



Technische Universität München
Fakultät für Mathematik
Lehrstuhl für Mathematische Optimierung

Optimal Control under Uncertainty: Theory and Numerical Solution with Low-Rank Tensors

Sebastian Raphael Garreis

Vollständiger Abdruck der von der Fakultät für Mathematik der Technischen Universität München zur Erlangung des akademischen Grades eines

Doktors der Naturwissenschaften (Dr. rer. nat.)

genehmigten Dissertation.

Vorsitzende: Prof. Dr. Caroline Lasser

Prüfer der Dissertation: 1. Prof. Dr. Michael Ulbrich
2. Prof. Dr. Christian Clason
Universität Duisburg-Essen
3. Prof. Dr. Matthias Heinkenschloss
Rice University, Houston, Texas, USA
(schriftliche Beurteilung)

Die Dissertation wurde am 19. September 2018 bei der Technischen Universität München eingereicht und durch die Fakultät für Mathematik am 5. Februar 2019 angenommen.

Abstract

In this work, a class of optimal control problems under uncertainty constrained by semilinear, elliptic partial differential equations is analyzed. An inexact trust-region algorithm with a suitable error control procedure is presented to solve such problems adaptively. The state and adjoint equations are formulated in a tensor Banach space and solved using low-rank tensor methods. Numerical results show how the algorithms adapt to the problem data. The dissertation is concluded by an outlook to alternative risk measures, which yield risk-averse controls.

Zusammenfassung

In dieser Arbeit wird eine Klasse von Optimalsteuerungsproblemen unter Unsicherheit analysiert, die semilineare, elliptische partielle Differentialgleichungen als Nebenbedingung haben. Ein inexaktes Trust-Region-Verfahren mit einer geeigneten Vorgehensweise zur Fehlerkontrolle wird vorgestellt, um solche Probleme adaptiv zu lösen. Die Zustands- und adjungierten Gleichungen werden in einem Tensor-Banachraum formuliert und mit Niedrigrangtensor-Methoden gelöst. Numerische Ergebnisse zeigen, wie sich die Algorithmen an die Problem-Data anpassen. Die Dissertation wird mit einem Ausblick auf alternative Risikomaße abgeschlossen, die risikoaverse Steuerungen liefern.

Notation

The following notation is used in this thesis:

Sets:

\emptyset	empty set
$ S $	cardinality of the set S
2^S	power set of the set S
$\text{int}(S)$	the topological interior of the set S
$\overline{S}^{\ \cdot\ } = \text{cl}(S)$	the topological closure of the set S w. r. t. the norm $\ \cdot\ $
$S \times T$	product set of the sets S and T
\mathbb{N}	natural numbers: $\mathbb{N} = \{1, 2, 3, \dots\}$
\mathbb{N}_0	non-negative integers: $\mathbb{N}_0 = \{0, 1, 2, \dots\}$
$[n]$	the first $n \in \mathbb{N}$ natural numbers: $[n] = \{1, 2, \dots, n\}$
\mathbb{R}	real numbers
\mathbb{C}	complex numbers
\mathbb{K}	general field
$\mathbb{R}_{>0}, \mathbb{R}_{\geq 0}, \mathbb{R}_{<0}, \mathbb{R}_{\leq 0}$	set of positive, non-negative, negative, non-positive real numbers, respectively
$\overline{\mathbb{R}}$	extended real line $\overline{\mathbb{R}} = \mathbb{R} \cup \{-\infty, \infty\} = [-\infty, \infty]$
$(a, b), [a, b], [a, b), (a, b]$	open, closed, half-open intervals, respectively, with endpoints $a, b \in \overline{\mathbb{R}}, a \leq b$, empty in the (half-)open case if $a = b$

Vectors/matrices/tensors:

S^n	set of column vectors with $n \in \mathbb{N}$ components from the set S
$S^{m \times n}$	set of $(m \times n)$ -matrices with components from the set S
$S^{n_1 \times \dots \times n_d}$	set of tensors of order $d \in \mathbb{N}$ (d -dimensional arrays) with dimensions $n_1, n_2, \dots, n_d \in \mathbb{N}$ and entries from S
$\mathbf{x}(K_1, K_2, \dots, K_d)$	indexing of a tensor $\mathbf{x} \in S^{n_1 \times \dots \times n_d}$: modes $i \in [d]$ indexed by a list K_i with elements from $[n_i]$ remain, modes indexed by a single index $K_i \in [n_i]$ are cut out
\cdot	shorthand for $[n_i]$ (the set/list of all indices) when indexing a tensor

$\mathbf{x} \otimes \mathbf{y}, \lambda \otimes \mu, Y \otimes Z$	Kronecker/outer product of the vectors/matrices/tensors \mathbf{x} and \mathbf{y} , product of the measures λ and μ , tensor product of the Banach spaces Y and Z
$\mathbf{x} \odot \mathbf{y}$	Hadamard (componentwise) product of the vectors/matrices/tensors \mathbf{x} and \mathbf{y}
$\mathbf{x} \oslash \mathbf{y}$	componentwise quotient of the vectors/matrices/tensors \mathbf{x} and \mathbf{y}
\mathbf{x}^λ	componentwise exponentiation of the vector/matrix/tensor \mathbf{x} with the exponent $\lambda \in \mathbb{R}$
$\langle \mathbf{x}, \mathbf{y} \rangle_{s,t}$	contraction of the tensors \mathbf{x} and \mathbf{y} along the modes s and t
$\langle \mathbf{x}, \mathbf{y} \rangle$	inner product of the tensors \mathbf{x} and \mathbf{y} of the same size
$\ \mathbf{x}\ _F$	Frobenius norm of the tensor/matrix \mathbf{x} : $\ \mathbf{x}\ _F = \sqrt{\langle \mathbf{x}, \mathbf{x} \rangle}$
$\mathbb{1}$	the vector/matrix/tensor only containing ones or the function which is constant one
I, \mathbb{I}	identity mapping, identity matrix
\mathbf{x}^\top	transpose of the vector or matrix \mathbf{x}

Banach spaces and operators:

\mathbb{R}^n	Euclidean space equipped with the inner product $\langle \mathbf{x}, \mathbf{y} \rangle = \mathbf{x}^\top \mathbf{y}$ and norm $\ \cdot\ \equiv \ \cdot\ _2$ (see below)
$\ \cdot\ _p$	p -norm on \mathbb{R}^n : $\ \mathbf{x}\ _p = \left(\sum_{k=1}^n x_k ^p\right)^{1/p}$ for $p \in [1, \infty)$, $\ \mathbf{x}\ _\infty = \max_{k \in [n]} x_k $
$Y, Y^*, \ \cdot\ _Y$	a general Banach space, its dual space, the norm on Y
$\langle \cdot, \cdot \rangle_{Y^*, Y}$	dual pairing between Y^* and Y , i. e., $\langle f, y \rangle_{Y^*, Y} = f(y)$
$\langle \cdot, \cdot \rangle_{Y, Y^*}$	dual pairing between Y and Y^* , i. e., $\langle y, f \rangle_{Y, Y^*} = f(y)$
$y_k \rightarrow y, y_k \rightharpoonup y, y_k \rightharpoonup^* y$	Convergence of the sequence $(y_k)_{k \in \mathbb{N}} \subset Y$ to $y \in Y$ w. r. t. the strong, weak, weak* topology, respectively, as $k \rightarrow \infty$. The weak* topology is defined if Y is the dual space of a given normed space.
$\mathcal{L}(Y, Z)$	space of linear, bounded operators mapping from the Banach space Y to the Banach space Z
$\ \cdot\ _{\mathcal{L}(Y, Z)}$	induced norm $\ A\ _{\mathcal{L}(Y, Z)} = \sup_{\ y\ _Y \leq 1} \ Ay\ _Z$
$A: Y \rightarrow Z$	linear, bounded operator A mapping from Y to Z
$A^*: Z^* \rightarrow Y^*$	adjoint operator A^* mapping from the dual space Z^* to the dual space Y^*
$Y \hookrightarrow Z$	continuous embedding of Banach spaces
$Y \hookrightarrow\hookrightarrow Z$	compact embedding of Banach spaces
$\iota: Y \hookrightarrow Z$	canonical embedding of Banach spaces $Y \subset Z$

$U, (\cdot, \cdot)_U$ $\ \cdot\ _U$	a general Hilbert space, the inner product on U norm induced by the inner product, i.e., $\ u\ _U = \sqrt{(u, u)_U}$
$N : Y \rightarrow Z$	a general function mapping from the Banach space Y to the Banach space Z
$N' : Y \rightarrow \mathcal{L}(Y, Z)$	its first derivative
$N_{y_1} : Y \rightarrow \mathcal{L}(Y_1, Z)$	partial derivative of N w. r. t. y_1 for $Y = Y_1 \times Y_2$
$\frac{\partial}{\partial y_1}$	partial derivative of an expression w. r. t. y_1
$\frac{d}{dy_1}$	total derivative of an expression w. r. t. y_1
$N'' : Y \rightarrow \mathcal{L}(Y, \mathcal{L}(Y, Z))$	second derivative of N (etc.)
$N_{y_2 y_1} : Y \rightarrow \mathcal{L}(Y_2, \mathcal{L}(Y_1, Z))$	partial second derivative for $Y = Y_1 \times Y_2 \times Y_3$
$J : U \rightarrow \mathbb{R}$	a general functional on the Hilbert space U
$\nabla J : U \rightarrow U$	its gradient (Riesz representative of $J'(u) \in U^* \cong U$)
$\nabla_{u_1} J : U \rightarrow U_1$	partial gradient for $U = U_1 \times U_2$
$\nabla^2 J : U \rightarrow \mathcal{L}(U, U)$	Hessian of J (linearization of the gradient)
$\nabla_{u_2 u_1}^2 J : U \rightarrow \mathcal{L}(U_2, U_1)$	partial Hessian of J for $U = U_1 \times U_2 \times U_3$

Spatial domain and functions:

Ω, x	spatial domain $\Omega \subset \mathbb{R}^n$, point $x \in \Omega$
$1_\Omega : \mathbb{R}^n \rightarrow \mathbb{R}$	indicator function of the set Ω : $1_\Omega(x) = 1$ if $x \in \Omega$ and $1_\Omega(x) = 0$ if $x \notin \Omega$
$\partial\Omega$	boundary of the domain Ω
λ	Lebesgue measure on Ω
$f : \Omega \rightarrow \mathbb{R}^m$	a function mapping from $\Omega \subset \mathbb{R}^n$ to \mathbb{R}^m
$D^j f : \Omega \rightarrow \mathbb{R}^m$	(weak) k -th partial derivative of f , where $j \in \mathbb{N}_0^n$ is a multi-index

Probability:

$(\Xi, \mathcal{F}, \mathbb{P}), \xi$	probability space Ξ with σ -algebra \mathcal{F} and probability measure \mathbb{P} , random value $\xi \in \Xi$
$X : \Xi \rightarrow \mathbb{R}$	a real-valued random variable
\mathbb{E}	expectation: $\mathbb{E}[X] := \int_{\Xi} X(\xi) d\mathbb{P}$
Var	variance: $\text{Var}[X] := \mathbb{E}[(X - \mathbb{E}[X])^2]$
Cov	covariance: $\text{Cov}[X, \tilde{X}] = \mathbb{E}[(X - \mathbb{E}[X])(\tilde{X} - \mathbb{E}[\tilde{X}])]$
CVaR_β	conditional value-at-risk (expected shortfall) with quan- tile parameter $\beta \in (0, 1)$

Function spaces:

$\mathcal{C}(\Omega)$	space of continuous functions on Ω
$\mathcal{C}^k(\Omega)$	space of k -times continuously differentiable functions on Ω , $k \in \{1, 2, \dots, \infty\}$
$\mathcal{C}_C^k(\Omega)$	space of k -times continuously differentiable functions on Ω with compact support
$L^p(\Omega)$	Banach space of equivalence classes of measurable, p -integrable ($p \in [1, \infty)$) or essentially bounded ($p = \infty$) functions on Ω
$\ \cdot\ _{L^p(\Omega)}$	norm on $L^p(\Omega)$: $\ f\ _{L^p(\Omega)} = (\int_{\Omega} f(x) ^p dx)^{1/p}$ for $p \in [1, \infty)$, $\ f\ _{L^\infty(\Omega)} = \text{ess sup}_{x \in \Omega} f(x) $
$(\cdot, \cdot)_{L^2(\Omega)}$	inner product on $L^2(\Omega)$: $(f, g)_{L^2(\Omega)} = \int_{\Omega} f(x)g(x) dx$
$W^{k,p}(\Omega)$	Sobolev space of k -times weakly differentiable $L^p(\Omega)$ -functions with $L^p(\Omega)$ -derivatives
$\ \cdot\ _{W^{k,p}(\Omega)}$	norm on $W^{k,p}(\Omega)$: $\ f\ _{W^{k,p}(\Omega)} = \left(\sum_{ j \leq k} \ D^j f\ _{L^p(\Omega)}^p \right)^{1/p}$ for $p \in [1, \infty)$ and $\ f\ _{W^{k,\infty}(\Omega)} = \sum_{ j \leq k} \ D^j f\ _{L^\infty(\Omega)}$
$H^k(\Omega)$	shorthand for the Hilbert space $W^{k,2}(\Omega)$
$(\cdot, \cdot)_{H^k(\Omega)}$	inner product on $H^k(\Omega)$: $(f, g)_{H^k(\Omega)} = \sum_{ j \leq k} (D^j f, D^j g)_{L^2(\Omega)}$
$H_0^1(\Omega) = \overline{\mathcal{C}_C^\infty(\Omega)}^{\ \cdot\ _{H^1(\Omega)}}$	Hilbert space of $H^1(\Omega)$ -functions with zero boundary data in the sense of traces
$(\cdot, \cdot)_{H_0^1(\Omega)}$	inner product on $H_0^1(\Omega)$ defined by $(f, g)_{H_0^1(\Omega)} = (\nabla f, \nabla g)_{L^2(\Omega)^n}$, inducing the norm $\ \cdot\ _{H_0^1(\Omega)}$ on $H_0^1(\Omega)$ and a seminorm on $H^1(\Omega)$
$H^{-1}(\Omega)$	dual space of $H_0^1(\Omega)$: $H^{-1}(\Omega) = H_0^1(\Omega)^*$
$L_{\mathbb{P}}^p(\Xi; Y)$	Bochner space of strongly \mathbb{P} -measurable, p -integrable ($p \in [1, \infty)$) or essentially bounded ($p = \infty$) functions from Ξ to the Banach space Y
$\ \cdot\ _{L_{\mathbb{P}}^p(\Xi; Y)}$	norm $\ \mathbf{y}\ _{L_{\mathbb{P}}^p(\Xi; Y)} = (\int_{\Xi} \ \mathbf{y}(\xi)\ _Y^p d\mathbb{P})^{1/p}$ for $p \in [1, \infty)$, $\ \mathbf{y}\ _{L_{\mathbb{P}}^\infty(\Xi; Y)} = \text{ess sup}_{\xi \in \Xi} \ \mathbf{y}(\xi)\ _Y$
$\mathcal{P}_{\tilde{d}}(\Xi)$	space of polynomials of <i>total</i> degree at most $\tilde{d} \in \mathbb{N}_0$ on $\Xi \subset \mathbb{R}^m$
$\mathcal{P}_d(\Xi)$	space of polynomials of <i>coordinate</i> degree at most $d \in \mathbb{N}_0^m$ on $\Xi \subset \mathbb{R}^m$

Contents

Abstract	i
Notation	iii
1. Introduction	1
2. Tensors	9
2.1. Finite-Dimensional Tensors and Low-Rank Formats	9
2.1.1. Basics of Finite-Dimensional Tensors and Notation	9
2.1.2. Low-Rank Tensor Formats	12
2.1.3. Available Algorithms	14
2.2. Tensor Products of Hilbert and Banach Spaces	21
3. A Class of Optimal Control Problems under Uncertainty	27
3.1. Problem Formulation	27
3.2. A Class of Semilinear, Elliptic PDEs with Uncertain Coefficients	29
3.3. Existence of a Solution to the Optimal Control Problem	37
3.4. Derivatives of the Reduced Objective Function	39
3.4.1. Discussion of the Example	39
3.4.2. Differentiability of the Control-to-State Mapping	43
3.4.3. Second Derivatives	44
4. An Inexact Trust-Region Algorithm for the Solution of Optimal Control Problems with Control Constraints	47
4.1. Formulation of the Algorithm	48
4.2. Convergence Proof	50
4.3. Satisfying the Conditions Required by the Algorithm	54
4.4. Solution of Subproblems by a Semismooth Newton Method	64
5. Realization of the Required Error Estimates for the Model Problem	67
5.1. Realization of the Error Estimates in the Deterministic Case	67
5.2. Realization of the Error Estimates in the Stochastic Case	74
6. Discretization of the Model Problem	79
6.1. Space Discretization: Finite Elements	79
6.2. Stochastic Discretization: Polynomial Chaos	82
6.3. Choice of the Trust-Region Model	90
6.4. Solution of the Discrete Semismooth Newton System	92

7. A Posteriori Error Estimation and Adaptive Solution of a Class of Parametric PDEs Using Low-Rank Tensors	95
7.1. Realization of the Deterministic a Posteriori Error Estimator	97
7.2. A Posteriori Error Estimation in $L^2_{\mathbb{P}}(\Xi)$	100
7.3. Combination of Both Error Estimators	102
7.4. Realization with Low-Rank Tensors	110
8. Implementation and Numerical Results	115
8.1. Implementation Details	117
8.2. Results for the Deterministic Problem	120
8.3. Results for the Problem with Uncertainties	129
9. Alternative Risk Measures	137
9.1. Definition and Properties of Risk Measures	137
9.2. Mean-Variance Risk Measure	139
9.2.1. Discussion of the Example	141
9.2.2. Possible Error Estimates for Reduced Objective Function and Gradient	143
9.3. Convex Combination of Mean and Conditional Value-at-Risk	145
9.3.1. Definition and Derivation of the Properties	145
9.3.2. Smoothing by a Log-Barrier Approach	148
9.3.3. Application to Optimal Control under Uncertainty	152
9.3.4. Implementation and Numerical Results	161
10. Conclusions and Perspectives	173
A. Appendix	177
A.1. Tensor Spaces	177
A.2. General L^p Spaces and Operator Theory	179
A.3. Superposition Operators between L^p Spaces	180
A.4. Superposition Operators from $L^p_{\mathbb{P}}(\Xi; H^1_0(\Omega))$ to Its Dual	183
A.5. Alternative Approach for the Discussion of the Semilinear, Elliptic PDE . . .	185
A.6. Concrete Computations for Section 9.3	188
Bibliography	193

1. Introduction

Many phenomena in physics or engineering applications, such as the flight of an aircraft, the distribution of temperature in a heated system, the diffusion of a liquid or a gas in a medium, or a car crash, can be modeled mathematically by differential equations, in particular partial differential equations (PDEs). If analytical solutions to these equations are not available, numerical methods relying on, e. g., finite element (FE) or wavelet discretizations can be used to simulate such systems. Based on these simulations, certain components of the system, for instance, the shape of the wings of the aircraft, the applied heat distribution, the locations of the points where the liquid is inserted into the medium, or the design of the car, can be optimized to meet certain requirements or to minimize a cost functional. In the mentioned examples, one could aim to maximize the lift of the aircraft while keeping the drag below a given threshold, to attain a desired heat or liquid distribution, or to minimize the deformation of the occupants' space in the car while having a small enough deceleration to not injure them heavily. The described problems can be formulated as *optimal control problems* or *optimization problems with PDE constraints*. The mentioned goals appear in the objective functions of them while further requirements may be posed as constraints.

Often, the input parameters of the simulations, such as the wind speed or material properties, are not known to high accuracy a priori or are uncertain by nature. Hence, they can only be assumed to follow a given or estimated probability distribution or to belong to a set. It can be important to investigate how this uncertainty influences output quantities of the simulations. For instance, one could be interested in estimating the distribution of the lift given the distribution of the wind speed or in computing a lower and an upper bound for the heat distribution knowing bounds for the thermal conductivity and the specific heat capacity of the material. *Uncertainty quantification* has become an important field of research in recent years, see, e. g., [106] for an overview or the SIAM/ASA Journal on Uncertainty Quantification, which appeared in 2013 for the first time. A natural next step is to control the systems under uncertain influences optimally, a field of research called *optimal control under uncertainty*, which is the focus of this thesis.

In particular, we consider the problem of selecting a deterministic control for a system described by a semilinear, elliptic PDE with uncertain input parameters, which follow given probability distributions. We develop necessary theory for a class of such problems and investigate efficient numerical algorithms for solving them. The stochastic space is discretized by a stochastic Galerkin method with a full tensor product basis, resulting in exponentially many unknowns in the discretized system. For example, if we have $m \in \mathbb{N}$ uncertain parameters and the dependence of a quantity on each of them is discretized using $n \in \mathbb{N}$ basis functions, respectively, the full tensor product basis consists of n^m functions. To overcome the curse of dimensionality when having many parameters, the respective coefficients in tensor form, i. e., in form of a multi-dimensional array, are represented in a modern low-rank tensor format. These formats can reduce the storage and computational complexity drastically as already

observed in our previous work [46], which is the basis for some parts of this dissertation, but they offer only a limited set of efficiently implementable operations. Therefore, it is a main motivation of this thesis to investigate and develop algorithms which can be implemented using low-rank tensors. The complexity of low-rank tensor arithmetics depends strongly on the required tensor ranks. Truncation to smaller ranks is often necessary to make the algorithms efficient. This yields rounding errors, which have to be controlled suitably during the optimization process to achieve global convergence of the algorithm. Additionally, we control the errors resulting from stochastic and FE discretization and balance all error contributions.

In recent years, the mentioned efficient low-rank tensor formats, namely the Tensor Train format [88] and the hierarchical Tucker format [56], and numerical algorithms for them with good complexity have been developed, see [54, 67, 51] for an overview. Similar ideas have been used in quantum physics before, see, e. g., [83, 64] and the references therein. These formats and algorithms are an important ingredient of the methods developed in this thesis. For instance, they allow the numerical solution of PDEs with many random inputs in a reasonable amount of time. A more detailed introduction and overview of low-rank tensor methods and applications is given in Chapter 2.

General Problem Setting

For formulating a general problem setting of optimal control under uncertainty matching the problems considered in this work, let $(\Xi, \mathcal{F}, \mathbb{P})$ be a complete probability space. An element $\xi \in \Xi$ of this space stands for uncertain parameters. Let the control space U , the state space Y , and the image space Z be Banach spaces. The system to be controlled is under the influence of uncertainty and is described by the state equation

$$E : Y \times U \times \Xi \rightarrow Z, \quad E(y(\xi), u, \xi) = 0 \text{ for almost every (a. e.) } \xi \in \Xi, \quad (1.1)$$

where $y(\xi) \in Y$ is the uncertain state. The control $u \in U$ is *deterministic*, i. e., it does not depend on the uncertain inputs and shall be chosen prior to the observation of the uncertainty. Additionally, it is required to belong to a set of admissible controls $U_{\text{ad}} \subset U$. We formulate the optimal control problem

$$\min_{y(\cdot) \in Y, u \in U} \mathcal{R}[J_1(y(\cdot), u, \cdot)] + J_2(u) \quad \text{s. t.} \quad E(y(\xi), u, \xi) = 0 \text{ for a. e. } \xi \in \Xi, \quad u \in U_{\text{ad}}, \quad (1.2)$$

where $J_1 : Y \times U \times \Xi \rightarrow \mathbb{R}$ is the state-dependent part of the objective function, e. g., a tracking-type term, and $J_2 : U \rightarrow \mathbb{R}$ is a purely deterministic part used to, e. g., regularize the control. We assume that $J_1(y(\cdot), u, \cdot)$ is always measurable w. r. t. ξ and hence a *random variable* with values in \mathbb{R} . A *risk measure* $\mathcal{R} : \mathcal{X} \rightarrow \mathbb{R}$ has to be applied to obtain a real-valued objective function which can be minimized. The domain $\mathcal{X} \subset \{f : \Xi \rightarrow \mathbb{R}\}$ is a suitable set of random variables, such as $\mathcal{X} = L^1_{\mathbb{P}}(\Xi)$. The risk measure \mathcal{R} should return a typical value of the random variable, e. g., the expectation or a quantile.

Provided that the state equation $E(y(\xi), u, \xi) = 0$ admits a unique solution $y(u)(\xi) \in Y$ for every $u \in U_{\text{ad}}$ and almost every $\xi \in \Xi$, (1.2) can be reduced to an optimization problem

over the deterministic control variable only:

$$\min_{u \in U} \mathcal{R}[\hat{J}_1(u, \cdot)] + J_2(u) \quad \text{s. t.} \quad u \in U_{\text{ad}}, \quad (1.3)$$

where $\hat{J}_1(u, \xi) := J_1(y(u)(\xi), u, \xi)$. An algorithm for solving problems of the form (1.3) is presented and analyzed in Chapter 4. It can deal with inexact objective function and gradient evaluations as well as an inexact projection onto the feasible set U_{ad} .

Modeling Optimal Control under Uncertainty

We want to stress that the formulation (1.2) does clearly not cover all imaginable formulations of optimal control problems under uncertainty. State constraints are excluded and the control may not be chosen dependent on the concrete realization of ξ . For an overview of different variants of modeling optimal control problems under uncertainty we refer to [4, Sec. 2]. The setup (1.2) yields “robust deterministic controls”, cf. [4, Sec. 2], i. e., controls which are chosen before observing the random inputs, but account for the uncertainty in the system. Further works following this approach are [69, 70, 71] and our paper [46].

In contrast to that, optimal controls for specific realizations of the random parameters ξ are considered in [25, 4]. Having computed the set of these controls in an offline phase allows to quickly select the optimal control online after observing the concrete values of the parameters. Alternatively, this can be understood as an uncertainty quantification task. Knowing how the optimal control is computed, one can estimate its distribution from the distribution of the parameters. Another alternative is to work with the distribution of the uncertain state $y(\xi)$ instead of only inserting it into a deterministic objective function. For instance, in [108] the stochastic moments of the state are fit to desired values. Similarly, the variance of the state is penalized in [20].

If we consider the formulation (1.2), an important modeling aspect is the choice of the risk measure \mathcal{R} . A natural choice is the expectation $\mathcal{R} \equiv \mathbb{E}$ used in [24, 69, 70, 68], our work [46], and in the main part of this thesis (Chapters 3, 5, 6, and 8). This approach is called *risk-neutral control*, see [103, Sec. 6.4] and [3], and does not account for rare but possibly harmful or costly events or the variability of the cost functional, but intends to have a small cost on average.

More general risk measures allowing for *risk-averse control*, see [103, Chap. 6] and [71, 3, 72], and the underlying theory are presented in Chapter 9. In particular, one can choose the smooth *mean-variance risk measure* $\mathcal{R} \equiv \mathbb{E} + \lambda \text{Var}$ for some $\lambda > 0$ [3], which is discussed in more detail in Section 9.2. Other risk-averse risk measures such as the *conditional value-at-risk (CVaR)*—also called tail expectation—enjoy desirable theoretical properties which the mean-variance risk measure lacks, but are nonsmooth. Hence, it is suitable to apply smoothing techniques which retain important properties of the risk measure and enable the use of gradient-based optimization methods [71]. An interior point solution technique for such risk measures is discussed in Section 9.3.

Numerical Solution of PDEs with Uncertain Inputs

A key ingredient of an algorithm for solving (1.2) is a suitable numerical solution method for the state equation $E(y(\xi), u, \xi) = 0$, which is a PDE with uncertain inputs in our context. Overviews over numerical methods for elliptic PDEs with random inputs are given in [102, 53], where the equations are formulated in the Bochner space $L^2_{\mathbb{P}}(\Xi; Y)$, where Y is a Hilbert space, and in [7], which describes stochastic collocation techniques for linear, elliptic PDEs. In contrast to that, we consider the optimal control of *semilinear*, elliptic PDEs under uncertainty, which requires a more sophisticated analysis. In particular, it is necessary to consider the state space $L^p_{\mathbb{P}}(\Xi; Y)$ with $p > 3$ in our setting, see Chapter 3.

One approach to discretize the parameter space is a sampling-based method. This means that certain realizations of the random parameters ξ are inserted and the respective deterministic PDE is solved. The ensemble of deterministic solutions is then used to approximate the full solution, its distribution function, or stochastic moments. An advantage of such methods is that they are often non-intrusive meaning that deterministic black-box solvers can be reused for the computation. Monte Carlo-type methods, e. g., multilevel Monte Carlo [16, 30, 107, 4] and quasi-Monte Carlo [84, 80], are very popular. They can be applied under fairly general circumstances and feature a typically rather slow convergence speed, which does not depend on the number of uncertain parameters. If the stochastic space Ξ has tensor product structure, adaptive sparse grids can be used for quadrature or interpolation [24, 25, 69, 68]. Under smoothness assumptions on the integrand, they feature exponential convergence [86, 47], but often need a number of collocation points which increases exponentially in the number of parameters [45]. This can lead to a high computational cost if systems with many parameters shall be solved. Furthermore, some negative quadrature weights may appear, which can make discretized optimization problems non-convex or ill-posed although their continuous counterparts do not have these properties. Additionally, it is not possible then to make a smooth reformulation of optimization problems with certain nonsmooth risk measures such as the CVaR, where the nonsmoothness is moved to pointwise bound constraints. The solution of many similar deterministic PDEs in every collocation point can be sped up by using reduced basis methods [42, 28], proper orthogonal decomposition [25], or multigrid techniques [24]. The references [4, 24, 25, 69, 68, 28] indeed focus on optimal control and not only on solution techniques for stochastic PDEs.

Another option is a stochastic Galerkin discretization of the PDE with uncertain inputs [38, 39, 22, 23], where polynomials are chosen to discretize the dependence on the uncertain parameters. To make this efficient, either one chooses a small enough, “sparse” polynomial basis or a large basis with a data-sparse representation of the coefficients, e. g., a tensor product basis with a low-rank tensor representation as in [40]. In the latter case, one can benefit from the good complexity properties of modern low-rank tensor formats w. r. t. the number of parameters, but has to take errors due to tensor truncation into account. The papers [38, 39, 22, 23, 40] deal with adaptive solution techniques for linear, elliptic PDEs with uncertain inputs based on a posteriori error estimates. In this work, we derive a mixture of those, extend it to semilinear, elliptic PDEs, and apply it in the context of optimal control. Hence, we do not focus on error estimates in the energy norm for example, but choose a suitable reference norm, see Chapter 7. To simplify the discretization of the nonlinearity, weighted Lagrange polynomials are chosen as bases of the spaces of polynomials instead of the

usual basis of polynomials with increasing degree, cf. [44] and see Chapter 6. This is related to a full-grid Gaussian quadrature formula, which provides high accuracy and only positive weights. As pointed out in the discussion about negative quadrature weights appearing in sparse grid quadrature formulas, the latter is helpful in an optimization context because smooth reformulations of certain nonsmooth problems are possible in the discrete setting. Stochastic Galerkin methods for control problems were already used in [63, 82, 78, 79, 46]. In combination with low-rank tensors—the option we follow in this work—they were also applied in [21], but no control constraints were posed and the paper puts more focus on efficient linear algebra and works with a fixed discretization, whereas we select the discretization adaptively during the optimization process.

