

Trust-Region RB Methods for PDE-Constrained Optimization and Optimal Input Design

Andrea Petrocchi* Matthias K. Scharrer** Stefan Volkwein*

* Department of Mathematics and Statistics, Universität Konstanz, Universitätsstr. 10, D-78457 Konstanz, Germany. (e-mail: andrea.petrocchi@uni-konstanz.de, stefan.volkwein@uni-konstanz.de)

** Virtual Vehicle Research GmbH, Inffeldgasse 21A, A-8010 Graz, Austria. (e-mail: matthias.scharrer@v2c2.at)

Abstract: In this paper we propose an algorithm for the bi-level optimal input design involving a parameter-dependent evolution problem. In the inner cycle a control is fixed and the parameter is optimized in order to minimize a cost function that measure the discrepancy from some data. In the outer cycle the found parameter is fixed and the control is now optimized in order to minimize a suitable measure of uncertainty of the parameters. The inner cycle uses a trust-region reduced basis approximation of the model with creation and enrichment of the reduced basis on-the-fly. Numerical examples illustrate the efficiency of the proposed approach.

Copyright © 2022 The Authors. This is an open access article under the CC BY-NC-ND license (<https://creativecommons.org/licenses/by-nc-nd/4.0/>)

Keywords: Reduced basis methods, a-posteriori error, trust-region optimization, evolution problems, optimal input design.

1. INTRODUCTION

Mathematical models based on partial differential equations (PDEs) are very useful in many fields, such as natural sciences, medicine and engineering. However, PDE models involve unknown parameters which have to be estimated from experiments, while we might have some control input variables that we are able to change in real time. Parameter estimation (PE) comprehends methods and algorithms are able to find or approximate the underlying parameters using empirical data or observations. It is often the case that the choice of a particular control speeds up or slows down the PE, hence being an important choice when the temporal or computational costs are limited.

In this work two problems are discussed. First, the step of PE is not trivial: usually such problems are ill-posed or non-convex, hence resulting in non-unique or non-global solutions. Even when well-posed, some applications deal with such high-dimensional discretizations that the solution is prohibitively costly. For this reason model order reduction methods are developed, where the “expensive” models are replaced by cheaper and less accurate surrogates. In our case the reduced basis (RB) is used, creating a reduced-order model on the optimization path in such a way that the local minimum found by our taylored optimization method in the reduced space is very close to a local minimum in the full space.

* The authors would like to acknowledge the financial support within the COMET K2 Competence Centers for Excellent Technologies from the Austrian Federal Ministry for Climate Action (BMK), the Austrian Federal Ministry for Digital and Economic Affairs (BMDW), the Province of Styria (Dept. 12) and the Styrian Business Promotion Agency (SFG).

The second problem belongs to the optimal input/experimental design; we want to find the “best” control, namely the one that gives us a “better” PE in a bi-level algorithm. In particular, assuming there are measurement error in the observations, the “best” control is the one that minimizes the uncertainty given by the parameter optimization.

2. PARAMETER ESTIMATION PROBLEM

We assume that all parameters are stacked in a vector $\mu \in \mathbb{R}^d$, and belong to the compact (admissible) set

$$\mathcal{P}_{\text{ad}} = \{\mu \in \mathbb{R}^d \mid \mu_i^a \leq \mu_i \leq \mu_i^b \text{ for } i = 1, \dots, d\}$$

The *state* variable y satisfies the evolution problem

$$\begin{aligned} \frac{d}{dt} \langle y(t), \varphi \rangle_H + a_\mu(y(t), \varphi) &= \langle f_\mu(t; u), \varphi \rangle_{V', V} \\ y(0) &= y_\circ \quad \text{in } H \end{aligned} \quad (1)$$

for all $\varphi \in V$ and $t \in (0, T]$, where $y_\circ \in H$, V , H are Hilbert spaces with $V \hookrightarrow H \hookrightarrow V'$ (Gelfand triple) and u denotes the *control* belonging to a convex, bounded, closed subset \mathcal{U}_{ad} of a Hilbert space \mathcal{U} . Further, the bilinear form $a_\mu : V \times V \rightarrow \mathbb{R}$ is continuous, coercive, symmetric and $f_\mu(\cdot; u) \in L^2(0; T; V')$ for any $(\mu, u) \in \mathcal{X}_{\text{ad}} = \mathcal{P}_{\text{ad}} \times \mathcal{U}_{\text{ad}}$.

It follows that (1) admits a unique solution $y = y_\mu \in \mathcal{Y} = W(0, T) = L^2(0, T; V) \cap H^1(0, T; V')$ for any $\mu \in \mathcal{P}_{\text{ad}}$; cf., e.g., Hinze et al. (2009). Furthermore, we assume that a_μ and f_μ depend affinely on the parameters:

$$a_\mu = \sum_{l=1}^{m_a} \vartheta_l^a(\mu) \hat{a}_l, \quad f_\mu(t; u) = \sum_{l=1}^{m_f} \vartheta_l^f(\mu) \hat{f}_l(t; u) \quad (2)$$

for any $(\mu, u) \in \mathcal{X}_{\text{ad}}$ and $t \in [0, T]$. Otherwise, we apply the *empirical interpolation method* to get approximations satisfying (2); cf., e.g., Hesthaven et al. (2016).

