t), u_i \rangle_W u_i \right\|_W^2 = \sum_{i=\ell+1}^m |\langle \dot{y}(t), u_i \rangle_W|^2$$

for all $t \in (0, T)$. Consequently, (2.10) and (2.7) imply

$$\begin{aligned} \frac{1}{2} \frac{d}{dt} \|\vartheta^\ell(t)\|_W^2 &\leq \frac{\|A\|}{2} \left(\|\varrho^\ell(t)\|_W^2 + \|\vartheta^\ell(t)\|_W^2 \right) + \|A\| \|\vartheta^\ell(t)\|_W^2 \\ &\quad + L_f \|\varrho^\ell(t) + \vartheta^\ell(t)\|_W \|\vartheta^\ell(t)\|_W \\ &\quad + \frac{1}{2} \left(\|\mathcal{P}^\ell \dot{y}(t) - \dot{y}(t)\|_W^2 + \|\vartheta^\ell(t)\|_W^2 \right) \\ &\leq \frac{\|A\|}{2} \|\varrho^\ell(t)\|_W^2 + \left(\frac{1}{2} \|A\| + \frac{1}{2} + L_f \right) \|\vartheta^\ell(t)\|_W^2 \\ &\quad + L_f \|\varrho^\ell(t)\|_W \|\vartheta^\ell(t)\|_W + \sum_{i=\ell+1}^m |\langle \dot{y}(t), u_i \rangle_W|^2 \\ &\leq \frac{\|A\| + L_f}{2} \|\varrho^\ell(t)\|_W^2 + \left(\frac{3}{2} (\|A\| + L_f) + \frac{1}{2} \right) \|\vartheta^\ell(t)\|_W^2 \\ &\quad + \sum_{i=\ell+1}^m |\langle \dot{y}(t), u_i \rangle_W|^2. \end{aligned}$$

Consequently,

$$\begin{aligned} \frac{d}{dt} \|\vartheta^\ell(t)\|_W^2 &\leq \left(3(\|A\| + L_f) + 1 \right) \|\vartheta^\ell(t)\|_W^2 + (\|A\| + L_f) \|\varrho^\ell(t)\|_W^2 \\ &\quad + \sum_{i=\ell+1}^m |\langle \dot{y}(t), u_i \rangle_W|^2. \end{aligned}$$

Using Gronwall's lemma (see Exercise 2.1)) and (2.8) we arrive at

$$(2.11) \quad \begin{aligned} \|\vartheta^\ell(t)\|_W^2 &\leq c_1 \left(\|\vartheta^\ell(0)\|_W^2 + (\|A\| + L_f) \int_0^t \|\varrho^\ell(s)\|_W^2 ds \right) \\ &\quad + c_1 \sum_{i=\ell+1}^m \int_0^t |\langle \dot{y}(s), u_i \rangle_W|^2 ds \\ &\leq c_2 \left(\|\vartheta^\ell(0)\|_W^2 + \sum_{i=\ell+1}^m \left(\lambda_i + \int_0^T |\langle \dot{y}(t), u_i \rangle_W|^2 dt \right) \right) \end{aligned}$$

where $c_1 = \exp(3(\|A\| + L_f) + 1)T$ and $c_2 = c_1 \max\{\|A\| + L_f, 1\}$.

Theorem 2.1. *Let $y \in C([0, T]; \mathbb{R}^m) \cap C^1(0, T; \mathbb{R}^m)$ be the unique solution to (1.30), $\ell \in \{1, \dots, m\}$ be fixed and $\{u_i\}_{i=1}^\ell$ a POD basis of rank ℓ solving (2.5). Let y^ℓ be the unique solution to the reduced-order model (2.4). Then*

$$\int_0^T \|y(t) - y^\ell(t)\|_W^2 dt \leq C \sum_{i=\ell+1}^m \left(\lambda_i + \int_0^T |\langle \dot{y}(t), u_i \rangle_W|^2 dt \right)$$

for a constant $C > 0$.

Proof. From (2.8), (2.11) and $\vartheta^\ell(0) = \mathcal{P}^\ell y_0 - y^\ell(0) = 0$ we find

$$\begin{aligned} \int_0^T \|y(t) - y^\ell(t)\|_W^2 dt &= \int_0^T \|\varrho^\ell(t) + \vartheta^\ell(t)\|_W^2 dt \\ &\leq 2 \int_0^T \|\varrho^\ell(t)\|_W^2 + \|\vartheta^\ell(t)\|_W^2 dt \\ &\leq 2 \sum_{i=\ell+1}^m \lambda_i + c_3 \sum_{i=\ell+1}^m \left(\lambda_i + \int_0^T |\langle \dot{y}(t), u_i \rangle_W|^2 dt \right) \end{aligned}$$

with $c_3 = 2c_2$. Setting $C = 2 + c_3$ the claim follows directly. \square

Remark 2.2. The term

$$\sum_{i=\ell+1}^m \int_0^T |\langle \dot{y}(t), u_i \rangle_W|^2 dt$$

can not be estimated by the sum over the eigenvalues $\lambda_{\ell+1}, \dots, \lambda_m$. If we replace (2.5) by

$$(2.12a) \quad \min \int_0^T \left\| y(t) - \sum_{i=1}^\ell \langle y(t), u_i \rangle_W u_i \right\|_W^2 + \left\| \dot{y}(t) - \sum_{i=1}^\ell \langle \dot{y}(t), u_i \rangle_W u_i \right\|_W^2 dt$$

subject to

$$(2.12b) \quad \langle u_j, u_i \rangle_W = \delta_{ij} \quad \text{for } 1 \leq i, j \leq \ell,$$

we end up with the estimate

$$\int_0^T \|y(t) - y^\ell(t)\|_W^2 dt \leq \tilde{C} \sum_{i=\ell+1}^m \tilde{\lambda}_i$$

for a constant $\tilde{C} > 0$. In this case the time derivatives are also included in the snapshot ensemble. Of course, the operator \mathcal{R} defined in (1.41) has to be replaced. It turns out that the POD basis $\{u_i\}_{i=1}^\ell$ is given by the eigenvalue problem

$$(2.13) \quad \tilde{\mathcal{R}} \tilde{u}_i = \tilde{\lambda}_i \tilde{u}_i \quad \text{for } 1 \leq i \leq m \quad \text{and} \quad \tilde{\lambda}_1 \geq \tilde{\lambda}_2 \geq \dots \geq \tilde{\lambda}_m \geq 0$$

where the operator $\tilde{\mathcal{R}} : \mathbb{R}^m \rightarrow \mathbb{R}^m$ is defined by

$$\tilde{\mathcal{R}} u = \int_0^T \langle y(t), u \rangle_W y(t) + \langle \dot{y}(t), u \rangle_W \dot{y}(t) dt$$

for $u \in \mathbb{R}^m$. \diamond

Remark 2.3. Suppose that we build the matrix $Y \in \mathbb{R}^{m \times (2n)}$ using the column vectors $y_j \approx y(t_j)$, $1 \leq j \leq n$, and $y_j \approx \dot{y}(t_{j-m})$, $m+1 \leq j \leq 2m$. Then, the discrete variant $\tilde{\mathcal{R}}^n$ of the operator $\tilde{\mathcal{R}}$ introduced in Remark 2.2 is given by

$$\begin{aligned} \tilde{\mathcal{R}}^n u &= \sum_{j=1}^n \alpha_j \langle y_j, u \rangle_W y_j + \alpha_j \langle y_{m+j}, u \rangle_W y_{m+j} \\ &= \sum_{j=1}^n \alpha_j \left(\left(\sum_{k=1}^m \sum_{\nu=1}^m Y_{kj} W_{k\nu} u_\nu \right) Y_{\cdot, j} + \left(\sum_{k=1}^m \sum_{\nu=1}^m Y_{k, m+j} W_{k\nu} u_\nu \right) Y_{\cdot, m+j} \right) \\ &= \sum_{j=1}^n \sum_{k=1}^m \sum_{\nu=1}^m \left(\left(Y_{\cdot, j} D_{jj} Y_{jk}^T + Y_{\cdot, m+j} D_{jj} Y_{m+j, k}^T \right) W_{k\nu} u_\nu \right) \\ &= Y \underbrace{\begin{pmatrix} D & 0 \\ 0 & D \end{pmatrix}}_{=: \tilde{D} \in \mathbb{R}^{2n \times 2n}} Y^T W u = Y \tilde{D} Y^T W u \end{aligned}$$

with non-negative weights introduced in $(\hat{\mathbf{P}}_W^{n, \ell})$ and the diagonal matrix $D = \text{diag}(\alpha_1, \dots, \alpha_n) \in \mathbb{R}^{n \times n}$. Thus, we have $\tilde{\mathcal{R}} = Y \tilde{D} Y^T W \in \mathbb{R}^{m \times m}$, which is of the same form as in (1.35). The discrete version to (2.13) is

$$(2.14) \quad Y \tilde{D} Y^T W \tilde{u}_i = \tilde{\lambda}_i \tilde{u}_i \text{ for } 1 \leq i \leq m \quad \text{and} \quad \tilde{\lambda}_1 \geq \tilde{\lambda}_2 \geq \dots \geq \tilde{\lambda}_m \geq 0$$