Optimization Algorithms

In context of optimal control of PDEs, different inexact trust-region type algorithms have been proposed to solve problems adaptively and control the appearing errors due to discretization or the inexact solution of equations in finite dimensions while guaranteeing global convergence. In [119] deterministic PDE-constrained optimization problems are solved by an inexact trust-region sequential quadratic programming (SQP) algorithm, which is applied to the non-reduced optimal control problem, where the state equation is viewed as an equality constraint, cf. (1.2). The sources of inexactness are an adaptive finite element discretization and the inexact solution of the linearized state equation by a conjugate gradient (CG) method. Error estimates are used to control the error which is required for the optimization method to converge. The paper [118] extends this approach to control constraints using an inexact projection onto the feasible set.

In [70] and its predecessor [69], a class of inexact trust-region algorithms formulated in Hilbert space is presented and applied to solve optimal control problems of PDEs under uncertainty. Since these algorithms are posed on the control space U , they aim for solving the reduced optimal control problem (1.3). The algorithms in [69] allow for inexact gradient evaluations. They are extended to inexact objective function evaluations in [70]. In the application, the inexactness due to the adaptive discretization of the stochastic space by sparse grids is controlled by the algorithm, but no further error sources are considered. Furthermore, the proposed algorithms cannot handle control-constrained problems.

We extend [70] to the control-constrained case in Chapter 4 by using a possibly inexact projection onto the feasible set as done in [118], but allow for arbitrary mappings approximating the exact projection instead of working with the exact projection on a discrete subspace. We then apply it to the optimal control problem under uncertainty and control all appearing errors simultaneously, namely the FE and the stochastic Galerkin discretization error as well as the algebraic error coming from a low-rank tensor solver used to solve the discretized stochastic PDE. In the paper [46], we applied a semismooth Newton method having fast local convergence properties to solve such problems with a fixed discretization. This method does not feature global convergence and the errors due to the used low-rank tensor arithmetics were not controlled in [46]. Now, we use it as a solver for the trust-region subproblems. The trust-region framework guarantees global convergence and controls all errors appropriately.

Outline of the Thesis

The rest of this thesis is structured as follows: An introduction to tensors including the used notation, operations, low-rank formats, and available algorithms is given in Chapter 2, where Section 2.1 focuses on finite-dimensional tensors. As already mentioned, one motivation and central question of the dissertation is how low-rank tensor methods can be used to solve optimal control problems under uncertainty. Therefore, we review tensor products of Hilbert spaces and certain Banach spaces in Section 2.2.

This concept is important for the formulation of the class of stochastic PDEs considered in Chapter 3 as well as its analysis and discretization in tensor form. We analyze semilinear, elliptic PDEs with uncertain inputs $\xi \in \Xi$, where the nonlinearity is a superposition operator induced by a \mathcal{C}^2 -function $\varphi : \mathbb{R} \rightarrow \mathbb{R}$ fulfilling suitable growth conditions. Two formulations are presented: Either a deterministic PDE can be solved for a. e. realization $\xi \in \Xi$ or a weak formulation w. r. t. the inputs ξ can be considered by integrating over the parameter space Ξ . The latter formulation is essential to apply a stochastic Galerkin discretization. The solution of the PDE is then represented as an element of a suitable Bochner space $L_{\mathbb{P}}^p(\Xi; Y)$ with $p \in (3, \infty)$. This requirement on p is necessary for well-definedness and differentiability properties of the appearing superposition operator. The connection between the Bochner space $L_{\mathbb{P}}^p(\Xi; Y)$ and a tensor product of Banach spaces is established to derive a discretization in tensor form later. We show the equivalence of the pointwise and the weak formulation. In fact, they both have the same unique solution. At this point, it is important to show that the solution constructed pointwise for a. e. ξ is indeed measurable w. r. t. ξ and has the required integrability properties due to a priori estimates, i. e., it is contained in the space $L_{\mathbb{P}}^p(\Xi; Y)$. We formulate a class of risk-neutral optimal control problems with a tracking-type objective function and analyze it. The integrability properties of the state are used to apply the dominated convergence theorem at some points, e. g., to show the existence of a solution to the optimal control problem (Theorem 3.19) and the differentiability of the control-to-state mapping (Theorem 3.26). An adjoint representation of the gradient and Hessian of the reduced objective function is derived. The theory from, e. g., [60, Sec. 1.6] cannot be applied at this point since the linearized state equation operator is not boundedly invertible as a mapping from the state space $L_{\mathbb{P}}^p(\Xi; Y)$ to its dual. It is shown that a pointwisely defined adjoint state can be used to compute the gradient and a meaningful interpretation of the adjoint equation in the space $L_{\mathbb{P}}^p(\Xi; Y)$ is presented.

The reduced optimal control problem, which is an optimization problem formulated in a Hilbert space U with a smooth objective function and a nonempty, closed, convex constraint set $U_{\text{ad}} \subset U$, shall be solved adaptively. Due to discretization errors and inexact state and adjoint equation solves, the objective function and its gradient can only be evaluated inexactly in many relevant cases. Possibly, additional errors appear if the projection onto the feasible set U_{ad} is computed inexactly. Tailored to this situation, an inexact trust-region algorithm suitable for smooth problems with control constraints is formulated in Chapter 4. It can deal with the mentioned inexact evaluations and features global convergence provided the errors are controlled as demanded by the algorithm. The respective required error bounds are up to fixed, but possibly unknown constants, which typically appear in error estimates due to discretization. Global convergence is proven and it is discussed in detail how the required error tolerances can be met and implemented. In particular, the computation of

a generalized Cauchy point by a projected linesearch with an inexact, but suitably refined projection is discussed. A semismooth Newton method in function space is proposed for the solution of the trust-region subproblems.

The presented trust-region algorithm requires to evaluate the objective function, gradient, and criticality measure of the reduced optimal control problem to a certain accuracy. In Chapter 5, these error bounds are transferred to error bounds on the state, adjoint state, and the projection onto the set of admissible controls. It turns out that—depending on the problem data—the state and adjoint state errors need not be controlled in the $L^2_{\mathbb{P}}(\Xi; Y)$ -norm, but a weaker norm such as the $L^2_{\mathbb{P}}(\Xi; Y)$ -norm can be sufficient in certain cases.

Chapter 6 describes the discretization of the model problem presented in Chapter 3 and in particular of the PDE with uncertain inputs. We use linear finite elements for the space discretization and a stochastic Galerkin discretization with weighted Lagrange polynomials in tensor product form for the parameter space. It is shown how the discretized equations are formulated within the tensor calculus and how they can be solved by low-rank tensor solvers. The choice of the trust-region model in the discrete setting and the implementation of the semismooth Newton method for box constraints is explained in detail.

As derived in Chapter 5, error control in $L^2_{\mathbb{P}}(\Xi; Y)$ of the state and adjoint state can be sufficient to obtain a globally convergent algorithm. Based on the presented discretization, an a posteriori error estimate is discussed in Chapter 7. It splits into error contributions stemming from the FE discretization error, the stochastic Galerkin error for each random parameter and the algebraic error coming from the iterative low-rank tensor solver. The implementation of the evaluation of this error estimate with low-rank tensors is described. This is crucial because the solutions of the discretized PDEs are represented by low-rank tensors, the full representation of which, i. e., a full array of real numbers should never be computed explicitly due to efficiency reasons. Based on this, the PDEs with uncertain inputs can be solved adaptively.

In Chapter 8, numerical results of the described solution technique are shown for different setups. It can be seen how the algorithms adapt to the problem data, especially in terms of mesh refinement and polynomial grades chosen for the discretization of each parameter. Furthermore, we compare the optimal control obtained from deterministic optimization with the robust one.

Chapter 9 shows how different, risk-averse risk measures can be incorporated into the optimal control problem since the previous parts of the thesis deal with risk-neutral control. A general introduction to risk measures is given. Then, including the mean-variance risk measure into the objective function is discussed in theory. At the end, we investigate a log-barrier reformulation for problems with risk-averse, nonsmooth risk measures which are convex combinations of the mean and the CVaR. It is examined how solving the respective barrier problem, which is an approximation of the original problem, affects the underlying risk measure. At the end, numerical results of an implementation with a fixed discretization and low-rank tensors are presented. It is observed that the control can be chosen such that very large cost function values can be excluded with high probability.

The thesis is concluded in Chapter 10, where we also point out future research perspectives and a few aspects not covered in this work.

Notation

In analogy to [46], the following notation is used throughout the thesis: Deterministic functions, e. g., the control u , and the corresponding function spaces and operators are written in italic font. If we want to emphasize that a function such as the state \mathbf{y} belongs to some Bochner space of the form $L_{\mathbb{P}}^p(\Xi; Y)$, i. e., it also depends on the uncertain parameters ξ , we use bold and italic font. Functions from finite-dimensional (discretized) spaces are written in roman font (u and \mathbf{y}), and the coefficients representing them, such as the vector \mathbf{u} and the tensor \mathbf{y} , are in sans-serif font. The same holds for the respective operators.

2. Tensors

One motivation of this work is the research on low-rank tensor methods carried out in the numerical linear algebra community in the past years. Tracing back to, e. g., even the 1920s, where one can find the roots of the canonical polyadic decomposition [61], low-rank tensor methods recently got attention within the scientific computing community because of their ability to break the curse of dimensionality. Especially new low-rank formats such as Tensor Train (TT) [88] and Hierarchical Tucker (HT) [56] have the properties of providing stable low-rank approximations and scaling well w. r. t. the tensor dimension. Older formats such as the canonical polyadic decomposition (CP) [93] and the Tucker format [109] lack either the first or the second property, respectively. Suitable numerical algorithms within the formats have been developed during the last two decades, approximately.

In this chapter, we give an introduction to finite-dimensional tensors and low-rank formats, which we will use in our numerical algorithms, and to tensor Hilbert spaces, which are suitable for the formulation of parametric problems in function spaces. Certain examples of tensor Banach spaces needed in Chapter 3 are also discussed.

2.1. Finite-Dimensional Tensors and Low-Rank Formats

After discretization, the state of the systems considered in this work can be represented as a real tensor, i. e., a multi-dimensional array of real numbers. To avoid the curse of dimensionality we will represent it in a low-rank format. Both—finite-dimensional tensors and the corresponding low-rank formats—are the topic of this section, which is very similar to the description in our paper [46].

2.1.1. Basics of Finite-Dimensional Tensors and Notation

In this thesis, we denote a d -dimensional array of numbers in the field \mathbb{R} by a *finite-dimensional tensor* $\mathbf{x} \in \mathbb{R}^{n_1 \times \dots \times n_d}$, cf. [76, Sec. 1], where \mathbb{C} is used as underlying field. The array dimension d is called the *order* of the tensor and the numbers $1, \dots, d$ are the *modes*. We denote the *dimension of the i -th mode* by n_i and write $\mathbf{x}(k_1, \dots, k_d)$ for the (k_1, \dots, k_d) -component of \mathbf{x} for readability purposes instead of using a subscript notation, where $k_i \in [n_i]$ and $i \in [d]$, writing $[m] := \{1, \dots, m\}$ here and throughout. The tensors of all ones will be denoted by $\mathbb{1}$ if it is clear which space is used. *Reshaping* a tensor as a vector in a certain order, e. g., reverse lexicographical as in [76], is denoted by $\text{vec}(\mathbf{x}) \in \mathbb{R}^{\prod_{i=1}^d n_i}$ and reshaping as a matrix by $\mathbf{x}^{(t)} \in \mathbb{R}^{(\prod_{i \in t} n_i) \times (\prod_{j \notin t} n_j)}$, where $t \subset [d]$ is the set of modes that become the rows of the resulting matrix. Analogously we write $\text{ten}(x)$ for reshaping a vector or matrix x back into a tensor, when the dimensions are clear. For the *extraction* of parts of a tensor we also allow indexing with ordered lists and write “ \bullet ” for all indices of a dimension in ascending order: $\mathbf{x}(\bullet, 4, (1, 5), \bullet, \dots, \bullet) \in \mathbb{R}^{n_1 \times 2 \times n_4 \times \dots \times n_d}$ is obtained from \mathbf{x} by

fixing the second index at 4 and taking the first and the fifth component of the third mode and all components in the rest of the modes. Note that modes that are indexed by a single index only are cut out in the result while modes indexed by a list remain. This is especially relevant when indexing by a list containing only one index: $\mathbf{x}(2, \bullet, \dots, \bullet) \in \mathbb{R}^{n_2 \times \dots \times n_d}$, but $\mathbf{x}((2), \bullet, \dots, \bullet) \in \mathbb{R}^{1 \times n_2 \times \dots \times n_d}$. Sometimes, the notation becomes shorter if sets instead of lists are used for indexing. The respective set then stands for the list containing all its elements in ascending order. In the i -th mode of a tensor, one could replace “ \bullet ” by “[n_i]” for example. For a compact notation, we also allow indexing by index vectors, i. e., $\mathbf{x}(k) = \mathbf{x}(k_1, \dots, k_d)$ for some $k \in \times_{i=1}^d [n_i]$.

As tensors form a vector space, all vector space operations are defined for tensors. Additionally, *componentwise operations* are useful, especially for function-related tensors: $\mathbf{x} \odot \mathbf{y}$ denotes multiplication, $\mathbf{x} \oslash \mathbf{y}$ division and \mathbf{x}^λ exponentiation by a scalar $\lambda \in \mathbb{R}$. The componentwise application of a function $f : \mathbb{R} \rightarrow \mathbb{R}$ is written as $f(\mathbf{x})$.

Definition 2.1 (*i -mode matrix product*). Let $\mathbf{x} \in \mathbb{R}^{n_1 \times \dots \times n_d}$ be a tensor, $i \in [d]$ a mode and $\mathbf{A} : \mathbb{R}^{n_i} \rightarrow \mathbb{R}^m$ a linear operator or a matrix.

Then the i -mode matrix product $\mathbf{A} \circ_i \mathbf{x} \in \mathbb{R}^{n_1 \times \dots \times n_{i-1} \times m \times n_{i+1} \times \dots \times n_d}$ is defined by

$$(\mathbf{A} \circ_i \mathbf{x})(k_1, \dots, k_d) := (\mathbf{A}\mathbf{x}(k_1, \dots, k_{i-1}, \bullet, k_{i+1}, \dots, k_d))(k_i)$$

for all $(k_1, \dots, k_d) \in [n_1] \times \dots \times [n_{i-1}] \times [m] \times [n_{i+1}] \times \dots \times [n_d]$.

Remark 2.2. Equivalently, the i -mode matrix product is given by $\mathbf{A} \circ_i \mathbf{x} = \text{ten}(\mathbf{A}\mathbf{x}^{\{i\}})$, cf. [76, Sec. 4.1], using reshaping and the standard product of matrices. The definition means that we take all vectors that result from fixing all indices of the tensor \mathbf{x} except the i -th one, apply \mathbf{A} and then use the resulting vectors to form $\mathbf{A} \circ_i \mathbf{x}$.

We view tensors of order 1 as column vectors and order-2-tensors as matrices, where the first index is for the rows. The contraction of two tensors is a further important operation. It generalizes inner and outer products as well as matrix multiplication.

Definition 2.3 (*tensor contraction*). Let $\mathbf{x} \in \mathbb{R}^{n_1 \times \dots \times n_d}$ and $\mathbf{z} \in \mathbb{R}^{\tilde{n}_1 \times \dots \times \tilde{n}_{\tilde{d}}}$ be tensors and let $s = (s_1, \dots, s_p)$, $t = (t_1, \dots, t_p)$ with $s_i \in [d]$, $t_i \in [\tilde{d}]$ be ordered lists of modes that shall be contracted. We require $n_{s_i} = \tilde{n}_{t_i}$ for all $i \in [p]$. Furthermore, let \bar{s}_i ($i \in [d-p]$) and \bar{t}_i ($i \in [\tilde{d}-p]$) be the remaining, untouched modes in ascending order.

Then the *contraction* $\langle \mathbf{x}, \mathbf{z} \rangle_{s,t} \in \mathbb{R}^{n_{\bar{s}_1} \times \dots \times n_{\bar{s}_{d-p}} \times \tilde{n}_{\bar{t}_1} \times \dots \times \tilde{n}_{\bar{t}_{\tilde{d}-p}}}$ of \mathbf{x} and \mathbf{z} along the modes s and t is defined as

$$\langle \mathbf{x}, \mathbf{z} \rangle_{s,t}(k_{\bar{s}_1}, \dots, k_{\bar{s}_{d-p}}, \ell_{\bar{t}_1}, \dots, \ell_{\bar{t}_{\tilde{d}-p}}) = \sum_{k_{s_1}, \dots, k_{s_p}=1}^{n_{s_1}, \dots, n_{s_p}} \mathbf{x}(k_1, \dots, k_d) \mathbf{z}(\ell_1, \dots, \ell_{\tilde{d}}) |_{\ell_{t_i}=k_{s_i} \forall i \in [p]}$$

componentwise for all indices $k_{\bar{s}_i} \in [n_{\bar{s}_i}]$ ($i \in [d-p]$) and $\ell_{\bar{t}_j} \in [\tilde{n}_{\bar{t}_j}]$ ($j \in [\tilde{d}-p]$). If s or t are sets rather than ordered lists, they stand for the lists of elements in ascending order as in the case of indexing.

Here, similar to a matrix product, a component of the resulting tensor is obtained by fixing indices in the untouched modes and computing an inner product of the resulting tensors of

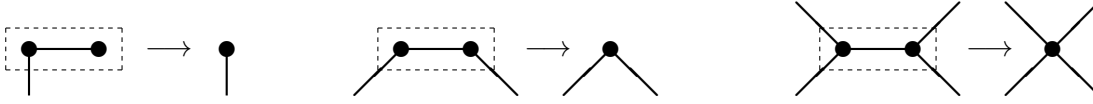


Figure 2.1.: Examples for tensor network diagrams: matrix-vector-product resulting in a vector (left), matrix-matrix-product resulting in a matrix (middle), two tensors of order 3 contracted along one mode resulting in a tensor of order 4 (right)

same size. Tensor contraction can be nicely visualized by tensor network diagrams [76, Sec. 5.1]. The resulting tensor of the contraction is represented by a network, where each node stands for a tensor in the contraction and the edges connected to it are the tensor modes. If two nodes are connected by an edge, the respective tensors are contracted along the respective modes. The “free” edges indicate the modes of the resulting tensor. Some simple examples of tensor networks are shown in Figure 2.1. Note that for a full description of the tensor contraction one should, e. g., display the mode numbers next to the edges, which we skip here as long as it is clear or not essential.

We sometimes need the following special cases of tensor contraction:

Definition 2.4 (special cases of tensor contraction). Let $\mathbf{x}, \mathbf{y} \in \mathbb{R}^{n_1 \times \dots \times n_d}$ and $\mathbf{z} \in \mathbb{R}^{\tilde{n}_1 \times \dots \times \tilde{n}_{\tilde{d}}}$ be tensors.

- The *outer product*¹ $\mathbf{x} \otimes \mathbf{z} := \langle \mathbf{x}, \mathbf{z} \rangle_{\emptyset, \emptyset} \in \mathbb{R}^{n_1 \times \dots \times n_d \times \tilde{n}_1 \times \dots \times \tilde{n}_{\tilde{d}}}$ of the tensors \mathbf{x} and \mathbf{z} is obtained as

$$(\mathbf{x} \otimes \mathbf{z})(k_1, \dots, k_d, \ell_1, \dots, \ell_{\tilde{d}}) = \mathbf{x}(k_1, \dots, k_d) \mathbf{z}(\ell_1, \dots, \ell_{\tilde{d}}).$$

- An *elementary tensor* or *rank-1-tensor* $\bigotimes_{i=1}^d \mathbf{u}^i \in \mathbb{R}^{n_1 \times \dots \times n_d}$ is an outer product of vectors:

$$\left(\bigotimes_{i=1}^d \mathbf{u}^i \right) (k_1, \dots, k_d) = \prod_{i=1}^d u_{k_i}^i,$$

where $u_{k_i}^i$ is the k_i -th component of the vector $\mathbf{u}^i \in \mathbb{R}^{n_i}$.

- The *inner product* $\langle \mathbf{x}, \mathbf{y} \rangle := \langle \mathbf{x}, \mathbf{y} \rangle_{[d], [d]} \in \mathbb{R}$ of two tensors of same size is

$$\langle \mathbf{x}, \mathbf{y} \rangle = \sum_{k_1=1}^{n_1} \dots \sum_{k_d=1}^{n_d} \mathbf{x}(k_1, \dots, k_d) \mathbf{y}(k_1, \dots, k_d).$$

- The *Frobenius norm* of a tensor \mathbf{x} is $\|\mathbf{x}\|_F := \sqrt{\langle \mathbf{x}, \mathbf{x} \rangle}$.

We conclude this section with the statement that the i -mode matrix product is also a special case of contraction since the tensor $\mathbf{A} \circ_i \mathbf{x}$ is obtained as contraction of the second order tensor \mathbf{A} with the tensor \mathbf{x} along the modes (2) and (i).

¹Note that the vectorization/matricization of the outer product of two tensors coincides with the standard Kronecker product (also denoted by " \otimes ") of vectorizations/matricizations of them.

2.1.2. Low-Rank Tensor Formats

As the amount of data of a full tensor is in general very high, namely of order $\mathcal{O}(n_{\max}^d)$ with $n_{\max} := \max_{i \in [d]} n_i$, it is important to find a representation or approximation of it that requires substantially less memory. One approach is to approximate a tensor by one of low rank. Different notions of tensor rank have been developed in the past. For instance, in the *canonical polyadic (CP) decomposition*, cf. [54, Chap. 7], a rank- R -tensor \mathbf{x} is written as the sum of R elementary tensors:

$$\mathbf{x} = \sum_{j=1}^R \bigotimes_{i=1}^d \mathbf{u}^{i,j} \quad \text{with} \quad \mathbf{u}^{i,j} \in \mathbb{R}^{n_i} \text{ for } i \in [d], j \in [R].$$

This representation requires storage of order $\mathcal{O}(Rn_{\max}d)$. We note that this is linear in d instead of exponential in the case of storing the full tensor. But one should take into account that R can grow if a good approximation of a given tensor is desired. Furthermore, this format has certain drawbacks as the possible ill-posedness of the best approximation of a given tensor by a tensor of rank at most R [33]. Concretely, this means that there exists a rank-3-tensor of order $d \geq 3$, which can be approximated arbitrarily well by a rank-2-tensor [54, Prop. 9.10].

A widely used approach for the matrix case $d = 2$ is a low rank approximation by a truncated singular value decomposition (SVD). The quadratic error (squared Frobenius norm) made by this approximation is bounded by the sum of the truncated, squared singular values [54, Lem. 2.30] and can thus be estimated easily. The SVD approximation is optimal w. r. t. this error in the set of matrices of a fixed maximum rank. A generalization of this approximation technique to d -dimensional tensors is therefore desirable, but cannot be available in the CP format due to the ill-posedness of the approximation problem. Thus, a different approach works with subspaces of the \mathbb{R}^{n_i} , as done in the *Tucker format* or *tensor subspace representation* [54, Chap. 8]: Here, for representing a tensor $\mathbf{x} \in \bigotimes_{i=1}^d \mathbb{R}^{n_i}$ (more on this tensor product of Hilbert spaces in Section 2.2), bases of r_i -dimensional subspaces of the \mathbb{R}^{n_i} are selected, stored as matrices $\mathbf{U}_i \in \mathbb{R}^{n_i \times r_i}$ and a basis $\bigotimes_{i=1}^d \mathbf{U}_i$ of the tensor product of subspaces is formed. One can combine the basis elements linearly with coefficients stored in the so-called *core tensor* $\mathbf{c} \in \mathbb{R}^{r_1 \times \dots \times r_d}$ to obtain the tensor \mathbf{x} . Equivalently, this representation of \mathbf{x} can be written as a contraction of the tensors \mathbf{U}_i and \mathbf{c} w. r. t. a tensor network as depicted in Figure 2.2. Actually, the CP format is a special case of this network with $r_i = R$ for all $i \in [d]$, \mathbf{U}_i containing exactly the vectors $\mathbf{u}^{i,j}$ and $\mathbf{c} \in \mathbb{R}^{R \times \dots \times R}$ being diagonal in the sense that only the entries $\mathbf{c}(k_1, \dots, k_d)$ with $k_1 = \dots = k_d$ are allowed to be non-zero. For the Tucker format the *HOSVD* (higher order SVD, which uses the standard SVD in substeps) is available for truncating a tensor to lower Tucker rank giving a quasi-optimal result² and offering error control similar to the one in standard SVD [81]. We give an impression of SVD based tensor truncation in Subsection 2.1.3, using the example of conversion from TT to HT tensors. Unfortunately, defining $r_{\max} := \max_{i \in [d]} r_i$, the required storage for the Tucker format is $\sum_{i=1}^d n_i r_i + \prod_{i=1}^d r_i = \mathcal{O}(n_{\max} r_{\max} d + r_{\max}^d)$ and grows exponentially in d .

²Best approximation in the Frobenius norm up to a factor that depends on the order d of the tensor (\sqrt{d} for the Tucker format, see, e. g., [51]).

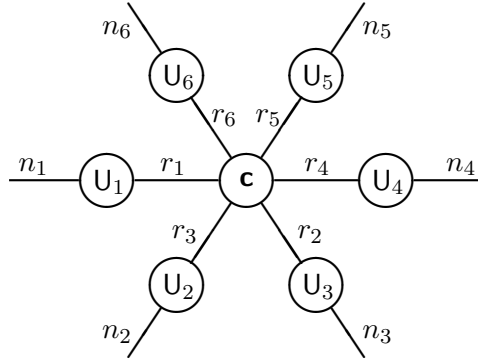


Figure 2.2.: Tensor network representing the Tucker format in the case $d = 6$. The numbers next to the edges are the respective mode dimensions.

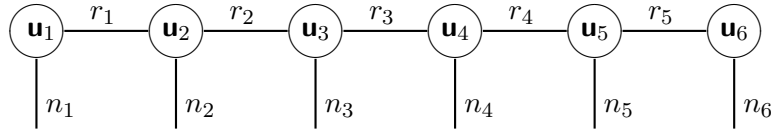


Figure 2.3.: Tensor network representing the Tensor Train format in the case $d = 6$ with the dimensions of the modes displayed next to the edges

Two more recent low-rank tensor formats, meeting the requirement of having a storage and arithmetic complexity scaling linearly in n_{\max} and d , and for which a variant of a higher-order SVD is available, are the *hierarchical Tucker (HT) format* [56], and the *tensor train (TT) format* [88], also known as *matrix product states (MPS)* from the quantum physics community [115]. The latter works as follows:

The *TT rank* of a tensor $\mathbf{x} \in \mathbb{R}^{n_1 \times \dots \times n_d}$ is a vector $r := (r_0, r_1, \dots, r_d)$ with $r_0 = r_d = 1$ and $r_i \in [n_i]$ for $i \in \{2, \dots, d-1\}$. In the TT format, $\mathbf{x} = \mathcal{TT}(\mathbf{u})$ is represented by a d -tuple $\mathbf{u} := (\mathbf{u}_1, \dots, \mathbf{u}_d) \in \times_{i=1}^d \mathbb{R}^{r_{i-1} \times n_i \times r_i}$ of order-3-tensors:

$$\mathbf{x}(k_1, \dots, k_d) = \mathbf{u}_1(\cdot, k_1, \cdot) \mathbf{u}_2(\cdot, k_2, \cdot) \cdots \mathbf{u}_d(\cdot, k_d, \cdot) \in \mathbb{R}^{1 \times 1} \cong \mathbb{R} \quad \forall k_i \in [n_i], i \in [d].$$

Note that this is a simple product of matrices, resulting in a scalar value because the first and the last matrix of the chain are a row and a column vector, respectively. The tensor \mathbf{x} can also be written as a contraction of the so-called core tensors \mathbf{u}_i as shown in Figure 2.3. The storage requirement of a tensor in the TT format is exactly $\sum_{i=1}^d r_{i-1} n_i r_i = \mathcal{O}(n_{\max} r_{\max}^2 d)$ with $r_{\max} := \max_{i \in \{0, 1, \dots, d\}} r_i$ and scales linearly in n_{\max} and d . As in the case of the Tucker format, the low-rank approximation problem is well-posed for TT tensors and an SVD-based truncation to smaller TT rank, the so-called *TT-SVD*, is available and gives a quasi-optimal result up to the factor $\sqrt{d-1}$ [88, Cor.2.4]. An error control mechanism based on the truncated singular values allows to choose the ranks such that a prescribed approximation accuracy can be guaranteed.