The goal is to estimate unknown model parameters $\mu \in \mathcal{P}_{\text{ad}}$, where the control input should be chosen in an optimal way explained later. Here, we suppose that $u \in \mathcal{U}_{\text{ad}}$ is fixed and consider

$$\min J(y, \mu) \text{ subject to } (y, \mu) \in \mathcal{Y} \times \mathcal{P}_{\text{ad}} \text{ satisfies (1) } \quad (\mathbf{P})$$

As (1) is uniquely solvable, we can define the *reduced cost* $\hat{J}(\mu) = J(y_\mu, \mu)$, where y_μ solves (1). Then, (\mathbf{P}) is equivalent to

$$\min \hat{J}(\mu) \text{ subject to (s.t.) } \mu \in \mathcal{P}_{\text{ad}} \quad (\hat{\mathbf{P}})$$

In our application the cost quantifies the discrepancy to a given desired (or observed) state $\hat{y} \in L^2(0, T; H)$:

$$\hat{J}(\mu) = 1 + \frac{1}{2} \int_0^T \|y_\mu(t) - \hat{y}(t)\|_H^2 dt + \frac{\sigma}{2} \|\mu - \hat{\mu}\|_2^2 \quad (3)$$

where $\|\cdot\|_2$ stands for the Euclidean norm, σ is a non-negative weight and $\hat{\mu} \in \mathbb{R}^d$ is a reference parameter. Existence of an optimal solution $\bar{\mu}$ can be ensured under continuity for $\mu \rightarrow a_\mu$ and $\mu \rightarrow f_\mu$, but – due to non-convexity – there are possibly many local solutions; cf., e.g., Hinze et al. (2009). A local optimal solution $\bar{\mu}$ to $(\hat{\mathbf{P}})$ is characterized by first-order necessary optimality conditions. Let $\bar{y} = y_{\bar{\mu}}$ be the optimal state associated with $\bar{\mu}$ and the *adjoint variable* $\bar{p} = p_{\bar{\mu}} \in \mathcal{Y}$ the solution of the *adjoint equation*

$$\begin{aligned} -\frac{d}{dt} \langle \bar{p}(t), \varphi \rangle_H + a_{\bar{\mu}}(\bar{p}(t), \varphi) &= \langle \hat{y}(t) - \bar{y}(t), \varphi \rangle_H \\ \bar{p}(T) &= 0 \quad \text{in } H \end{aligned}$$

for all $\varphi \in V$ and $t \in [0, T]$. Using the *adjoint approach* the gradient of the cost functional at $\bar{\mu}$ is given as

$$\begin{aligned} \nabla \hat{J}(\bar{\mu}) &= \sigma(\bar{\mu} - \hat{\mu}) + \int_0^T \sum_{i=1}^{m_a} \nabla \vartheta_i^a(\mu) \hat{a}_q(\bar{y}(t), \bar{p}(t)) dt \\ &\quad - \int_0^T \sum_{i=1}^{m_f} \nabla \vartheta_i^f(\mu) \langle \hat{f}_i(t; u), p(t) \rangle_{V', V} dt \in \mathbb{R}^d \end{aligned}$$

3. DISCRETIZATION

Next we introduce a high-dimensional discretization, called *full-order model* (FOM), which we assume to be accurate enough, but however expensive to solve. To reduce significantly the computational costs a further approximation is applied, the *reduced-order model* (ROM), faster to solve but less accurate.

3.1 FOM

Let $\varphi_1, \dots, \varphi_N \in V$ be given linearly independent functions and $V^N = \text{span}\{\varphi_1, \dots, \varphi_N\} \subset V$. The FOM for (1) reads: for given $\mu \in \mathcal{P}_{\text{ad}}$ the function $y_\mu^N(t) \in V^N$ solves

$$\begin{aligned} \frac{d}{dt} \langle y_\mu^N(t), \varphi \rangle_H + a_\mu(y_\mu^N(t), \varphi) &= \langle f_\mu(t; u), \varphi \rangle_{V', V} \\ y_\mu^N(0) &= \mathcal{P}^N y_\circ \end{aligned} \quad (4)$$

for all $\varphi \in V^N$ and $t \in (0, T]$, where $\mathcal{P}^N : H \rightarrow V^N$ is a projection. Due to $y_\mu^N(t) \in V^N$ we have

$$y_\mu^N(t) = \sum_{i=1}^N y_{\mu i}(t) \varphi_i \quad \text{for } t \in [0, T]$$

so that (4) reduces into finding the coefficient vector $y_\mu(t) = (y_{\mu i}(t))_{1 \leq i \leq N}$ solving

$$\begin{aligned} M \dot{y}_\mu(t) + A_\mu y_\mu(t) &= f_\mu(t; u), \quad t \in (0, T] \\ M y_\mu(0) &= y_\circ \end{aligned} \quad (5)$$

for $f_\mu(t; u) = ((f_\mu(t; u), \varphi_i)_{V', V}) \in \mathbb{R}^N$, $y_\circ = ((y_\circ, \varphi_i)_H) \in \mathbb{R}^N$, $A_\mu = ((a_\mu(\varphi_j, \varphi_i))) \in \mathbb{R}^{N \times N}$, $M = ((\langle \varphi_j, \varphi_i \rangle_H)) \in \mathbb{R}^{N \times N}$.

Remark 1. Due to (2), both A_μ and f_μ satisfy

$$A_\mu = \sum_{l=1}^{m_a} \vartheta_l^a(\mu) \hat{A}_l, \quad f_\mu(t; u) = \sum_{l=1}^{m_f} \vartheta_l^f(\mu) \hat{f}_l(t; u)$$

for $\hat{A}_l = ((\hat{a}_l(\varphi_j, \varphi_i)))$ and $\hat{f}_l(t; u) = ((\hat{f}_l(t; u), \varphi_i)_{V', V})$. \diamond

For solving (5) we apply the implicit Euler method for the time integration; cf., e.g., Quarteroni (2017). To simplify the presentation we utilize an equidistant time grid $t_k = (k-1)\Delta t$, $k = 1, \dots, K$ and $\Delta t = T/(K-1)$. Then, the problem is to find $\{y_\mu^k\}_{k=1}^K \subset \mathbb{R}^N$ solving

$$\begin{aligned} (M + \Delta t A_\mu) y_\mu^k &= M y_\mu^{k-1} + \Delta t f_\mu^k(u) \\ M y_\mu^1 &= y_\circ \end{aligned} \quad (6)$$

for $k = 2, \dots, K$ and $\mu \in \mathcal{P}_{\text{ad}}$ with $f_\mu^k(u) = f_\mu(t^k; u) \in \mathbb{R}^N$. Note that y_μ^k is the coefficient vector of the approximation $y_\mu^k \in V^N$ of $y_\mu^N(t_k) \in V^N$; cf. (8).