Setting $\tilde{u}_i = W^{-1/2} \bar{u}_i$ in (2.14) and multiplying by $W^{1/2}$ from the left yield

$$(2.15) \quad W^{1/2} Y \tilde{D} Y^T W^{1/2} \bar{u}_i = \lambda_i \bar{u}_i.$$

Let $\bar{Y} = W^{1/2} Y \tilde{D}^{1/2} \in \mathbb{R}^{m \times 2n}$. Using $W^T = W$ as well as $\tilde{D}^T = \tilde{D}$ we infer from (2.15) that the solution $\{\bar{u}_i\}_{i=1}^\ell$ is given by the symmetric $m \times m$ eigenvalue problem

$$\bar{Y} \bar{Y}^T \bar{u}_i = \lambda_i \bar{u}_i, \quad 1 \leq i \leq \ell \quad \text{and} \quad \langle \bar{u}_i, \bar{u}_j \rangle_{\mathbb{R}^m} = \delta_{ij}, \quad 1 \leq i, j \leq \ell$$

and $\tilde{u}_i = W^{-1/2} \bar{u}_i$. Note that

$$\bar{Y}^T \bar{Y} = \tilde{D}^{1/2} Y^T W Y \tilde{D}^{1/2} \in \mathbb{R}^{2n \times 2n}.$$

Thus, the POD basis of rank ℓ can also be computed by the methods of snapshots as follows: First solve the symmetric $2n \times 2n$ eigenvalue problem

$$\bar{Y}^T \bar{Y} \bar{v}_i = \lambda_i \bar{v}_i, \quad 1 \leq i \leq \ell \quad \text{and} \quad \langle \bar{v}_i, \bar{v}_j \rangle_{\mathbb{R}^{2n}} = \delta_{ij}, \quad 1 \leq i, j \leq \ell.$$

Then we set (by SVD)

$$\tilde{u}_i = W^{-1/2} \bar{u}_i = \frac{1}{\sqrt{\lambda_i}} W^{-1/2} \bar{Y} \bar{v}_i = \frac{1}{\sqrt{\lambda_i}} Y \tilde{D}^{1/2} \bar{v}_i$$

for $1 \leq i \leq \ell$. ◇

From a practical point of view we do not have the information on the whole trajectory in $[0, T]$. Therefore, let $\Delta t = T/(n-1)$ be a fixed time step size and $t_j = (j-1)\Delta t$ for $1 \leq j \leq n$ a given time grid in $[0, T]$. To simplify the presentation we choose an equidistant grid. Of course, non-equidistant meshes can be treated analogously [8]. We compute a POD basis $\{u_i^n\}_{i=1}^\ell$ of rank ℓ by solving the constrained minimization problem $(\hat{\mathbf{P}}_W^{n, \ell})$. After the POD basis has been determined,

we derive the reduced-order model as described in Section 2.1. Thus,

$$y^\ell(t) = \sum_{i=1}^{\ell} y_j^\ell(t) u_i^n, \quad t \in [0, T],$$

solves the POD Galerkin projection of (1.30)

$$(2.16a) \quad \langle \dot{y}^\ell(t), u_i^n \rangle_W = \langle Ay^\ell(t) + f(t, y^\ell(t)), u_i^n \rangle_W \quad \text{for } i = 1 \dots, \ell \text{ and } t \in (0, T],$$

$$(2.16b) \quad \langle y^\ell(0), u_i^n \rangle_W = \langle y_0, u_i^n \rangle_W \quad \text{for } i = 1 \dots, \ell.$$

To solve (2.16) we apply the implicit Euler method. By Y_j we denote an approximation for y^ℓ at the time t_j , $1 \leq j \leq n$. Then, the discrete system for the sequence $\{Y_j\}_{j=1}^n$ in $V_n^\ell = \text{span}\{u_1^n, \dots, u_\ell^n\}$ looks like

$$(2.17a) \quad \left\langle \frac{Y_j - Y_{j-1}}{\Delta t}, u_i^n \right\rangle_W = \langle AY_j + f(t, Y_j), u_i^n \rangle_W \quad \text{for } i = 1 \dots, \ell, \quad 2 \leq j \leq n,$$

$$(2.17b) \quad \langle Y_1, u_i^n \rangle_W = \langle y_0, u_i^n \rangle_W \quad \text{for } i = 1 \dots, \ell.$$

We are interested in estimating

$$\sum_{j=1}^n \alpha_j \|y(t_j) - Y_j\|_W^2.$$

Let us introduce the projection $\mathcal{P}_n^\ell : \mathbb{R}^m \rightarrow V_n^\ell$ by

$$(2.18) \quad \mathcal{P}_n^\ell = \sum_{i=1}^{\ell} \langle u, u_i^n \rangle_W u_i^n \quad \text{for } u \in \mathbb{R}^m.$$

It follows that \mathcal{P}_n^ℓ is linear and bounded (and therefore continuous). In particular, $\|\mathcal{P}_n^\ell\|_{L(\mathbb{R}^m)} = 1$.

We shall make use of the decomposition

$$y(t_j) - Y_j = y(t_j) - \mathcal{P}_n^\ell y(t_j) + \mathcal{P}_n^\ell y(t_j) - Y_j = \varrho_j^\ell + \vartheta_j^\ell,$$

where $\varrho_j^\ell = y(t_j) - \mathcal{P}_n^\ell y(t_j)$ and $\vartheta_j^\ell = \mathcal{P}_n^\ell y(t_j) - Y_j$. Note that

$$\sum_{j=1}^n \alpha_j \left\| y(t_j) - \sum_{i=1}^{\ell} \langle y(t_j), u_i^n \rangle_W u_i^n \right\|_W^2 = \sum_{j=1}^n \alpha_j \|y(t_j) - \mathcal{P}_n^\ell y(t_j)\|_W^2 = \sum_{j=1}^n \alpha_j \|\varrho_j^\ell\|_W^2.$$

Since $\{u_i^n\}_{i=1}^{\ell}$ is the POD basis of rank ℓ , we have

$$(2.19) \quad \sum_{j=1}^n \alpha_j \|\varrho_j^\ell\|_W^2 = \sum_{i=\ell+1}^m \lambda_i^n.$$

Next we estimate the terms ϑ_j^ℓ . Using the notation $\bar{\partial}\vartheta_j^\ell = (\vartheta_j^\ell - \vartheta_{j-1}^\ell)/\Delta t$ for $2 \leq j \leq n$ we obtain by (1.30a) and (2.17a)

$$\begin{aligned}
(2.20) \quad \langle \bar{\partial}\vartheta_j^\ell, u_i \rangle &= \left\langle \mathcal{P}_n^\ell \left(\frac{y(t_j) - y(t_{j-1})}{\Delta t} \right) - \frac{Y_j - Y_{j-1}}{\Delta t}, u_i^n \right\rangle_W \\
&= \langle \dot{y}(t_j) - (AY_j + f(t_j, Y_j)), u_i^n \rangle_W \\
&\quad + \left\langle \mathcal{P}_n^\ell \left(\frac{y(t_j) - y(t_{j-1})}{\Delta t} \right) - \dot{y}(t_j), u_i^n \right\rangle_W \\
&= \langle A(y(t_j) - Y_j) + f(t_j, y(t_j)) - f(t_j, Y_j), u_i^n \rangle_W \\
&\quad + \left\langle \mathcal{P}_n^\ell \left(\frac{y(t_j) - y(t_{j-1})}{\Delta t} \right) - \frac{y(t_j) - y(t_{j-1})}{\Delta t}, u_i^n \right\rangle_W \\
&\quad + \left\langle \frac{y(t_j) - y(t_{j-1})}{\Delta t} - \dot{y}(t_j), u_i^n \right\rangle_W \\
&= \langle A(y(t_j) - Y_j) + f(t_j, y(t_j)) - f(t_j, Y_j) + z_j^\ell + w_j^\ell, u_i^n \rangle_W
\end{aligned}$$

for $1 \leq i \leq \ell$ and $2 \leq j \leq n$, where

$$z_j^\ell = \mathcal{P}_n^\ell \left(\frac{y(t_j) - y(t_{j-1})}{\Delta t} \right) - \frac{y(t_j) - y(t_{j-1})}{\Delta t}, \quad w_j^\ell = \frac{y(t_j) - y(t_{j-1})}{\Delta t} - \dot{y}(t_j).$$