In the related *hierarchical Tucker (HT) decomposition* [56], the idea of writing a tensor as a contraction of smaller ones is generalized in a certain way: The contraction is performed according to a more general, binary dimension tree \mathcal{T} , which has one root node $[d]$. Each node

$t \in \mathcal{T}$ containing more than one index splits up into two children $t_l, t_r \neq \emptyset$ with $t = t_l \cup t_r$, where we assume w.l.o.g. that t_l contains only smaller indices than the ones in t_r . The leaves of the tree are sets of only one index so that we have d leaves and $2d - 1$ nodes overall. To simplify the notation we call the leaf nodes t_1, \dots, t_d , the parent node of t_1 and t_2 would be called t_{12} and the root node is $t_{12\dots d}$. The tree tells us how to form a tensor starting with subspaces of the \mathbb{R}^{n_i} and recursively selecting subspaces of tensor products of subspaces. For each leaf node t_i ($i \in [d]$) a matrix $U_i \in \mathbb{R}^{n_i \times r_i}$ is stored, which contains vectors that span a linear subspace V_i of \mathbb{R}^{n_i} as in the Tucker format. One can use an orthogonal basis here for example. Now, the full tensor is represented by these leaf matrices and additional transfer matrices $B_t \in \mathbb{R}^{r_{t_l} \times r_{t_r}}$ ($r_t \leq r_l r_r$, $r_{[d]} = 1$) assigned to the non-leaf nodes t , as follows: For computing the full tensor, which is normally never done explicitly, one would recursively go through the tree from the leaves to the root. At each parent node t a generating system of the tensor product of the two subspaces at the children t_l and t_r would be built combining all vectors, i. e., the matrix $U_r \otimes U_l \in \mathbb{R}^{n_l n_r \times r_l r_r}$ would be computed. This can also be viewed as a tensor; the Kronecker product is a possible matricization of it. Then, by multiplying with the transfer matrix B_t , which can also be viewed as a three-dimensional tensor $\mathbf{b}_t \in \mathbb{R}^{r_l \times r_r \times r_t}$, a generating system of a subspace would be obtained. It can be represented by the matrix $U_t = (U_r \otimes U_l) B_t \in \mathbb{R}^{n_l n_r \times r_t}$, which would contain vectors generating a subspace of $\bigotimes_{i \in t} V_i$. Arriving at the root node one would obtain a $\mathbb{R}^{\prod_{i=1}^d n_i}$ -vectorization of the represented tensor. On the other hand we can view the represented tensor to be a contraction of the leaf matrices U_i and transfer tensors \mathbf{b}_t according to a tensor network defined by the tree \mathcal{T} as shown in Figure 2.4, where a typical balanced tree is used. The required storage for the hierarchical Tucker representation is

- $\mathcal{O}(n_{\max} r_{\max} d)$ for the basis matrices U_i at each of the d leaf nodes and
- $\mathcal{O}(r_{\max}^3 d)$ for the transfer matrices B_t at the non-leaf nodes with $r_{\max} := \max_{t \in \mathcal{T}} r_t$,

and sums to $\mathcal{O}(n_{\max} r_{\max} d + r_{\max}^3 d)$. This is linear in d , but again, for a good approximation the hierarchical Tucker rank, i. e., the vector of ranks at the nodes of \mathcal{T} , can be required to be chosen dependent on d and n_{\max} . Truncation to a lower rank is well-posed and can be done by the SVD-based *hierarchical SVD (HSVD)* [50] with error control and quasi-optimal approximation up to the factor $\sqrt{2d - 3}$.

2.1.3. Available Algorithms

Due to their good complexity, we want to use TT or HT tensors for all high-dimensional computations. Many relevant algorithms are available within (MATLAB) toolboxes, which we present in the following. In addition, we discuss some iterative algorithms used for computing quantities which cannot be computed directly with low-rank tensors in a reasonable amount of time. For instance, the componentwise application of a function $f : \mathbb{R} \rightarrow \mathbb{R}$ to a low-rank tensor could *in theory* be computed by building the full array, applying f to its components and truncating the result to obtain a low-rank tensor again. But this approach is highly prohibitive because even the storage required for the full tensor could easily exceed any available storage, at least if the order d is large enough and $n_i \geq 2$ for all $i \in [d]$. Therefore, it is important to compute or approximate every tensor within a computation in a low-rank

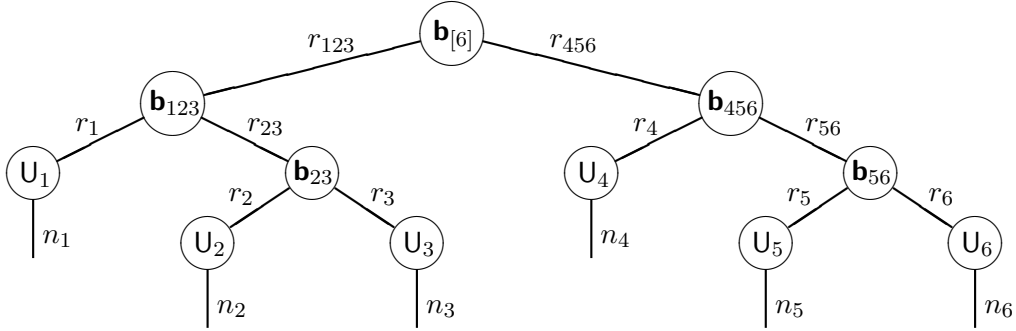


Figure 2.4.: Tensor network representing the Hierarchical Tucker format in the case $d = 6$ with a balanced tree. The dimensions of the modes are displayed next to the edges.

format without ever forming a full tensor explicitly. If there are no direct algorithms available, iterative algorithms operating on the low-rank tensor representation or inexact versions of well-known iterative algorithms formulated on the tensor space $\mathbb{R}^{n_1 \times \dots \times n_d}$ with truncation in between should be chosen.

Toolboxes and Conversion between Formats

The `TT-Toolbox` [89] provides efficient implementations of TT tensors and important operations in MATLAB and PYTHON. These are, e. g., elementwise addition and multiplication, extraction of elements of the tensor, implementations of linear operators, and truncation to a lower rank by TT-SVD. Truncation is crucial in practice to avoid infeasible rank growth because when using the standard implementations for summing or multiplying two TT tensors elementwise, the ranks sum or multiply, respectively. Clearly, truncation comes with the drawback of additional errors introduced in the computation. Therefore, it is an important aspect of this thesis how those errors have to be controlled to obtain convergence of the final algorithms, which use tensor computations in substeps.

Due to the similar structure, comparable algorithms are available for HT tensors. The `htucker` toolbox [76] implements tensors in this format and—similarly to the `TT-Toolbox`—efficient versions of the most important operations, e. g., extraction of parts of the tensor, application of linear operators to the i -th mode, contraction, tensor orthogonalization and truncation to a lower rank by HSVD [50]. For elementwise addition and multiplication, exact and special, truncated versions are available such that the explosion of ranks can be avoided. HT tensors can be orthogonalized such that the matrices \mathbf{U}_t at all nodes t except for the root node form an orthogonal basis of the respective tensor subspace. This procedure can simplify certain computations and make some tensor algorithms more stable. All mentioned operations can be performed in a reasonable amount of time, which is linear in n_{\max} and d and polynomial in r_{\max} , in particular at most $\mathcal{O}(n_{\max} d r_{\max}^2 + d r_{\max}^4)$, but often less. The same holds for tensors in TT format.

For the algorithms developed in this thesis, operations and subsolvers that currently are not all available within a single format or toolbox are needed. Thus, we use HT tensors with a linear, TT-like dimension tree as provided by `htucker`, see Figure 2.5.

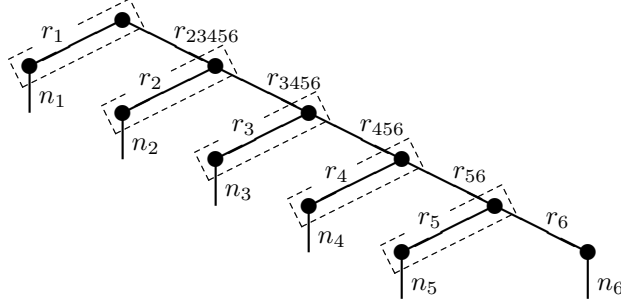


Figure 2.5.: Tensor network representing the Hierarchical Tucker format in the case $d = 6$ with a linear tree, which makes the conversion to TT tensors simple. By performing the marked contractions, a TT network is obtained.

This allows to convert between the HT and TT formats, cf. [54, Sec.12.2.2], so that both the `TT-Toolbox` and the `htucker` toolbox can be used simultaneously. By performing the marked contractions in Figure 2.5, a TT network is obtained. Conversely, the TT cores $\mathbf{u}_1, \dots, \mathbf{u}_{d-1}$ can be split into the basis matrices $\mathbf{U}_1, \dots, \mathbf{U}_{d-1}$ and the respective transfer tensors in order to convert from the TT to the HT format: In fact, we aim for finding $r_i \in [n_i]$, $\mathbf{U}_i \in \mathbb{R}^{n_i \times r_i}$ and $\mathbf{b}_{i,\dots,d} \in \mathbb{R}^{r_i \times r_{i+1}, \dots, d \times r_{i,\dots,d}}$ such that $\langle \mathbf{U}_i, \mathbf{b}_{i,\dots,d} \rangle_{2,1} = \tilde{\mathbf{u}}_i \in \mathbb{R}^{n_i \times r_{i+1}, \dots, d \times r_{i,\dots,d}}$ holds, where $\tilde{\mathbf{u}}_i$ is a reshaped version of \mathbf{u}_i . Writing this in a matricized form, we have

$$\mathbf{U}_i \mathbf{b}_{i,\dots,d}^{(r_i)} = \tilde{\mathbf{u}}_i^{(n_i)} \in \mathbb{R}^{n_i \times r_{i+1}, \dots, d r_{i,\dots,d}}$$

with $\mathbf{b}_{i,\dots,d}^{(r_i)} \in \mathbb{R}^{r_i \times r_{i+1}, \dots, d r_{i,\dots,d}}$. Clearly, one could choose $r_i = n_i$, $\mathbf{U}_i = \mathbf{I}_{n_i}$ (identity matrix), and $\mathbf{b}_{i,\dots,d} = \tilde{\mathbf{u}}_i$ to obtain the result, but this approach is not suitable due to the possibly (for large n_i) large ranks $r_i = n_i$ so that the large matrices $\mathbf{U}_i \in \mathbb{R}^{n_i \times n_i}$ have to be stored. Low-rank tensor implementations use dense linear algebra because the representing small tensors and matrices cannot be expected to be sparse so that one cannot benefit from the sparsity of the matrix \mathbf{I}_{n_i} here. A smaller rank r_i can be obtained by using an SVD of $\tilde{\mathbf{u}}_i^{(n_i)} = \mathbf{V}_i \Sigma_i \mathbf{W}_i^\top$, where $\mathbf{V}_i \in \mathbb{R}^{n_i \times R_i}$ and $\mathbf{W}_i \in \mathbb{R}^{r_{i+1}, \dots, d r_{i,\dots,d} \times R_i}$ are orthogonal matrices, and $\Sigma = \text{diag}((\sigma_1^i, \dots, \sigma_{R_i}^i)) \in \mathbb{R}^{R_i \times R_i}$, where $R_i = \min\{n_i, r_{i+1}, \dots, d r_{i,\dots,d}\}$, i. e., we use the economic variant of SVD or *thin SVD* [49, Sec.2.5.4]. The singular values shall be ordered: $\sigma_1^i \geq \sigma_2^i \geq \dots \geq \sigma_{R_i}^i$. Truncating this decomposition to rank $r_i \leq R_i$, we set $\mathbf{U}_i := \mathbf{V}_i(\cdot, [r_i])$ and $\mathbf{b}_{i,\dots,d} := \text{ten}(\Sigma_i([r_i], [r_i]) \mathbf{W}_i(\cdot, [r_i])^\top)$ to obtain an approximate decomposition of the TT core tensor \mathbf{u}_i . The accuracy of this approximation can be controlled easily since the squared Frobenius error is given by the sum of the truncated, squared singular values.

Additional Arithmetic Operations

When designing optimization algorithms for constrained problems in tensor space, one often needs more componentwise operations than only addition, subtraction, and multiplication. Standard algorithms for constrained problems typically require the computation of special componentwise functions such as penalty terms, projections, indicators, e. g., for a semi-smooth Newton method, or the reciprocal \mathbf{x}^{-1} , which is needed to compute the derivative of a log-barrier term in an interior point method. One option is to approximate such quantities

by applying Newton's method to a suitable equation. In the case of the elementwise reciprocal $\mathbf{y} = \mathbf{x}^{-1}$, one can consider the equation $\mathbf{x} - \mathbf{y}^{-1} = 0$ and solve it for \mathbf{y} , see also [43, Sec. 4.4]. The resulting Newton equation given the current iterate \mathbf{y}_k is $\mathbf{y}_k^{-2}(\mathbf{y}_{k+1} - \mathbf{y}_k) = -\mathbf{x} + \mathbf{y}_k^{-1}$, giving the Newton update $\mathbf{y}_{k+1} = \mathbf{y}_k \odot (2 \cdot \mathbb{1} - \mathbf{x} \odot \mathbf{y}_k)$. This iteration can be performed with available operations within the low-rank format; especially no componentwise division is needed. It can be shown to converge if one chooses the first iterate \mathbf{y}_0 such that $\text{sgn}(\mathbf{y}_0) = \text{sgn}(\mathbf{x})$ and $|\mathbf{y}_0| \leq |\mathbf{x}|^{-1}$ (componentwise). One can choose, e. g., $\frac{1}{\|\mathbf{x}\|_{\mathbb{F}}^2} \mathbf{x}$ for general \mathbf{x} or $\frac{1}{\|\mathbf{x}\|_{\mathbb{F}}} \mathbb{1}$ for $\mathbf{x} > 0$. This iteration was already proposed to approximate the inverse of structured matrices [55, Sec. 4.1] and is called Newton-Schulz method [101]. When implementing it with low-rank tensors—as done in the example section of `htucker`—the ranks of the iterates have to be bounded by truncation. This causes the implementation to be an *inexact* version of Newton's method, which could only be expected to converge if the error caused by the truncation was controlled suitably, see [55, Sec. 2]. In practice, this approach is sometimes unfeasible because the required ranks are too large. If truncation is performed always to a fixed rank, one typically obtains only a solution of a certain accuracy, which can be quantified by $\|\mathbb{1} - \mathbf{x} \odot \mathbf{y}\|_{\mathbb{F}}$. In order to guarantee convergence of this truncated iteration for $\mathbf{x} > 0$ one would have to ensure that all iterates fulfill $\mathbf{y}_k > 0$ elementwise. By error control in the maximum norm $\|\cdot\|_{\infty}$ this goal could be achieved. Such a rounding procedure is currently not available to our knowledge. Using the equivalence of norms in finite-dimensional spaces is not suitable here because of the magnitude of the equivalence constant: We have $\|\mathbb{1}\|_{\mathbb{F}} = \sqrt{n_1 \cdots n_d}$ for $\mathbf{x} \in \mathbb{R}^{n_1 \times \cdots \times n_d}$, whereas $\|\mathbb{1}\|_{\infty} = 1$, but can only ensure $\|\mathbf{x}\|_{\infty} \leq \|\mathbf{x}\|_{\mathbb{F}}$ in general, i. e., it could be necessary to decrease the Frobenius norm of the error further although the maximum norm is small enough already.

A similar iteration, namely $\mathbf{y}_{k+1} = \frac{1}{2} \mathbf{y}_k \odot (3 \cdot \mathbb{1} - \mathbf{y}_k^2 \odot \mathbf{x})$, can be employed to compute $\mathbf{x}^{-1/2}$ for $\mathbf{x} > 0$. This gives then access to $|\mathbf{x}| \approx \mathbf{x}^2 \odot (\mathbf{x}^2 + \varepsilon \cdot \mathbb{1})^{-1/2}$, where the addition of $\varepsilon \cdot \mathbb{1}$ with a small $\varepsilon > 0$ makes the algorithm more stable. Using the componentwise absolute value, quantities such as quadratic penalty terms or projections onto boxes can be approximated with low-rank tensors and used in an optimization algorithm.

In our experiments [46] it was found that the Newton-Schulz iteration performs often well in practice, but does not give satisfying results sometimes if the entries of the positive tensor $\mathbf{x} > 0$ are too close to 0. Then we computed $\mathbf{y} = \mathbf{x}^{-1}$ by solving the linear system $\mathbf{x} \odot \mathbf{y} = \mathbb{1}$ by an iterative method working on the low-rank tensor representation, such as ALS, AMEN or optimization on manifolds, see below.

The iterative methods for computing componentwise functions of low-rank tensors in a low-rank format described above are all designed for a specific problem. This has the advantage that they can be analyzed very well, but comes with the drawback that a separate algorithm has to be designed for each function. A more general approach for computing componentwise functions is given by sampling some entries of the resulting \mathbf{y} , i. e., computing $\mathbf{y}(k) = f(\mathbf{x}(k))$ for some indices k in some index set $K \subset \times_{i=1}^d [n_i]$. Then the full tensor in the low-rank format is computed by some tensor completion method. Such methods aim to find a tensor of a prescribed or adaptively selected low rank which contains the sampled entries, sometimes at least approximately. Suitable algorithms are cross approximation [91, 15, 14] or tensor completion using optimization algorithms such as ALS [62] or optimization on the manifold of tensors of fixed TT or HT rank [73, 104, 32].

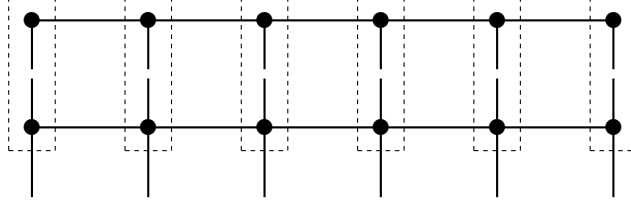


Figure 2.6.: Product of a TT matrix (below) and a TT tensor/vector (above) represented as a contraction. If the modes along which one contracts in the resulting network are “vectorized”, a TT network is obtained. The TT ranks of the result are thus obtained by multiplying the TT ranks of the TT vector and the TT matrix. See also [12].

Solution of Linear Equations

In many applications, the solution of linear systems involving low-rank tensors is a key tool. They arise, e. g., when discretizing PDEs in high space or stochastic dimension. It is then important that the corresponding linear operators can at least be applied efficiently to low-rank tensors. For instance, they can be of the following forms:

- Linear operators in the CP format: Given $d \in \mathbb{N}$ and $R \in \mathbb{N}$, let $\mathbf{A}_i^{(j)} : \mathbb{R}^{n_i} \rightarrow \mathbb{R}^{m_i}$ be linear operators or matrices for $i \in [d]$, $j \in [R]$. Consider the operator

$$\mathbf{A} : \mathbb{R}^{n_1 \times \dots \times n_d} \rightarrow \mathbb{R}^{m_1 \times \dots \times m_d}, \quad \mathbf{A} = \sum_{j=1}^R \bigotimes_{i=1}^m \mathbf{A}_i^{(j)}, \quad \mathbf{A}(\mathbf{x}_1 \otimes \dots \otimes \mathbf{x}_m) = \sum_{j=1}^R \bigotimes_{i=1}^m \mathbf{A}_i^{(j)} \mathbf{x}_i,$$

defined by its action on elementary tensors. Operators of this type can be efficiently applied to low-rank tensors of any type using i -mode matrix products and summation:

$$\mathbf{A}\mathbf{x} = \sum_{j=1}^R \mathbf{A}_m^{(j)} \circ_m \dots \mathbf{A}_1^{(j)} \circ_1 \mathbf{x}.$$

- Linear operators in special low-rank formats, such as HTD for linear operators [76] or `tt_matrix` [89]. Operators of this type can easily be applied to HT or TT tensors giving HT or TT tensors as a result, respectively, essentially by performing a suitable contraction. This is visualized in Figure 2.6 for the case of a product of a TT matrix and a TT tensor (vector).
- Componentwise multiplication by a tensor $\mathbf{y} \in \mathbb{R}^{n_1 \times \dots \times n_d}$: $\mathbf{x} \mapsto \mathbf{y} \odot \mathbf{x}$ is linear and can be computed by the respective algorithms.

In general, such linear operators between tensor spaces cannot be inverted easily except for rank-1-operators: If the matrices $\mathbf{A}_i \in \mathbb{R}^{n_i \times n_i}$ ($i \in [d]$) are invertible, it holds that

$$\left(\bigotimes_{i=1}^m \mathbf{A}_i \right)^{-1} = \bigotimes_{i=1}^m \mathbf{A}_i^{-1}$$

so that the inverse of a rank-1-operator can be applied to low-rank tensors via i -mode matrix products. This fact can be used to construct good preconditioners for operators which are perturbations of rank-1-operators. For general linear systems, however, iterative methods have to be applied to compute the solution $\mathbf{x} \in \mathbb{R}^{n_1 \times \dots \times n_d}$ of $\mathbf{A}\mathbf{x} = \mathbf{b}$ with $\mathbf{b} \in \mathbb{R}^{n_1 \times \dots \times n_d}$ and $\mathbf{A} \in \mathcal{L}(\mathbb{R}^{n_1 \times \dots \times n_d}, \mathbb{R}^{n_1 \times \dots \times n_d})$, which we assume to be a symmetric operator for simplicity. There are several options to do this:

- Similarly to the Newton-Schulz method, one can apply standard iterative methods such as the preconditioned conjugate gradient (PCG) method [75] formulated on the space $\mathbb{R}^{n_1 \times \dots \times n_d}$ but implemented with low-rank tensors only and round the iterates in between to avoid infeasible rank growth. See also [11] for a review of such methods, and [9, 10, 13] for concrete examples. To obtain provable convergence of such an inexact scheme, the errors caused by truncation typically have to be controlled suitably so that a guaranteed rank bound for the iterates of the method cannot be provided even though the solution \mathbf{x} of the system may have low rank. If a suitable soft thresholding procedure is used for rounding to a fixed tensor rank, convergence results of fixed point methods can be established [11].
- Alternative approaches guess the rank r of the solution \mathbf{x} and formulate the equation $\mathbf{A}\mathbf{x} = \mathbf{b}$ as a least squares problem, which is then solved on the set of tensors of rank at most r by optimization over the representing tensors. In the case of the TT format, $\mathbf{x} = \mathcal{TT}(\mathbf{u}_1, \dots, \mathbf{u}_d)$ holds and one would solve

$$\min_{\mathbf{u}} \frac{1}{2} \|\mathbf{A} \mathcal{TT}(\mathbf{u}) - \mathbf{b}\|_{\mathbb{F}}^2.$$

Since the map $\mathbf{u} \mapsto \mathcal{TT}(\mathbf{u})$ is only multilinear and therefore *nonlinear* in \mathbf{u} , this is a general nonlinear and especially non-quadratic problem. It can be solved locally by block-wise optimization over the TT cores for example. If all cores except \mathbf{u}_i are fixed, a quadratic function in \mathbf{u}_i is obtained, a stationary point of which can be found by solving a linear system. This is known as alternating linear scheme (ALS) [62, 98]. The idea of block-wise optimization can also be applied to more general, nonlinear optimization problems in tensor space.

- When optimizing over the TT cores, the problem appears that the TT representation of a rank- r -tensor is not unique. In fact, certain rescalings by orthogonal matrices do not change the represented tensor [98]. Along those rescalings, the objective function will be constant and its curvature will be zero so that typical convergence results requiring second-order sufficient conditions cannot be applied. One way to overcome this issue is to factor out the ambiguity of the representation and work on the manifold of tensors of fixed rank r . Riemannian optimization algorithms on manifolds are well understood [1] and were also developed especially for low-rank tensor manifolds [100, 73, 74, 104, 32, 113].
- Estimating the rank of a solution of a linear system is a very hard task in general. Therefore, ALS or methods operating on a fixed manifold may not give accurate enough results. Modifications of ALS such as MALS [62], DMRG [87, 90], or AMEn [36,

35] allow the adaption of ranks by essentially optimizing over the contraction of two “neighboring” TT cores at once and then splitting this optimized tensor into two cores again as described above when discussing the conversion from TT to HT tensors. At this point, a rank can be chosen such that, e. g., a prescribed accuracy is met.

Initially, we used the HT PCG method in our implementations with an adequate rank-1-preconditioner, which is a viable approach. In numerical experiments we found that AMEn yielded more accurate results with smaller ranks and better computing times [46]. One has to mention that—to our knowledge—the convergence theory of this method is still limited, but the method performs well in practice. In principle, it supports operators \mathbf{A} in the `tt_matrix` format and was extended to operators in the so-called $\{d, R\}$ -format in [35], which is essentially the CP format with sparse matrices. We extended it further to also have componentwise multiplication operators³ and a rank-1-preconditioner in decomposed form, i. e., a rank-1-operator $\mathbf{T} \approx \mathbf{A}^{-1}$, which has to be given in the form $\mathbf{T} = \bigotimes_{i=1}^d \mathbf{P}_i \mathbf{P}_i^*$, where $\mathbf{P}_i \in \mathcal{L}(\mathbb{R}^{n_i}, \mathbb{R}^{n_i})$. There, we also allowed for function handles R_i so that also, e. g., inverses of triangular matrices or certain linear transformations, such as the fast Fourier transform, could be implemented efficiently.

Determination of Maximum/Minimum Entries

For determining stepsizes or feasibility in constrained optimization, it can be necessary to compute or estimate the maximum or minimum entry of a given low-rank tensor. In [43] it is proposed to use the fact that the linear mapping $\mathbf{y} \mapsto \mathbf{x} \odot \mathbf{y}$ has exactly the components of \mathbf{x} as eigenvalues. A normed basis of eigenvectors is given by the tensors $\bigotimes_{i=1}^d \mathbf{e}_{k_i}^i$ ($k \in \times_{i=1}^m [n_i]$), where $\mathbf{e}_{k_i}^i \in \mathbb{R}^{n_i}$ denotes the k_i -th unit vector. The element of \mathbf{x} with maximum absolute value can thus be found by performing a power iteration with initial tensor $\mathbf{y}_0 := \frac{1}{\sqrt{n_1 \cdots n_d}} \mathbb{1}$. The iterates are given by $\mathbf{y}_j = \frac{1}{\|\mathbf{x}^j\|_F} \mathbf{x}^j$. Hence, the iteration can be sped up by squaring each iterate so that only $\frac{1}{\|\mathbf{x}^{2^j}\|_F} \mathbf{x}^{2^j}$ is computed for each $j \in \mathbb{N}$. This is implemented in the example section of `htucker`. The maximum entry of \mathbf{x} is approximated by the sequence $\lambda_j := \frac{\langle \mathbf{x} \odot \mathbf{x}^{2^j}, \mathbf{x}^{2^j} \rangle}{\langle \mathbf{x}^{2^j}, \mathbf{x}^{2^j} \rangle}$ (Rayleigh quotient). It can be shown that $|\lambda_j| \leq \max_k |\mathbf{x}(k)|$ holds for all j . Since the ranks square in each iteration, truncation is necessary. Therefore, one can only expect to be able to compute a lower bound on the absolute value of the maximum entry with this procedure. Therefore, we also experimented with an upper bound, which can be obtained from this iteration: It holds $\|\mathbf{x}\|_\infty \leq \|\mathbf{x}\|_{2^j} = \langle \mathbf{x}^{2^{j-1}}, \mathbf{x}^{2^{j-1}} \rangle^{2^{-j}}$ for all $j \in \mathbb{N}$ and $\lim_{j \rightarrow \infty} \|\mathbf{x}\|_{2^j} = \|\mathbf{x}\|_\infty$. The approximation $\|\mathbf{x}\|_{2^j}$ can be computed from the repeatedly squared tensor if one handles arithmetic overflow suitably.

Numerical experiments [46] showed us that rounding can be a severe issue in this iteration so that we did not obtain accurate enough results. For tensors describing a function on a product set, we therefore chose the multilevel coordinate search (MCS) method [65] a non-rigorous global optimization routine for box-constrained problems where only function values are available, to optimize the represented function globally to approximate the maximum element of the tensor. This approach yielded better results in practice but is limited to function-related tensors.

³The implementation of this task was essentially done by Prof. Dr. Michael Ulbrich.

2.2. Tensor Products of Hilbert and Banach Spaces

In this subsection we introduce the basic tools and results for tensor products of Hilbert spaces, following [54] and [116, Sec. 3.4]. We will need this concept to formulate parametric PDE problems. Note that we only collect the most important facts about tensor products of Hilbert spaces and point out only a few generalizations to Banach spaces as far as it is necessary for the functional analytic setting of the semilinear, elliptic PDE with uncertain inputs discussed in Chapter 3.

Definition 2.5 (Algebraic and topological tensor product of vector spaces⁴). Let V and W be vector spaces over a field $\mathbb{K} \in \{\mathbb{R}, \mathbb{C}\}$. The *algebraic tensor product space* $V \otimes_a W$ of V and W is defined as the quotient vector space

$$V \otimes_a W := \text{span}\{(v, w) : v \in V, w \in W\} / \mathcal{N}$$

with

$$\mathcal{N} := \text{span}\left\{ \sum_{i=1}^m \sum_{j=1}^n \alpha_i \beta_j (v_i, w_j) - \left(\sum_{i=1}^m \alpha_i v_i, \sum_{j=1}^n \beta_j w_j \right) : \right. \\ \left. m, n \in \mathbb{N}, \alpha_i, \beta_j \in \mathbb{K}, v_i \in V, w_j \in W \right\},$$

where the span of the set is taken in the algebraic sense, meaning that it consists of *finite* formal linear combinations of pairs (v_i, w_i) . This first set of formal linear combinations is normally referred to as the *free vector space* over the set $V \times W$. The equivalence class of a pair (v, w) is denoted by $v \otimes w$ and called *elementary tensor*.

If a norm $\|\cdot\|$ is given on $V \otimes_a W$, the *topological tensor space* is defined as the closure of the algebraic tensor space with respect to this norm:

$$V \otimes_{\|\cdot\|} W := \overline{V \otimes_a W}^{\|\cdot\|}.$$

Proposition 2.6 (Characterization of equivalence classes⁵). *Let V and W be vector spaces over a field \mathbb{K} and let $v, \tilde{v} \in V$ and $w, \tilde{w} \in W$ be given. Consider the algebraic tensor space $V \otimes_a W$. The pair (\tilde{v}, \tilde{w}) belongs to the equivalence class $v \otimes w$ if and only if the following holds:*

$$(v = 0 \vee w = 0) \wedge (\tilde{v} = 0 \vee \tilde{w} = 0) \tag{2.1}$$

or

$$(v \neq 0 \wedge w \neq 0) \vee (\tilde{v} \neq 0 \wedge \tilde{w} \neq 0) \quad \text{and} \quad \exists c \in \mathbb{K} \setminus \{0\} \text{ s. t. } \tilde{v} = cv \text{ and } \tilde{w} = c\tilde{w}.$$

Case (2.1) holds exactly for pairs belonging to the equivalence class $0 \otimes 0$.

Proof. The proof of this proposition is given in the appendix (Proposition A.1). □

⁴This definition follows [116, Sec. 3.4] and [54, Sec. 3.2.2].

⁵This proposition follows [116, Exercise 3.11].