Remark 2. Similarly, we get the Galerkin approximations for the adjoint state, the reduced cost, and the reduced cost gradient. More precisely, the discrete adjoint sequence $\{p_\mu^k\}_{k=1}^K \subset \mathbb{R}^N$ satisfies:

$$\begin{aligned} (M + \Delta t A_\mu) p_\mu^k &= M(p_\mu^{k+1} + \Delta t(\hat{y}^k - y_\mu^k)) \\ p_\mu^K &= 0 \end{aligned}$$

where $\{y_\mu^k\}_{k=1}^K$ solves (6) and $\hat{y}^k = \sum_{i=1}^N \hat{y}_i^k \varphi_i \approx \hat{y}(t_k)$ holds. Further, the reduced cost \hat{J} in (3) is approximated by

$$\begin{aligned} \hat{J}^N(\mu) &= 1 + \frac{1}{2} \sum_{k=1}^K \alpha_k (y_\mu^k - \hat{y}^k)^\top M (y_\mu^k - \hat{y}^k) \\ &\quad + \frac{\sigma}{2} \|\mu - \hat{\mu}\|_2^2 \end{aligned} \quad (7)$$

where the α_k 's are trapezoidal weights. Then,

$$\begin{aligned} \nabla \hat{J}^N(\mu) &= \sum_{k=1}^K \alpha_k \left[\sum_{l=1}^{m_a} \nabla \vartheta_l^a(\mu) \hat{A}_l y_\mu^k - \sum_{l=1}^{m_f} \nabla \vartheta_l^f(\mu) \hat{f}_l^k(u) \right]^\top p_\mu^k \\ &\quad + \sigma(\mu - \hat{\mu}) \end{aligned}$$

is the gradient of the reduced cost. \diamond

3.2 ROM

For the moment we assume that for $\ell \ll N$ we are given an ℓ -dimensional reduced-order space $V^\ell = \text{span}\{\psi_i\}_{i=1}^\ell \subset V^N$, whose construction will be explained in Section 4. We define the coefficient matrix $\Psi \in \mathbb{R}^{N \times \ell}$ satisfying $\psi_j = \sum_{i=1}^N \Psi_{ij} \varphi_i$ for $j = 1, \dots, \ell$. Then, the ROM for (4) reads as follows: find $\{y_\mu^{k, \ell}\}_{k=1}^K \subset \mathbb{R}^\ell$ such that

$$\begin{aligned} (M^\ell + \Delta t A_\mu^\ell) y_\mu^{k, \ell} &= M^\ell y_\mu^{k-1, \ell} + \Delta t f_\mu^{k, \ell}(u) \\ M^\ell y_\mu^{1, \ell} &= y_\circ^\ell \end{aligned}$$

with $y_\circ^\ell = \Psi^\top y_\circ$, $M^\ell = \Psi^\top M \Psi$, $A_\mu^\ell = \sum_{l=1}^{m_a} \vartheta_l^a(\mu) \Psi^\top \hat{A}_l \Psi$ and $f_\mu^{k, \ell}(u) = \sum_{l=1}^{m_f} \vartheta_l^f(\mu) \Psi^\top \hat{f}_l^k(u)$. The solution $\{y_\mu^{k, \ell}\}_{k=1}^K$ uniquely exists (see Kunisch and Volkwein (2001)) and it is interpreted as a reduced-order approximation for $\{y_\mu^k\}_{k=1}^K$:

$$y_\mu^k \approx \tilde{y}_\mu^{k, \ell} := \Psi y_\mu^{k, \ell} \in \mathbb{R}^N$$

For $1 \leq k \leq K$ we also use the additional notations

$$y_\mu^k = \sum_{i=1}^N y_{\mu i}^k \varphi_i \in V^N \quad \text{and} \quad y_\mu^{k, \ell} = \sum_{i=1}^\ell y_{\mu i}^{k, \ell} \psi_i \in V^\ell \quad (8)$$

Similarly, setting $\hat{y}^{k,\ell} = \Psi^\top M \hat{y}^k$ we construct a reduced-order adjoint approximation $\{\mathbf{p}_\mu^{k,\ell}\}_{k=1}^K \subset \mathbb{R}^\ell$ by solving

$$\begin{aligned} (M^\ell + \Delta t A_\mu^\ell) \mathbf{p}_\mu^{k,\ell} &= M^\ell \mathbf{p}_\mu^{k+1,\ell} + \Delta t (\hat{y}^{k,\ell} - M^\ell y_\mu^{k,\ell}) \\ \mathbf{p}^{K,\ell} &= 0 \end{aligned}$$

Finally, the reduced cost function is approximated by

$$\begin{aligned} \hat{J}^\ell(\mu) &= 1 + \frac{1}{2} \sum_{k=1}^K \alpha_k \left[(y_\mu^{k,\ell})^\top M^\ell y_\mu^{k,\ell} - 2(y_\mu^{k,\ell})^\top \hat{y}^{k,\ell} \right] \\ &\quad + \frac{1}{2} \sum_{k=1}^K \alpha_k (\hat{y}^k)^\top M \hat{y}^k + \frac{\sigma}{2} \|\mu - \hat{\mu}\|_2^2 \end{aligned} \quad (9)$$

and its gradient is defined similarly.

3.3 A-posteriori RB error estimates

Now we present a-posteriori error estimates for the error of the ROM without evaluating the full-order solution. For our trust-region (TR) algorithm we will need estimates for the state and for the cost function.