Multiplying (2.20) by $\langle \vartheta_j^\ell, u_i^n \rangle_W$ and adding all ℓ equations we arrive at

$$(2.21) \quad \langle \bar{\partial}\vartheta_j^\ell, \vartheta_j^\ell \rangle = \langle A(y(t_j) - Y_j) + f(t_j, y(t_j)) - f(t_j, Y_j) + z_j^\ell + w_j^\ell, \vartheta_j^\ell \rangle_W$$

for $j = 2, \dots, n$. Note that

$$\begin{aligned}
2 \langle u - \tilde{u}, u \rangle_W &= 2 \|u\|_W^2 - 2 \langle \tilde{u}, u \rangle_W \\
&= \|u\|_W^2 + \|u\|_W^2 - 2 \langle \tilde{u}, u \rangle_W + \|\tilde{u}\|_W^2 - \|\tilde{u}\|_W^2 \\
&= \|u\|_W^2 - \|\tilde{u}\|_W^2 + \|u - \tilde{u}\|_W^2
\end{aligned}$$

for all $u, \tilde{u} \in \mathbb{R}^m$. Choosing $u = \vartheta_j^\ell$ and $\tilde{u} = \vartheta_{j-1}^\ell$ we infer from (2.21)

$$(2.22) \quad 2 \langle \bar{\partial}\vartheta_j^\ell, \vartheta_j^\ell \rangle = \frac{1}{\Delta t} \left(\|\vartheta_j^\ell\|_W^2 - \|\vartheta_{j-1}^\ell\|_W^2 + \|\vartheta_j^\ell - \vartheta_{j-1}^\ell\|_W^2 \right).$$

Inserting (2.22) into (2.21) and using the Cauchy-Schwarz inequality we obtain

$$\begin{aligned}
\|\vartheta_j^\ell\|_W^2 &\leq \|\vartheta_{j-1}^\ell\|_W^2 + \Delta t \|A\| (\|\varrho_j^\ell\|_W + \|\vartheta_j^\ell\|_W) \|\vartheta_j^\ell\|_W \\
&\quad + \Delta t \left(\|f(t_j, y(t_j)) - f(t_j, Y_j)\|_W + \|z_j^\ell\|_W + \|w_j^\ell\|_W \right) \|\vartheta_j^\ell\|_W.
\end{aligned}$$

Suppose that f is Lipschitz-continuous with respect to the second argument. Then there exists a constant $L_f \geq 0$ such that

$$\|f(t_j, y(t_j)) - f(t_j, Y_j)\|_W \leq L_f \|y(t_j) - Y_j\|_W \quad \text{for } j = 2, \dots, n.$$

Hence, by Young's inequality we find

$$\|\vartheta_j^\ell\|_W^2 \leq \|\vartheta_{j-1}^\ell\|_W^2 + \Delta t \left(c_1 \|\varrho_j^\ell\|_W^2 + c_2 \|\vartheta_j^\ell\|_W^2 + \|z_j^\ell\|_W^2 + \|w_j^\ell\|_W^2 \right),$$

where $c_1 = \max\{\|A\|, L_f\}$ and $c_2 = \max\{3\|A\|, 3L_f, 2\}$. Suppose that

$$(2.23) \quad 0 < \Delta t \leq \frac{1}{2c_2}$$

holds. With (2.23) holding we have

$$0 \leq 1 - 2c_2\Delta t < 1 - c_2\Delta t \quad \text{and} \quad 1 - c_2\Delta t \geq 1 - \frac{1}{2} = \frac{1}{2}.$$

Thus,

$$(2.24) \quad \frac{1}{1 - c_2\Delta t} = \frac{1 - c_2\Delta t + c_2\Delta t}{1 - c_2\Delta t} = 1 + \frac{c_2\Delta t}{1 - c_2\Delta t} \leq 1 + 2c_2\Delta t$$

Using (2.24) we infer that

$$\|\vartheta_j^\ell\|_W^2 \leq (1 + 2c_2\Delta t) \left(\|\vartheta_{j-1}^\ell\|_W^2 + \Delta t (\|z_j^\ell\|_W^2 + \|w_j^\ell\|_W^2 + c_1 \|\varrho_j^\ell\|_W^2) \right).$$

Summation on j yields

$$\|\vartheta_j^\ell\|_W^2 \leq (1 + 2c_2\Delta t)^j \left(\|\vartheta_0^\ell\|_W^2 + \Delta t \sum_{k=1}^j (\|z_k^\ell\|_W^2 + \|w_k^\ell\|_W^2 + c_1 \|\varrho_k^\ell\|_W^2) \right).$$

Note that

$$(1 + 2c_2\Delta t)^j = \left(1 + \frac{2c_2j\Delta t}{j} \right)^j \leq e^{2c_2j\Delta t}.$$

Thus,

$$\|\vartheta_j^\ell\|_W^2 \leq e^{2c_2j\Delta t} \left(\|\vartheta_0^\ell\|_W^2 + \Delta t \sum_{k=1}^j (\|z_k^\ell\|_W^2 + \|w_k^\ell\|_W^2 + c_1 \|\varrho_k^\ell\|_W^2) \right).$$

We next estimate the term involving w_k^ℓ :

$$\begin{aligned} \Delta t \sum_{k=1}^j \|w_k^\ell\|_W^2 &= \Delta t \sum_{k=1}^j \left\| \frac{y(t_k) - y(t_{k-1})}{\Delta t} - \dot{y}(t_k) \right\|_W^2 \\ &= \frac{1}{\Delta t} \sum_{k=1}^j \|y(t_k) - y(t_{k-1}) - \Delta t \dot{y}(t_k)\|_W^2 \\ &= \frac{1}{\Delta t} \sum_{k=1}^j \left\| \int_{t_{k-1}}^{t_k} (t_{k-1} - s) \ddot{y}(s) \, ds \right\|_W^2 \\ &\leq \frac{1}{\Delta t} \sum_{k=1}^j \int_{t_{k-1}}^{t_k} |t_{k-1} - s|^2 \, ds \int_{t_{k-1}}^{t_k} \|\ddot{y}(s)\|_W^2 \, ds \\ &\leq \frac{(\Delta t)^2}{3} \sum_{k=1}^j \|\ddot{y}\|_{L^2(t_{k-1}, t_k; \mathbb{R}^m)}^2 = \frac{(\Delta t)^2}{3} \|\ddot{y}\|_{L^2(0, t_j; \mathbb{R}^m)}^2. \end{aligned}$$

The term z_k^ℓ can be estimated as follows:

$$\begin{aligned}
\|z_k^\ell\|_W^2 &= \left\| \mathcal{P}_n^\ell \left(\frac{y(t_k) - y(t_{k-1})}{\Delta t} \right) - \frac{y(t_k) - y(t_{k-1})}{\Delta t} \right\|_W^2 \\
&= \left\| \mathcal{P}_n^\ell \left(\frac{y(t_k) - y(t_{k-1})}{\Delta t} \right) - \mathcal{P}_n^\ell \dot{y}(t_k) + \mathcal{P}_n^\ell \dot{y}(t_k) - \frac{y(t_k) - y(t_{k-1})}{\Delta t} \right\|_W^2 \\
&\leq 2 \|\mathcal{P}_n^\ell\|_{L(\mathbb{R}^m)}^2 \left\| \frac{y(t_k) - y(t_{k-1})}{\Delta t} - \dot{y}(t_k) \right\|_W^2 \\
&\quad + 2 \left\| \mathcal{P}_n^\ell \dot{y}(t_k) - \dot{y}(t_k) + \dot{y}(t_k) - \frac{y(t_k) - y(t_{k-1})}{\Delta t} \right\|_W^2 \\
&\leq 2 \|w_k^\ell\|_W^2 + 4 \|\mathcal{P}_n^\ell \dot{y}(t_k) - \dot{y}(t_k)\|_W^2 + 4 \left\| \dot{y}(t_k) - \frac{y(t_k) - y(t_{k-1})}{\Delta t} \right\|_W^2 \\
&= 4 \|\mathcal{P}_n^\ell \dot{y}(t_k) - \dot{y}(t_k)\|_W^2 + 6 \|w_k^\ell\|_W^2.
\end{aligned}$$