Remark 2.7. More basic facts and notations about tensor spaces are given in [54, Sec. 3.2]. We collect only some of them:

- Tensor spaces are again vector spaces [54, Def. 3.9].
- It holds that $V \otimes_a W = \text{span}\{v \otimes w : v \in V, w \in W\}$ [54, Eq. (3.11)].
- The tensor product of two finite-dimensional vector spaces V and W is again finite-dimensional. We have $V \otimes_a W = V \otimes_{\|\cdot\|} W =: V \otimes W$ [54, Notation 3.8(a)].
- The map $\otimes : V \times W \rightarrow V \otimes W$ is bilinear [54, Lem. 3.10].
- If B is a basis of V and C is a basis of W , then $\{b \otimes c : b \in B, c \in C\}$ is a basis of $V \otimes_a W$. We obtain $\dim(V \otimes_a W) = \dim(V) \dim(W)$ [54, Lem. 3.11].
- If V, W and X are vector spaces over \mathbb{K} and a bilinear map $\otimes : V \times W \rightarrow X$ fulfills that $\text{span}\{v \otimes w : v \in V, w \in W\} = X$ and that sets $B \subset V$ and $C \subset W$ of linearly independent vectors are mapped to a set $\{b \otimes c : b \in B, c \in C\} \subset X$ of linearly independent vectors, then X is isomorphic to $V \otimes_a W$ [54, Prop. 3.12].

Example 2.8. Let $\mathbb{K} \in \{\mathbb{R}, \mathbb{C}\}$ and let S and T be some suitable sets, e. g., $[n]$ for some $n \in \mathbb{N}$ to represent the vector space \mathbb{K}^n , \mathbb{N} to represent vector spaces of sequences or some set $S \subset \mathbb{R}^n$ to represent more general function spaces. Let $V \subset \{v : S \rightarrow \mathbb{K}\}$ and $W \subset \{w : T \rightarrow \mathbb{K}\}$ be vector spaces of \mathbb{K} -valued functions on the given sets. Now we can consider the algebraic tensor product space $V \otimes_a W$ and make the identification of elementary tensors $v \otimes w \cong vw$, i. e., we say that the elementary tensor $v \otimes w$ is a pointwise product of functions. We get an isomorphism between the sets $X_1 := V \otimes_a W$ and $X_2 := \text{span}\{vw, v \in V, w \in W\} \subset \{x : S \times T \rightarrow \mathbb{K}\}$ with the usual equality that $x = y$ for $x, y \in X_2$ if and only if $x(s, t) = y(s, t)$ for all $s \in S$ and all $t \in T$. Therefore, we will frequently use the outer product $\otimes : V \times W \rightarrow V \otimes_a W, (v \otimes w)(s, t) := v(s)w(t)$ for function-related tensors, see also [34, Chap. I, Sec. 2.4].

Analogously, if U is a Banach space over \mathbb{K} and we consider $U \otimes_a V$, we can identify $u \otimes v \cong uv(\cdot)$, i. e.,

$$U \otimes_a V \cong \{y : S \rightarrow U \text{ s. t. } \exists n \in \mathbb{N}, (u_k, v_k) \in U \times V \text{ with } y(s) = \sum_{k=1}^n v_k(s)u_k \forall s \in S\},$$

see [34, Chap. I, Sec. 2.4]). This means that we can identify elements of the tensor product space $U \otimes_a V$ by U -valued functions on the set S .

We briefly discuss how linear maps from one tensor space to another can look like. It is important to observe that linear maps are uniquely determined by their action on a basis, in our case on certain elementary tensors. When having linear maps $\varphi : V \rightarrow W$ and $\psi : \tilde{V} \rightarrow \tilde{W}$ between \mathbb{K} -vector spaces V, \tilde{V}, W , and \tilde{W} , the *tensor product of these maps* is defined by

$$(\varphi \otimes \psi)(v \otimes \tilde{v}) := \varphi(v) \otimes \psi(\tilde{v}). \quad (2.2)$$

This tensor product is an elementary tensor in the space $\mathcal{L}(V, W) \otimes_a \mathcal{L}(\tilde{V}, \tilde{W})$. It holds that $\mathcal{L}(V, W) \otimes_a \mathcal{L}(\tilde{V}, \tilde{W}) \subset \mathcal{L}(V \otimes_a \tilde{V}, W \otimes_a \tilde{W})$ with equality if the spaces V and \tilde{V} are finite dimensional [54, Prop. 3.49].

Let now V and W be Hilbert spaces over $\mathbb{K} \in \{\mathbb{R}, \mathbb{C}\}$ with the inner products $(\cdot, \cdot)_V$ and $(\cdot, \cdot)_W$, respectively. Elements of their algebraic tensor product $X_a = V \otimes_a W$ can be written as finite linear combination $x = \sum_{i=1}^m \alpha_i v_i \otimes w_i$ of elementary tensors. Therefore, an inner product on the algebraic tensor product space can be defined in a natural way [116, Sec. 3.4]: For $x = \sum_{i=1}^m \alpha_i v_i \otimes w_i \in X_a$ and $y = \sum_{j=1}^n \beta_j \tilde{v}_j \otimes \tilde{w}_j \in X_a$ the inner product is given by

$$(x, y)_X := \sum_{i=1}^m \sum_{j=1}^n \overline{\alpha_i} \beta_j (v_i, \tilde{v}_j)_V (w_i, \tilde{w}_j)_W, \quad (2.3)$$

where $\overline{\alpha_i}$ denotes the complex conjugate of α_i . It can be shown that this is well-defined meaning that the value of the inner product does not depend on the representation of x and y and that this defines a scalar product, the so-called *induced scalar product* [54, Lem. 4.124]. The completion X of X_a with respect to the norm induced by this scalar product is called the *complete tensor product* of the Hilbert spaces V and W and is itself a Hilbert space. We denote it by $V \otimes W$.

We collect some facts about the complete tensor product $V \otimes W$ of two Hilbert spaces V and W and the induced inner product $(\cdot, \cdot)_{V \otimes W}$:

- If B is an orthonormal basis of V and C is an orthonormal basis of W , the set $\{b \otimes c : b \in B, c \in C\}$ is an orthonormal basis of $V \otimes W$ [116, Thm. 3.12(b)].
- The norm $\|x\|_{V \otimes W} := \sqrt{(x, x)_{V \otimes W}}$ is a *reasonable crossnorm* [54, Def. 4.66, Lem. 4.67], i. e., $\|v \otimes w\|_{V \otimes W} = \|v\|_V \|w\|_W$ holds for all $v \in V, w \in W$ and $\|\varphi \otimes \psi\|_{(V \otimes W)^*} = \|\varphi\|_{V^*} \|\psi\|_{W^*}$ holds for all $\varphi \in V^*, \psi \in W^*$ [54, Prop. 4.127].
- As a Hilbert space, the space $V \otimes W$ is reflexive. Its dual space is $(V \otimes W)^* = V^* \otimes W^*$ [54, Lem. 4.75].
- The norm $\|\cdot\|_{V \otimes W}$ is a *uniform crossnorm* [54, Def. 4.77], i. e., for any operators $A \in \mathcal{L}(V, V)$ and $B \in \mathcal{L}(W, W)$, the operator $A \otimes B$ belongs to $\mathcal{L}(V \otimes W, V \otimes W)$ and has the operator norm $\|A \otimes B\|_{\mathcal{L}(V \otimes W, V \otimes W)} = \|A\|_{\mathcal{L}(V, V)} \|B\|_{\mathcal{L}(W, W)}$ [54, Prop. 4.127]. The same holds for the norm $\|\cdot\|_{(V \otimes W)^*}$, which is also a uniform and reasonable crossnorm [54, Prop. 4.80].

Note that we can extend all statements in this section to tensor products $\bigotimes_{i=1}^m V_i$ of $m \in \mathbb{N}$ vector or Hilbert spaces V_i , since the tensor product is associative. $\otimes : \times_{i=1}^m V_i \rightarrow \bigotimes_{i=1}^m V_i$ is then a multilinear map and the properties hold in an analogous sense, see also [54, Sec. 3.2.4].

Example 2.9 (Tensor product of L^2 -spaces). Let $m \in \mathbb{N}$ open sets $\Xi_i \subset \mathbb{R}$ ($i \in [m]$) be given and equipped with the Borel σ -algebras and given probability measures \mathbb{P}_i , respectively. We consider the spaces $L^2_{\mathbb{P}_i}(\Xi_i)$ of \mathbb{R} -valued random variables with finite variance and their tensor product $\bigotimes_{i=1}^m L^2_{\mathbb{P}_i}(\Xi_i)$. Equipped with the scalar product $(v, \tilde{v})_{L^2_{\mathbb{P}_i}(\Xi_i)} := \int_{\Xi_i} v \tilde{v} d\mathbb{P}_i$ these spaces are Hilbert spaces. According to Example 2.8 we identify elementary tensors

by the product of univariate functions: $(v^1 \otimes v^2 \otimes \dots \otimes v^m)(\xi) := v^1(\xi_1)v^2(\xi_2) \cdots v^m(\xi_m)$ for every $\xi \in \Xi := \times_{i=1}^m \Xi_i \subset \mathbb{R}^m$. We compute the inner product of two elementary tensors:

$$\left(\bigotimes_{i=1}^m v^i, \bigotimes_{i=1}^m \tilde{v}^i \right)_{\bigotimes_{i=1}^m L_{\mathbb{P}_i}^2(\Xi_i)} = \prod_{i=1}^m (v^i, \tilde{v}^i)_{L_{\mathbb{P}_i}^2(\Xi_i)} = \prod_{i=1}^m \left(\int_{\Xi_i} v^i \tilde{v}^i d\mathbb{P}_i \right) = \int_{\Xi} \bigotimes_{i=1}^m v^i \bigotimes_{i=1}^m \tilde{v}^i d\mathbb{P}.$$

This holds due to Fubini's theorem and $\mathbb{P} := \otimes_{i=1}^m \mathbb{P}_i$ is the product measure of the measures \mathbb{P}_i ($i \in [m]$). Due to linearity we see that this scalar product is exactly the $L_{\mathbb{P}}^2(\Xi)$ -inner product $(\mathbf{v}, \tilde{\mathbf{v}})_{L_{\mathbb{P}}^2(\Xi)} = \int_{\Xi} \mathbf{v} \tilde{\mathbf{v}} d\mathbb{P}$. Since any algebraic tensor has a finite norm in $L_{\mathbb{P}}^2(\Xi)$, we get $\bigotimes_{i=1}^m L_{\mathbb{P}_i}^2(\Xi_i) \subset L_{\mathbb{P}}^2(\Xi)$ with the used identification of the tensor product using that the closure of a subset of a set is contained in its closure. We want to show that the sets are in fact equal. Note that the simple functions on Ξ that are obtained as finite linear combinations of tensor products of simple functions in one variable are dense in $L_{\mathbb{P}}^2(\Xi)$. Since they are also square integrable and therefore contained in the tensor product space ${}_a \bigotimes_{i=1}^m L_{\mathbb{P}_i}^2(\Xi_i)$, this algebraic tensor product is actually dense in $L_{\mathbb{P}}^2(\Xi)$ and its closure therefore equal to $L_{\mathbb{P}}^2(\Xi)$, cf. [34, Chap. I, Sec. 7].

This example raises the question if it is generalizable to L^p -spaces with $p \in [1, \infty]$. If it comes to tensor products of general Banach spaces instead of Hilbert spaces, the construction (2.3) is no longer valid since no inner product exists and a similar construction with norms is not possible. Instead, one has to define a suitable norm on the algebraic tensor product and complete it w. r. t. this norm. In general, one can always define the so-called *projective norm* [54, Sec. 4.2.4], which is the strongest norm⁶ on $V \otimes_a W$ ensuring continuity of the tensor product mapping $V \times W \ni (v, w) \mapsto v \otimes w \in V \otimes_a W$, and the *injective norm* [54, Sec. 4.7.2], which is the weakest norm on $V \otimes_a W$ yielding a dual norm on $(V \otimes_a W)^*$ rendering the tensor product mapping $V^* \times W^* \ni (\varphi, \psi) \mapsto \varphi \otimes \psi \in V^* \otimes_a W^* \subset (V \otimes_a W)^*$ continuous. However, in the case of L^p spaces, it is more suitable to work with the L^p norm on the product set:

Example 2.10 (Tensor product of L^p -spaces).⁷ Let $p \in [1, \infty)$ and consider the spaces $L_{\mathbb{P}_i}^p(\Xi_i)$ of \mathbb{R} -valued, p -integrable random variables and the algebraic tensor product space ${}_a \bigotimes_{i=1}^m L_{\mathbb{P}_i}^p(\Xi_i)$ in analogy to Example 2.9. Identifying functions in this space with functions on Ξ , we equip ${}_a \bigotimes_{i=1}^m L_{\mathbb{P}_i}^p(\Xi_i)$ with the $L_{\mathbb{P}}^p(\Xi)$ norm. Indeed, this norm is finite for any algebraic tensor since it is finite for any elementary tensor:

$$\left\| \bigotimes_{i=1}^m v^i \right\|_{L_{\mathbb{P}}^p(\Xi)}^p = \int_{\Xi} \left| \prod_{i=1}^m v^i \right|^p d\mathbb{P} = \prod_{i=1}^m \int_{\Xi_i} |v^i|^p d\mathbb{P}_i = \prod_{i=1}^m \|v^i\|_{L_{\mathbb{P}_i}^p(\Xi_i)}^p < \infty.$$

Therefore, we can complete the algebraic tensor space w. r. t. to the $L_{\mathbb{P}}^p(\Xi)$ norm to obtain $\bigotimes_{i=1}^m L_{\mathbb{P}_i}^p(\Xi_i) \subset L_{\mathbb{P}}^p(\Xi)$. By the simple functions argument from Example 2.9 we even obtain equality of the sets. Note that this does not work any longer for $p = \infty$ so that in general $\bigotimes_{i=1}^m L_{\mathbb{P}_i}^{\infty}(\Xi_i) \subsetneq L_{\mathbb{P}}^{\infty}(\Xi)$ holds, see [34, Chap. I, Exercise 7.2].

⁶See [54, Prop. 4.46] for the precise meaning of this statement.

⁷cf. [54, Example 4.40] and [34, Chap. I, Sec. 7]

Switching back to Hilbert spaces, we introduce a space, which is typically used for the functional analytic setting of linear, elliptic PDEs with uncertain coefficients.

Example 2.11 (An anisotropic Sobolev space). Let for some $n \in \mathbb{N}$ an open domain $\Omega \subset \mathbb{R}^n$ be given and consider the Sobolev space $H^1(\Omega)$ of weakly differentiable $L^2(\Omega)$ functions with $L^2(\Omega)$ derivatives. Moreover, let $(\Xi, \mathcal{F}, \mathbb{P})$ be a probability space, cf. Example 2.9. Now we consider the tensor product $H^1(\Omega) \otimes L^2_{\mathbb{P}}(\Xi)$. Again, we compute the induced inner product of two elementary tensors with $v, \tilde{v} \in H^1(\Omega)$, $w, \tilde{w} \in L^2_{\mathbb{P}}(\Xi)$:

$$\begin{aligned} (v \otimes w, \tilde{v} \otimes \tilde{w})_{H^1(\Omega) \otimes L^2(\Xi)} &= \left(\int_{\Omega} v \tilde{v} + \nabla v \cdot \nabla \tilde{v} \, dx \right) \left(\int_{\Xi} w \tilde{w} \, d\mathbb{P} \right) = \\ &= \int_{\Xi} \int_{\Omega} (v \otimes w) (\tilde{v} \otimes \tilde{w}) + \nabla_x (v \otimes w) \cdot \nabla_x (\tilde{v} \otimes \tilde{w}) \, dx \, d\mathbb{P} \end{aligned}$$

Similarly as above this extends to the whole space, which can also be characterized analogously to an anisotropic Sobolev space [54, Sec. 4.2.3]. We identify $H^1(\Omega) \otimes L^2_{\mathbb{P}}(\Xi) = \{\mathbf{y} \in L^2(\Omega \times \Xi) : \partial_{x_i} \mathbf{y} \in L^2(\Omega \times \Xi) \forall i \in [n]\}$. The partial derivatives w. r. t. the space variables are weak derivatives. This space is also isomorphic to the Bochner space $L^2_{\mathbb{P}}(\Xi; H^1(\Omega))$ of $H^1(\Omega)$ -valued random variables with finite variance [102, Thm. B.17], [34, Chap. 1, Sec. 7].

Again, we want to extend this example to Bochner spaces with general $p \in [1, \infty)$ to obtain the space of choice for the semilinear, elliptic PDE with uncertain coefficients discussed in Chapter 3. For a given Banach space Y and a measure space $(\Omega, \mathcal{A}, \mu)$, the Bochner space $L^p_{\mu}(\Omega; Y)$ ($p \in [1, \infty]$) is the space of equivalence classes of strongly measurable functions $\mathbf{y} : \Xi \rightarrow Y$ such that the norm

$$\begin{aligned} \|\mathbf{y}\|_{L^p_{\mathbb{P}}(\Xi; Y)} &:= \left(\int_{\Omega} \|\mathbf{y}(\omega)\|_Y^p \, d\mu(\omega) \right)^{\frac{1}{p}} \quad (p \in [1, \infty)), \\ \|\mathbf{y}\|_{L^{\infty}_{\mathbb{P}}(\Xi; Y)} &:= \operatorname{ess\,sup}_{\omega \in \Omega} \|\mathbf{y}(\omega)\|_Y \end{aligned}$$

is finite. If Y is a Hilbert space, $L^2_{\mu}(\Omega; Y)$ is a Hilbert space with inner product

$$(\mathbf{y}, \mathbf{v})_{L^2_{\mu}(\Omega; Y)} := \int_{\Omega} (\mathbf{y}(\omega), \mathbf{v}(\omega))_Y \, d\mu(\omega).$$

More information about Bochner spaces can be found in, e. g., [66, Chaps. 1 and 2].

Example 2.12 (Connection to Bochner spaces). Let Y be a Banach space, let $(\Xi, \mathcal{F}, \mathbb{P})$ be a probability space and let $p \in [1, \infty)$. Consider the algebraic tensor product $Y \otimes_a L^p_{\mathbb{P}}(\Xi)$. As seen in Example 2.8, we can identify an elementary tensor $y \otimes w$ with a Y -valued function \mathbf{y} on Ξ via $\mathbf{y}(\xi) = y w(\xi)$. Using this identification, the space $Y \otimes_a L^p_{\mathbb{P}}(\Xi)$ can be endowed with the $L^p_{\mathbb{P}}(\Xi; Y)$ -norm, which is finite on the algebraic tensor product:

$$\|y \otimes w\|_{L^p_{\mathbb{P}}(\Xi; Y)}^p = \int_{\Xi} \|y\|_Y^p |w(\xi)|^p \, d\mathbb{P} = \|y\|_Y^p \|w\|_{L^p_{\mathbb{P}}(\Xi)}^p < \infty.$$

By the argument that simple functions of the form $\sum_{k=1}^n y_k 1_{A_k}(\cdot)$ with $y_k \in Y$, $A_k \in \mathcal{F}$

are dense in $L_{\mathbb{P}}^p(\Xi; Y)$, see [34, Chap. I, Sec. 7.2], we obtain equality of the completion of the algebraic tensor product and the Bochner space: $Y \otimes L_{\mathbb{P}}^p(\Xi) \cong L_{\mathbb{P}}^p(\Xi; Y)$.

We have seen in (2.2) that tensor products of linear maps $\varphi_i \in \mathcal{L}(V_i, W_i)$ are defined by their action on elementary tensors. By linearity, this definition extends to the algebraic tensor product space: $\bigotimes_{i=1}^m \varphi_i : {}_a \bigotimes_{i=1}^m V_i \rightarrow {}_a \bigotimes_{i=1}^m W_i \subset \bigotimes_{i=1}^m W_i$. If $\bigotimes_{i=1}^m \varphi_i$ is bounded, it can be uniquely extended to the space $\bigotimes_{i=1}^m V_i$ as it is defined on a dense subset. Note that it is important to know the norms on the tensor products of spaces to derive continuity of the tensor product of operators even though each operator φ_i itself might be bounded.

3. A Class of Optimal Control Problems under Uncertainty

In this chapter, we present a class of risk-neutral optimal control problems under uncertainty with a semilinear, elliptic PDE as state equation. A suitable functional analytic setting in a Bochner space is developed such that all appearing functions are well-defined and the optimal control problem admits a solution. The connection to tensor products of Banach spaces as discussed in Example 2.12 is drawn. Furthermore, adjoint formulations of the derivatives of the reduced objective function are derived to enable the efficient, gradient-based solution of problems of this class.

3.1. Problem Formulation

Let $\xi \in \mathbb{R}^m$ be a vector of $m \in \mathbb{N}$ independently distributed, real-valued random variables. Note that we do not consider a more general probability space here, but involve directly the finite noise assumption as done in [70]. These finitely many random variables can also originate from a truncated Karhunen-Loève expansion [106, Thm. 11.4] of some random field. Then they are only uncorrelated, and independence has to be assumed or shown additionally; for example, it follows from uncorrelatedness if they are Gaussian. They will act as uncertain parameters in the systems we consider. For each parameter ξ_i ($i \in [m]$) we have a corresponding sample space $\Xi_i \subset \mathbb{R}$ equipped with the Borel σ -algebra and a probability measure \mathbb{P}_i . Due to the independence of the random variables, the random vector ξ gives rise to the product measure $\mathbb{P} := \otimes_{i=1}^m \mathbb{P}_i$ on $\Xi := \times_{i=1}^m \Xi_i$, cf. Example 2.9. We define the mean as $\bar{\xi} := \int_{\Xi} \xi \, d\mathbb{P}$. For a function $\mathbf{v} : \Xi \rightarrow \mathbb{R}$, the expectation is defined by $\mathbb{E}[\mathbf{v}] := \int_{\Xi} \mathbf{v} \, d\mathbb{P}$.

Now let Y and U be Hilbert spaces of deterministic functions and denote by Y^* the dual space of Y . We consider parametrized, nonlinear PDEs of the form

$$A(\xi)y(\xi) + N(y(\xi), \xi) = B(\xi)u + b(\xi) \tag{3.1}$$

with strongly measurable functions $A : \Gamma \rightarrow \mathcal{L}(Y, Y^*)$, $N : Y \times \Xi \rightarrow Y^*$, $B : \Xi \rightarrow \mathcal{L}(U, Y^*)$ and $b : \Xi \rightarrow Y^*$, cf. [70]. Moreover, $y(\xi) \in Y$ is the parameter-dependent state and $u \in U$ the deterministic control. Writing

$$E(y, u, \xi) := A(\xi)y + N(y, \xi) - B(\xi)u - b(\xi)$$

we have a state equation of the form (1.1). We want to control the system prior to the observation of the parameters, only knowing their distribution. Typically, one wants to minimize an objective function depending on the state and the control, like the squared deviation $J_1(y(\xi), u, \xi) = \frac{1}{2} \|Q(\xi)y(\xi) - \hat{q}(\xi)\|_H^2$ of the state from a desired state $\hat{q}(\xi) \in H$,

skipping regularization of the control here, where H is a Hilbert space and $Q : \Xi \rightarrow \mathcal{L}(Y, H)$ is strongly measurable. This quantity is a random variable. Therefore, we have to incorporate a risk measure into the objective function to handle the uncertainty in the system. This risk measure can be the expectation or the conditional value-at-risk (CVaR) for example.

In [69, 70] the authors view all parameter-dependent quantities in (3.1) as Banach space-valued random variables and formulate the problem in a suitable Bochner space. The state is a function $y : \Xi \rightarrow Y$, which is required to have finite variance, i. e., to be contained in the Hilbert space $L_{\mathbb{P}}^2(\Xi; Y)$. They assume the continuous dependence of the solution y on the parameters ξ to be able to use a stochastic collocation approach and to approximate stochastic integrals with sparse grids. In contrast to that we formulate (3.1) in the Bochner space $L_{\mathbb{P}}^p(\Xi; Y)$ and apply a stochastic Galerkin discretization later, where $p \in [2, \infty)$ has to be chosen appropriately, and with the identification $\mathbf{y}(x, \xi) = y(\xi)(x)$ as described in Section 2.2. The state space is $\mathbf{Y} := L_{\mathbb{P}}^p(\Xi; Y)$ and the image space is identified with $\mathbf{Y}^* = L_{\mathbb{P}}^{p^*}(\Xi; Y^*)$ with $p^* := \frac{p}{p-1} \in (1, 2]$. Then, for some $\mathbf{b} \in \mathbf{Y}^*$, $\mathbf{b}(\cdot, \xi)$ belongs to Y^* and \mathbf{b} is applied to $\mathbf{y} \in \mathbf{Y}$ via $\langle \mathbf{b}, \mathbf{y} \rangle_{\mathbf{Y}^*, \mathbf{Y}} = \int_{\Xi} \langle \mathbf{b}(\cdot, \xi), \mathbf{y}(\cdot, \xi) \rangle_{Y^*, Y} d\mathbb{P}$. Problem (3.1) is formulated equivalently as

$$\mathbf{A}\mathbf{y} + \mathbf{N}(\mathbf{y}) = \mathbf{B}u + \mathbf{b} \quad (3.2)$$

with $\mathbf{y} \in \mathbf{Y}$, $\mathbf{A} \in \mathcal{L}(\mathbf{Y}, \mathbf{Y}^*)$, $\mathbf{N} : \mathbf{Y} \rightarrow \mathbf{Y}^*$, $\mathbf{B} \in \mathcal{L}(U, \mathbf{Y}^*)$ and $\mathbf{b} \in \mathbf{Y}^*$:

$$\begin{aligned} (\mathbf{A}\mathbf{y})(\cdot, \xi) &:= A(\xi)y(\xi), & (\mathbf{N}(\mathbf{y}))(\cdot, \xi) &:= N(y(\xi), \xi), \\ (\mathbf{B}u)(\cdot, \xi) &:= B(\xi)u, & \mathbf{b}(\cdot, \xi) &:= b(\xi). \end{aligned}$$

The space of the desired state $\hat{\mathbf{q}}$ is $\mathbf{H} := L_{\mathbb{P}}^2(\Xi, H)$, $\hat{\mathbf{q}}(\cdot, \xi) = \hat{q}(\xi)$, and the operator $\mathbf{Q} \in \mathcal{L}(\mathbf{Y}, \mathbf{H})$ is the tensor version of Q , i. e., $(\mathbf{Q}\mathbf{y})(\cdot, \xi) := Q(\xi)y(\xi)$.

Throughout, we assume that $\hat{\mathbf{q}} \in L_{\mathbb{P}}^2(\Xi; H)$ holds. If the operators A , B , the functional b given after (3.1), and the operator Q are well-defined and in particular regular enough w. r. t. the parameters, meaning that they are q -integrable with a large enough $q \in [1, \infty)$ or even essentially bounded, the tensor operators \mathbf{A} , \mathbf{B} , \mathbf{Q} , and the functional \mathbf{b} are also well-defined:

Proposition 3.1. *Let Y and Z be Banach spaces and let $\Xi \subset \mathbb{R}^m$ be measurable and equipped with the probability measure \mathbb{P} . Let $p_Y, p_Z \in [1, \infty]$, $p_Y \geq p_Z$ be given and define $\mathbf{Y} := L_{\mathbb{P}}^{p_Y}(\Xi; Y)$ and $\mathbf{Z} := L_{\mathbb{P}}^{p_Z}(\Xi; Z)$. Assume that $A \in L_{\mathbb{P}}^{r_A}(\Xi, \mathcal{L}(Y, Z))$ holds for $r_A \in [1, \infty]$ such that $\frac{1}{r_A} + \frac{1}{p_Y} = \frac{1}{p_Z}$, i. e., $\int_{\Xi} \|A(\xi)\|_{\mathcal{L}(Y, Z)}^{r_A} d\mathbb{P} \leq C^{r_A} < \infty$ for $r_A < \infty$ and $\|A(\xi)\|_{\mathcal{L}(Y, Z)} \leq C$ for \mathbb{P} -a. e. $\xi \in \Xi$ in the case $r_A = \infty$.*

Then, the operator \mathbf{A} defined by $(\mathbf{A}\mathbf{y})(\cdot, \xi) = A(\xi)y(\xi)$ for all $\xi \in \Xi$ belongs to $\mathcal{L}(\mathbf{Y}, \mathbf{Z})$ and in particular it holds that $\|\mathbf{A}\mathbf{y}\|_{\mathbf{Z}} \leq C\|\mathbf{y}\|_{\mathbf{Y}}$ for all $\mathbf{y} \in \mathbf{Y}$.

Proof. This result is obtained by some standard estimates. In the case $q_Z < \infty$ we have

$$\begin{aligned} \|\mathbf{A}\mathbf{y}\|_{\mathbf{Z}}^{p_Z} &= \int_{\Xi} \|(\mathbf{A}\mathbf{y})(\cdot, \xi)\|_{\mathbf{Z}}^{p_Z} d\mathbb{P} = \int_{\Xi} \|A(\xi)y(\xi)\|_{\mathbf{Z}}^{p_Z} d\mathbb{P} \\ &\stackrel{\text{H\"older's inequality}}{\leq} \left\| \|A(\cdot)\|_{\mathcal{L}(Y, Z)}^{p_Z} \right\|_{L_{\mathbb{P}}^{r_A/p_Z}(\Xi)} \cdot \left\| \|y(\cdot)\|_{\mathbf{Y}}^{p_Z} \right\|_{L_{\mathbb{P}}^{p_Y/p_Z}(\Xi)} \leq C^{p_Z} \|\mathbf{y}\|_{\mathbf{Y}}^{p_Z} \end{aligned}$$

using $\frac{p_Z}{r_A} + \frac{p_Z}{p_Y} = 1$. In the case $p_Z = \infty$, it follows that $p_Y = \infty$ and $r_A = \infty$. Thus,

$$\begin{aligned} \|\mathbf{A}\mathbf{y}\|_Z &= \operatorname{ess\,sup}_{\xi \in \Xi} \|(\mathbf{A}\mathbf{y})(\cdot, \xi)\|_Z = \operatorname{ess\,sup}_{\xi \in \Xi} \|A(\xi)y(\xi)\|_Z \\ &\leq \left(\operatorname{ess\,sup}_{\xi \in \Xi} \|A(\xi)\|_{\mathcal{L}(Y, Z)} \right) \left(\operatorname{ess\,sup}_{\xi \in \Xi} \|y(\xi)\|_Y \right) \leq C \|\mathbf{y}\|_{\mathbf{Y}}. \end{aligned}$$

□

This proposition can be applied to the operator \mathbf{A} from (3.2) by setting $Z = Y^*$, $p_Y = p \in [2, \infty)$, and $p_Z = p^* = \frac{p}{p-1} \in (1, 2]$. We then need $r_A = \frac{p}{p-2}$ to obtain a bounded linear operator $A \in \mathcal{L}(Y, Y^*)$.