Let us define the μ -dependent space-time energy norm for the sequence $\{y^k\}_{k=1}^K \subset V^N$ as

$$\|y^j\|_\mu = \left(\|y^j\|_H^2 + \Delta t \sum_{k=2}^j a_\mu(y^k, y^k) \right)^{1/2}$$

for $j = 2, \dots, K$, and $\|y^1\|_\mu = \|y^1\|_H$. Furthermore, let $\underline{\alpha}_\mu > 0$ be the μ -dependent coercitivity constant of the bilinear form satisfying $a_\mu(\varphi, \varphi) \geq \underline{\alpha}_\mu \|\varphi\|_V^2$ for all $\varphi \in V$ and $\mu \in \mathcal{P}_{\text{ad}}$. Then we can define the error estimates in the following prepositions.

Proposition 3. Let $\{y_\mu^k\}_{k=1}^K$ and $\{y_\mu^{k,\ell}\}_{k=1}^K$ be given by (8) and $\text{err}_\mu^k = y_\mu^k - y_\mu^{k,\ell} \in V^N$ for $k = 1, \dots, K$. For $\varphi \in V^N$ we define the residual

$$\begin{aligned} \langle \text{res}_\mu^k, \varphi \rangle_{(V^N)', V^N} &= \left\langle \frac{y_\mu^{k-1,\ell} - y_\mu^{k,\ell}}{\Delta t}, \varphi \right\rangle_H - a_\mu(y_\mu^{k,\ell}, \varphi) \\ &\quad + \langle f_\mu(t_k; u), \varphi \rangle_{V', V} \end{aligned}$$

and $\varepsilon_\mu^k = \|\text{res}_\mu^k\|_{(V^N)'}$. Then $\|\text{err}_\mu^j\|_\mu \leq \Delta_\mu^j$ for $j \in \{1, \dots, K\}$, where

$$\Delta_\mu^j = \left(\|\text{err}_\mu^1\|_H^2 + \frac{\Delta t}{\underline{\alpha}_\mu} \sum_{k=2}^j |\varepsilon_\mu^k|^2 \right)^{1/2} \quad (10)$$

Proof. We refer to Grepl and Patera (2005) for a proof.

Proposition 4. Let $\{y_\mu^k\}_{k=1}^K$ and $\{y_\mu^{k,\ell}\}_{k=1}^K$ be given by (8), the reduced cost functions $\hat{J}^N(\mu)$ and $\hat{J}^\ell(\mu)$ be given by (7) and (9) respectively, the estimator Δ_μ^j be defined by (10), and the cost error defined by $\text{err}_\mu^j = \hat{J}^N(\mu) - \hat{J}^\ell(\mu)$. Then, we have

$$|\text{err}_\mu^j| \leq \Delta_\mu^j \quad (11)$$

where

$$\Delta_\mu^j = \frac{1}{2} \sum_{k=1}^K \alpha_k \left((\Delta_\mu^k)^2 + 2\Delta_\mu^k \|y_\mu^{k,\ell} - \hat{y}^k\|_H \right), \quad (12)$$

which, again, does not require the evaluation of the full-order solution $\{y_\mu^k\}_{k=1}^K$.

Proof. We skip the proof, available in Petrocchi et al. (2022).

4. TR-RB APPROXIMATION

It is fairly easy to find a reduced basis for one fixed parameter, as we see in the next subsection. But when the parameter changes, the basis found before is not necessarily a good reduced basis for the new state. For this reason, a Greedy algorithm is often used (in the offline phase) to find a ‘‘common’’ reduced basis for all the parameters in the admissible set \mathcal{P}_{ad} so that in the online phase a fast computation of a state for a new admissible parameter is possible. This leads to the well-known offline/online decomposition (see, e.g., Grepl and Patera (2005), Binev et al. (2011), Haasdonk (2013)). In the context of optimization the offline computation might not be very suitable, because during the optimization method only parameters from a small (but a-priorily unknown) subset of \mathcal{P}_{ad} are required. In Qian et al. (2017), Keil et al. (2021), Banholzer et al. (2020)) the RB space is built during a TR optimization process, where – following the optimization path – the RB space is enriched if it is necessary. We will utilize these ideas to develop our algorithm for the optimal experimental design. Since we cannot recall all details of the algorithm we refer to Petrocchi et al. (2022) and Banholzer et al. (2020) for further information. See also the recent paper Banholzer et al. (2022), where the TR approach is used in a multi-objective parameter optimization problem.

4.1 Creating and enriching a reduced basis

For a fixed $\mu \in \mathcal{P}_{\text{ad}}$ we compute a reduced-order approximation of y^k by the proper orthogonal decomposition (POD) method (cf., e.g., Kunisch and Volkwein (2001)). Given some relevant snapshots $\{z^k\}_{k=1}^K \subset \mathbb{R}^N$, we use POD to compute a coefficient matrix $\Psi \in \mathbb{R}^{N \times \ell}$ by solving

$$\min \left\{ \sum_{k=1}^K \alpha_k \|z_\mu^k - \Psi \Psi^\top W z_\mu^k\|_W^2 \mid \Psi^\top W \Psi = I \right\}$$

with the positive, symmetric matrix $W = (\langle \varphi_j, \varphi_i \rangle_V)$. It turns out that we have

$$\psi_i = \sum_{i=1}^N \Psi_{ij} \varphi_i \in V^N, \quad \langle \psi_j, \psi_i \rangle_V = \delta_{ij}, \quad 1 \leq i, j \leq \ell, \quad (13)$$

so the reduced space can be characterized by the matrix Ψ . In general, the dimension of the reduced space ℓ may vary; for example, it can be fixed a-priorily or it can be chosen based on the eigenvalue decay in the POD method; see, e.g., Gubisch and Volkwein (2017). Details about our choice for the RB dimension and about our choice of snapshots are discussed in Petrocchi et al. (2022).

If, on the other hand, we already have computed an RB space $V^\ell = \text{span}\{\psi_1, \dots, \psi_\ell\}$ characterized by the matrix Ψ (see (13)), let us assume that we want to enrich the basis in the parameter $\mu_+ \in \mathcal{P}_{\text{ad}}$. Then, we can find another matrix $\Psi_+ \in \mathbb{R}^{N \times \ell_+}$ using POD on the subspace that is orthogonal to the reduced space, and then we merge the matrices (using Gram-Schmidt W -orthonormalization, if necessary). Again, details are available in Petrocchi et al. (2022).