Recall that $\Delta t \leq 2\alpha_k$ for $1 \leq k \leq n$. Hence,

$$\Delta t \sum_{k=1}^j \|z_k^\ell\|_W^2 \leq 8 \sum_{k=1}^n \alpha_k \|\mathcal{P}_n^\ell \dot{y}(t_k) - \dot{y}(t_k)\|_W^2 + 2(\Delta t)^2 \|\ddot{y}\|_{L^2(0, t_j; \mathbb{R}^m)}^2.$$

Further, $\vartheta_0^\ell = \mathcal{P}_n^\ell y_0 - Y_1 = 0$ and $0 \leq j\Delta t \leq T$ for $j = 0, \dots, n-1$. Summarizing

$$\|\vartheta_j^\ell\|_W^2 \leq c_3 \left(\sum_{k=1}^n 8\alpha_k \left(\|\mathcal{P}_n^\ell \dot{y}(t_k) - \dot{y}(t_k)\|_W^2 + 2c_1 \|\varrho_k^\ell\|_W^2 \right) + \frac{7}{3} (\Delta t)^2 \|\ddot{y}\|_{L^2(0, t_j; \mathbb{R}^m)}^2 \right),$$

where $c_3 = e^{2c_2 T} \max\{7/3, 2c_1, 8\}$ is independent of ℓ and $\{t_j\}_{j=1}^n$. From $\sum_{k=1}^n \alpha_k = T$ and (2.19) we infer

$$\begin{aligned}
\sum_{j=1}^n \alpha_j \|\vartheta_j^\ell\|_W^2 &\leq c_3 T \left(\sum_{j=1}^n \alpha_j \left(\|\mathcal{P}_n^\ell \dot{y}(t_j) - \dot{y}(t_j)\|_W^2 + \|\varrho_j^\ell\|_W^2 \right) \right. \\
(2.25) \quad &\quad \left. + (\Delta t)^2 \|\ddot{y}\|_{L^2(0, T; \mathbb{R}^m)}^2 \right) \\
&\leq c_4 \left(\sum_{i=\ell+1}^m \left(\lambda_i^n + \sum_{j=1}^n \alpha_j |\langle \dot{y}(t_j), u_i^n \rangle_W|^2 \right) + (\Delta t)^2 \right)
\end{aligned}$$

with $c_4 = c_3 T \max\{1, \|\ddot{y}\|_{L^2(0, T; \mathbb{R}^m)}^2\}$.

Theorem 2.4. *Let $y \in C([0, T]; \mathbb{R}^m) \cap C^1(0, T; \mathbb{R}^m)$ be the unique solution to (1.30) satisfying $\ddot{y} \in L^2(0, T; \mathbb{R}^m)$ and $\ell \in \{1, \dots, m\}$ be fixed. Suppose that $\{u_i^n\}_{i=1}^\ell$ is a POD basis of rank ℓ solving $(\hat{\mathbf{P}}_W^{n, \ell})$. Assume that (2.17) possesses a unique solution $\{Y_j\}_{j=1}^n$. Then there exists a constant $C > 0$ such that*

$$\sum_{j=1}^n \alpha_j \|y(t_j) - Y_j\|_W^2 \leq C \left((\Delta t)^2 + \sum_{i=\ell+1}^m \left(\lambda_i^n + \sum_{j=1}^n \alpha_j |\langle \dot{y}(t_j), u_i^n \rangle_W|^2 \right) \right)$$

provided Δt is sufficiently small and f is Lipschitz-continuous with respect to the second argument.

Proof. The claim follows directly from (2.19), (2.25), and

$$\begin{aligned} \sum_{j=1}^n \alpha_j \|y(t_j) - Y_j\|_W^2 &\leq 2 \sum_{j=1}^n \alpha_j \left(\|\vartheta_j^\ell\|_W^2 + \|\varrho_j^\ell\|_W^2 \right) \\ &\leq 2c_4 \left(\sum_{i=\ell+1}^m \left(\lambda_i^n + \sum_{j=1}^n |\langle \dot{y}(t_j), u_i^n \rangle_W|^2 \right) + (\Delta t)^2 \right) \\ &\quad + 2 \sum_{i=\ell+1}^m \lambda_i^n \end{aligned}$$

provided Δt is sufficiently small and f is Lipschitz-continuous with respect to the second argument. \square

Remark 2.5. Compared to the estimate in Theorem 2.1 we observe the term

$$(2.26) \quad \sum_{j=1}^n \alpha_j |\langle \dot{y}(t_j), u_i^n \rangle_W|^2$$

instead of the term

$$(2.27) \quad \int_0^T |\langle \dot{y}(t), u_i \rangle_W|^2 dt.$$

Note that (2.26) is the trapezoidal approximation of (2.27). Furthermore, the error $O((\Delta t)^2)$ appears in the estimate of Theorem 2.4 due to the Euler method. \diamond

Next we address the fact that the eigenvalues $\{\lambda_i^n\}_{i=1}^m$ and the associated eigenvectors $\{u_i^n\}$ (i.e., the POD basis) depend on the chosen time grid $\{t_j\}_{j=1}^n$. We apply the asymptotic theory presented in Section 1.3. Then, it follows from Theorem 1.14 that there exists a number $\bar{n} \in \mathbb{N}$ satisfying

$$\begin{aligned} \sum_{i=\ell+1}^m \lambda_i^n &\leq 2 \sum_{i=\ell+1}^m \lambda_i, \\ \sum_{i=\ell+1}^m \sum_{j=1}^n \alpha_j |\langle \dot{y}(t_j), u_i^n \rangle_W|^2 &\leq 2 \sum_{i=\ell+1}^m \int_0^T |\langle \dot{y}(t), u_i \rangle_W|^2 dt \end{aligned}$$

for $n \geq \bar{n}$ provided $\sum_{i=\ell+1}^m \lambda_i \neq 0$ and $\int_0^T |\langle \dot{y}(t), u_i \rangle_W|^2 dt \neq 0$ hold. Thus, we infer from Theorems 2.1 and 2.4 the following result.

Theorem 2.6. *Let all hypothesis of Theorems 1.14, 2.1 and 2.4 be satisfied. If $\int_0^T |\langle \dot{y}(t), u_i \rangle_W|^2 dt \neq 0$, then there exists a constant $C > 0$ and a number $\bar{n} \in \mathbb{N}$ such that*

$$\sum_{j=1}^n \alpha_j \|y(t_j) - Y_j\|_W^2 \leq C \left((\Delta t)^2 + \sum_{i=\ell+1}^m \left(\lambda_i + \int_0^T |\langle \dot{y}(t), u_i \rangle|^2 dt \right) \right)$$

for all $n \geq \bar{n}$.

2.3. Exercises.

- 2.1) Prove the *Gronwall lemma*: For $T > 0$ let $\eta : [0, T] \rightarrow \mathbb{R}$ be a non-negative, differentiable function satisfying

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

where φ and ψ are real-valued, non-negative, integrable functions on $[0, T]$. Then

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

In particular, if

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

show that $\eta = 0$ holds in $[0, T]$.

- 2.2) Show that the operator \mathcal{P}_n^ℓ defined in (2.18) is linear, bounded and satisfies $\|\mathcal{P}_n^\ell\|_{L(\mathbb{R}^m)} = 1$.
- 2.3) Prove that the first-order necessary optimality condition for (2.12) is given by $\tilde{\mathcal{R}}\tilde{u}_i = \tilde{\lambda}_i\tilde{u}_i$, $1 \leq i \leq \ell$.
- 2.4) Show that $\tilde{\mathcal{R}}$ is linear, bounded, self-adjoint and non-negative provided $y \in H^1(0, T; \mathbb{R}^m)$, i.e.,

$$\int_0^T \|y(t)\|_W^2 + \|\dot{y}(t)\|_W^2 dt < \infty$$

holds.

3. The linear-quadratic control problem

In this section we introduce the optimal state-feedback and the linear-quadratic regulator (LQR) problem. Utilizing dynamic programming necessary optimality conditions are derived. It turns out that for the LQR problem the state-feedback solution can be determined by solving a differential matrix Riccati equation. The presented theory is taken from the book [2].