Remark 3.2.

- Note that in the case $p = 2$ we require the uniform boundedness of the operators $A(\xi)$ due to $r_A = \infty$; 2-integrability w. r. t. the parameters would not be sufficient.
- The same proposition applies in a similar manner to the operators \mathbf{B} and \mathbf{Q} and the functional \mathbf{b} . The operator Q , for example, must belong to $L_{\mathbb{P}}^{2p/(p-2)}(\Xi; \mathcal{L}(Y, H))$ for the operator $\mathbf{Q} \in \mathcal{L}(Y, H)$ to be well-defined since we have $p_Y = p \in [2, \infty)$, $p_H = 2$.
- For the nonlinear operator \mathbf{N} , a similar estimation cannot be done; therefore it will be part of the assumptions that it is well-defined.

Using the expectation as risk measure, we obtain the regularized optimal control problem

$$\min_{\mathbf{y} \in \mathbf{Y}, u \in U} \mathbf{J}(\mathbf{y}, u) := \frac{1}{2} \|\mathbf{Q}\mathbf{y} - \hat{\mathbf{q}}\|_H^2 + \frac{\gamma}{2} \|u\|_U^2 \quad \text{s. t.} \quad \mathbf{E}(\mathbf{y}, u) = 0, \quad u \in U_{\text{ad}} \quad (3.3)$$

with the definition $\mathbf{E}(\mathbf{y}, u) := \mathbf{A}\mathbf{y} + \mathbf{N}(\mathbf{y}) - \mathbf{B}u - \mathbf{b}$ and a nonempty, closed and convex set of admissible controls $U_{\text{ad}} \subset U$. Note that the \mathbf{H} -inner product is the expectation of the H -inner product: $(\hat{\mathbf{q}}, \mathbf{v})_{\mathbf{H}} = \int_{\Xi} (\hat{\mathbf{q}}(\cdot, \xi), \mathbf{v}(\cdot, \xi))_H \, d\mathbb{P}$. Hence, the objective function can be written as the expectation of some deterministic objective function:

$$\mathbf{J}(\mathbf{y}, u) = \int_{\Xi} J[\xi](\mathbf{y}(\xi), u) \, d\mathbb{P} \quad \text{with} \quad J[\xi](y, u) := \frac{1}{2} \|Q(\xi)y - \hat{q}(\xi)\|_H^2 + \frac{\gamma}{2} \|u\|_U^2. \quad (3.4)$$

Having

$$J_1(y, u, \xi) := \frac{1}{2} \|Q(\xi)y - \hat{q}(\xi)\|_H^2, \quad J_2(u) := \frac{\gamma}{2} \|u\|_U^2, \quad (3.5)$$

and $\mathcal{R} \equiv \mathbb{E}$ this fits into the general setting (1.2).

3.2. A Class of Semilinear, Elliptic PDEs with Uncertain Coefficients

In this section, we concretize the model problem and discuss a class of parameter-dependent, elliptic PDEs. For this purpose let $\Omega \subset \mathbb{R}^n$ ($n \in \{2, 3\}$) be an open, bounded domain with Lipschitz boundary $\partial\Omega$. We choose $Y := H_0^1(\Omega)$ as state space, i. e., the Sobolev space of weakly differentiable $L^2(\Omega)$ functions with $L^2(\Omega)$ derivatives and zero boundary data in the

sense of traces. In the light of Poincaré's inequality [60, Thm. 1.13], we equip this space with the inner product $(v, \tilde{v})_{H_0^1(\Omega)} = \int_{\Omega} \nabla v \cdot \nabla \tilde{v} \, dx$, which is not the standard $H^1(\Omega)$ -inner product given by $(v, \tilde{v})_{H^1(\Omega)} = \int_{\Omega} v \tilde{v} + \nabla v \cdot \nabla \tilde{v} \, dx$. $U := L^2(\Omega_u)$ shall be the control space, where Ω_u can be a measurable subset of Ω for example or—for finite-dimensional controls— $\Omega_u = [n_u]$ for some $n_u \in \mathbb{N}$. The control will act as distributed control on the system via a linear, bounded operator $D \in \mathcal{L}(L^2(\Omega_u), L^2(\Omega))$. Boundary control with, e. g., Ω_u being a subset of $\partial\Omega$ can also be handled by the algorithm presented later, but is not discussed here. The parameter-dependence of the PDE is due to an uncertain, space-dependent function $\kappa \in L^\infty(\Omega \times \Xi)$ and due to a parameter-dependent right-hand side offset $f(\xi) \in L^2(\Omega)$. The considered semilinear, elliptic PDE is

$$-\operatorname{div}(\kappa(\cdot, \xi)\nabla y) + \varphi(y) = Du + f(\xi) \quad (\text{in } \Omega), \quad y = 0 \quad (\text{on } \partial\Omega), \quad (3.6)$$

where $\varphi : \mathbb{R} \rightarrow \mathbb{R}$ is possibly nonlinear.

Assumption 3.3. We require the following additional properties of the problem data:

- The function $\kappa \in L^\infty(\Omega \times \Xi)$ is uniformly bounded and positive, i. e., there exist constants $\underline{\kappa}, \bar{\kappa} \in (0, \infty)$, $\underline{\kappa} \leq \bar{\kappa}$, such that $\underline{\kappa} \leq \kappa(x, \xi) \leq \bar{\kappa}$ holds for almost every $x \in \Omega$ and $\xi \in \Xi$, cf. [44].
- The function $\varphi : \mathbb{R} \rightarrow \mathbb{R}$ is twice continuously differentiable and monotonically increasing. Its second derivative fulfills the growth condition

$$|\varphi''(t)| \leq a''_{\varphi''} + c''_{\varphi''}|t|^{p-3} \quad (3.7)$$

for all $t \in \mathbb{R}$ with constants $a''_{\varphi''}, c''_{\varphi''} \geq 0$ and the exponent $p \in (3, \infty)$ if $n = 2$ and $p \in (3, 6]$ if $n = 3$.

- It holds that $f \in L_{\mathbb{P}}^{r_f}(\Xi; L^2(\Omega))$ for some $r_f \in [p, \infty]$ with p from (3.7).
- $D : U \rightarrow L^2(\Omega)$ is linear and bounded.

Remark 3.4.

- The growth condition (3.7) implies similar conditions for the function φ itself and its first derivative by Lemma A.10 with the exponents $p - 1$ and $p - 2$, respectively.
- An additional dependence of φ on the parameters ξ and the space variable x could be included easily (see Section A.4), but is skipped here to keep the presentation compact.
- The assumption on φ is used to show that the induced Nemytskii operator is well-defined and twice continuously differentiable.

The weak formulation of (3.6) is

$$(\kappa(\cdot, \xi)\nabla y, \nabla v)_{L^2(\Omega)^n} + \int_{\Omega} \varphi(y)v \, dx = (Du, v)_{L^2(\Omega)} + (f(\xi), v)_{L^2(\Omega)} \quad \forall v \in Y. \quad (3.8)$$

Proposition 3.5. *Under Assumption 3.3, the state equation (3.8) has a unique solution $y(\xi) = S[\xi](u) \in Y \cap \mathcal{C}(\bar{\Omega})$. It satisfies the estimate*

$$\|S[\xi](u)\|_{H_0^1(\Omega)} \leq \frac{C_\Omega}{\kappa} \|Du + f(\xi) - \varphi(0)\|_{L^2(\Omega)}, \quad (3.9)$$

where C_Ω is the constant from Poincaré's inequality depending only on Ω .

Proof. Existence is shown in, e. g., [60, Thm. 1.25, Remark 1.12]. The Sobolev embedding $H^1(\Omega) \hookrightarrow L^p(\Omega)$ implies that the nonlinear part of the equation is well-defined. The given a priori estimate can be proven by inserting $v = y$ into (3.8) and using $(\varphi(y) - \varphi(0))y \geq 0$ because of the monotonicity of φ . \square

The map $S[\xi] : U \rightarrow Y$ is the parametrized control-to-state mapping. We want to emphasize that the uncertain parameters enter only in the realization $f(\xi)$ of the right-hand side in the estimate (3.9).

With $Y^* = H_0^1(\Omega)^* =: H^{-1}(\Omega)$ and the definitions

$$\begin{aligned} A(\xi) : Y &\rightarrow Y^*, \quad \langle A(\xi)y, v \rangle_{Y^*, Y} = (\kappa(\cdot, \xi) \nabla y, \nabla v)_{L^2(\Omega)^n}, \\ N : Y \times \Xi &\rightarrow Y^*, \quad \langle N(y, \xi), v \rangle_{Y^*, Y} = \int_{\Omega} \varphi(y) v \, dx, \\ B(\xi) : U &\rightarrow Y^*, \quad \langle B(\xi)u, v \rangle_{Y^*, Y} = (Du, v)_{L^2(\Omega)}, \\ b(\xi) &\in Y^*, \quad \langle b(\xi), v \rangle_{Y^*, Y} = (f(\xi), v)_{L^2(\Omega)} \end{aligned} \quad (3.10)$$

for every $v \in H_0^1(\Omega)$, (3.8) can be written in the form of (3.1). Parameter dependence of the operator B can be handled, but we consider only this fixed form for simplicity reasons.

Now, writing $\mathbf{y}(\cdot, \xi) = y(\xi)(\cdot)$, we use the state space $\mathbf{Y} = L_{\mathbb{P}}^p(\Xi; H_0^1(\Omega))$. We require p to be chosen according to Assumption 3.3. Especially, if φ is a truly nonlinear function, $p > 3$ is required to get a twice continuously differentiable Nemytskii operator. In the linear case $\varphi \equiv 0$, $p = 2$ —as also allowed in Section 3.1—is sufficient and a better choice since Hilbert space theory can be applied then, see, e. g., [46]. We write (3.10) as an equation in \mathbf{Y} , in a weak sense also w. r. t. the parameters: $\mathbf{A}\mathbf{y} + \mathbf{N}(\mathbf{y}) = \mathbf{B}u + \mathbf{b}$ as in (3.2), where

$$\begin{aligned} \mathbf{A} : \mathbf{Y} &\rightarrow \mathbf{Y}^*, \quad \langle \mathbf{A}\mathbf{y}, \mathbf{v} \rangle_{\mathbf{Y}^*, \mathbf{Y}} = \int_{\Xi} (\kappa(\cdot, \xi) \nabla_x \mathbf{y}(\cdot, \xi), \nabla_x \mathbf{v}(\cdot, \xi))_{L^2(\Omega)^n} \, d\mathbb{P}, \\ \mathbf{N} : \mathbf{Y} &\rightarrow \mathbf{Y}^*, \quad \langle \mathbf{N}(\mathbf{y}), \mathbf{v} \rangle_{\mathbf{Y}^*, \mathbf{Y}} = \int_{\Xi} \int_{\Omega} \varphi(\mathbf{y}(x, \xi)) \mathbf{v}(x, \xi) \, dx \, d\mathbb{P} \\ \mathbf{B} : U &\rightarrow \mathbf{Y}^*, \quad \langle \mathbf{B}u, \mathbf{v} \rangle_{\mathbf{Y}^*, \mathbf{Y}} = \int_{\Xi} (Du, \mathbf{v}(\cdot, \xi))_{L^2(\Omega)} \, d\mathbb{P} \\ \mathbf{b} &\in \mathbf{Y}^*, \quad \langle \mathbf{b}, \mathbf{v} \rangle_{\mathbf{Y}^*, \mathbf{Y}} = \int_{\Xi} (f(\xi), \mathbf{v}(\cdot, \xi))_{L^2(\Omega)} \, d\mathbb{P} \quad \forall \mathbf{v} \in \mathbf{Y}. \end{aligned} \quad (3.11)$$

This equation has a unique solution in \mathbf{Y} , which we prove using the theory of monotone operators and pointwise considerations. First note that the operators in (3.11) are well-

defined. For the linear operator \mathbf{A} this follows from Proposition 3.1, because A belongs to $L_{\mathbb{P}}^{\infty}(\Xi; \mathcal{L}(Y, Y^*)) \subset L_{\mathbb{P}}^{p/(p-2)}(\Xi; \mathcal{L}(Y, Y^*))$:

$$\|A(\xi)\|_{\mathcal{L}(Y, Y^*)} = \sup_{y, v \in Y, \|y\|_Y = \|v\|_Y = 1} (\kappa(\cdot, \xi) \nabla y, \nabla v)_{L^2(\Omega)^n} \leq \|\kappa(\cdot, \xi)\|_{L^{\infty}(\Omega)} \leq \bar{\kappa}$$

holds for a. e. $\xi \in \Xi$ by, e. g., Hölder's inequality and the Cauchy-Schwarz inequality on \mathbb{R}^n . This gives that $\int_{\Xi} \|A(\xi)\|_{\mathcal{L}(Y, Y^*)}^{p/(p-2)} d\mathbb{P} \leq \bar{\kappa}^{p/(p-2)}$ for any $p \in [2, \infty)$ and especially $p \in (3, \infty)$ by Proposition A.2 and thus that \mathbf{A} is bounded: $\|\mathbf{A}\|_{\mathcal{L}(\mathbf{Y}, \mathbf{Y}^*)} \leq \bar{\kappa}$. Similarly, we compute bounds for the norms

$$\|B(\xi)\|_{\mathcal{L}(U, Y^*)} = \sup_{u \in U, v \in Y, \|u\|_U = 1, \|v\|_Y = 1} (Du, v)_{L^2(\Omega)} \leq \|D\|_{\mathcal{L}(U, L^2(\Omega))}$$

resulting in $\|\mathbf{B}\|_{\mathcal{L}(U, \mathbf{Y}^*)} \leq \|D\|_{\mathcal{L}(U, L^2(\Omega))}$, and

$$\|b(\xi)\|_{Y^*} = \sup_{v \in Y, \|v\|_Y = 1} (f(\xi), v)_{L^2(\Omega)} \leq \|f(\xi)\|_{L^2(\Omega)}$$

implying $\|\mathbf{b}\|_{\mathbf{Y}^*} \leq \|f\|_{L_{\mathbb{P}}^* (\Xi; L^2(\Omega))}$ as in Proposition 3.1. As long as $r_f \geq p^*$ holds, which can be deduced from Assumption 3.3, this is bounded by Proposition A.2. Therefore, all appearing linear operators are bounded and also twice continuously differentiable.

Proposition 3.6. *Under the conditions on φ and p stated in Assumption 3.3, the nonlinear superposition operator $\mathbf{N} : \mathbf{Y} \rightarrow \mathbf{Y}^*$ as defined in (3.11) is well-defined, monotone and twice continuously Fréchet-differentiable. The derivatives are given by*

$$\begin{aligned} \langle \mathbf{N}'(\mathbf{y})\mathbf{v}, \tilde{\mathbf{v}} \rangle_{\mathbf{Y}^*, \mathbf{Y}} &= \int_{\Xi} \int_{\Omega} \varphi'(\mathbf{y}(x, \xi)) \mathbf{v}(x, \xi) \tilde{\mathbf{v}}(x, \xi) dx d\mathbb{P}, \\ \langle [\mathbf{N}''(\mathbf{y})\mathbf{v}]\mathbf{w}, \tilde{\mathbf{v}} \rangle_{\mathbf{Y}^*, \mathbf{Y}} &= \int_{\Xi} \int_{\Omega} \varphi''(\mathbf{y}(x, \xi)) \mathbf{v}(x, \xi) \mathbf{w}(x, \xi) \tilde{\mathbf{v}}(x, \xi) dx d\mathbb{P}. \end{aligned}$$

Proof. This statement is proven in Section A.4. Note that Lemma A.10 shows that Assumption A.11 is satisfied for the function φ . \square

For a discussion of the more concrete case $\varphi(t) = t^3$ and $p = 4$ see Section A.5 in the appendix.

Unique Solvability of the State Equation

Now we show that the equation

$$\tilde{\mathbf{N}}(\mathbf{y}) := \mathbf{A}\mathbf{y} + \mathbf{N}(\mathbf{y}) = \tilde{\mathbf{b}}, \quad \mathbf{y} \in \mathbf{Y} \tag{3.12}$$

has a unique solution $\mathbf{y} \in \mathbf{Y}$ for all $\tilde{\mathbf{b}} \in \mathbf{Y}^*$ with $\tilde{\mathbf{b}} = \mathbf{B}u + \mathbf{b}$, where $u \in U$ is given. For this purpose we use the following theorem, which is obtained by combining Theorem 26.A, Definition 25.2, and Definition 26.1 in [117]:

Theorem 3.7 (Excerpt from the Minty-Browder theorem about monotone operators). *Let Y be a real, reflexive Banach space. Let $\tilde{N} : Y \rightarrow Y^*$ be*

- *strictly monotone, i. e., $\langle \tilde{N}(y) - \tilde{N}(\tilde{y}), y - \tilde{y} \rangle_{Y^*, Y} > 0$ holds for all $y \neq \tilde{y} \in Y$,*
- *coercive, i. e., $\lim_{\|y\|_Y \rightarrow \infty} \frac{\langle \tilde{N}(y), y \rangle_{Y^*, Y}}{\|y\|_Y} = \infty$, and*
- *hemicontinuous, i. e., the real function $t \mapsto \langle \tilde{N}(y + t\tilde{y}), v \rangle_{Y^*, Y}$ is continuous on $[0, 1]$ for all $y, \tilde{y}, v \in Y$.*

Then, the equation

$$\tilde{N}(y) = b, \quad y \in Y$$

has a unique solution $y \in Y$ for all right-hand sides $b \in Y^*$. The inverse operator $\tilde{N}^{-1} : Y^* \rightarrow Y$ exists and is strictly monotone, i. e., $\langle \tilde{N}^{-1}(b) - \tilde{N}^{-1}(\tilde{b}), b - \tilde{b} \rangle_{Y, Y^*} > 0$ for all $b \neq \tilde{b} \in Y^*$ and bounded. If \tilde{N} is additionally

- *strongly monotone, i. e., there exists $c > 0$ such that $\langle \tilde{N}(y) - \tilde{N}(\tilde{y}), y - \tilde{y} \rangle_{Y^*, Y} \geq c\|y - \tilde{y}\|_Y^2$ holds for all $y, \tilde{y} \in Y$,*

the inverse operator \tilde{N}^{-1} is Lipschitz continuous with Lipschitz constant c^{-1} .

Remark 3.8. Note that a strongly monotone operator \tilde{N} is always

- monotone, i. e., $\langle \tilde{N}(y) - \tilde{N}(\tilde{y}), y - \tilde{y} \rangle_{Y^*, Y} \geq 0$ for all $y, \tilde{y} \in Y$,
- coercive, and
- strictly monotone.

Corollary 3.9. *Under Assumption 3.3, the state equation (3.8) has a unique solution $y(\xi) = S[\xi](u) \in Y$ for almost every $\xi \in \Xi$. The parametrized control-to-state mapping $S[\xi] : U \rightarrow Y$ is Lipschitz continuous with constant $\frac{C_\Omega \|D\|_{\mathcal{L}(U, L^2(\Omega))}}{\underline{\kappa}}$.*

Proof. The operator $A(\xi)$ defined in (3.10) is strongly monotone on $Y = H_0^1(\Omega)$ with constant $\underline{\kappa}$ for a. e. $\xi \in \Xi$. The nonlinear operator N is well-defined by the growth condition (3.7) and the Sobolev embedding $H_0^1(\Omega) \hookrightarrow L^p(\Omega)$, which implies $\varphi(y) \in L^{p/(p-1)}(\Omega)$. It is monotone due to $(\varphi(y(x)) - \varphi(\tilde{y}(x)))(y(x) - \tilde{y}(x)) \geq 0$ for all $x \in \Omega$ and arbitrary $y, \tilde{y} \in Y$ because of the monotonicity of φ . Thus, the operator $\tilde{N}(\xi) \equiv A(\xi) + N : Y \rightarrow Y^*$ is strongly monotone with constant $\underline{\kappa}$. It is continuous (and therefore hemicontinuous) because $A(\xi)$ is bounded and N is continuous as discussed in Section A.3. Hence, Theorem 3.7 can be applied to deduce the statements. \square

Remark 3.10. Proposition 3.5 provides the same statement as Corollary 3.9 with the additional fact that the solution $y(\xi)$ is continuous on $\bar{\Omega}$.

When trying to verify the prerequisites from Theorem 3.7 for (3.12) it turns out that strict monotonicity and hemicontinuity of \tilde{N} on $\mathbf{Y} = L_{\mathbb{P}}^p(\Xi; H_0^1(\Omega))$ can be shown quite simply. Hemicontinuity follows from well-definedness and continuity, and strict monotonicity follows

from the strict monotonicity of \mathbf{A} and the monotonicity of \mathbf{N} . Strict monotonicity of the operator \mathbf{A} is equivalent to $\langle \mathbf{A}\mathbf{y}, \mathbf{y} \rangle_{\mathbf{Y}^*, \mathbf{Y}} > 0$ for all $\mathbf{y} \in \mathbf{Y} \setminus \{0\}$. This is true due to

$$\begin{aligned} \langle \mathbf{A}\mathbf{y}, \mathbf{y} \rangle_{\mathbf{Y}^*, \mathbf{Y}} &= \int_{\Xi} (\kappa(\cdot, \xi) \nabla_x \mathbf{y}(\cdot, \xi), \nabla_x \mathbf{y}(\cdot, \xi))_{L^2(\Omega)^n} \, d\mathbb{P} \\ &\geq \underline{\kappa} \int_{\Xi} \|\mathbf{y}(\cdot, \xi)\|_{H_0^1(\Omega)}^2 \, d\mathbb{P} \geq 0. \end{aligned}$$

The last term is 0 if and only if $\|\mathbf{y}(\cdot, \xi)\|_{H_0^1(\Omega)} = 0$ for almost every $\xi \in \Xi$ such that for $\mathbf{y} \neq 0$ we get strict positivity. But the operator $\tilde{\mathbf{N}}$ is not coercive on \mathbf{Y} . This issue can be overcome by considering a larger space endowed with a weaker norm, see Section A.5 for a concrete choice of such space in the case $\varphi(t) = t^3$.

In the general case, we show existence and uniqueness of a solution by constructing it pointwise for a. e. $\xi \in \Xi$.

Definition 3.11 (Control-to-state mapping). Let $\mathbf{S} : U \rightarrow \mathbf{Y}$ defined by $\mathbf{S}(u)(x, \xi) := S[\xi](u)(x)$ for almost every $\xi \in \Xi$, where $S[\xi](u)$ is the weak solution of (3.8), i. e.,

$$\langle A(\xi)S[\xi](u) + N(S[\xi](u), \xi) - \tilde{b}(\xi), v \rangle_{\mathbf{Y}^*, \mathbf{Y}} = 0 \quad \forall v \in \mathbf{Y} \quad (3.13)$$

with $\tilde{b}(\xi) := B(\xi)u + b(\xi)$.

We show that \mathbf{S} is well-defined. Measurability of the weak solution $\mathbf{S}(u)$ w. r. t. ξ can be shown using the following theorem:

Theorem 3.12. *Let $A : \Xi \rightarrow \mathcal{L}(Y, Y^*)$ be Bochner measurable and such that $A(\xi)$ is strongly monotone with constant $\underline{\kappa} > 0$ (independent of ξ) for almost every $\xi \in \Xi$. Let $b : \Xi \rightarrow Y^*$ also be Bochner measurable and let $N : Y \rightarrow Y^*$ be a deterministic, monotone and continuous operator.*

Then, the function $\xi \mapsto y(\xi)$, which is defined almost everywhere on Ξ , where $y(\xi)$ is the unique solution of the equation $A(\xi)y + N(y) = b(\xi)$, is also Bochner measurable.

Proof. W. l. o. g. we assume $N(0) = 0$. If this is not the case, we can subtract the value $N(0)$ from both sides of the equation without changing the required properties of the nonlinearity and the right-hand side. For the rest of the proof, we consider only random vectors $\xi \in \Xi_{\underline{\kappa}} := \{\xi \in \Xi : A(\xi) \text{ is strongly monotone with constant } \underline{\kappa}\}$. By assumption, the set $\Xi_{\underline{\kappa}} \subset \Xi$ has measure 1. Observe that the equation $A(\xi)y + N(y) = b(\xi)$ has a unique solution $y(\xi)$ fulfilling $\|y(\xi)\|_Y \leq \frac{1}{\underline{\kappa}} \|b(\xi)\|_{Y^*}$ by Theorem 3.7 in analogy to Corollary 3.9 for every $\xi \in \Xi_{\underline{\kappa}}$.

Due to the Bochner measurability, there exist sequences $(A_n)_{n \in \mathbb{N}}$, $A_n : \Xi \rightarrow \mathcal{L}(Y, Y^*)$ and $(b_n)_{n \in \mathbb{N}}$, $b_n : \Xi \rightarrow Y^*$ of simple functions with $\lim_{n \rightarrow \infty} \|A(\xi) - A_n(\xi)\|_{\mathcal{L}(Y, Y^*)} = 0$ for every $\xi \in \Xi_A \subset \Xi$ and $\lim_{n \rightarrow \infty} \|b(\xi) - b_n(\xi)\|_{Y^*} = 0$ for every $\xi \in \Xi_b \subset \Xi$, where $\mathbb{P}(\Xi_A) = 1 = \mathbb{P}(\Xi_b)$. Overall, we obtain

$$\lim_{n \rightarrow \infty} \|A(\xi) - A_n(\xi)\|_{\mathcal{L}(Y, Y^*)} + \|b(\xi) - b_n(\xi)\|_{Y^*} = 0$$

for all $\xi \in \Xi_{\underline{\kappa}} \cap \Xi_A \cap \Xi_b$, which is a set of measure 1. Due to convergence we can assume w. l. o. g. that also $A_n(\xi)$ is strongly monotone with constant $\frac{\underline{\kappa}}{2}$ for all $\xi \in \Xi_{\underline{\kappa}} \cap \Xi_A \cap \Xi_b$ and all $n \in \mathbb{N}$ by possibly restricting to a subsequence.

Now define the sequence $(y_n)_{n \in \mathbb{N}}$, $y_n : \Xi \rightarrow Y$, where $y_n(\xi)$ is the unique solution of the equation $A_n(\xi)y + N(y) = b_n(\xi)$ for $\xi \in \Xi_{\underline{\kappa}} \cap \Xi_A \cap \Xi_b$ and $y_n(\xi) = 0$ for $\xi \in \Xi \setminus (\Xi_{\underline{\kappa}} \cap \Xi_A \cap \Xi_b)$. The elements of this sequence fulfill $\|y_n(\xi)\|_Y \leq \frac{2}{\underline{\kappa}} \|b(\xi)\|_{Y^*}$ for all $\xi \in \Xi$. They are also simple functions. If A_n takes $N_A < \infty$ many values on $\mathcal{L}(Y, Y^*)$ and b_n takes $N_b < \infty$ many values on Y^* , then $y_n(\xi)$ admits at most $N_A N_b + 1 < \infty$ many values.

Now we show that indeed $\lim_{n \rightarrow \infty} \|y_n(\xi) - y(\xi)\|_Y = 0$ holds for almost every $\xi \in \Xi$, namely for all $\xi \in \Xi_{\underline{\kappa}} \cap \Xi_A \cap \Xi_b$. By strong monotonicity of $A(\xi)$ and monotonicity of N we get, skipping the argument ξ :

$$\begin{aligned} 0 &\leq \underline{\kappa} \|y - y_n\|_Y^2 \\ &\leq \langle A(y - y_n), y - y_n \rangle_{Y^*, Y} + \langle N(y) - N(y_n), y - y_n \rangle_{Y^*, Y} \\ &= \langle Ay + N(y), y - y_n \rangle_{Y^*, Y} - \langle Ay_n + N(y_n), y - y_n \rangle_{Y^*, Y} \\ &= \langle b, y - y_n \rangle_{Y^*, Y} + \langle (A_n - A)y_n, y - y_n \rangle_{Y^*, Y} - \langle b_n, y - y_n \rangle_{Y^*, Y} \\ &= \langle b - b_n + (A_n - A)y_n, y - y_n \rangle_{Y^*, Y} \\ &\leq (\|b - b_n\|_{Y^*} + \frac{2}{\underline{\kappa}} \|A_n - A\|_{\mathcal{L}(Y, Y^*)} \|b\|_{Y^*}) \|y - y_n\|_Y. \end{aligned}$$

This results in

$$\|y(\xi) - y_n(\xi)\|_Y \leq \frac{1}{\underline{\kappa}} \|b(\xi) - b_n(\xi)\|_{Y^*} + \frac{2}{\underline{\kappa}^2} \|A_n(\xi) - A(\xi)\|_{\mathcal{L}(Y, Y^*)} \|b(\xi)\|_{Y^*}$$

proving the limit $\lim_{n \rightarrow \infty} \|y_n(\xi) - y(\xi)\|_Y = 0$ for all $\xi \in \Xi_{\underline{\kappa}} \cap \Xi_A \cap \Xi_b$. Thus, the function $y : \Xi \rightarrow Y$ is Bochner measurable. \square

Remark 3.13. The proof of Theorem 3.12 shows that the solution y of the equation $Ay + N(y) = b$ depends continuously on the operator A and the right hand side b .

Corollary 3.14. *The function $\xi \mapsto S[\xi](u) \in Y$, where the parametrized control-to-state mapping $S[\xi]$ is implicitly defined by (3.13), is measurable.*

Proof. Obviously, all prerequisites from Theorem 3.12 are fulfilled. \square

Corollary 3.15. *Under Assumption 3.3, the control-to-state mapping from Definition 3.11 is well-defined. It holds that $\mathbf{S}(u) \in L^{r_f}(\Xi; Y)$, i. e., $\mathbf{S}(u)$ inherits its ξ -regularity (integrability or essential boundedness) from $f(\cdot)$.*

Proof. By (3.9) we have the estimate

$$\|\mathbf{S}(u)(\cdot, \xi)\|_{H_0^1(D)} \leq \frac{C_\Omega}{\underline{\kappa}} \|Du + f(\xi) - \varphi(0)\|_{L^2(\Omega)} \quad \text{for a. e. } \xi \in \Xi.$$

Since $f \in L_{\mathbb{P}}^{r_f}(\Xi; L^2(\Omega))$, this shows, together with the Bochner measurability (Corollary 3.14), that $\mathbf{S}(u) \in L_{\mathbb{P}}^{r_f}(\Xi; Y) \subset L_{\mathbb{P}}^p(\Xi; Y) = \mathbf{Y}$. The inclusion of Bochner spaces holds true by Proposition A.2 if $r_f \geq p$. \square

Note that this regularity can be reduced if the function $\xi \mapsto \frac{1}{\underline{\kappa}(\xi)}$ with a parameter-dependent coercivity constant $\underline{\kappa}(\xi) > 0$ for the operator $A(\xi)$ does not belong to $L_{\mathbb{P}}^\infty(\Xi)$. We do not consider this case because we need the uniform coercivity for error estimation later.