4.2 TR framework

The TR optimization method computes iteratively a first-order critical point of $(\hat{\mathbf{P}})$. At each iteration $k \geq 0$ of

the optimization algorithm, we call such approximated optimal parameter $\mu^{(k)}$. We consider a cheaply computable model $m^{(k)}$ (approximation of the reduced cost) that can be trusted to accurately represent the function \hat{J} in a reasonable neighborhood of $\mu^{(k)}$, called *trust region* $\mathcal{T}(\delta^{(k)}) = \{\mu : \|\mu - \mu^{(k)}\|_2 \leq \delta^{(k)}\}$, where $\delta^{(k)}$ is called *TR radius*. The TR method finds $\mu^{(k+1)}$ by solving the problem

$$\min_{s \in \mathbb{R}^d} m^{(k)}(s) \quad \text{s.t.} \quad \|s\|_2 \leq \delta^{(k)}, \mu^{(k)} + s \in \mathcal{P}_{\text{ad}} \quad (14)$$

Setting $\tilde{\mu} = \mu^{(k)} + s$ the RB version of (14) is

$$\min_{\tilde{\mu} \in \mathcal{P}_{\text{ad}}} \hat{J}^{\ell, (k)}(\tilde{\mu}) \quad \text{s.t.} \quad q^{(k)}(\tilde{\mu}) = \frac{\Delta_{\tilde{\mu}}^{\hat{J}, (k)}}{\hat{J}^{\ell, (k)}(\tilde{\mu})} \leq \delta^{(k)} \quad (15)$$

Here and whenever some quantity depends on the iteration k , we show it in the superscript, like the RB cost $\hat{J}^{\ell, (k)}$. The so-called efficiency $q^{(k)}$ helps us to quantify the accuracy of the RB and to define the TR.

How the model behaves in the TR tells us if we need to enrich the basis, or reduce the TR radius (in case the model is not accurate enough), or if we can enlarge the RB basis and even skip the enrichment process, in case the model is already accurate enough in the trust region.

4.3 TR subproblem

Let us suppose that we have fixed a TR radius $\delta^{(k)}$. Then a proposed parameter, solution of (15), is evaluated using a projected Armijo-BFGS algorithm (we refer to Kelley (1999) for the details). The iterations of the BFGS algorithm are indicated in the second superscript of the parameter: $\{\mu^{(k, j)}\}_{j=1}^{m_k}$ is the sequence of BFGS iterates and $\tilde{\mu} = \mu^{(k, m_k)}$ is the result of the BFGS algorithm. The maximum number of BFGS iterates is fixed at 400, and the algorithm finishes automatically when one of the termination criterion is satisfied:

$$q^{(k)}(\mu^{(k, j)}) \geq \beta_1 \delta^{(k)} \quad (16a)$$

$$\|\mu^{(k, j)} - \mathcal{P}_{\mathcal{P}_{\text{ad}}}(\mu^{(k, j)} + \nabla \hat{J}^{\ell, (k)}(\mu^{(k, j)}))\|_2 \leq \varepsilon_{\text{sub}}, \quad (16b)$$

where $\beta_1 \in (0, 1]$ and ε_{sub} is the tolerance of the subproblem. Equation (16a) tells us if we are too close to the border of the trust region (where the RB model is less accurate) and $\mathcal{P}_{\mathcal{P}_{\text{ad}}}$ is the projection onto \mathcal{P}_{ad} .

Once that the projected BFGS method evaluates the descent direction $d^{(k, j)}$ for the parameter $\mu^{(k, j)}$, the Armijo backtracking finds $\mu^{(k, j+1)}$ as:

$$\mu^{(k, j+1)} = \mathcal{P}_{\mathcal{P}_{\text{ad}}}(\mu^{(k, j)} + \alpha_{(k, j)} d^{(k, j)}) \in \mathcal{P}_{\text{ad}}$$

where $\alpha_{(k, j)} = 0.5^\beta$ and the power $\beta = \beta(k, j)$ is the smallest integer such that the sufficient decrease and the TR constraint are satisfied:

$$\begin{aligned} \hat{J}^{\ell, (k)}(\mu^{(k, j+1)}) - \hat{J}^{\ell, (k)}(\mu^{(k, j)}) \\ \leq -\frac{\alpha_o}{\alpha_{(k, j)}} \|\mu^{(k, j+1)} - \mu^{(k, j)}\|_2^2 \end{aligned} \quad (17a)$$

with $\alpha_o = 10^{-4}$ and

$$q^{(k)}(\mu^{(k, j+1)}) \leq \delta^{(k)} \quad (17b)$$

4.4 Modification of the trust region

Once that the subproblem has found a suitable candidate $\tilde{\mu} := \mu^{(k, m_k)}$, there are different possibilities: for example,

we might realize that the candidate is on the border of the trust region, and since the model is not accurate there, we could shrink the TR radius and find a new candidate or just enrich the basis; on the other hand, if we notice that some “predicted sufficient reduction” is satisfied, we can enlarge the TR radius.