3.1. The LQR problem. The goal is to find a state-feedback control law of the form

$$u(t) = -Kx(t) \quad \text{for } t \in [0, T]$$

with $u : [0, T] \rightarrow \mathbb{R}^{m_u}$, $x : [0, T] \rightarrow \mathbb{R}^{m_x}$, $K \in \mathbb{R}^{m_u \times m_x}$ so that u minimizes the quadratic cost functional

$$(3.1a) \quad J(x, u) = \int_0^T x(t)^T Qx(t) + u(t)^T Ru(t) dt + x(T)^T Mx(T),$$

where the state x and the control u are related by the linear initial value problem

$$(3.1b) \quad \dot{x}(t) = Ax(t) + Bu(t) \text{ for } t \in (0, T] \quad \text{and} \quad x(0) = x_0.$$

In (3.1a) the matrices $Q, M \in \mathbb{R}^{m_x \times m_x}$ are symmetric, positive semi-definite, $R \in \mathbb{R}^{m_u \times m_u}$ is symmetric, positive definite and in (3.1b) we have $A \in \mathbb{R}^{m_x \times m_x}$, $B \in \mathbb{R}^{m_x \times m_u}$ and $x_0 \in \mathbb{R}^{m_x}$. The final time T is fixed, but the final state $x(T)$ is free. Thus, we aim to track the state to the state $\bar{x} = 0$ as good as possible. The terms $x(t)^T Qx(t)$ and $x(T)^T Mx(T)$ are measures for the control accuracy and the term $u(t)^T Ru(t)$ measures the control effort. Problem (3.1) is called the *linear-quadratic regulator problem (LQR problem)*.

3.2. The Hamilton-Jacobi-Bellman equation. In this section we derive first-order necessary optimality conditions for the LQR problem. Since generalizing the problem to a non-linear problem does not cause more difficulties in the deviation, we consider the problem to find a state-control feedback control law

$$u(t) = \Phi(x(t), t), \quad t \in [0, T],$$

such that the cost-functional

$$(3.2a) \quad J_t(x, u) = \int_t^T L(x(s), u(s), s) ds + g(x(T))$$

is minimized subject to the non-linear system dynamics

$$(3.2b) \quad \dot{x}(s) = F(x(s), u(s), s) \text{ for } s \in (0, T] \quad \text{and} \quad x(t) = x_t.$$

We suppose that the functions $L : \mathbb{R}^{m_x} \times \mathbb{R}^{m_u} \times [0, T] \rightarrow [0, \infty)$ and $g : \mathbb{R}^{m_x} \rightarrow [0, \infty)$ satisfy

$$L(0, 0, s) = 0 \text{ for } s \in [0, T] \quad \text{and} \quad g(0) = 0$$

Moreover, let $F : \mathbb{R}^{m_x} \times \mathbb{R}^{m_u} \times [0, T] \rightarrow \mathbb{R}^{m_x}$ be continuous and locally Lipschitz-continuous with respect to the variable x . Moreover, $x_t \in \mathbb{R}^{m_x}$ holds. To derive optimality conditions we use the so-called *Bellman principle* (or *dynamic programming principle*). The essential assumption is that the system can be characterized by its state $x(t)$ at the time $t \in [0, T]$ which completely summarizes the effect of all $u(s)$ for $0 \leq s \leq t$. The dynamic programming principle was first proposed by Bellman [1].

Theorem 3.1 (Bellman principle). *Let $t \in [0, T]$. If $u^*(s)$ is optimal for $s \in [t, T]$ and x^* is the associated optimal state, starting at the state $x_t \in \mathbb{R}^{m_x}$, then $u^*(s)$ is also optimal over the subinterval $[t + \Delta t, T]$ for any $\Delta t \in [0, T - t]$ starting at $x_{t+\Delta t} = x^*(t + \Delta t)$.*

Proof. We show Theorem 3.1 by contradiction. Suppose that there exists a control u^{**} so that

$$(3.3) \quad \begin{aligned} & \int_{t+\Delta t}^T L(x^{**}(s), u^{**}(s), s) ds + g(x^{**}(T)) \\ & < \int_{t+\Delta t}^T L(x^*(s), u^*(s), s) ds + g(x^*(T)), \end{aligned}$$

where

$$\dot{x}^*(s) = F(x^*(s), u^*(s), s) \quad \text{and} \quad \dot{x}^{**}(s) = F(x^{**}(s), u^{**}(s), s)$$

hold for $s \in [t + \Delta t, T]$. We define the control

$$(3.4) \quad u(s) = \begin{cases} u^*(s) & \text{if } s \in [t, t + \Delta t], \\ u^{**}(s) & \text{if } s \in (t + \Delta t, T]. \end{cases}$$

By $x(s)$ we denote the state satisfying $\dot{x}(s) = F(x(s), u(s), s)$ for $s \in [t, T]$ and $x(t) = x_t$. Then we derive from (3.3) and (3.4) that

$$\begin{aligned}
 & \int_t^T L(x(s), u(s), s) \, ds + g(x(T)) \\
 (3.5) \quad &= \int_t^{t+\Delta t} L(x^*(s), u^*(s), s) \, ds + \int_{t+\Delta t}^T L(x^{**}(s), u^{**}(s), s) \, ds + g(x^{**}(T)) \\
 &< \int_t^{t+\Delta t} L(x^*(s), u^*(s), s) \, ds + \int_{t+\Delta t}^T L(x^*(s), u^*(s), s) \, ds + g(x^*(T)) \\
 &= \int_t^T L(x^*(s), u^*(s), s) \, ds + g(x^*(T)).
 \end{aligned}$$

Recall that $u^*(s)$ is optimal for $s \in [t, T]$ by assumption. From (3.5) it follows that the control u given by (3.4) yields a smaller value of the cost functional. This is a contradiction. \square

Next we derive the Hamilton-Jacobi-Bellman equation for (3.2). Let $V^* : \mathbb{R}^{m_x} \times [0, T] \rightarrow \mathbb{R}$ denote the minimal value function given by

$$\begin{aligned}
 & V^*(x_t, t) \\
 (3.6) \quad &= \min_{u: [t, T] \rightarrow \mathbb{R}^{m_u}} \left\{ J_t(x, u) \mid \dot{x}(s) = F(x(s), u(s), s), \, s \in (t, T] \text{ and } x(t) = x_t \right\}
 \end{aligned}$$

for $(x_t, t) \in \mathbb{R}^{m_x} \times [0, T]$, where

$$J_t(x, u) = \int_t^T L(x(s), u(s), s) \, ds + g(x(T)).$$

From the linearity of the integral and (3.6) we conclude

$$\begin{aligned}
 & V^*(x_t, t) \\
 (3.7) \quad &= \min_{u: [t, t+\Delta t] \rightarrow \mathbb{R}^{m_u}} \left\{ \int_t^{t+\Delta t} L(x(s), u(s), s) \, ds + V^*(x(t+\Delta t), t+\Delta t) \mid \right. \\
 & \quad \left. \dot{x}(s) = F(x(s), u(s), s), \, s \in (t, t+\Delta t] \text{ and } x(t) = x_t \right\}
 \end{aligned}$$

for $(x_t, t) \in \mathbb{R}^{m_x} \times [0, T - \Delta t]$, where we have used the Bellman principle. Thus, by using the Bellman principle the problem of finding an optimal control over the interval $[t, T]$ has been reduced to the problem of finding an optimal control over the interval $[t, t + \Delta t]$.

Now we replace the integral in (3.7) by $L(x(t), u(t), t)\Delta t$, perform a Taylor approximation for $V^*(x(t+\Delta t), t+\Delta t)$ about the point $(x_t, t) = (x(t), t)$ and approximate $x(t+\Delta t) - x(t)$ by $F(x(t), u(t), t)\Delta t$. Then we find

$$\begin{aligned} V^*(x_t, t) &= \min_{u_t \in \mathbb{R}^{m_u}} \left\{ L(x_t, u_t, t)\Delta t + V^*(x_t, t) + \frac{\partial V^*}{\partial t}(x_t, t)\Delta t \right. \\ &\quad \left. + \nabla V^*(x_t, t)^T F(x_t, u_t, t)\Delta t + o(\Delta t) \right\} \\ &= V^*(x_t, t) + \frac{\partial V^*}{\partial t}(x_t, t)\Delta t \\ &\quad + \Delta t \min_{u_t \in \mathbb{R}^{m_u}} \left\{ L(x_t, u_t, t) + \nabla V^*(x_t, t)^T F(x_t, u_t, t) + \frac{o(\Delta t)}{\Delta t} \right\} \end{aligned}$$

for any $\Delta t > 0$. Thus,

$$-\frac{\partial V^*}{\partial t}(x_t, t) = \min_{u_t \in \mathbb{R}^{m_u}} \left\{ L(x_t, u_t, t) + \nabla V^*(x_t, t)^T F(x_t, u_t, t) + \frac{o(\Delta t)}{\Delta t} \right\}.$$

Taking the limit $\Delta t \rightarrow 0$ and using $V^*(x_t, T) = g(x_t)$ we obtain

$$(3.8a) \quad -\frac{\partial V^*}{\partial t}(x_t, t) = \min_{u_t \in \mathbb{R}^{m_u}} \left\{ L(x_t, u_t, t) + \nabla V^*(x_t, t)^T F(x_t, u_t, t) \right\}$$

for all $(x_t, t) \in \mathbb{R}^{m_x} \times [0, T)$ and

$$(3.8b) \quad V^*(x_t, T) = g(x_t)$$

for all $x_t \in \mathbb{R}^{m_x}$.