Proposition 3.16. *Under Assumption 3.3, the function $\mathbf{y} = \mathbf{S}(u)$ (Definition 3.11) is the unique solution of the equation $\mathbf{A}\mathbf{y} + \mathbf{N}(\mathbf{y}) = \mathbf{B}u + \mathbf{b}$ as defined in (3.11) and (3.12).*

Proof. The equation can be written as

$$\int_{\Xi} \langle A(\xi)\mathbf{y}(\cdot, \xi) + N(\mathbf{y}(\cdot, \xi), \xi) - \tilde{b}(\xi), \mathbf{v}(\cdot, \xi) \rangle_{Y^*, Y} d\mathbb{P} = 0 \quad \forall \mathbf{v} \in \mathbf{Y}.$$

The integrand vanishes for almost every $\xi \in \Xi$ if $\mathbf{y} = \mathbf{S}(u)$ is inserted. Hence, $\mathbf{S}(u) \in L_{\mathbb{P}}^p(\Xi; H_0^1(\Omega))$ (Corollary 3.15) is a solution of the equation.

Uniqueness of the solution in \mathbf{Y} can be proven as follows: Let $\mathbf{y}, \tilde{\mathbf{y}} \in \mathbf{Y}$ both solve (3.12). Then, subtracting the state equation for $\tilde{\mathbf{y}}$ from the one with \mathbf{y} and testing with $\mathbf{y} - \tilde{\mathbf{y}}$ yields

$$0 = \langle \mathbf{A}(\mathbf{y} - \tilde{\mathbf{y}}), \mathbf{y} - \tilde{\mathbf{y}} \rangle_{Y^*, Y} + \langle \mathbf{N}(\mathbf{y}) - \mathbf{N}(\tilde{\mathbf{y}}), \mathbf{y} - \tilde{\mathbf{y}} \rangle_{Y^*, Y} \geq \kappa \|\mathbf{y} - \tilde{\mathbf{y}}\|_{L_{\mathbb{P}}^2(\Xi; H_0^1(\Omega))}^2 \geq 0,$$

giving $\mathbf{y} = \tilde{\mathbf{y}}$ almost everywhere. \square

Example 3.17. To be more concrete, let us consider a specific form of the coefficient function κ . We define a mean function $\kappa_0(x)$ and some functions $\eta_i(x)$ ($i \in [m]$) describing the amount of influence of the parameter ξ_i over the domain Ω . We set

$$\kappa(x, \xi) := \kappa_0(x) \left(1 + \sum_{i=1}^m \xi_i \eta_i(x) \right).$$

Defining $\kappa_i(x) := \kappa_0(x)\eta_i(x)$ we have $\kappa(x, \xi) = \kappa_0(x) + \sum_{i=1}^m \xi_i \kappa_i(x)$, which is the form for random fields originating from a truncated Karhunen-Loève expansion. If we have $\kappa_0, \eta_i \in L^\infty(\Omega)$ for all $i \in [m]$, κ_0 uniformly positive ($0 < \underline{\kappa}_0 \leq \kappa_0(x) \leq \bar{\kappa}_0 < \infty$), bounded random variables, i. e., $\exists C > 0$ s. t. $\mathbb{P}(\{\xi \in \Xi : |\xi| \leq C\}) = 1$, and η_i also suitably bounded (e. g., $|\eta_i(x)| \leq \bar{\eta}_i < \frac{1}{C_m}$ for a. e. $x \in \Omega$ and all $i \in [m]$), then the prerequisites on κ from above are fulfilled and we get an $L_{\mathbb{P}}^2(\Xi; H_0^1(\Omega))$ -elliptic operator \mathbf{A} .

With κ having the specific form from Example 3.17, the operator \mathbf{A} has a useful structure:

Lemma 3.18. *The operator \mathbf{A} defined in (3.11) with $\kappa(x, \xi) = \kappa_0(x)(1 + \sum_{i=1}^m \xi_i \eta_i(x))$ as in Example 3.17 can be written as*

$$\mathbf{A} = A_0 \otimes \left(\bigotimes_{j=1}^m I \right) + \sum_{i=1}^m A_i \otimes \left(\bigotimes_{j=1}^m \tilde{S}_{ij} \right)$$

with the Kronecker product of operators as in Section 2.2, identifying $\mathbf{Y}^* = Y^* \otimes L_{\mathbb{P}_1}^{p^*}(\Xi_1) \otimes \dots \otimes L_{\mathbb{P}_m}^{p^*}(\Xi_m)$, and with $A_0, A_i \in \mathcal{L}(Y, Y^*)$ and $\tilde{S}_{ij} \in \mathcal{L}(L_{\mathbb{P}_j}^p(\Xi_j), L_{\mathbb{P}_j}^{p^*}(\Xi_j))$ ($i, j \in [m]$) defined by

$$\begin{aligned} \langle A_0 y, v \rangle_{Y^*, Y} &:= (\kappa_0 \nabla y, \nabla v)_{L^2(\Omega)^n}, & \langle A_i y, v \rangle_{Y^*, Y} &:= (\eta_i \kappa_0 \nabla y, \nabla v)_{L^2(\Omega)^n}, \\ \tilde{S}_{ij} &= I \text{ for } i, j \in [m], i \neq j, & (\tilde{S}_{jj} v_j)(\xi_j) &:= \xi_j \cdot v_j(\xi_j) \text{ for } j \in [m]. \end{aligned}$$

Proof. At first, we split up the parametrized operator $A(\xi)$ using the definitions of $\kappa(\cdot, \xi)$, A_0 and A_i :

$$\begin{aligned} \langle A(\xi)y, v \rangle_{Y^*, Y} &= (\kappa(\cdot, \xi) \nabla y, \nabla v)_{L^2(\Omega)^n} \\ &= (\kappa_0 \nabla y, \nabla v)_{L^2(\Omega)^n} + \sum_{i=1}^m \xi_i (\kappa_0 \eta_i \nabla y, \nabla v)_{L^2(\Omega)^n} \\ &= \left\langle \left(A_0 + \sum_{i=1}^m \xi_i A_i \right) y, v \right\rangle_{Y^*, Y} \end{aligned}$$

Next we consider an elementary tensor $\mathbf{y} = y \otimes v_1 \otimes \dots \otimes v_m \in \mathbf{Y}$, i. e., $\mathbf{y}(x, \xi) = y(x) \cdot \prod_{j=1}^m v_j(\xi_j)$, and compute for fixed ξ :

$$\begin{aligned} A(\xi)\mathbf{y}(\cdot, \xi) &= \left(A_0 + \sum_{i=1}^m \xi_i A_i \right) \left(y(\cdot) \cdot \prod_{j=1}^m v_j(\xi_j) \right) = \\ &= A_0 y(\cdot) \cdot \prod_{j=1}^m v_j(\xi_j) + \sum_{i=1}^m A_i y(\cdot) \cdot \xi_i \cdot \prod_{j=1}^m v_j(\xi_j) = \\ &= \left((A_0 y) \otimes \left(\bigotimes_{j=1}^m v_j \right) \right) (\cdot, \xi) + \sum_{i=1}^m \left((A_i y) \otimes \left(\bigotimes_{j=1}^m \tilde{S}_{ij} v_j \right) \right) (\cdot, \xi). \end{aligned}$$

Thus, we can write the operator \mathbf{A} as $\mathbf{A} = A_0 \otimes \left(\bigotimes_{j=1}^m I \right) + \sum_{i=1}^m A_i \otimes \left(\bigotimes_{j=1}^m \tilde{S}_{ij} \right)$. Because of linearity and continuity this extends to algebraic and topological tensors \mathbf{y} as seen at the end of Section 2.2. The operator \mathbf{A} is continuous from $L_{\mathbb{P}}^p(\Xi; Y)$ to $L_{\mathbb{P}}^{p^*}(\Xi; Y^*)$ as already noted. \square

The structure of the operator \mathbf{A} provided in Lemma 3.18 uses the identification of the state space \mathbf{Y} with a tensor product of Banach spaces and is useful because it makes the discrete version of the operator \mathbf{A} easily applicable to tensors in low-rank formats that implement i -mode matrix multiplication and componentwise sums. Formally, we can write $\mathbf{A}\mathbf{y} = A_0 \circ_1 \mathbf{y} + \sum_{i=1}^m A_i \circ_1 \tilde{S}_{ii} \circ_{i+1} \mathbf{y}$ using the application of a linear operator to a certain mode of a tensor, which we have defined only for finite dimensional tensors so that this will make more sense in the discretized setting (Example 6.5). In our work [46] it was already demonstrated that also more difficult operators as resulting from domain parametrizations can be efficiently handled with low-rank tensors. In this thesis we present an adaptive approach and its convergence theory and base our numerical tests only on this example since reliable a posteriori error estimates, which we use for the adaptive solution of the state equation, can be derived in this setting, see Chapter 7.

3.3. Existence of a Solution to the Optimal Control Problem

A natural question is whether the optimal control problem (3.3) with the state equation from Section 3.2 admits a solution. Typically, continuity under weak convergence is used in existence proofs. In nonlinear cases, this is often done by means of compact embeddings,

which yield strong convergence (w. r. t. a weaker norm) of weakly converging sequences. In the Bochner space setting considered here, only the continuous embedding $L^p \hookrightarrow L^q$ is available for $p \geq q$ (see Proposition A.2), but this embedding is not compact. Therefore, pointwise considerations for a. e. $\xi \in \Xi$ are used to prove the following existence theorem.

Theorem 3.19. *Let Assumption 3.3 hold with $p \in (3, \infty)$ in the case $n = 2$ and $p \in (3, 6)$ in the case $n = 3$. Furthermore, let $\gamma > 0$, $\hat{\mathbf{q}} \in L^2_{\mathbb{P}}(\Xi; H)$, $\mathbf{Q} \in L^{r_Q}_{\mathbb{P}}(\Xi; \mathcal{L}(Y, H))$ with $\frac{1}{r_Q} + \frac{1}{r_f} = \frac{1}{2}$, and let $Q(\xi) : Y \rightarrow H$ be a compact operator for a. e. ξ .*

Then, Problem (3.3) with the state equation discussed in Section 3.2 has a solution.

Proof. The statement is proven analogously to [111, Lem. 9.4]. Let $(\mathbf{y}_k, u_k)_{k \in \mathbb{N}}$ be a feasible, infimizing sequence, where $\mathbf{y}_k = \mathbf{S}(u_k)$ is uniquely given due to Proposition 3.16. Because of the regularization term in the objective function with $\gamma > 0$ and the non-negative tracking term, the sequence (u_k) is bounded in $L^2(\Omega_u) = U$. By (3.9), we have for a. e. $\xi \in \Xi$ that the corresponding, pointwise states fulfill

$$\|\mathbf{y}_k(\xi)\|_{H_0^1(\Omega)} \leq \frac{C_{\Omega}}{\kappa} (\|D\|_{\mathcal{L}(U, L^2(\Omega))} \|u_k\|_U + \|f(\xi) - \varphi(0)\|_{L^2(\Omega)}). \quad (3.14)$$

Note that for a. e. $\xi \in \Xi$ the state equation has a unique solution for *arbitrary* u by Assumption 3.3. By boundedness, and since U_{ad} is a convex, closed subset of a reflexive Banach space, we can extract subsequences for a. e. $\xi \in \Xi$, again denoted by (u_k) and $(\mathbf{y}_k(\xi))$, which converge weakly to some limits $u \in U_{\text{ad}}$ and $\mathbf{y}(\xi) \in H_0^1(\Omega)$. The compact embedding $H_0^1(\Omega) \hookrightarrow L^p(\Omega)$ implies that the sequences $(\mathbf{y}_k(\xi))$ converge strongly in $L^p(\Omega)$ to $\mathbf{y}(\xi)$ for a. e. $\xi \in \Xi$. The growth condition and the continuity of φ as well as the continuous embedding $L^{p^*}(\Omega) \hookrightarrow H^{-1}(\Omega)$ yield that $\varphi(\mathbf{y}_k(\xi)) \rightarrow \varphi(\mathbf{y}(\xi))$ strongly in $H^{-1}(\Omega)$. It follows that

$$\begin{aligned} A(\xi)\mathbf{y}_k(\xi) + \varphi(\mathbf{y}_k(\xi)) &\rightharpoonup A(\xi)\mathbf{y}(\xi) + \varphi(\mathbf{y}(\xi)) \text{ in } H^{-1}(\Omega), \\ Du_k + f(\xi) &\rightharpoonup Du + f(\xi) \text{ in } L^2(\Omega) \subset H^{-1}(\Omega) \end{aligned}$$

so that $\mathbf{y}(\xi)$ solves (3.8) for a. e. ξ . Thus, we have that $\mathbf{y} = \mathbf{S}(u) \in L^{r_f}(\Xi; Y)$ by Corollary 3.15 giving that (\mathbf{y}, u) is feasible for (3.3). The regularization term of the objective function (3.4) is convex and continuous under strong convergence in $L^2(\Omega_u)$, and thus weakly lower semicontinuous. Furthermore, the tracking term converges: Since $\mathbf{y}_k(\xi) \rightarrow \mathbf{y}(\xi)$ in Y and $Q(\xi)$ maps Y to H compactly, the sequence $(Q(\xi)\mathbf{y}_k(\xi))$ converges strongly to $Q(\xi)\mathbf{y}(\xi)$ in H for a. e. ξ . By continuity, $\|Q(\xi)\mathbf{y}_k(\xi) - \hat{q}(\xi)\|_H^2 \rightarrow \|Q(\xi)\mathbf{y}(\xi) - \hat{q}(\xi)\|_H^2$ for a. e. ξ as $k \rightarrow \infty$. Because of the boundedness of $\|u_k\|_U$ by some constant $C_u > 0$ and (3.14), we can bound

$$\|Q(\xi)\mathbf{y}_k(\xi) - \hat{q}(\xi)\|_H^2 \leq \left(\|Q(\xi)\|_{\mathcal{L}(Y, H)} \frac{C_{\Omega}}{\kappa} (\|D\|_{\mathcal{L}(U, L^2(\Omega))} C_u + \|f(\xi) - \varphi(0)\|_{L^2(\Omega)}) + \|\hat{q}(\xi)\|_H \right)^2$$

for every k . The ξ -regularities of $Q(\cdot)$, $f(\cdot)$ and $\hat{q}(\cdot)$ are chosen exactly such that the estimate on the right-hand side is integrable w. r. t. \mathbb{P} . Therefore, the dominated convergence theorem can be applied to conclude

$$\int_{\Xi} \|Q(\xi)\mathbf{y}_k(\xi) - \hat{q}(\xi)\|_H^2 d\mathbb{P} \longrightarrow \int_{\Xi} \|Q(\xi)\mathbf{y}(\xi) - \hat{q}(\xi)\|_H^2 d\mathbb{P} \quad \text{as } k \rightarrow \infty.$$

It follows from the derived continuity properties of \mathbf{J} that (\mathbf{y}, u) solves (3.3). \square

3.4. Derivatives of the Reduced Objective Function

We will apply an adaptive, inexact, gradient-based, nonlinear optimization method to solve problems of type (3.3). This algorithm will be formulated in the function space to be able to apply adaptive solution techniques. Therefore, we need expressions for the derivatives in function space and make the following assumption:

Assumption 3.20. The state equation (3.2) has a unique weak solution $\mathbf{y} = \mathbf{S}(u)$ for every control $u \in U$. The nonlinear operator $\mathbf{N} : \mathbf{Y} \rightarrow \mathbf{Z}$ is twice continuously Fréchet-differentiable.

Then we can reduce problem (3.3) to the control and get

$$\min_{u \in U} \hat{\mathbf{J}}(u) := \mathbf{J}(\mathbf{S}(u), u) \quad \text{s. t.} \quad u \in U_{\text{ad}}. \quad (3.15)$$

We use the adjoint approach for the computation of the gradient of the reduced objective function. The formal adjoint equation for the adjoint state $\mathbf{z} \in \mathbf{Y}$ is

$$\mathbf{E}_{\mathbf{y}}(\mathbf{S}(u), u)^* \mathbf{z} = -\mathbf{J}_{\mathbf{y}}(\mathbf{S}(u), u)$$

and more concretely

$$\langle \mathbf{z}, \mathbf{A}\mathbf{v} + \mathbf{N}'(\mathbf{S}(u))\mathbf{v} \rangle_{\mathbf{Y}, \mathbf{Y}^*} = -(\mathbf{Q}\mathbf{S}(u) - \hat{\mathbf{q}}, \mathbf{Q}\mathbf{v})_{\mathbf{H}} \quad \forall \mathbf{v} \in \mathbf{Y}.$$

Typically, one would assume that the partial derivative $\mathbf{E}_{\mathbf{y}}(\mathbf{S}(u), u)$ is boundedly invertible for every $u \in U$. Then, the adjoint equation has a unique solution $\mathbf{T}(u)$, the control-to-state mapping $\mathbf{S} : U \rightarrow \mathbf{Y}$ is continuously F-differentiable and the adjoint representation of the first derivative of the reduced objective function is given by

$$\hat{\mathbf{J}}'(u) = \mathbf{E}_u(\mathbf{S}(u), u)^* \mathbf{T}(u) + \mathbf{J}_u(\mathbf{y}(u), u)$$

and in our concrete case

$$\langle \hat{\mathbf{J}}'(u), w \rangle_{U^*, U} = \langle \mathbf{T}(u), -\mathbf{B}w \rangle_{\mathbf{Y}, \mathbf{Y}^*} + \gamma(u, w)_U \quad \forall w \in U$$

or in short notation $\nabla \hat{\mathbf{J}}(u) = -\mathbf{B}^* \mathbf{T}(u) + \gamma u$.

We will see that in our case the operator $\mathbf{E}_{\mathbf{y}}(\mathbf{S}(u), u)$ is not necessarily boundedly invertible. Therefore, we will verify that a pointwisely defined adjoint state is sufficient for the computation of the gradient and discuss this for the example presented beforehand.

3.4.1. Discussion of the Example

For the example from Section 3.2, Assumption 3.20 is fulfilled. We have proven that the state equation has a unique weak solution for every control (Proposition 3.16) and that the nonlinear operator is twice continuously differentiable (Proposition 3.6). Next, we show that the formal adjoint equation admits a unique weak solution, which is also constructed pointwise. Again, its regularity w. r. t. ξ is discussed using standard a priori estimates.

First note that the partial derivative $\mathbf{E}_y(\mathbf{S}(u), u) = \mathbf{A} + \mathbf{N}'(\mathbf{S}(u))$ does not necessarily have a bounded inverse for every u . Assume that $\mathbf{N}'(\mathbf{S}(u)) \equiv 0$, then $\mathbf{E}_y(\mathbf{S}(u), u) = \mathbf{A}$. This operator is indeed boundedly invertible from $L_{\mathbb{P}}^2(\Xi; Y)$ to $L_{\mathbb{P}}^2(\Xi; Y^*)$. For $p > 2$ we have $\mathbf{A}(L_{\mathbb{P}}^p(\Xi; Y)) = L_{\mathbb{P}}^p(\Xi; Y^*) \subsetneq L_{\mathbb{P}}^{p^*}(\Xi; Y^*)$. The first equality follows from the fact that $A(\xi) : Y \rightarrow Y^*$ is boundedly invertible and $A(\cdot) \in L_{\mathbb{P}}^\infty(\Xi; \mathcal{L}(Y, Y^*))$ as well as $A(\cdot)^{-1} \in L_{\mathbb{P}}^\infty(\Xi; \mathcal{L}(Y^*, Y))$. Hence, $\mathbf{A} : \mathbf{Y} \rightarrow \mathbf{Y}^*$ is not surjective and thus not invertible if $p > 2$.

In the deterministic case, the standard adjoint approach [60, Sec. 1.6] can be applied since the operator $A(\xi) + N_y(y, \xi)$ is boundedly invertible for every $y \in Y$ and almost every $\xi \in \Xi$. This follows from Theorem 3.7 using the strong monotonicity of $A(\xi)$ and the monotonicity of $N_y(y, \xi)$ (Lemma A.4). We can therefore conclude that the pointwise adjoint equation

$$A(\xi)z(\xi) + N_y(y(\xi), \xi)z(\xi) = -Q(\xi)^*(Q(\xi)y(\xi) - \hat{q}(\xi)) \text{ for a. e. } \xi \in \Xi \quad (3.16)$$

with $y(\xi) = S[\xi](u)$ has a unique solution $z(\xi) = T[\xi](u)$ for a. e. $\xi \in \Xi$ and that the parametrized control-to-state mapping $u \mapsto S[\xi](u)$ is continuously differentiable by the implicit function theorem. The derivative of the parametric reduced objective function $\hat{J}[\xi](u) := J[\xi](S[\xi](u), u)$, see the definition of $J[\xi]$ in (3.4), is then given by

$$\langle \hat{J}[\xi]'(u), w \rangle_{U^*, U} = \langle T[\xi](u), -B(\xi)w \rangle_{Y^*, Y} + \gamma(u, w)_U \quad \forall w \in U. \quad (3.17)$$

For $\mathbf{y} \in \mathbf{Y}$, we now consider the formal adjoint equation

$$\mathbf{A}z + \mathbf{N}'(\mathbf{y})z = -\mathbf{Q}^*(\mathbf{Q}\mathbf{y} - \hat{\mathbf{q}}), \quad (3.18)$$

where we use that the operators \mathbf{A} and $\mathbf{N}'(\mathbf{y})$ are self-adjoint in our example.

Definition 3.21. Let $\mathbf{T} : U \rightarrow L_{\mathbb{P}}^{p/(p-1)}(\Xi; Y)$, $\mathbf{T}(u)(x, \xi) := T[\xi](u)(x)$ be defined for almost every $\xi \in \Xi$, where $T[\xi](u)$ is the weak solution of (3.16), i. e.,

$$\langle A(\xi)T[\xi](u) + N_y(S[\xi](u), \xi)T[\xi](u) - \tilde{b}(\xi), v \rangle_{Y^*, Y} = 0 \quad \forall v \in Y \quad (3.19)$$

with $\tilde{b}(\xi) := -Q(\xi)^*(Q(\xi)S[\xi](u) - \hat{q}(\xi))$ and $S[\xi]$ as in Definition 3.11.

Proposition 3.22. Assume that Assumption 3.3 holds true and $\mathbf{Q} \in L_{\mathbb{P}}^{2p/(p-2)}(\Xi; \mathcal{L}(Y, H))$ as well as $\hat{\mathbf{q}} \in L_{\mathbb{P}}^2(\Xi; H)$. Then, the operator \mathbf{T} from Definition 3.21 is well-defined.

Proof. Since $z \mapsto N_y(y, \xi)z$ is monotone (Proposition A.4) and $A(\xi)$ is strongly monotone with constant $\underline{\kappa}$, (3.19) has a unique solution $T[\xi](u)$ for a. e. $\xi \in \Xi$ due to Theorem 3.7. The estimate

$$\begin{aligned} \|\mathbf{T}(u)(\cdot, \xi)\|_Y &\leq \frac{1}{\underline{\kappa}} \|Q(\xi)^*(Q(\xi)S[\xi](u) - \hat{q}(\xi))\|_{Y^*} \\ &\leq \frac{1}{\underline{\kappa}} \|Q(\xi)\|_{\mathcal{L}(Y, H)} \|Q(\xi)S[\xi](u) - \hat{q}(\xi)\|_H \end{aligned} \quad (3.20)$$

can be derived in analogy to (3.9). Therefore, the adjoint state inherits its regularity w. r. t. ξ from the ξ -regularity of $Q^*(\cdot)(Q(\cdot)S[\cdot](u) - \hat{q}(\cdot))$. Since $S[\cdot](u) \in L_{\mathbb{P}}^p(\Xi; Y)$ by Corollary 3.15, Hölder's inequality implies that $Q(\cdot)S[\cdot](u) \in L_{\mathbb{P}}^2(\Xi; H)$, hence $Q(\cdot)^*(Q(\cdot)S[\cdot](u) - \hat{q}(\cdot)) \in L_{\mathbb{P}}^{p/(p-1)}(\Xi; Y^*)$.

The function $\xi \mapsto T[\xi](u)$ is measurable due to Theorem 3.12 as long as the mapping $\xi \mapsto N_y(S[\xi](u), \xi) \in \mathcal{L}(Y, Y^*)$ is measurable. This is true because N_y is deterministic and continuous and $\xi \mapsto S[\xi](u)$ is measurable by Corollary 3.14. Measurability and the regularity estimate (3.20) yield $\mathbf{T}(u) \in L_{\mathbb{P}}^{p/(p-1)}(\Xi; Y)$. \square

Lemma 3.23. *Let Assumption 3.3 hold and assume $\mathbf{y} \in L_{\mathbb{P}}^{r_f}(\Xi; Y)$ for some $r_f \in [p, \infty]$, $\mathbf{Q} \in L_{\mathbb{P}}^{r_Q}(\Xi; \mathcal{L}(Y, H))$ for some $r_Q \in [\frac{2p}{p-2}, \infty]$, and $\hat{\mathbf{q}} \in L_{\mathbb{P}}^{r_{\hat{q}}}(\Xi; H)$ for some $r_{\hat{q}} \in [2, \infty]$. Then, $\mathbf{T}(u) \in L_{\mathbb{P}}^{r_z}(\Xi; Y)$ (Definition 3.21) holds, where $r_z \in [\frac{p}{p-1}, \infty]$ depends on r_f , r_Q , and $r_{\hat{q}}$ and is defined as in the column “Exponents” of Table 3.1.*

Proof. In Proposition 3.22, $\mathbf{T}(u) \in L_{\mathbb{P}}^{p/(p-1)}(\Xi; Y)$ is shown. The regularity exponent is deduced from (3.20). In the same way, the regularity exponents for $\mathbf{Q}\mathbf{y}$, $\mathbf{Q}\mathbf{y} - \hat{\mathbf{q}}$, and $\mathbf{T}(u)$ depicted in Table 3.1 follow from Hölder’s inequality. \square

Proposition 3.24. *Let Assumption 3.3 hold and assume $\mathbf{y} \in L_{\mathbb{P}}^{r_f}(\Xi; Y)$ for some $r_f \in [p, \infty]$, $\mathbf{Q} \in L_{\mathbb{P}}^{r_Q}(\Xi; \mathcal{L}(Y, H))$ and $\hat{\mathbf{q}} \in L_{\mathbb{P}}^{r_{\hat{q}}}(\Xi; H)$, where $r_Q \in [2p, \infty]$ and $r_{\hat{q}} \in [p, \infty]$ fulfill the properties in columns “Estimation 1” and “Estimation 2” of Table 3.1, respectively. Then, the formal adjoint equation (3.18) has the unique solution $\mathbf{T}(u) \in \mathbf{Y}$ (Definition 3.21).*

Proof. This can be proven analogously to Proposition 3.16.

Equation (3.18) can be written as

$$\int_{\Xi} \langle A(\xi)z(\cdot, \xi) + N_y(\mathbf{y}(\cdot, \xi), \xi)z(\cdot, \xi) - \tilde{b}(\xi), \mathbf{v}(\cdot, \xi) \rangle_{Y^*, Y} d\mathbb{P} = 0 \quad \forall \mathbf{v} \in \mathbf{Y}.$$

It is well-posed in \mathbf{Y} since Hölder’s inequality and the growth of φ' (Lemma A.10) yield $\mathbf{A}\mathbf{z} \in L_{\mathbb{P}}^p(\Xi; Y^*)$ with $p \geq p^*$, $\varphi'(\mathbf{y}) \in L_{\mathbb{P}}^{r_f/(p-2)}(\Xi; L^{p/(p-2)}(\Omega))$, $\mathbf{N}'(\mathbf{y})\mathbf{z} \in L_{\mathbb{P}}^{r_{N'z}}(\Xi; Y^*)$ with $r_{N'z} = \frac{pr_f}{r_f + p(p-2)} \geq \frac{p}{p-1} = p^*$ if $r_f < \infty$ and $r_{N'z} = p$ if $r_f = \infty$, and $\tilde{b}(\cdot) = -Q(\cdot)^*(Q(\cdot)S[\cdot](u) - \hat{\mathbf{q}}(\cdot)) \in L_{\mathbb{P}}^p(\Xi; Y^*)$ for $\mathbf{z} \in \mathbf{Y}$ and the assumed ξ -regularities of \mathbf{y} , \mathbf{Q} , and $\hat{\mathbf{q}}$. The integrand vanishes for almost every $\xi \in \Xi$ if $\mathbf{z} = \mathbf{T}(u)$ is inserted. Hence, $\mathbf{T}(u) \in L_{\mathbb{P}}^p(\Xi; H_0^1(\Omega)) = \mathbf{Y}$ is a solution of the equation. Lemma 3.23 provides the additional ξ -regularity.

Since $\mathbf{N}'(\mathbf{y})$ is monotone by Proposition A.4, the operator $\mathbf{A} + \mathbf{N}'(\mathbf{y})$ is strictly monotone on \mathbf{Y} . This proves uniqueness of a solution as in Proposition 3.16. \square

Table 3.1 shows how the ξ -regularity of \mathbf{y} , \mathbf{Q} , and $\hat{\mathbf{q}}$ affects the ξ -regularity of the adjoint state \mathbf{z} as deduced from (3.20). The column “Exponents” depicts the minimum requirement on the exponents for every function in (3.3) to be well-defined. Then, $L_{\mathbb{P}}^{p/(p-1)}$ -regularity of the adjoint state can be shown (Proposition 3.22). In the columns “Estimation 1” and “Estimation 2” it is listed which values of the exponent r_Q and $r_{\hat{q}}$ are required for the adjoint state to belong to $\mathbf{Y} = L_{\mathbb{P}}^p(\Xi, Y)$ if $\mathbf{y} \in L_{\mathbb{P}}^{r_f}(\Xi; Y)$ holds (e. g., by Corollary 3.15). In the case $r_Q = \infty$ and $r_{\hat{q}} \geq r_f$ (column “Example”), we get $\mathbf{z} \in L_{\mathbb{P}}^{r_f}(\Xi; Y)$, i. e., the adjoint state has the same ξ -regularity as the state itself. In the following, we will restrict the discussion to this case, i. e., $\mathbf{Q} \in L_{\mathbb{P}}^{\infty}(\Xi; \mathcal{L}(Y, H))$ and $\hat{\mathbf{q}} \in L_{\mathbb{P}}^{r_f}(\Xi; L^2(\Omega))$, since then the state, the desired state, and the adjoint state all enjoy $L_{\mathbb{P}}^{r_f}$ -regularity.