The approximated generalized Cauchy point (AGC) is defined as $\mu_{\text{AGC}}^{(k)} = \mathcal{P}_{\mathcal{P}_{\text{ad}}}(\mu^{(k)} - \alpha^{(k, 0)} \nabla \hat{J}^{\ell, (k)}(\mu^{(k)}))$, with $\alpha^{(k, 0)}$ chosen such that (17a) and (17b) are satisfied; cf. Yue and Meerbergen (2013). The cost of the evaluation of this parameter is “free”, since it is evaluated in the first iteration of the BFGS algorithm. Then, an error-aware sufficient decrease condition is introduced (Qian et al. (2017)):

$$\hat{J}^{\ell, (k+1)}(\mu^{(k+1)}) \leq \hat{J}^{\ell, (k)}(\mu_{\text{AGC}}^{(k)}) \quad (18)$$

where we highlight that $\hat{J}^{\ell, (k)}$ refers to the reduced model at iteration k , while $\hat{J}^{\ell, (k+1)}$ refers instead to the model after the $(k+1)$ -th (eventual) enrichment. This condition is central in the proof of convergence and the “sufficient” and “necessary” conditions we will see in the algorithm refer to the verification of it. Let us say explicitly that this condition is not straightforward to verify. Indeed, the left-hand side requires the evaluation of the cost function post-enhancement on the new parameter. We want to avoid this evaluation if some other conditions are not satisfied. For this reason, we postpone the evaluation of the FOM solution until we have to check the termination criterion or we want to enrich the RB.

To ensure that the candidate is a good parameter, a sufficient condition is analyzed:

$$\hat{J}^{\ell, (k)}(\tilde{\mu}) + \Delta_{\tilde{\mu}}^{\hat{J}, (k)} < \hat{J}^{\ell, (k)}(\mu_{\text{AGC}}^{(k)}) \quad (19)$$

If (19) is satisfied, the candidate is accepted, $\mu^{(k+1)} = \tilde{\mu}$ and the model is “updated” (that means, the basis is enriched) there. Then, the TR radius is doubled if the predicted sufficient reduction (of model $m^{(k)}$) is realized, namely if

$$\rho^{(k)} = \frac{\hat{J}^N(\mu^{(k)}) - \hat{J}^N(\mu^{(k+1)})}{\hat{J}^{\ell, (k)}(\mu^{(k)}) - \hat{J}^{\ell, (k)}(\mu^{(k+1)})} \geq \beta_2, \quad (20)$$

where $\beta_2 \in (0, 1)$.

If the sufficient condition (19) does not hold, we check a necessary condition:

$$\hat{J}^{\ell, (k)}(\tilde{\mu}) - \Delta_{\tilde{\mu}}^{\hat{J}, (k)} < \hat{J}^{\ell, (k)}(\mu_{\text{AGC}}^{(k)})$$

If this also fails, it means that the point $\tilde{\mu}$ probably needs an enhancement too big to satisfy the error-aware sufficient decrease condition (18). Then it must be rejected, the model must be enriched and the TR radius shrunken. If, on the other hand, the sufficient condition fails and the necessary condition holds, we enrich the model and check (18) on the candidate parameter:

$$\hat{J}^{\ell, (k+1)}(\tilde{\mu}) \leq \hat{J}^{\ell, (k)}(\mu_{\text{AGC}}^{(k)})$$

If this holds, the candidate and the enrichment are accepted (and again, if (20) holds the radius is doubled), while if it fails both the parameter and the enrichment are rejected and the radius is shrunken.

Defining

$$g^N(\mu^{(k+1)}) = \|\mu^{(k+1)} - \mathcal{P}_{\mathcal{P}_{\text{ad}}}(\mu^{(k+1)} - \nabla \hat{J}^N(\mu^{(k+1)}))\|_2$$

the stopping criterion for the parameter estimation is $g^N(\mu^{(k+1)}) \leq \varepsilon_{tr}$, where ε_{tr} is the overall tolerance.

4.5 Skipping the enrichment

Enriching the basis often leads to too many basis elements, hence wasting the ROM purpose. For this reason, in Banholzer et al. (2020) the possibility of skipping the enrichment was included. In particular, this happens when all these three conditions hold:

$$q^{(k)}(\mu^{(k+1)}) \leq \beta_3 \delta^{(k+1)} \quad (21a)$$

$$\frac{|g^N(\mu^{(k+1)}) - g^{\ell, (k)}(\mu^{(k+1)})|}{g^{\ell, (k)}(\mu^{(k+1)})} \leq \tau_g \quad (21b)$$

$$\frac{\|\nabla \hat{J}^N(\mu^{(k+1)}) - \nabla \hat{J}^{\ell, (k)}(\mu^{(k+1)})\|_2}{\|\nabla \hat{J}^N(\mu^{(k+1)})\|_2} \leq \min \{ \tau_{grad}, \beta_3 \delta^{(k+1)} \} \quad (21c)$$

with $\beta_3 \in (0, 1)$, $g^{\ell, (k)}(\mu) = \|\mu - \mathcal{P}_{\mathcal{P}_{ad}}(\mu - \nabla \hat{J}^{\ell, (k)}(\mu))\|_2$, and $\tau_g, \tau_{grad} > 0$ are arbitrary accuracies. Inequality (21a) indicates how much the model $m^{(k+1)}$ is trustworthy, inequality (21b) checks the convergence criterion of the reduced model, and inequality (21c) checks the RB accuracy of the cost gradient.

Convergence results for this algorithm are showed in Keil et al. (2021) and Banholzer et al. (2020) for an elliptic system. The proof can be adapted to our case, as discussed in Petrocchi et al. (2022).

5. OPTIMAL INPUT DESIGN

The result of the parameter estimation can depend strongly on the chosen control, even leading to convergence errors. In order to gain confidence in the accuracy of the parameter estimation, we want to find a “optimal” control, with the purpose of minimizing the uncertainty of the estimated parameter.

Following standard ideas of the optimal experimental/input design (Goodwin and Payne (1977), Atkinson et al. (2007)), we use an adaptive algorithm storing past data, such that the cost function will measure the discrepancy will all data stored. The superscript $[n]$ on objects like data vectors, state variables and cost functions specifies the experiment index, namely what control function it is being considered. The algorithm starts from initial guesses $u^{[1]}$ and $\mu^{[1]}$ and finds a new parameter and a new control at each iteration. Then, for $n > 1$, the algorithm iterates on three phases:

- The first phase is the simulation of some data $\hat{y}^{[n-1]}$ using the control $u^{[n-1]}$ and the true parameter μ^* .
- The second phase is the TR-RB approximation, from which we obtain the parameter $\mu^{[n]}$.
- The third phase is the optimal input design, whose purpose is to find a new control $u^{[n]}$. The algorithm stops when the control just found is very close to the previous one.