To solve (3.8) we proceed in two steps. First we compute a solution u_t to

$$u^*(t) = \operatorname{argmin}_{u_t \in \mathbb{R}^{m_u}} \left\{ L(x_t, u_t, t) + \nabla V^*(x_t, t)^T F(x_t, u_t, t) \right\}$$

and set

$$(3.9) \quad \Psi(\nabla V^*(x_t, t), x_t, t) = u^*(t),$$

which gives us a control law. Then we insert (3.9) into (3.8a) and solve

$$\begin{aligned} -\frac{\partial V^*}{\partial t}(x_t, t) &= L(x_t, \Psi(\nabla V^*(x_t, t), x_t, t), t) \\ &\quad + \nabla V^*(x_t, t)^T F(x_t, \Psi(\nabla V^*(x_t, t), x_t, t), t) \end{aligned}$$

for all $(x_t, t) \in \mathbb{R}^{m_x} \times [0, T)$. Finally, we can compute the gradient $\nabla V^*(x_t, t)$ and deduce the state-feedback law

$$u^*(t) = \Phi(x_t, t) = \Psi(\nabla V^*(x_t, t), x_t, t) \quad \text{for all } (x_t, t) \in \mathbb{R}^{m_x} \times [0, T).$$

Remark 3.2. 1) In general, it is not possible to solve (3.8) analytically. However, for the LQR problem we can derive an explicit solution for the state-feedback law.

2) Note that the Hamilton-Jacobi-Bellman equation are only necessary optimality conditions. \diamond

3.3. The state-feedback law for the LQR problem. For the LQR problem we have

$$L(x, u, t) = x^T Q x + u^T R u, \quad g(x) = x^T M x, \quad F(x, u, t) = A x + B u$$

for $(x, u, t) \in \mathbb{R}^{m_x} \times \mathbb{R}^{m_u} \times [0, T]$. For brevity, we focus on the situation, where the matrices A, B, Q, R are time-invariant. However, most of the presented theory also holds for the time-varying case.

First we minimize

$$x^T Q x + u^T R u + \nabla V^*(x, t)^T (A x + B u)$$

with respect to u . First-order necessary optimality conditions are given by

$$u^T R \tilde{u} + \tilde{u}^T R u + \nabla V^*(x, t)^T B \tilde{u} = 0 \quad \text{for all } \tilde{u} \in \mathbb{R}^{m_u}.$$

By assumption, R is symmetric and positive definite. Then we find

$$(2R u + B^T \nabla V^*(x, t))^T \tilde{u} = 0 \quad \text{for all } \tilde{u} \in \mathbb{R}^{m_u}$$

and

$$(3.10) \quad u^* = -\frac{1}{2} R^{-1} B^T \nabla V^*(x, t).$$

For the minimal value function V^* we make the quadratic ansatz

$$(3.11) \quad V^*(x, t) = x^T P(t) x, \quad P(t) \in \mathbb{R}^{m_x \times m_x} \text{ symmetric.}$$

Then, we have $\nabla V^*(x, t) = 2P(t)x$ so that

$$u^* = -R^{-1} B^T P(t) x.$$

Note that

$$\begin{aligned} \frac{\partial V^*}{\partial t}(x_t, t) &= x_t^T \dot{P}(t) x_t, \\ L(x_t, -R^{-1} B^T P(t) x_t, t) &= x_t^T Q x_t + x_t^T P(t) B R^{-1} B^T P(t) x_t \\ &= x_t^T (Q + P(t) B R^{-1} B^T P(t)) x_t, \\ F(x_t, -R^{-1} B^T P(t) x_t, t) &= A x_t - B R^{-1} B^T P(t) x_t = (A - B R^{-1} B^T P(t)) x_t, \\ \nabla V^*(x_t, t) &= 2P(t) x_t. \end{aligned}$$

Consequently,

$$\begin{aligned} -x_t^T \dot{P}(t) x_t &= -\frac{\partial V^*}{\partial t}(x_t, t) \\ &= x_t^T (Q + P(t) B R^{-1} B^T P(t)) x_t + (2P(t) x_t)^T (A - B R^{-1} B^T P(t)) x_t \end{aligned}$$

for all $x_t \in \mathbb{R}^{m_x}$, which yields

$$\begin{aligned} -x_t^T \dot{P}(t) x_t &= x_t^T (Q + P(t) B R^{-1} B^T P(t) + 2P(t) A - 2P(t) B R^{-1} B^T P(t)) x_t \\ &= x_t^T (2P(t) A + Q - P(t) B R^{-1} B^T P(t)) x_t \end{aligned}$$

for all $x_t \in \mathbb{R}^{m_x}$. From $P(t) = P(t)^T$ we deduce that

$$2x_t^T P(t) A x_t = x_t^T P(t) A x_t + x_t^T A^T P(t) x_t = x_t^T (A^T P(t) + P(t) A) x_t.$$

Using $V^*(x_t, T) = x_t^T P(T)x_t$ and (3.8b) we get

$$(3.12a)$$

$$-x_t^T \dot{P}(t)x_t = x_t^T (A^T P(t) + P(t)A + Q - P(t)BR^{-1}B^T P(t))x_t, \quad t \in [0, T)$$

$$(3.12b)$$

$$x_t^T P(T)x_t = x_t^T Mx_t.$$

Since (3.12) holds for all $x_t \in \mathbb{R}^{m_x}$ we obtain the following *matrix Riccati equation*

$$(3.13a) \quad -\dot{P}(t) = A^T P(t) + P(t)A + Q - P(t)BR^{-1}B^T P(t), \quad t \in [0, T)$$

$$(3.13b) \quad P(T) = M.$$

Finally, the optimal state-feedback is given by

$$u^*(t) = -K(t)x(t) \quad \text{and} \quad K(t) = R^{-1}B^T P(t).$$

Example 3.3. Let us consider the problem

$$\min \int_0^T |x(t)|^2 + |u(t)|^2 dt \quad \text{s.t.} \quad \dot{x}(t) = u(t) \quad \text{for } t \in (0, T].$$

Choosing $m_x = m_u = 1$, $A = M = 0$ and $B = Q = R = 1$ the matrix Riccati equation has the form

$$-\dot{P}(t) = 1 - P(t)^2 \quad \text{for } t \in [0, T) \quad \text{and} \quad P(T) = 0.$$

This scalar ordinary differential equation can be solved by separation of variables. Its solution is

$$P(t) = \frac{1 - e^{-2(T-t)}}{1 + e^{-2(T-t)}}$$

with the optimal control $u^*(t) = -P(t)x(t)$. ◇

3.4. Exercises. Let us consider the one-dimensional heat equation

$$(3.14a) \quad \theta_t(t, x) = \theta_{xx}(t, x) + u(t)\chi(x) \quad \text{for all } (t, x) \in Q = (0, T) \times \Omega,$$

$$(3.14b) \quad \theta_x(t, 0) = \theta_x(t, 1) = 0 \quad \text{for all } t \in (0, T),$$

$$(3.14c) \quad \theta(0, x) = \theta_0(x) \quad \text{for all } x \in \Omega = (0, 1) \subset \mathbb{R},$$

where $\theta = \theta(t, x)$ is the temperature, $u = u(t)$ the control input, $\chi = \chi(x)$ a given control shape function and $\theta_0 = \theta_0(x)$ a given initial condition.