Function and space	Exponents	Estimation 1	Estimation 2	Example
$\mathbf{y} \in L_{\mathbb{P}}^{r_f}(\Xi; Y)$	$r_f \geq p \in (3, \infty)$	$r_f \in [p, \infty)$	$r_f = \infty$	$r_f \geq p$
$\mathbf{Q} \in L_{\mathbb{P}}^{r_Q}(\Xi; \mathcal{L}(Y, H))$	$r_Q \geq \frac{2p}{p-2}$	$r_Q \geq \frac{2pr_f}{r_f-p}$	$r_Q \geq 2p$	$r_Q = \infty$
$\hat{\mathbf{q}} \in L_{\mathbb{P}}^{r_{\hat{q}}}(\Xi; H)$	$r_{\hat{q}} \geq 2$	$r_{\hat{q}} \geq \frac{2pr_f}{r_f+p}$	$r_{\hat{q}} \geq r_Q$	$r_{\hat{q}} \geq r_f$
$\mathbf{Qy} \in L_{\mathbb{P}}^{r_{Qy}}(\Xi; H)$	$\frac{1}{r_{Qy}} = \frac{1}{r_Q} + \frac{1}{r_f}, r_{Qy} \geq 2$	$r_{Qy} \geq \frac{2pr_f}{r_f+p}$	$r_{Qy} = r_Q$	$r_{Qy} = r_f$
$\mathbf{Qy} - \hat{\mathbf{q}} \in L_{\mathbb{P}}^{\hat{r}}(\Xi; H)$	$\hat{r} = \min\{r_{Qy}, r_{\hat{q}}\} \geq 2$	$\hat{r} \geq \frac{2pr_f}{r_f+p}$	$\hat{r} = r_Q$	$\hat{r} = r_f$
$\mathbf{z} \in L_{\mathbb{P}}^{r_z}(\Xi; Y)$	$\frac{1}{r_z} = \frac{1}{\hat{r}} + \frac{1}{r_Q}, r_z \geq \frac{p}{p-1}$	$r_z \geq p$	$r_z = \frac{r_Q}{2} \geq p$	$r_z = r_f$

Table 3.1.: Overview of the integrability inheritance

Remark 3.25. The adjoint equation can be well-defined in a larger space than \mathbf{Y} . If, e.g., $\mathbf{y} \in L_{\mathbb{P}}^{\infty}(\Xi; H_0^1(\Omega) \cap L^{\infty}(\Omega))$ ($r_f = \infty$), we have $\varphi'(\mathbf{y}) \in L_{\lambda \otimes \mathbb{P}}^{\infty}(\Omega \times \Xi)$ by the continuity of φ' . Note that $L_{\mathbb{P}}^{\infty}(\Xi; L^{\infty}(\Omega)) \subset L_{\lambda \otimes \mathbb{P}}^{\infty}(\Omega \times \Xi)$. It follows that

$$\varphi'(\mathbf{y})\mathbf{v} \in L_{\lambda \otimes \mathbb{P}}^2(\Omega \times \Xi) \cong L_{\mathbb{P}}^2(\Xi; L^2(\Omega))$$

and thus by $Y^* \subset L^2(\Omega) \subset Y$ that $\mathbf{N}'(\mathbf{y}) \in \mathcal{L}(L_{\mathbb{P}}^2(\Xi; Y), L_{\mathbb{P}}^2(\Xi; Y^*))$. The adjoint equation is then even well-posed in $L_{\mathbb{P}}^2(\Xi; Y) \supset \mathbf{Y}$.

From $\mathbf{y} \in L_{\mathbb{P}}^{r_f}(\Xi; H_0^1(\Omega) \cap L^{r_f}(\Omega))$ with $r_f \in [p, \infty)$ we can conclude that

$$\varphi'(\mathbf{y}) \in L_{\mathbb{P}}^{r_f/(p-2)}(\Xi; L^{r_f/(p-2)}(\Omega))$$

and that $\mathbf{N}'(\mathbf{y}) \in \mathcal{L}(L_{\mathbb{P}}^{\tilde{r}}(\Xi; Y), L_{\mathbb{P}}^{\tilde{r}*}(\Xi; Y^*))$ holds with $\tilde{r} = \frac{2r_f}{r_f-p+2} \leq p$ by Proposition 3.1. The adjoint equation is then even well-defined in $L_{\mathbb{P}}^{\tilde{r}}(\Xi; Y) \supset \mathbf{Y}$.

Now we use that the derivative of the parametric, reduced objective function is given by

$$\hat{J}[\xi]'(u) = -B(\xi)^*T[\xi](u) + \gamma(u, \cdot)_U. \quad (3.21)$$

It is measurable as a composition of measurable functions and has the same ξ -regularity as the adjoint state because $B(\cdot) \in L_{\mathbb{P}}^{\infty}(\Xi; \mathcal{L}(U, Y^*))$, especially $\|\hat{J}[\cdot]'(u)\|_{U^*} \in L_{\mathbb{P}}^{p/(p-1)}(\Xi) \subset L_{\mathbb{P}}^1(\Xi)$ by Propositions 3.22 and A.2. Furthermore, the function $\xi \mapsto \hat{J}[\xi](u)$ belongs to $L_{\mathbb{P}}^{p/2}(\Xi) \subset L_{\mathbb{P}}^1(\Xi)$ if $\mathbf{Q} \in L_{\mathbb{P}}^{\infty}(\Xi; \mathcal{L}(Y, H))$. Therefore, we can apply the chain rule and get

$$\frac{d}{du} \hat{\mathbf{J}}(u) = \frac{d}{du} \int_{\Xi} \hat{J}[\xi](u) d\mathbb{P} = \int_{\Xi} \hat{J}[\xi]'(u) d\mathbb{P} = -\mathbf{B}^*\mathbf{T}(u) + \gamma(u, \cdot)_U$$

since $\mathbb{E} : L_{\mathbb{P}}^1(\Xi) \rightarrow \mathbb{R}$ is linear and bounded. The pointwise adjoint state can thus be used for the computation of the gradient of the reduced objective function $\hat{\mathbf{J}}$.

3.4.2. Differentiability of the Control-to-State Mapping

By the argumentation above, we see that $\hat{\mathbf{J}}$ is continuously differentiable, but have not used differentiability properties of the control-to-state mapping \mathbf{S} (Definition 3.11) in a concrete way. As discussed in Corollary 3.15, the control-to-state mapping \mathbf{S} is well-defined. It is even continuously differentiable:

Theorem 3.26. *Under Assumption 3.3 and using the Definitions (3.10), the control-to-state mapping \mathbf{S} (Definition 3.11) is continuously differentiable with the derivative $\mathbf{S}' : U \rightarrow \mathcal{L}(U, \mathbf{Y})$, $[\mathbf{S}'(u)w](\cdot, \xi) := S[\xi]'(u)w$.*

Proof. The pointwise derivative applied to $w \in U$ is the solution of

$$[A(\xi) + N_y(S[\xi](u), \xi)](S[\xi]'(u)w) = -B(\xi)w.$$

In analogy to the discussion of the adjoint state in the proof of Proposition 3.22 we conclude that $S[\xi]'(u)w$ admits the same ξ -regularity as $B(\cdot)$ by

$$\|S[\xi]'(u)w\|_Y \leq \frac{1}{\kappa} \|B(\xi)w\|_{Y^*} \leq \frac{1}{\kappa} \|B(\xi)\|_{\mathcal{L}(U, Y^*)} \|w\|_U \quad (3.22)$$

giving $\|S[\xi]'(u)\|_{\mathcal{L}(U, Y)} \leq \frac{1}{\kappa} \|B(\xi)\|_{\mathcal{L}(U, Y^*)}$. In our case B is constant and thus it holds that $B \in L_{\mathbb{P}}^{\infty}(\Xi; \mathcal{L}(U, Y^*))$ and therefore $S[\xi]'(u) \in L_{\mathbb{P}}^{\infty}(\Xi; \mathcal{L}(U, Y)) \subset \mathbf{Y}$. Linearity w. r. t. w is obvious; the given operator \mathbf{S}' is well-defined. Since $S[\xi]$ is F-differentiable for a. e. ξ by the implicit function theorem, we have that

$$\lim_{\|w\|_U \rightarrow 0} \frac{\|S[\xi](u+w) - S[\xi](u) - S[\xi]'(u)w\|_Y}{\|w\|_U} = 0$$

for a. e. $\xi \in \Xi$.

Next, the following holds true, cf. the proof of [111, Prop. A.11]:

$$\begin{aligned} & \lim_{\|w\|_U \rightarrow 0} \frac{\|\mathbf{S}(u+w) - \mathbf{S}(u) - \mathbf{S}'(u)w\|_Y}{\|w\|_U} \\ &= \lim_{\|w\|_U \rightarrow 0} \frac{\left(\int_{\Xi} \|S[\xi](u+w) - S[\xi](u) - S[\xi]'(u)w\|_Y^p \, d\mathbb{P} \right)^{1/p}}{\|w\|_U} \\ &= \left(\lim_{\|w\|_U \rightarrow 0} \int_{\Xi} \frac{\|S[\xi](u+w) - S[\xi](u) - S[\xi]'(u)w\|_Y^p}{\|w\|_U^p} \, d\mathbb{P} \right)^{1/p}. \\ &= \left(\lim_{\|w\|_U \rightarrow 0} \int_{\Xi} \frac{\| \int_0^1 [S[\xi]'(u+\tau w) - S[\xi]'(u)]w \, d\tau \|_Y^p}{\|w\|_U^p} \, d\mathbb{P} \right)^{1/p} = 0. \end{aligned}$$

This limit is zero by the dominated convergence theorem because the integrand can be bounded by

$$\sup_{\tau \in [0,1]} \|S[\xi]'(u+\tau w) - S[\xi]'(u)\|_{\mathcal{L}(U, Y)}^p \leq \frac{2^p}{\kappa^p} \|B(\xi)\|_{\mathcal{L}(U, Y^*)}^p,$$

where the upper bound is an $L_{\mathbb{P}}^1(\Xi)$ -function w. r. t. ξ for all $w \in U$. We obtain that \mathbf{S} is F-differentiable with the given derivative.

Continuity of \mathbf{S}' can be shown as follows: We know that $S[\xi]'$ is continuous for a. e. $\xi \in \Xi$. Moreover, we have $\|S[\xi]'(u)\|_{\mathcal{L}(U,Y)} \leq \frac{1}{\kappa} \|B(\xi)\|_{\mathcal{L}(U,Y^*)}$ independently of $u \in U$ by (3.22). For $u, \tilde{u}, w \in U$ it holds that

$$\begin{aligned} \|(\mathbf{S}'(u) - \mathbf{S}'(\tilde{u}))w\|_{\mathbf{Y}}^p &= \int_{\Xi} \|(S[\xi]'(u) - S[\xi]'(\tilde{u}))w\|_{\mathbf{Y}}^p d\mathbb{P} \\ &\leq \int_{\Xi} \|S[\xi]'(u) - S[\xi]'(\tilde{u})\|_{\mathcal{L}(U,Y)}^p d\mathbb{P} \cdot \|w\|_U^p, \end{aligned}$$

giving

$$0 \leq \|(\mathbf{S}'(u) - \mathbf{S}'(\tilde{u}))\|_{\mathcal{L}(U,\mathbf{Y})}^p \leq \int_{\Xi} \|S[\xi]'(u) - S[\xi]'(\tilde{u})\|_{\mathcal{L}(U,Y)}^p d\mathbb{P} \longrightarrow 0 \quad \text{as } \tilde{u} \xrightarrow{U} u.$$

This limit is obtained because the integrand converges to zero almost everywhere and can be bounded by $\frac{2^p}{\kappa^p} \|B(\cdot)\|_{\mathcal{L}(U,Y^*)}^p$, which is an $L^1_{\mathbb{P}}(\Xi)$ -function if at least $B \in L^p_{\mathbb{P}}(\Xi; \mathcal{L}(U, Y^*))$. \square

Overall, we have shown that the control-to-state mapping \mathbf{S} is well-defined and continuously differentiable. Therefore, using the chain rule and the boundedness and linearity of the expectation operator, we obtain that the reduced objective function $\hat{\mathbf{J}}$ is well-defined and continuously differentiable with the given derivative. Other risk measures can also be included as long as they are at least once continuously differentiable, cf. Section 9.2.

3.4.3. Second Derivatives

We can follow a similar strategy as in Subsection 3.4.1 to derive an expression for the second derivative of the reduced objective function applied to a given direction $s \in U$. Following [60, Sec. 1.6.5], we see that

$$\nabla^2 \hat{\mathbf{J}}[\xi](u)s = \iota B(\xi)^* h(\xi) + \gamma s \in U, \quad (3.23)$$

where $h(\xi) \in Y$ solves

$$[A(\xi) + N_y(S[\xi](u), \xi)]^* h(\xi) = Q(\xi)^* Q(\xi) \delta(\xi) + \langle T[\xi](u), [N_{yy}(S[\xi](u), \xi) \delta(\xi)] \cdot \rangle_{Y, Y^*} \quad (3.24)$$

with $\delta(\xi) = S[\xi]'(u)s \in Y$ solving

$$[A(\xi) + N_y(S[\xi](u), \xi)] \delta(\xi) = B(\xi)s. \quad (3.25)$$

The mapping $\iota : U^* \rightarrow U$ is the Riesz representation operator. It is used that N is twice continuously differentiable w. r. t. y . Both $\delta(\xi)$ and $h(\xi)$ are unique as discussed in Subsection 3.4.1. Since $A(\xi)$ and $N_y(S[\xi](u), \xi)$ are self-adjoint, the results for the adjoint equation are also applicable to the linearized state equation. As in Proposition 3.22 we have

$$\|\delta(\xi)\|_Y \leq \frac{1}{\kappa} \|B(\xi)\|_{\mathcal{L}(U,Y^*)} \|s\|_U \leq \|D\|_{\mathcal{L}(U, L^2(\Omega))} \|s\|_U$$

and therefore $\delta \in L^\infty_{\mathbb{P}}(\Xi; Y)$. In analogy to (3.20), it follows that

$$\|h(\xi)\|_Y \leq \frac{1}{\kappa} (\|Q(\xi)\|_{\mathcal{L}(Y,H)} \|Q(\xi) \delta(\xi)\|_H + \|T[\xi](u)\|_Y \|N_{yy}(S[\xi](u), \xi) \delta(\xi)\|_{\mathcal{L}(Y, Y^*)}) \quad (3.26)$$

and thus $h \in L_{\mathbb{P}}^{r_h}(\Xi; Y)$ with $r_h = \min\{\frac{r_Q}{2}, (\frac{1}{r_z} + \frac{p-3}{r_f})^{-1}\} \geq \frac{p}{p-2} > 1$ with the prerequisites and notation from Proposition 3.24, especially $r_z \geq p$. We obtain $\nabla^2 \hat{J}[\cdot](u)_s \in L_{\mathbb{P}}^{r_h}(\Xi; U) \subset L_{\mathbb{P}}^1(\Xi; U)$.

Again, the parametrized linearized state equation and adjoint equation can be written in weak form also w. r. t. the parameters. As in the case of the first derivative, the regularity estimates and measurability yield

$$\nabla^2 \hat{J}(u)_s = \int_{\Xi} \nabla^2 \hat{J}[\xi](u)_s d\mathbb{P}.$$

4. An Inexact Trust-Region Algorithm for the Solution of Optimal Control Problems with Control Constraints

We present and discuss an inexact and projection-based trust-region algorithm, which can be used for the solution of optimal control problems of the form

$$\min_{(y,u) \in Y \times U} J(y, u) \quad \text{s. t.} \quad E(y, u) = 0, \quad u \in U_{\text{ad}}. \quad (4.1)$$

We want to use this algorithm later to solve problem (3.3) adaptively and to control all arising errors caused by discretization, inexact solution of linear systems and tensor truncation during this procedure.

One possible solution algorithm would be an inexact, trust-region-based SQP algorithm such as [118], which extends [119, 57] to the case of control constraints handled by projections onto U_{ad} . Here it is necessary to minimize the residual of the linearized state equation approximately. If the space Z of the residual is not a Hilbert space as in the example from Section 3.2, this is possibly a nonsmooth problem so that this approach is not preferable in our setting. To circumvent this difficulty, we can work with the reduced problem (4.3) and use an extension of the algorithm presented in [70], which is based on [69]. The algorithm in [70] can handle inexact gradients as well as inexact objective function evaluations. Earlier versions of inexact trust-region algorithms formulated in Hilbert space [26, 27] have the disadvantage that the relative errors in the gradient and the objective function evaluation have to be bounded by fixed, prescribed constants of a certain magnitude, whereas [70] allows for fixed, but possibly unknown and arbitrarily large constants. This fits very well to PDE applications, where error estimates often contain unknown multiplicative constants. We extend [70, Algorithm 4.1] to the case of control constraints using an inexact projection (inspired by [118]) and present our version in the following.

Assumption 4.1. We make the following assumptions on problem (4.1):

- The space U is a Hilbert space; Y, Z are reflexive Banach spaces.
- The feasible control set $U_{\text{ad}} \subset U$ is nonempty, convex and closed.
- For each control $u \in \tilde{U} \subset U$, where $\tilde{U} \supset U_{\text{ad}}$ is an open set, there exists a unique state $y = S(u)$ fulfilling $E(S(u), u) = 0$. $S : \tilde{U} \rightarrow Y$ is the control-to-state mapping.
- The control-to-state mapping $S : \tilde{U} \rightarrow Y$ and the objective function $J : Y \times \tilde{U} \rightarrow \mathbb{R}$ are such that the reduced objective function $\hat{J} : \tilde{U} \rightarrow \mathbb{R}$ in (4.3) is continuously differentiable.

If Assumption 4.1 is fulfilled, Problem (4.1) can be reduced to the control:

$$\min_{u \in U} \hat{J}(u) := J(S(u), u) \quad \text{s. t.} \quad u \in U_{\text{ad}}. \quad (4.2)$$

This motivates that we describe an inexact trust-region algorithm for the solution of general problems of the form

$$\min_{u \in U} \hat{J}(u) \quad \text{s. t.} \quad u \in U_{\text{ad}}. \quad (4.3)$$

in this chapter and prove its convergence. Later, in Chapter 5, we return to the original problem (4.2), where the objective function \hat{J} is of reduced form.

Using parts of Assumption 4.1, we make the following assumption for problem (4.3):

Assumption 4.2.

- U is a Hilbert space.
- The feasible set $U_{\text{ad}} \subset U$ is nonempty, closed and convex.
- The objective function $\hat{J} : \tilde{U} \rightarrow \mathbb{R}$ is continuously differentiable on an open set \tilde{U} , $U_{\text{ad}} \subset \tilde{U} \subset U$, and bounded from below on U_{ad} . The Fréchet approximation condition holds uniformly on every level set, i. e.,

$$\sup_{u \in U_{\text{ad}} : \hat{J}(u) \leq \hat{J}(\tilde{u})} |\hat{J}(u + s) - \hat{J}(u) - (\nabla \hat{J}(u), s)_U| = o(\|s\|_U) \quad (s \rightarrow 0), \quad (4.4)$$

for every $\tilde{u} \in U_{\text{ad}}$.

4.1. Formulation of the Algorithm

For a comprehensive introduction to trust-region algorithms we refer to [31]. In the algorithm presented here, we use a typically, but not necessarily quadratic model $m_k(s)$ of $\hat{J}(u^k + s) - \hat{J}(u^k)$ with the current control u^k for the computation of the current step $s^k \in U$. The step computation approximately solves

$$\min_{s \in U} m_k(s) \quad \text{s. t.} \quad u^k + s \in U_{\text{ad}}, \quad \|s\|_U \leq \Delta_k \quad (4.5)$$

with the current trust region radius $\Delta_k > 0$. For the acceptance of the step, we allow for inexact evaluations of \hat{J} by using an approximation \hat{J}_k instead of \hat{J} . We define the actual, computed (as introduced in [26]), and predicted reduction, respectively, as

$$\text{ared}_k := \hat{J}(u^k) - \hat{J}(u^k + s^k), \quad (4.6a)$$

$$\text{cred}_k := \hat{J}_k(u^k) - \hat{J}_k(u^k + s^k), \quad (4.6b)$$

$$\text{pred}_k := m_k(0) - m_k(s^k). \quad (4.6c)$$

Furthermore, we define a criticality measure for the original problem (4.3), namely

$$\chi : \tilde{U} \rightarrow \mathbb{R}_{\geq 0}, \quad \chi(u) := \|u - P_{U_{\text{ad}}}(u - \tau \nabla \hat{J}(u))\|_U \quad (4.7)$$

with a fixed parameter $\tau > 0$ and the projection $P_{U_{\text{ad}}}$ onto the feasible set. The function χ is continuous, and $\chi(\bar{u}) = 0$ holds if and only if \bar{u} is a first order critical point for (4.3). In addition, a criticality measure for problem (4.5) without the trust-region constraint is defined:

$$\tilde{\chi}_k : U \rightarrow \mathbb{R}_{\geq 0}, \quad \tilde{\chi}_k(s) := \|u^k + s - P_{U_{\text{ad}}}(u^k + s - \tau \nabla m_k(s))\|_U.$$

The condition $\tilde{\chi}_k(\bar{s}) = 0$ holds if and only if \bar{s} is first order critical for the problem

$$\min_{s \in U} m_k(s) \quad \text{s. t.} \quad u^k + s \in U_{\text{ad}}. \quad (4.8)$$

Typically, the projection is also not computed exactly. Therefore, we introduce the *approximate criticality measure* for problem (4.8),

$$\chi_k : U \rightarrow \mathbb{R}_{\geq 0}, \quad \chi_k(s) := \|u^k + s - \hat{P}_{U_{\text{ad}}}(u^k + s - \tau \nabla m_k(s))\|_U, \quad (4.9)$$

now with an *approximate* projection $\hat{P}_{U_{\text{ad}}} : U \rightarrow U_{\text{ad}}$ onto the feasible set. This can be *any* mapping approximating the exact projection. Especially, it does not have to fulfill the variational inequality defining the projection on some discrete subspace $U \subset U$. In contrast to that, the properties of the discrete projection are used in [118, Lem. 5.3, Lem. 5.5] when proving the Cauchy decrease condition although an approximate version is used in the final implementation there [118, Sec. 5.2].

Remark 4.3. Consider the case that $U = L^2(\Omega)$. The gradient of the reduced objective function at $u \in U$ (discrete subspace of U) is, e. g., of the form $\nabla \hat{J}(u) = -B^*z + \gamma u$ with $B^* \equiv \iota : H_0^1(\Omega) \hookrightarrow L^2(\Omega)$ (canonical embedding). An approximation of it is obtained using an approximate adjoint state $z \in Y$, computed, e. g., by a finite element method with piecewise linear, continuous ansatz functions. The function $B^*z \in U$ then belongs to in the finite element space Y . To avoid further approximations, it makes sense that the current control u belongs to the space U of linear FE functions (not necessarily with zero boundary data). Then, the approximate gradient and any control obtained by a linesearch in its direction are also in U . This means that we do not explicitly discretize the control space, cf. the variational discretization concept of optimal control problems [59], but its discretization is implied by the one of the adjoint state.

In general, one could equip the discrete space U of linear FE functions with the inner product induced by the lumped mass matrix. The obtained discrete projection onto a box-constrained feasible set $U_{\text{ad}} \cap U$ can then be computed node-wise. As described, this is not possible in our case since the algorithm uses the exact inner product of U . The resulting discrete projection would have to be computed by solving a high-dimensional quadratic program with box-constraints. In order to use the node-wise projection instead as in [118, Sec. 5.2], we have to allow for an inexact projection and cannot use any projection properties of it.

Alternatively, one could use piecewise constant functions for the control. Then, the exact L^2 -projection onto a box (with constant bounds) can be computed simply, but an additional error in the approximate gradient occurs since the operator B^* cannot be applied exactly then.

To ensure global convergence of the algorithm we have to control the following quantities:

- The inexactness of the model gradient:

$$\|\nabla m_k(0) - \nabla \hat{J}(u^k)\|_U \leq \varrho_g(\Delta_k), \quad (4.10)$$

where $\varrho_g : \mathbb{R}_{>0} \rightarrow \mathbb{R}_{\geq 0}$ is a function satisfying $\lim_{t \rightarrow 0^+} \varrho_g(t) = 0$, e. g., $\varrho_g(t) = \mathbf{c}_g t$, $\mathbf{c}_g > 0$.

- The inexactness of the approximate criticality measure:

$$|\chi_k(0) - \chi(u^k)| \leq \varrho_c(\chi_k(0)), \quad (4.11)$$

where $\varrho_c : \mathbb{R}_{\geq 0} \rightarrow \mathbb{R}_{\geq 0}$ is a function satisfying $\lim_{t \rightarrow 0^+} \varrho_c(t) = 0$ and $\varrho_c(0) = 0$, e. g., $\varrho_c(t) = \mathbf{c}_c t$, $\mathbf{c}_c > 0$. Note that we do not have to use the exact criticality measure for the model problem here since it is not required in the convergence proof. It will only be used for the control of the inexactness.

- The quality of the computed reduction:

$$|\text{ared}_k - \text{cred}_k| \leq \varrho_r(\eta_3 \min\{\text{pred}_k, \mathbf{r}_k\}), \quad (4.12)$$

with $\eta_3 < \min\{\eta_1, 1 - \eta_2\}$, where $0 < \eta_1 < \eta_2 < 1$ are a priori chosen parameters for assessing the quality of the model, and with a forcing sequence $(\mathbf{r}_k)_{k \in \mathbb{N}_0} \subset \mathbb{R}_{>0}$ fulfilling $\lim_{k \rightarrow \infty} \mathbf{r}_k = 0$ and a function $\varrho_r : \mathbb{R}_{>0} \rightarrow \mathbb{R}_{\geq 0}$ fulfilling $\varrho_r(t) \leq t$ for all $t \in (0, \bar{t}]$ with some fixed $\bar{t} > 0$, e. g., $\varrho_r(t) = \mathbf{c}_r t^{\mathbf{e}_r}$, $\mathbf{c}_r > 0$, $\mathbf{e}_r > 1$. Note that $\varrho_r(t) = t$ would also be possible, but then it is not sufficient to know the error in (4.12) up to an unknown, multiplicative constant.

A trial step $s^k \in U_{\text{ad}} - u^k$, $\|s^k\|_U \leq \Delta_k$, has to fulfill the decrease condition

$$\text{pred}_k = m_k(0) - m_k(s^k) \geq \varrho_{t1}(\chi_k(0)) \cdot \min\{\varrho_{t2}(\chi_k(0)), \Delta_k\} \quad (4.13)$$

with monotonically increasing functions $\varrho_{t1}, \varrho_{t2} : \mathbb{R}_{>0} \rightarrow \mathbb{R}_{>0}$. These functions must be chosen such that (4.13) is satisfiable by, e. g., the generalized Cauchy point, see Section 4.3. A possible example is $\varrho_{t1}(t) = \mathbf{c}_{t1} t$, $\varrho_{t2}(t) = \mathbf{c}_{t2} t$ with $\mathbf{c}_{t1}, \mathbf{c}_{t2} > 0$.

The complete method is listed in Algorithm 1. Note that all iterates u^k belong to U_{ad} since $s^k \in U_{\text{ad}} - u^k$ is required for all trial steps. Therefore, it is sufficient to assume differentiability of \hat{J} only in an open neighborhood \tilde{U} of U_{ad} , see Assumption 4.2.

Remark 4.4. In case of an unsuccessful step, the formulation of Algorithm 1 allows to choose $\Delta_{k+1} \in (0, \nu_1 \|s^k\|_U] \subset (0, \nu_1 \Delta_k]$ if $\|s^k\|_U > 0$, which is a suitable strategy to avoid that s^k is feasible for the trust-region subproblem in the next iteration.

4.2. Convergence Proof

Provided all conditions in Algorithm 1 can be satisfied, we prove its convergence. This means that we assume for now that an adequate model, approximate projection, trial step, and inexact objective function exist in each iteration. We prove this in Section 4.3.

Algorithm 1: Inexact Trust-Region Method for Solving Problem (4.3)

Input: Initial iterate $u^0 \in U_{\text{ad}}$

Parameters : $\tau > 0$, error control functions $\varrho_c, \varrho_g, \varrho_{t1}, \varrho_{t2}, \varrho_r$,
 forcing sequence $(\mathbf{r}_k)_{k \in \mathbb{N}_0} \subset \mathbb{R}_{>0}$, $\lim_{k \rightarrow \infty} \mathbf{r}_k = 0$.
 $\Delta_{\max} \in (0, \infty]$, $\Delta_0 \in \mathbb{R}_{>0}$ s. t. $\Delta_0 \leq \Delta_{\max}$,
 $0 < \eta_1 < \eta_2 < 1$ and $0 < \eta_3 \leq \min\{\eta_1, 1 - \eta_2\}$,
 $0 < \nu_1 < 1 \leq \nu_2 < \nu_3$.

Output: Sequences $(u^k)_{k \in \mathbb{N}_0} \subset U_{\text{ad}}$, $(\Delta_k)_{k \in \mathbb{N}_0} \subset \mathbb{R}_{>0}$, $(\chi_k(0))_{k \in \mathbb{N}_0} \subset \mathbb{R}_{\geq 0}$

for $k := 0, 1, 2, \dots$ **do**

Choose a model $m_k : U \rightarrow \mathbb{R}$ and an approximate projection $\hat{P}_{U_{\text{ad}}} : U \rightarrow U_{\text{ad}}$ such that (4.10) and (4.11) hold. Compute $\chi_k(0)$ using m_k and $\hat{P}_{U_{\text{ad}}}$.

if $\chi_k(0) = 0$, **then**

Set $u^\ell := u^k$, $\Delta_\ell := \Delta_k$, and $\chi_\ell(0) = 0$ for all $\ell \geq k + 1$ and STOP.

end

Compute a trial step $s^k \in U_{\text{ad}} - u^k$, $\|s^k\|_U \leq \Delta_k$, fulfilling (4.13) with the computed $\chi_k(0)$, see Section 4.3.