Simulating measurement errors, for $n \geq 1$ (experiment index) we assume that the data is a random variable presenting some additive noise, i.e.

$$\hat{y}_i^{k, [n]} = y_{\mu^* i}^{k, [n]} + \eta_{n, k, i}, \quad (22)$$

where $\eta_{n, k, i} \sim \mathcal{N}(0, \sigma_d^2)$ for all $k = 1, \dots, K$ (temporal index), $i = 1, \dots, N$ (coefficient index in the FOM basis).

So, in a way, we assume we are able to measure the FOM coefficient of the state variable corresponding to the true parameter at each time step, plus some error. The data collected at each iteration, $\hat{y}^{[n]} \in \mathbb{R}^{NK}$, is a vector stacked with all entries $\hat{y}_i^{k, [n]}$ for $k = 1, \dots, K$ and $i = 1, \dots, N$.

Hence, the solution of the second phase $\mu^{[n]}$ will be an estimator of the true parameter μ^* . To evaluate the quality of such estimator we use the Fisher’s information matrix, which quantifies how informative an experiment is.

Skipping the details, the Fisher’s information matrix $M_{\mu^{[n]}} \in \mathbb{R}^{d \times d}$ of the estimator $\mu^{[n]}$ is in our case

$$(M_{\mu^{[n]}})_{j, l} = \frac{1}{\sigma_d^2} \int_0^T s_{\mu^{[n]}, j}^{[n]}(t)^\top s_{\mu^{[n]}, l}^{[n]}(t) dt$$

for any $j, l \in \{1, \dots, d\}$, where $s_{\mu^{[n]}, j}^{[n]}(t)$ is the coefficient vector (in the FOM basis) of the sensitivity of the output with respect to the j -th parameter, evaluated in $\mu^{[n]}$ and using control $u^{[n]}$. More information about how this is evaluated (in particular in the RB framework) can be found in Petrocchi et al. (2022).

Then, the measure of uncertainty is chosen as the trace of the inverse of the Fisher’s information matrix:

$$\phi(u) = \text{trace}(M_{\mu^{[n]}}^{-1})$$

6. NUMERICAL EXAMPLES

We consider examples in a one-dimensional spatial interval $\Omega = (0, 4)$ with $H = L^2(\Omega)$ and $V = H^1(\Omega)$. In both cases the initial value is $y_o(x) \equiv 1$ and the time horizon is $T = 2$. The chosen initial control is $u(t) = \frac{1}{2} \cos(10t) \in \mathcal{U}_{ad} = \{u \in L^2(0, T) \mid -3 \leq u(t) \leq 3, \forall t \in (0, T)\} \subset \mathcal{U} = L^2(0, T)$, and data is always evaluated with the formula (22), where $y_{\mu^*}^{k, [n]}$ is evaluated with a FOM solver and $\sigma_b^2 = 10^{-3}$. The initial reference parameter $\hat{\mu}$ in (3) is the middle point of \mathcal{P}_{ad} , and at any following iteration $n > 1$ the reference parameter is the optimal parameter $\mu^{[n-1]}$ found in the previous iteration. The constant σ is fixed to 10^{-8} .

In the discretized setting, the control $u = u(t)$ is a step function and it is determined by the vector $\mathbf{u} = (u_k) \in \mathbb{R}^K$ with $u(t_k) = u_k$ for $k \in \{1, \dots, K\}$. Furthermore, in the numerical example the optimal input optimization problem

$$\mathbf{u}^{[n]} = \arg \min \{ \phi(\mathbf{u}) \mid \mathbf{u} \in \mathcal{U}_{ad} \}$$

$$\mathcal{U}_{ad} = \{ \mathbf{u} \in \mathbb{R}^K \mid -3 \leq u_k \leq 3 \text{ for } k = 1, \dots, K \}$$

is solved using the function `fmin_l_bfgs_b` from the Python library `scipy.optimize` without specifying the gradient of $\phi(u)$ and approximating it numerically.

In the next tables we are going to analyze the number ℓ_y, ℓ_p of bases generated by the TR optimization method (respectively for the state and adjoint variables), the error $\text{err}^* = \|\mu^{[n]} - \mu^*\|_2$ and $\phi(u^{[n]}) = \min_u \phi(u)$, where the optimal solutions $\mu^{[n]}$ and $u^{[n]}$ are computed by the strategy explained in Section 5.

Run 1 (2 parameters). We consider $d = 2$ parameters, $\mu \in \mathcal{P}_{ad} = [0.24, 3.4] \times [0.11, 5.0]$ and an equation in strong form

$$\begin{cases} y_t(x, t) - \mu_1 y_{xx}(x, t) + \frac{1}{2}y(x, t) = 0 \\ y_x(0, t) = 0, \quad \mu_1 y_x(1, t) = \mu_2 u(t) \\ y(x, 0) = y_0(x) \end{cases}$$

so that $a_\mu(y, \varphi) = \mu_1 \int_\Omega y'(x)\varphi'(x) dx + \frac{1}{2} \int_\Omega y(x)\varphi(x) dx$ and $\langle f_\mu(t; u), \varphi \rangle_{V', V} = \mu_2 u(t)\varphi(L)$. The true parameter is $\mu^* = (1, 2)$ and the initial parameter is $\mu^{[1]} = (3.4, 0.11)$.

Table 1. Run 1.

| | Starting values | $n = 2$ | $n = 3$ |
|-------------------|-----------------|----------|----------|
| ℓ_y, ℓ_p | | 12, 30 | 18, 50 |
| err* | 3.05 | 3.14e-05 | 1.02e-06 |
| Second phase time | | 11 s | 9 s |
| $\phi(u^{[n]})$ | 4.5e-06 | 1.8e-09 | 1.8e-09 |
| Third phase time | | 288 s | 31 s |

We can observe how the second phase converges very fast to a parameter very close to the true parameter. Furthermore, the third phase reaches the “best” control in just one iteration. With this control the second phase in the next iteration finds a more accurate parameter and then the third phase does not find a better control, therefore stopping.