- 3.1) Apply a classical finite difference approximation for the spatial variable x (compare Example 1.10) and derive the finite-dimensional initial value problem for the finite difference approximations.
- 3.2) Utilizing the trapezoidal rule deduce a discretization for the quadratic cost functional

$$J(\theta, u) = \frac{1}{2} \int_{\Omega} |\theta(T, x) - \theta_T(x)|^2 dx + \frac{\kappa}{2} \int_0^T |u(t)|^2 dt,$$

where $\theta_T = \theta_T(x)$ is a given desired terminal state and $\kappa > 0$ denotes a fixed regularization parameter.

- 3.3) Formulate the matrix Riccati equation for the discretized quadratic cost functional — see part 3.2) — and the discretized heat equation — see part 3.1).

- 3.4) What is the matrix Riccati equation in the case if we apply a POD Galerkin approximation instead of a finite difference discretization? How can we solve the matrix Riccati equation numerically?

4. Balanced truncation

Let us consider the linear time-invariant system

$$(4.1a) \quad \dot{x}(t) = Ax(t) + Bu(t) \text{ for } t \in (0, \infty) \quad \text{and} \quad x(0) = x_0,$$

$$(4.1b) \quad y(t) = Cx(t) \quad \text{for } t \in [0, \infty)$$

where $x(t) \in \mathbb{R}^{m_x}$ is called the system state, $x_0 \in \mathbb{R}^{m_x}$ is the initial condition of the system, $u(t) \in \mathbb{R}^{m_u}$ is said to be the system input and $y(t) \in \mathbb{R}^{m_y}$ is called the system output. The matrices A , B and C are assumed to have appropriate sizes.

It is helpful to analyze the linear system (4.1) through the Laplace transform.

Definition 4.1. *Let $f(t)$ be a time-varying vector. Then its Laplace transform is defined by*

$$(4.2) \quad \mathcal{L}[f](s) = \int_0^\infty e^{-st} f(t) dt \quad \text{for } s \in \mathbb{R}.$$

The Laplace transform is defined for those values of s , for which (4.2) converges.

The Laplace transforms of $u(t)$ and $y(t)$ are given by

$$\mathcal{L}[u](s) = \int_0^\infty e^{-st} u(t) dt \quad \text{and} \quad \mathcal{L}[y](s) = \int_0^\infty e^{-st} y(t) dt = C\mathcal{L}[x](s),$$

where we have used (4.1b). Note that

$$\begin{aligned} \mathcal{L}[\dot{x}](s) &= \int_0^\infty e^{-st} \dot{x}(t) dt = - \int_0^\infty (-s)e^{-st} x(t) dt + (e^{-st} x(t)) \Big|_{s=0}^{s=\infty} \\ &= s\mathcal{L}[x](s) - x_0. \end{aligned}$$

Therefore, the Laplace transform of the dynamical system (4.1a) yields

$$s\mathcal{L}[x](s) - x(0) = A\mathcal{L}[x](s) + B\mathcal{L}[u](s),$$

which gives

$$\mathcal{L}[x](s) = (sI - A)^{-1}x(0) + (sI - A)^{-1}B\mathcal{L}[u](s).$$

Thus,

$$(4.3) \quad \mathcal{L}[y](s) = C\mathcal{L}[x](s) = C(sI - A)^{-1}x(0) + C(sI - A)^{-1}B\mathcal{L}[u](s).$$

For $x(0) = 0$ the expression (4.3) reduces to

$$(4.4) \quad \mathcal{L}[y](s) = G(s)\mathcal{L}[u](s)$$

where

$$(4.5) \quad G(s) = C(sI - A)^{-1}B$$

is called the *transfer matrix* of the system.

Given the initial state x_0 and the input $u(t)$, the dynamical system response $x(t)$ and $y(t)$ for $t \in [0, T]$ satisfy

$$x(t) = e^{tA}x_0 + \int_0^t e^{(t-s)A}Bu(s) ds \quad \text{and} \quad y(t) = Cx(t).$$

If $u(t) = 0$ holds for all $t \in [0, T]$, we infer that

$$x(t) = e^{(t-t_1)A}x(t_1)$$

for any $t_1, t \in [0, T]$. The matrix $e^{(t-t_1)A}$ acts as a transformation from one state to another. Therefore, $\Phi(t, t_1) = e^{(t-t_1)A}$ is often called the *state transition matrix*.

Definition 4.2. *The dynamical system (4.1a) or the pair (A, B) are called controllable if for any $x_0 \in \mathbb{R}^{m_x}$ and final state $x_T \in \mathbb{R}^{m_x}$ there exists a (piecewise continuous) input u such that the solution to (4.1a) satisfies $x(T) = x_T$. Otherwise, (A, B) is said to be uncontrollable.*

Controllability can be verified as stated in the next theorem. For a proof we refer to [14].

Theorem 4.3. *The following claims are equivalent:*

- 1) (A, B) are controllable.
- 2) The controllability gramian

$$W_c(t) = \int_0^t e^{sA} B B^T e^{sA^T} ds$$

is positive definite for every $t > 0$.

- 3) The controllability matrix

$$C = [B \ AB \ A^2B \ \dots \ A^{m_x-1}B] \in \mathbb{R}^{m_x \times (m_x m_u)}$$

has full rank.

Definition 4.4. 1) *The unforced system $\dot{x}(t) = Ax(t)$ is called stable, if the eigenvalues of A are in the open left half plane, i.e., $\Re\lambda < 0$ for every eigenvalue λ . A matrix with this property is said to be stable or Hurwitz.*
 2) *The dynamical system (4.1a) or (A, B) are called stabilizable if there exists a state-feedback $u(t) = -Kx(t)$ so that $A - BK$ is stable.*

The next result, which is proved in [14], is a consequence of Theorem 4.3.

Theorem 4.5. *The following claims are equivalent:*

- 1) (A, B) are stabilizable.
- 2) The matrix $[A - \lambda I \ B] \in \mathbb{R}^{m_x \times (m_x + m_u)}$ has full row rank for all $\lambda \in \mathbb{C}$ with a negative real part, i.e., $\Re\lambda < 0$.

Let us now consider the dual notions of observability.

Definition 4.6. *The dynamical system (4.1) or (A, C) are called observable if for any $t_1 \in (0, T]$, the initial condition $x_0 \in \mathbb{R}^{m_x}$ can be determined from the time history of the input $u(t)$ and the output $y(t)$ in the interval $[0, t_1] \subset [0, T]$. Otherwise, the system or (A, C) is said to be unobservable.*

For a proof of the next theorem we refer the reader to [14].

Theorem 4.7. *The following claims are equivalent:*

- 1) (A, C) is observable.
- 2) The observability gramian

$$W_o(t) = \int_0^t e^{sA^T} C^T C e^{sA} ds$$

is positive definite for every $t > 0$.

(3) *The observability matrix*

$$\mathcal{O} = \begin{pmatrix} C \\ CA \\ \vdots \\ CA^{m_x-1} \end{pmatrix} \in \mathbb{R}^{(m_x m_y) \times m_x}$$

has full rank.

We set

$$W_c = \int_0^\infty e^{sA} B B^T e^{sA^T} ds \quad \text{and} \quad W_o = \int_0^\infty e^{sA^T} C^T C e^{sA} ds.$$

It can be proved that W_c and W_o can be determined numerically by solving the *Lyapunov equations*

$$(4.6a) \quad A W_c + W_c A^T + B B^T = 0 \in \mathbb{R}^{n_x \times n_x},$$

$$(4.6b) \quad A^T W_o + W_o A + C^T C = 0 \in \mathbb{R}^{n_x \times n_x}.$$

The controllability gramian is a measure to what degree each state is excited by an input. Suppose that $x_1, x_2 \in \mathbb{R}^{n_x}$ are two states with $\|x_1\|_{\mathbb{R}^{n_x}} = \|x_2\|_{\mathbb{R}^{n_x}}$. If $x_1^T W_c x_1 > x_2^T W_c x_2$ holds, then we say that the state x_1 is more controllable than x_2 . This means, it takes a smaller input to drive the system from x_0 to x_1 than to x_2 . It can be proved that the gramian W_c is positive definite if and only if all states are reachable with some input u . On the other hand, the observability gramian W_o is a measure to what degree each state excites future outputs y . Let x_0 be an initial state. If $u = 0$ holds, we have

$$\begin{aligned} \|y\|_{L^2(0, \infty; \mathbb{R}^{m_y})}^2 &= \int_0^\infty y(s)^T y(s) ds = \int_0^\infty x(s)^T C^T C x(s) ds \\ &= \int_0^\infty x_0^T e^{sA^T} C^T C e^{sA} x_0 ds = x_0^T W_o x_0. \end{aligned}$$

We say that the state x_1 is *more observable* than another state x_2 if the corresponding output $y_1 = C x_1$ yields a larger value of the L^2 -norm than for $y_2 = C x_2$