Run 2 (4 parameters) Let $d = 4$ and $\mathcal{P}_{\text{ad}} = [0.24, 3.4] \times [0.11, 5] \times [0.013, 4] \times [0.97, 2.22]$ and

$$a_\mu(y, \varphi) = \mu_1 \int_0^{1.5} y'(x)\varphi'(x) dx + \mu_2 \int_{1.5}^3 y'(x)\varphi'(x) dx + \mu_3 \int_3^4 y'(x)\varphi'(x) dx + \frac{1}{2} \int_\Omega y(x)\varphi(x) dx$$

and $\langle f_\mu(t; u), \varphi \rangle_{V', V} = \mu_4 u(t)\varphi(1)$. The real parameter is $\mu^* = (1, 1.3, 0.8, 2)$ and the initial parameter is $(3, 5, 0.1, 2)$. The result are showed in the next table.

Table 2. Run 2.

| | Starting values | $n = 2$ | $n = 3$ |
|-------------------|-----------------|---------|---------|
| ℓ_y, ℓ_p | | 24, 30 | 30, 40 |
| err* | 4.26 | 3.23 | 6.9e-05 |
| Second phase time | | 180 s | 20 s |
| $\phi(u^{[n]})$ | 19 | 1.3e-06 | 2e-07 |
| Third phase time | | 540 s | 15 s |

In this example we can see how computational times and errors rise with bigger sizes of the parameter space. Let us observe that the TR algorithm works way faster than a standard weak greedy and we can appreciate the full potential of such algorithm when applied to higher dimensions. With the first control the algorithm fails to approximate the true parameter, but as soon as a new control is generated, the trust-region approximation gives a good estimator.

REFERENCES

- Atkinson, A.C., Donev, A.N., and Tobias, R.D. (2007). *Optimum experimental designs, with SAS*, volume 34 of *Oxford Statistical Science Series*. Oxford University Press, Oxford.
- Banholzer, S., Keil, T., Mechelli, L., Ohlberger, M., Schindler, F., and Volkwein, S. (2020). An adaptive projected newton non-conforming dual approach for trust-region reduced basis approximation of PDE-constrained parameter optimization. arXiv:2012.11653. To appear in *Pure and Applied Functional Analysis*, issue 5, 2022.
- Banholzer, S., Mechelli, L., and Volkwein, S. (2022). A trust region reduced basis Pascoletti-Serafini algorithm for multi-objective PDE-constrained parameter optimization. arXiv:2201.07744.
- Binev, P., Cohen, A., Dahmen, W., DeVore, R., Petrova, G., and Wojtaszczyk, P. (2011). Convergence rates for greedy algorithms in reduced basis methods. *SIAM J. Math. Anal.*, 43(3), 1457–1472.
- Goodwin, G.C. and Payne, R.L. (1977). *Dynamic system identification*, volume 136 of *Mathematics in Science and Engineering*. Academic Press, Inc. [Harcourt Brace Jovanovich, Publishers], New York-London. Experiment design and data analysis.
- Grepl, M.A. and Patera, A.T. (2005). A posteriori error bounds for reduced-bias approximations of parametrized parabolic partial differential equations. *M2AN Math. Model. Numer. Anal.*, 39(1), 157–181.
- Gubisch, M. and Volkwein, S. (2017). Proper orthogonal decomposition for linear-quadratic optimal control. In P. Benner, A. Cohen, M. Ohlberger, and K. Willcox (eds.), *Model Reduction and Approximation: Theory and Algorithms*, 5–66. SIAM, Philadelphia, PA.
- Haasdonk, B. (2013). Convergence rates of the POD-greedy method. *ESAIM Math. Model. Numer. Anal.*, 47(3), 859–873.
- Hesthaven, J., Rozza, G., and Stamm, B. (2016). *Certified Reduced Basis Methods for Parametrized Partial Differential Equations*. SpringerBriefs in Mathematics. Springer Cham.
- Hinze, M., Pinnau, R., Ulbrich, M., and Ulbrich, S. (2009). *Optimization with PDE constraints*, volume 23 of *Mathematical Modelling: Theory and Applications*. Springer, New York.
- Keil, T., Mechelli, L., Ohlberger, M., Schindler, F., and Volkwein, S. (2021). A non-conforming dual approach for adaptive trust-region reduced basis approximation of PDE-constrained parameter optimization. *ESAIM Math. Model. Numer. Anal.*, 55(3), 1239–1269.
- Kelley, C.T. (1999). *Iterative methods for optimization*, volume 18 of *Frontiers in Applied Mathematics*. SIAM, Philadelphia, PA.
- Kunisch, K. and Volkwein, S. (2001). Galerkin proper orthogonal decomposition methods for parabolic problems. *Numer. Math.*, 90(1), 117–148.
- Petrocchi, A., Scharrer, M.K., and Volkwein, S. (2022). Adaptive reduced basis methods for PDE-constrained optimization and optimal input design. Technical report, Universität Konstanz, Konstanzer Schriften in Mathematik. URL <http://nbn-resolving.de/urn:nbn:de:bsz:352-2-119b7uuht7b1s0>.
- Qian, E., Grepl, M., Veroy, K., and Willcox, K. (2017). A certified trust region reduced basis approach to PDE-constrained optimization. *SIAM J. Sci. Comput.*, 39(5), S434–S460.
- Quarteroni, A. (2017). *Numerical models for differential problems*, volume 16 of *MSE&A. Modeling, Simulation and Applications*. Springer, Cham.
- Yue, Y. and Meerbergen, K. (2013). Accelerating optimization of parametric linear systems by model order reduction. *SIAM J. Optim.*, 23(2), 1344–1370.