The gramians depend on the coordinates. Suppose that

$$(4.7) \quad x = T z$$

where $T \in \mathbb{R}^{n_x \times n_x}$ is a regular matrix. Then we obtain instead of (4.1) the system

$$(4.8a) \quad \dot{z}(t) = \tilde{A} z(t) + \tilde{B} u(t) \quad \text{for } t \in (0, \infty) \quad \text{and} \quad z(0) = z_0,$$

$$(4.8b) \quad y(t) = \tilde{C} z(t) \quad \text{for } t \in [0, \infty)$$

with

$$\tilde{A} = T^{-1} A T, \quad \tilde{B} = T^{-1} B, \quad \tilde{C} = C T, \quad z_0 = T^{-1} x_0.$$

Let W_c solve (4.6a). The controllability gramian \tilde{W}_c for (4.8) satisfies

$$\tilde{A} \tilde{W}_c + \tilde{W}_c \tilde{A}^T + \tilde{B} \tilde{B}^T = 0$$

i.e.,

$$(4.9) \quad T^{-1} A T \tilde{W}_c + \tilde{W}_c T^T A^T T^{-T} + T^{-1} B B^T T^{-T} = 0.$$

Multiplying (4.9) by \mathcal{T} from the left and by \mathcal{T}^T from the right yields

$$(4.10) \quad A\mathcal{T}\tilde{W}_c\mathcal{T}^T + \mathcal{T}\tilde{W}_c\mathcal{T}^T A^T + B B^T = 0.$$

From (4.6a) and (4.10) we infer that $W_c = \mathcal{T}\tilde{W}_c\mathcal{T}^T$ holds. Thus, the coordinate transformation (4.7) implies that the controllability gramian W_c is transformed as

$$W_c \mapsto \tilde{W}_c = \mathcal{T}^{-1}W_c\mathcal{T}^{-T}.$$

Now we suppose that W_o solves (4.6b). The observability gramian \tilde{W}_o for (4.8) satisfies

$$\tilde{A}^T\tilde{W}_o + \tilde{W}_o\tilde{A} + \tilde{C}^T\tilde{C} = 0$$

i.e.,

$$(4.11) \quad \mathcal{T}^T A^T \mathcal{T}^{-T} \tilde{W}_o + \tilde{W}_o \mathcal{T}^{-1} A \mathcal{T} + \mathcal{T}^T C^T C \mathcal{T} = 0.$$

Multiplying (4.9) by \mathcal{T}^{-T} from the left and by \mathcal{T}^{-1} from the right yields

$$(4.12) \quad A^T \mathcal{T}^{-T} \tilde{W}_o \mathcal{T}^{-1} + \mathcal{T}^{-T} \tilde{W}_o \mathcal{T}^{-1} A + C^T C = 0.$$

From (4.6b) and (4.12) we infer that $W_o = \mathcal{T}^{-T} \tilde{W}_o \mathcal{T}^{-1}$ holds. Thus, the coordinate transformation (4.7) implies that the observability gramian W_o is transformed as

$$W_o \mapsto \tilde{W}_o = \mathcal{T}^T W_o \mathcal{T}.$$

The goal is to find a transformation \mathcal{T} such that

$$(4.13) \quad \mathcal{T}^{-1}W_c\mathcal{T}^{-T} = \mathcal{T}^T W_o \mathcal{T} = \Sigma = \text{diag}(\sigma_1, \dots, \sigma_{m_x}).$$

The elements $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_{m_x}$ are called *Hankel singular values* of the system. They are independent of the coordinate system. It can be shown that a regular matrix \mathcal{T} which satisfies (4.13) exists if the system is controllable and observable, i.e., the matrices W_c and W_o are positive definite. The coordinate transformation \mathcal{T} is said to be a *balancing transformation*. Computing appropriately scaled eigenvalues of the product $W_c W_o$, the matrix \mathcal{T} can be determined. In the balanced coordinates, the states which are least influenced by the input u also have least influence on the output y . In *balanced truncation* the least controllable and observable states having little effect on the input-output performance are truncated.

Instead of (4.8) we only consider the system for the first $\ell \in \{1, \dots, m_x\}$ components of z :

$$(4.14a) \quad \dot{z}_\ell(t) = \tilde{A}_\ell z_\ell(t) + \tilde{B}_\ell u(t) \quad \text{for } t \in (0, \infty) \quad \text{and} \quad z_\ell(0) = z_{0\ell},$$

$$(4.14b) \quad y_\ell(t) = \tilde{C}_\ell z_\ell(t) \quad \text{for } t \in [0, \infty),$$

where

$$\tilde{A} = \left(\begin{array}{c|c} \tilde{A}_\ell & * \\ \hline * & * \end{array} \right), \quad \tilde{B} = \left(\begin{array}{c} \tilde{B}_\ell \\ * \end{array} \right), \quad \tilde{C} = (\tilde{C}_\ell \mid *), \quad z_{0\ell} = \left(\begin{array}{c} \tilde{z}_{0\ell} \\ * \end{array} \right),$$

and $\tilde{A}_\ell \in \mathbb{R}^{\ell \times \ell}$, $\tilde{B}_\ell \in \mathbb{R}^{\ell \times m_u}$, $\tilde{C}_\ell \in \mathbb{R}^{m_y \times \ell}$ and $z_{0\ell} \in \mathbb{R}^\ell$.

One big advantage of balanced truncation is that a-priori error bounds are known. These bounds are formulated for the transfer function. Suppose that $G(s) = C(sI - A)^{-1}B \in \mathbb{R}^{m_y \times m_u}$ is the transfer function of the system (4.1) and $G_\ell(s) = C_\ell(sI - A_\ell)^{-1}B_\ell \in \mathbb{R}^{m_y \times m_u}$ is the transfer function of the reduced system (4.14). Then we have

$$\|G - G_\ell\| = \max \left\{ \|(G - G_\ell)u\|_{L^2(0, \infty; \mathbb{R}^{m_y})} : \|u\|_{L^2(0, \infty; \mathbb{R}^{m_u})} = 1 \right\} > \sigma_{\ell+1}$$

and

$$\|G - G_\ell\| < 2 \sum_{i=\ell+1}^{m_x} \sigma_i.$$

Acknowledgement. The author wants to thank B. Gotthardt, M. Kahlbacher, M. Kanitsar, and M. Müller for their careful reading of the scriptum and their comments which improve the scriptum significantly.

REFERENCES

- [1] R.E. Bellman. The theory of dynamic programming. *Proc. Nat. Acad. Sci.*, USA, 38:716-719, 1952.
- [2] P. Dorato, C. Abdallah, and V. Cerone. *Linear-Quadratic Control. An Introduction*. Prentice Hall, Englewood Cliffs, New Jersey 07632, 1995.
- [3] P. Holmes, J.L. Lumley, and G. Berkooz. *Turbulence, Coherent Structures, Dynamical Systems and Symmetry*. Cambridge Monographs on Mechanics, Cambridge University Press, 1996.
- [4] 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.
- [5] T. Kato. *Perturbation Theory for Linear Operators*. Springer-Verlag, Berlin, 1980.
- [6] 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.
- [7] K. Kunisch and S. Volkwein. Galerkin proper orthogonal decomposition methods for parabolic problems. *Numerische Mathematik*, 90:117-148, 2001.
- [8] 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.
- [9] 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.
- [10] B. Noble. *Applied Linear Algebra*. Englewood Cliffs, NJ : Prentice-Hall, 1969.
- [11] M. Reed and B. Simon. *Methods of Modern Mathematical Physics I: Functional Analysis*. Academic Press, New York, 1980.
- [12] C.W. Rowley. Model reduction for fluids, using balanced proper orthogonal decomposition. *Int. J. on Bifurcation and Chaos*, 15:997-1013, 2005.
- [13] L. Sirovich. Turbulence and the dynamics of coherent structures, parts I-III. *Quarterly of Applied Mathematicss*, XLV:561-590, 1987.
- [14] K. Zhou, J.C. Doyle, and K. Glover. *Robust and Optimal Control*. Prentice Hall, Upper Saddle River, New Jersey, 07458, 1996.

S. VOLKWEIN, KARL-FRANZENS-UNIVERSITÄT GRAZ, INSTITUT FÜR MATHEMATIK UND WISSENSCHAFTLICHES RECHNEN, HEINRICHSTRASSE 36, A-8010 GRAZ, AUSTRIA
E-mail address: stefan.volkwein@uni-graz.at