Compute pred_k by (4.6c) and cred_k by (4.6b) with \hat{J}_k such that (4.12) holds.

if $\frac{\text{cred}_k}{\text{pred}_k} < \eta_1$ (*unsuccessful step*), **then**

$u^{k+1} := u^k$, choose $\Delta_{k+1} \in (0, \nu_1 \Delta_k]$.

else if $\frac{\text{cred}_k}{\text{pred}_k} \in [\eta_1, \eta_2)$ (*successful step*), **then**

$u^{k+1} := u^k + s^k$, choose $\Delta_{k+1} \in [\nu_1 \Delta_k, \min\{\nu_2 \Delta_k, \Delta_{\max}\}]$.

else if $\frac{\text{cred}_k}{\text{pred}_k} \geq \eta_2$ (*very successful step*), **then**

$u^{k+1} := u^k + s^k$, choose $\Delta_{k+1} \in [\min\{\nu_2 \Delta_k, \Delta_{\max}\}, \min\{\nu_3 \Delta_k, \Delta_{\max}\}]$.

end

end

We require the following assumption in addition to Assumption 4.2 to prove the convergence result given in Theorem 4.6.

Assumption 4.5. Each model $m_k : U \rightarrow \mathbb{R}$ is continuously differentiable. The Fréchet approximation condition holds uniformly over all models:

$$\sup_{k \in \mathbb{N}_0} |m_k(s) - m_k(0) - (\nabla m_k(0), s)_U| = o(\|s\|) \quad (s \rightarrow 0). \quad (4.14)$$

The model gradients are Lipschitz continuous on the sets of feasible search directions, i. e.,

$$\|\nabla m_k(s) - \nabla m_k(\hat{s})\|_U \leq \mathbf{c}_{m_k} \|s - \hat{s}\|_U$$

holds for all $s, \hat{s} \in U_{\text{ad}} - u^k$, $\|s\|_U \leq \Delta_k$, $\|\hat{s}\|_U \leq \Delta_k$ with some $\mathbf{c}_{m_k} > 0$. The Lipschitz constants shall be bounded uniformly: $\mathbf{c}_{m_k} \leq \mathbf{c}_m$ for some $\mathbf{c}_m > 0$ and all $k \in \mathbb{N}_0$.

Theorem 4.6. Let Assumptions 4.2 and 4.5 hold and let the sequence $(u^k)_{k \in \mathbb{N}_0} \subset U$ be generated by Algorithm 1. Then,

$$\liminf_{k \rightarrow \infty} \chi(u^k) = 0$$

holds with the criticality measure χ defined in (4.7).

4. An Inexact Trust-Region Algorithm

We apply the following two lemmas to prove Theorem 4.6:

Lemma 4.7. *Let Assumptions 4.2 and 4.5 hold and let the sequence of inexact criticality measures (as defined in (4.9)) generated by Algorithm 1 satisfy*

$$\chi_k(0) \geq \varepsilon > 0 \quad \text{for all } k \geq K_1 \in \mathbb{N}_0 \quad (4.15)$$

for some fixed $\varepsilon > 0$.

Then, $\lim_{k \rightarrow \infty} \Delta_k = 0$ holds for the sequence of corresponding trust-region radii.

Proof. First observe that the termination criterion “ $\chi_k(0) = 0$ ” in the algorithm is not met for any $k \in \mathbb{N}_0$ by assumption, because $\chi_k(0) = 0$ for some $k \in \mathbb{N}_0$ would yield $\chi_\ell(0) = 0$ for all $\ell \geq k$, which contradicts (4.15). Moreover, $\text{pred}_k > 0$ holds for all $k \in \mathbb{N}_0$ due to (4.13), the positivity property of ϱ_{t1} and ϱ_{t2} , and $\chi_k(0) > 0$, $\Delta_k > 0$.

Due to $\lim_{k \rightarrow \infty} \tau_k = 0$ and $\varrho_r(t) \leq t$ for small enough t , it holds that

$$\varrho_r(\eta_3 \min\{\text{pred}_k, \tau_k\}) \leq \eta_3 \text{pred}_k \quad \text{for all } k \geq K_2 \in \mathbb{N}_0$$

with some $K_2 \geq K_1$. By (4.12) we thus get

$$|\text{ared}_k - \text{cred}_k| \leq \eta_3 \text{pred}_k \quad \text{for all } k \geq K_2.$$

This implies

$$\text{ared}_k = \text{cred}_k + \text{ared}_k - \text{cred}_k \geq \text{cred}_k - \eta_3 \text{pred}_k = \left(\frac{\text{cred}_k}{\text{pred}_k} - \eta_3\right) \text{pred}_k \quad (4.16)$$

for all $k \geq K_2$. This is well-defined due to $\text{pred}_k > 0$.

Now we show that $\sum_{k=0}^{\infty} \Delta_k < \infty$ follows from (4.15).

For all unsuccessful steps $k \in \mathcal{I}_u \subset \mathbb{N}_0$ we have $\Delta_{k+1} \leq \nu_1 \Delta_k$ with the parameter $\nu_1 \in (0, 1)$. Thus, if there are only finitely many (very) successful steps, the sequence $(\Delta_k)_{k \in \mathbb{N}_0}$ is summable since then $\Delta_k \leq \nu_1^{k-K_3} \Delta_{K_3}$ for all $k \geq K_3$ for some $K_3 \in \mathbb{N}_0$. In the following, we consider the case of infinitely many (very) successful steps.

For a (very) successful step $k \in \mathcal{I}_s = \mathbb{N}_0 \setminus \mathcal{I}_u$, i. e.,

$$\frac{\text{cred}_k}{\text{pred}_k} \geq \eta_1 \quad \text{and} \quad \Delta_{k+1} \leq \min\{\nu_3 \Delta_k, \Delta_{\max}\} \leq \nu_3 \Delta_k,$$

we deduce from (4.16) and (4.13):

$$\begin{aligned} \text{ared}_k &\geq \left(\frac{\text{cred}_k}{\text{pred}_k} - \eta_3\right) \text{pred}_k \geq (\eta_1 - \eta_3) \text{pred}_k \\ &\geq (\eta_1 - \eta_3) \cdot \varrho_{t1}(\chi_k(0)) \cdot \min\{\varrho_{t2}(\chi_k(0)), \Delta_k\} \\ &\geq (\eta_1 - \eta_3) \cdot \varrho_{t1}(\varepsilon) \cdot \min\{\varrho_{t2}(\varepsilon), \Delta_k\} > 0, \end{aligned} \quad (4.17)$$

for all $k \in \mathcal{I}_s$, $k \geq K_2$, where we used in the last estimate that ϱ_{t1} and ϱ_{t2} are increasing. Since the sequence $(\hat{J}(u^k))_{k \in \mathcal{I}_s}$ is bounded from below by Assumption 4.2, we get

$$0 \leq \sum_{k \in \mathcal{I}_s, k \geq K_2} \text{ared}_k = \sum_{k \in \mathcal{I}_s, k \geq K_2} (\hat{J}(u^k) - \hat{J}(u^{k+1})) = \sum_{k=K_2}^{\infty} (\hat{J}(u^k) - \hat{J}(u^{k+1})) < \infty \quad (4.18)$$

using that $u^{k+1} = u^k + s^k$ in the case of a (very) successful step and $u^{k+1} = u^k$ in the case of an unsuccessful step. Due to $\eta_1 > \eta_3$, $\varrho_{t1}(\varepsilon) > 0$, $\varrho_{t2}(\varepsilon) > 0$ and $\lim_{\mathcal{I}_s \ni k \rightarrow \infty} \text{ared}_k = 0$ (by (4.18)), it follows from (4.17) that $\Delta_k \leq \frac{\text{ared}_k}{(\eta_1 - \eta_3) \cdot \varrho_{t1}(\varepsilon)}$ for all $k \in \mathcal{I}_s$, $k \geq K_4$, with some sufficiently large $K_4 \in \mathbb{N}_0$, $K_4 \geq K_2$, and thus

$$0 \leq \sum_{k \in \mathcal{I}_s} \Delta_k < \infty \quad (4.19)$$

by (4.18).

Now we consider two (very) successful steps $\tilde{k}, \hat{k} \in \mathcal{I}_s$, $\hat{k} \geq \tilde{k} + 2$ with only unsuccessful steps $k \in \{\tilde{k} + 1, \dots, \hat{k} - 1\}$ in between. Hence, we have $\Delta_k \leq \nu_3 \nu_1^{k - \tilde{k} - 1} \Delta_{\tilde{k}}$ for $\tilde{k} + 1 \leq k \leq \hat{k} - 1$ and thus (using the geometric series with $\nu_1 \in (0, 1)$)

$$\Sigma(\tilde{k}) := \sum_{k=\tilde{k}}^{\hat{k}-1} \Delta_k \leq \Delta_{\tilde{k}} \left(1 + \nu_3 \sum_{\ell=0}^{\hat{k}-\tilde{k}-2} \nu_1^\ell \right) \leq \Delta_{\tilde{k}} \left(1 + \frac{\nu_3}{1-\nu_1} \right).$$

Additionally, for $\tilde{k} \in \mathcal{I}_s$ such that $\tilde{k} + 1 \in \mathcal{I}_s$ we set $\Sigma(\tilde{k}) = \Delta_{\tilde{k}}$ and in the case $0 \notin \mathcal{I}_s$ we set and estimate

$$\Sigma(0) := \sum_{k=0}^{\hat{k}-1} \Delta_k \leq \Delta_0 \left(\sum_{\ell=0}^{\hat{k}-1} \nu_1^\ell \right) \leq \Delta_0 \cdot \frac{1}{1-\nu_1},$$

where $\hat{k} = \min \mathcal{I}_s$. Therefore,

$$0 \leq \sum_{k=0}^{\infty} \Delta_k = \sum_{\tilde{k} \in \mathcal{I}_s \cup \{0\}} \Sigma(\tilde{k}) \leq \left(1 + \frac{\nu_3}{1-\nu_1} \right) \cdot \sum_{\tilde{k} \in \mathcal{I}_s \cup \{0\}} \Delta_{\tilde{k}} < \infty$$

follows with $\Sigma(\tilde{k}) \leq \Delta_{\tilde{k}} \left(1 + \frac{\nu_3}{1-\nu_1} \right)$ for all $\tilde{k} \in \mathcal{I}_s \cup \{0\}$ (using $\nu_3 \geq 1$) and (4.19). We see that $(\Delta_k)_{k \in \mathbb{N}_0}$ is summable and thus $\lim_{k \rightarrow \infty} \Delta_k = 0$. \square

Lemma 4.8. *Under Assumptions 4.2 and 4.5*

$$\liminf_{k \rightarrow \infty} \chi_k(0) = 0 \quad (4.20)$$

holds true for every sequence generated by Algorithm 1, where the inexact criticality measure χ_k is defined as in (4.9).

Proof. For a proof by contradiction, assume that (4.20) is false, giving that (4.15) is true for some fixed $\varepsilon > 0$. By Lemma 4.7 we have that $\lim_{k \rightarrow \infty} \Delta_k = 0$ and thus $\lim_{k \rightarrow \infty} \|s^k\|_U = 0$.

In analogy to (4.16), we can estimate

$$\text{cred}_k \geq \text{ared}_k - |\text{ared}_k - \text{cred}_k| \geq \text{ared}_k - \eta_3 \text{pred}_k \quad (4.21)$$

for all $k \geq K_2 \geq K_1$.

As in (4.17) we infer

$$\text{pred}_k \geq \varrho_{t1}(\varepsilon) \cdot \Delta_k$$

for all $k \geq K_5$ with some sufficiently large $K_5 \in \mathbb{N}_0$, $K_5 \geq K_2$, due to (4.13) and $\lim_{k \rightarrow \infty} \Delta_k = 0$, i. e., $\Delta_k \leq \varrho_{t2}(\varepsilon)$ for all $k \geq K_5$ because $\varrho_{t2}(\varepsilon) > 0$. Thus, we can bound

$$|o(\Delta_k)| \leq (1 - \eta_3 - \eta_2)\text{pred}_k \quad (4.22)$$

for all $k \geq K_6$ with $K_6 \geq K_5 \geq K_2$ sufficiently large, since $(1 - \eta_3 - \eta_2) > 0$. Using this and the bounds indicated below, we estimate for $k \geq K_6$ and using $s_k \rightarrow 0$:

$$\begin{aligned} \text{cred}_k &\stackrel{(4.21)}{\geq} \text{ared}_k - \eta_3 \text{pred}_k \stackrel{(4.4)}{\geq} -(\nabla \hat{J}(u^k), s^k)_U - \eta_3 \text{pred}_k - |o(\|s^k\|)| \\ &\stackrel{(4.14)}{\geq} \text{pred}_k + (\nabla m_k(0), s^k)_U - (\nabla J(u^k), s^k)_U - \eta_3 \text{pred}_k - |o(\|s^k\|)| \\ &\geq (1 - \eta_3)\text{pred}_k - \|\nabla m_k(0) - \nabla \hat{J}(u^k)\|_U \|s^k\|_U - |o(\|s^k\|)| \\ &\stackrel{(4.10),}{\geq} \underset{\|s^k\|_U \leq \Delta_k}{(1 - \eta_3)\text{pred}_k - \varrho_g(\Delta_k)\Delta_k - |o(\Delta_k)|} \stackrel{\varrho_g(t \rightarrow 0)}{\underset{(t \rightarrow 0^+)}{\geq}} (1 - \eta_3)\text{pred}_k - |o(\Delta_k)| \\ &\stackrel{(4.22)}{\geq} (1 - \eta_3)\text{pred}_k - (1 - \eta_3 - \eta_2)\text{pred}_k = \eta_2 \text{pred}_k \end{aligned}$$

Note that $u^k \in \{u \in U_{\text{ad}} : \hat{J}(u) \leq \hat{J}(u^{K_2})\}$ holds for all $k \geq K_2$ due to (4.17) so that (4.4) is applicable. In fact, the objective function values $(\hat{J}(u^k))_{k \geq K_2}$ are non-increasing since (4.17) holds for (very) successful steps and the function values do not change for unsuccessful steps. Using $\text{pred}_k > 0$ as in the proof of Lemma 4.7, it follows that $\frac{\text{cred}_k}{\text{pred}_k} \geq \eta_2$ for all $k \geq K_6$, i. e., all steps $k \geq K_6$ are successful giving $\Delta_{k+1} \geq \min\{\nu_2 \Delta_k, \Delta_{\max}\} \geq \Delta_k > 0$ due to $\nu_2 \geq 1$. This contradicts $\lim_{k \rightarrow \infty} \Delta_k = 0$, proving (4.20). \square

Using Lemma 4.8, the proof of Theorem 4.6 is very short:

Proof of Theorem 4.6. Due to (4.11) we have

$$\chi(u^k) \leq \chi_k(0) + |\chi_k(0) - \chi(u^k)| \leq \chi_k(0) + \varrho_c(\chi_k(0)).$$

This is also true if the algorithm is stopped due to $\chi_k(0) = 0$, because then $\chi(u^k) = 0$ follows from (4.11). Therefore, $\chi(u^\ell) = \chi(u^k) = 0 = \chi_\ell(0)$ for all $\ell \geq k$. The bound on $\chi(u^k)$ and $\lim_{t \rightarrow 0^+} \varrho_c(t) = 0$ as well as $\varrho_c(0) = 0$ show

$$0 \leq \liminf_{k \rightarrow \infty} \chi(u^k) \leq \liminf_{k \rightarrow \infty} \chi_k(0) + \varrho_c(\chi_k(0)) = 0.$$

\square

4.3. Satisfying the Conditions Required by the Algorithm

We show that Algorithm 1 is realizable, i. e., that all conditions can be satisfied under certain assumptions. This includes the computation of a generalized Cauchy point satisfying (4.13). In addition to Assumptions 4.2 and 4.5 we require:

Assumption 4.9.

- For every iterate $u^k \in U_{\text{ad}}$, one can compute $\nabla m_k(0) \in U$ such that

$$\|\nabla m_k(0) - \nabla \hat{J}(u^k)\|_U \leq c_g \varepsilon_g \quad (4.23)$$

holds with some fixed but possibly unknown constant $c_g > 0$ and a still to be specified $\varepsilon_g > 0$.

- For every $w^k(t) := u^k - t\nabla m_k(0) \in U$ with the values of $t > 0$ specified later, one can compute $\hat{P}_{U_{\text{ad}}}(w^k(t)) \in U_{\text{ad}}$ such that

$$\|P_{U_{\text{ad}}}(w^k(t)) - \hat{P}_{U_{\text{ad}}}(w^k(t))\|_U \leq c_p \varepsilon_p \quad (4.24)$$

holds with some fixed but possibly unknown constant $c_p > 0$ and a still to be specified $\varepsilon_p > 0$.

- For every iterate $u^k \in U_{\text{ad}}$ and every trial step $s^k \in U$, one can compute $\hat{J}_k(u) \in \mathbb{R}$ such that

$$|\hat{J}_k(u) - \hat{J}(u)| \leq c_o \varepsilon_o$$

holds with some fixed but possibly unknown constant $c_o > 0$ and a still to be specified $\varepsilon_o > 0$, where $u \in \{u^k, u^k + s^k\}$.

Remark 4.10. It is necessary to require (4.24) not only for $t = \tau$ because an Armijo type linesearch, which tests different values of t , is employed for computing the generalized Cauchy point later.

Model and Approximate Projection

To ensure (4.10) and (4.11), the following result is helpful:

Proposition 4.11. *Let $\mathbf{c}_s \in [0, 1]$ be given and let the following be fulfilled:*

$$\|\nabla m_k(0) - \nabla \hat{J}(u^k)\|_U \leq \min\left\{\frac{1-\mathbf{c}_s}{\tau} \varrho_c(\chi_k(0)), \varrho_g(\Delta_k)\right\}, \quad (4.25)$$

$$\|[P_{U_{\text{ad}}} - \hat{P}_{U_{\text{ad}}}] (u^k - \tau \nabla m_k(0))\|_U \leq \mathbf{c}_s \varrho_c(\chi_k(0)). \quad (4.26)$$

Then, (4.10) and (4.11) hold, where χ and χ_k are defined as in (4.7) and (4.9), respectively.

Proof. The estimate (4.10) follows directly from (4.25). Using the definitions (4.7) and (4.9), we estimate

$$\begin{aligned} |\chi_k(0) - \chi(u^k)| &\leq |\chi_k(0) - \tilde{\chi}_k(0)| + |\tilde{\chi}_k(0) - \chi(u^k)| \\ &\leq \|P_{U_{\text{ad}}}(u^k - \tau \nabla m_k(0)) - \hat{P}_{U_{\text{ad}}}(u^k - \tau \nabla m_k(0))\|_U \\ &\quad + \|P_{U_{\text{ad}}}(u^k - \tau \nabla \hat{J}(u^k)) - P_{U_{\text{ad}}}(u^k - \tau \nabla m_k(0))\|_U \\ &\leq \|[P_{U_{\text{ad}}} - \hat{P}_{U_{\text{ad}}}] (u^k - \tau \nabla m_k(0))\|_U + \tau \|\nabla m_k(0) - \nabla \hat{J}(u^k)\|_U \\ &\leq \mathbf{c}_s \varrho_c(\chi_k(0)) + (1 - \mathbf{c}_s) \varrho_c(\chi_k(0)) = \varrho_c(\chi_k(0)), \end{aligned}$$

where the second inequality is established using $|\|u\|_U - \|w\|_U| \leq \|u - w\|_U$ for any $u, w \in U$ and the last one follows from (4.26) and (4.25). This proves that condition (4.11) is satisfied. \square

The two conditions (4.25) and (4.26) can be satisfied as follows:

Lemma 4.12. *Let (4.23) and (4.24) with $t = \tau$ be fulfilled with ε_g and ε_p such that*

$$(1 + (1 - \mathbf{c}_s)\mathbf{c}_c)c_g\varepsilon_g + \frac{1-\mathbf{c}_s}{\tau}\mathbf{c}_c c_p \varepsilon_p \leq \frac{1-\mathbf{c}_s}{\tau}\mathbf{c}_c \cdot \chi(u^k), \quad (4.27a)$$

$$c_g\varepsilon_g \leq \mathbf{c}_g \cdot \Delta_k, \quad (4.27b)$$

$$\mathbf{c}_s\mathbf{c}_c\tau c_g\varepsilon_g + (1 + \mathbf{c}_s\mathbf{c}_c)c_p\varepsilon_p \leq \mathbf{c}_s\mathbf{c}_c \cdot \chi(u^k) \quad (4.27c)$$

hold with some constants $\mathbf{c}_c, \mathbf{c}_g > 0$.

Then, the estimates (4.25) and (4.26) hold true with the choices $\varrho_c(t) = \mathbf{c}_c t$ and $\varrho_g(t) = \mathbf{c}_g t$.

Proof. From

$$|\chi_k(0) - \chi(u^k)| \leq \|[P_{U_{\text{ad}}} - \hat{P}_{U_{\text{ad}}}] (u^k - \tau \nabla m_k(0))\|_U + \tau \|\nabla m_k(0) - \nabla \hat{J}(u^k)\|_U$$

(see the proof of Proposition 4.11), (4.23), and (4.24) it follows that

$$\chi_k(0) \geq \chi(u^k) - c_p\varepsilon_p - \tau c_g\varepsilon_g$$

holds. Computing $\nabla m_k(0)$, $\hat{P}_{U_{\text{ad}}}(u^k - \tau \nabla m_k(0))$, and $\chi_k(0)$ (in this order) according to Assumption 4.9, with ε_g and ε_p as in (4.27), we get

$$\begin{aligned} \|\nabla m_k(0) - \nabla \hat{J}(u^k)\|_U &\leq c_g\varepsilon_g \\ &\stackrel{(4.27a)}{\leq} \frac{1-\mathbf{c}_s}{\tau}\mathbf{c}_c (\chi(u^k) - c_p\varepsilon_p - \tau c_g\varepsilon_g) \leq \frac{1-\mathbf{c}_s}{\tau}\mathbf{c}_c \cdot \chi_k(0), \end{aligned}$$

$$\|\nabla m_k(0) - \nabla \hat{J}(u^k)\|_U \leq c_g\varepsilon_g \stackrel{(4.27b)}{\leq} \mathbf{c}_g \cdot \Delta_k,$$

and

$$\begin{aligned} \|[P_{U_{\text{ad}}} - \hat{P}_{U_{\text{ad}}}] (u^k - \tau \nabla m_k(0))\|_U &\leq c_p\varepsilon_p \\ &\stackrel{(4.27c)}{\leq} \mathbf{c}_s\mathbf{c}_c (\chi(u^k) - c_p\varepsilon_p - \tau c_g\varepsilon_g) \leq \mathbf{c}_s\mathbf{c}_c \cdot \chi_k(0). \end{aligned}$$

The choices $\varrho_c(t) = \mathbf{c}_c t$ and $\varrho_g(t) = \mathbf{c}_g t$ yield (4.25) and (4.26). \square

We see that (4.27) can always be fulfilled by choosing ε_g and ε_p small enough: Given some parameters $\mathbf{a}_1, \mathbf{a}_2 \in [0, 1]$, we require

$$\varepsilon_g \leq \min\left\{ \frac{\mathbf{a}_1(1-\mathbf{c}_s)\mathbf{c}_c}{\tau c_g(1+(1-\mathbf{c}_s)\mathbf{c}_c)} \chi(u^k), \frac{\mathbf{c}_g}{c_g} \Delta_k, \frac{\mathbf{a}_2}{\tau c_g} \chi(u^k) \right\} \quad (4.28)$$

as well as

$$\varepsilon_p \leq \min\left\{ \frac{1-\mathbf{a}_1}{c_p} \chi(u^k), \frac{(1-\mathbf{a}_2)\mathbf{c}_s\mathbf{c}_c}{(1+\mathbf{c}_s\mathbf{c}_c)c_p} \chi(u^k) \right\} \quad (4.29)$$

to fulfill (4.27).

If we choose $\mathbf{a}_1 := \frac{1+(1-\mathbf{c}_s)\mathbf{c}_c}{1+\mathbf{c}_c} \in (0, 1]$ and $\mathbf{a}_2 := \frac{(1-\mathbf{c}_s)\mathbf{c}_c}{1+\mathbf{c}_c} \in [0, 1)$, the bound (4.28) becomes

$$\varepsilon_g \leq \min\left\{(1 - \mathbf{c}_s) \frac{\mathbf{c}_c}{1+\mathbf{c}_c} \cdot \frac{\chi(u^k)}{\tau c_g}, \frac{\mathbf{c}_g}{c_g} \Delta_k\right\}$$

and the bound (4.29) becomes

$$\varepsilon_p \leq \mathbf{c}_s \frac{\mathbf{c}_c}{1+\mathbf{c}_c} \cdot \frac{\chi(u^k)}{c_p},$$

i. e., we balance the bounds involving $\chi(u^k)$ by this choice.

As long as $\chi(u^k) > 0$ and $\mathbf{c}_s \in (0, 1)$, the respective bounds on ε_g and ε_p are positive. From $\varepsilon_g > 0$ and $\varepsilon_p > 0$ it then follows that $\chi_k(0) > 0$ must hold at the end of the refinement procedure by (4.25) and (4.26).

In the case $\chi(u^k) = 0$ it could happen that an adaptive algorithm for computing $\nabla m_k(0)$ and the projection keeps increasing the accuracy, i. e., decreasing ε_g and ε_p towards 0, without being able to fulfill (4.25) and (4.26). For theoretical considerations, we assume then, e. g., that the exact gradient and projection are used in this case. In a practical implementation, the refinement procedure should be stopped if $\chi_k(0) \leq c_\chi \varepsilon_{\text{tol}}$, $\varepsilon_g \leq \frac{\varepsilon_{\text{tol}}}{\tau}$, and $\varepsilon_p \leq \varepsilon_{\text{tol}}$ holds for some tolerance $\varepsilon_{\text{tol}} > 0$ and a constant $c_\chi > 0$, because then

$$\begin{aligned} \chi(u^k) &\leq \chi_k(0) + |\chi(u^k) - \chi_k(0)| \leq c_\chi \varepsilon_{\text{tol}} + c_p \varepsilon_p + \tau c_g \varepsilon_g \\ &\leq (c_\chi + c_p + c_g) \varepsilon_{\text{tol}}. \end{aligned}$$

Exact computation of the projection: When we consider a box-constrained problem in $L^2(\Omega_u)$ and a discretization of the control space by linear finite elements for example, we have to evaluate the exact projection to compute the error in $\hat{P}_{U_{\text{ad}}}$. Therefore, it makes sense to compute the inexact criticality measure directly with the exact projection, i. e., $\chi_k(0) = \tilde{\chi}_k(0)$, but not to refine the U -grid in the following to save computational costs. The bound (4.26) is then always satisfied and it remains to fulfill (4.25) with $\chi_k \equiv \tilde{\chi}_k$ and $\mathbf{c}_s = 0$. By setting $\varepsilon_p = 0$, we arrive at the bound

$$\varepsilon_g \leq \min\left\{\frac{\mathbf{c}_c}{\mathbf{c}_c+1} \cdot \frac{\chi(u^k)}{\tau c_g}, \frac{\mathbf{c}_g}{c_g} \Delta_k\right\}$$

for ε_g , which yields that (4.25) holds with $\mathbf{c}_s = 0$. Again, the bound $\frac{\mathbf{c}_c}{\mathbf{c}_c+1} \cdot \frac{\chi(u^k)}{\tau c_g}$ on ε_g cannot be computed explicitly, but is positive as long as u^k is not stationary. If the model gradient is not computed exactly, $\chi_k(0) > 0$ holds at the end of the refinement procedure.

Often, we do not have access to the constant c_g in (4.23), but only can compute the error estimate ε_g . Therefore, we set $\mathbf{c}_c := c_g \tau \tilde{\mathbf{c}}_c$ and $\mathbf{c}_g := c_g \tilde{\mathbf{c}}_g$ for two constants $\tilde{\mathbf{c}}_c, \tilde{\mathbf{c}}_g > 0$ and require

$$\varepsilon_g \leq \min\{\tilde{\mathbf{c}}_c \chi_k(0), \tilde{\mathbf{c}}_g \Delta_k\}$$

to ensure (4.25) with $\mathbf{c}_s = 0$, $\varrho_c(t) = \mathbf{c}_c t$, and $\varrho_g(t) = \mathbf{c}_g t$.

Generalized Cauchy Point

Condition (4.13) will be satisfied by a generalized Cauchy point $s_C^k \in U_{\text{ad}} - u^k$, $\|s_C^k\|_U \leq \Delta_k$, which is computed by some type of linesearch with an inexact, and possibly refined projection $\hat{P}_{U_{\text{ad}}}$. It is very important to permit an inexact projection in this procedure because then U -grid refinement may not be necessary in every iteration. In this way, computational cost can be saved and the quality of the FE grid can be preserved by a suitable refinement method instead of adapting the refinement exactly to the projection which has to be computed. Based on the generalized Cauchy point, the decrease condition (4.13) can be evaluated for improved trial steps. One can use an Armijo type linesearch on a line segment or a projected linesearch for example. We describe the latter here, following [118].

Lemma 4.13. *Given $u^k \in U_{\text{ad}}$, $\nabla m_k(0) \in U$ and $t > 0$, the direction*

$$p^k = p^k(t) := \hat{P}_{U_{\text{ad}}}(w^k(t)) - u^k$$

with $w^k = w^k(t) := u^k - t\nabla m_k(0)$ is a descent direction for m_k in 0 in the sense that $(\nabla m_k(0), p^k)_U \leq -\frac{\mathbf{c}_i}{t} \|p^k\|_U^2$, provided the inexact projection $\hat{P}_{U_{\text{ad}}}$ satisfies

$$(\hat{P}_{U_{\text{ad}}}(w^k(t)) - w^k(t), \hat{P}_{U_{\text{ad}}}(w^k(t)) - u^k)_U \leq (1 - \mathbf{c}_i) \|p^k(t)\|_U^2 \quad (4.30)$$

for some arbitrary, but fixed constant $\mathbf{c}_i \in (0, 1]$ and $\|p^k\|_U > 0$.

Proof. We estimate

$$\begin{aligned} t(\nabla m_k(0), p^k)_U &= (u^k - w^k, \hat{P}_{U_{\text{ad}}}(w^k) - u^k)_U \\ &= -\|p^k\|_U^2 + (\hat{P}_{U_{\text{ad}}}(w^k) - w^k, \hat{P}_{U_{\text{ad}}}(w^k) - u^k)_U \leq -\mathbf{c}_i \|p^k\|_U^2. \end{aligned}$$

Due to $t > 0$, $\mathbf{c}_i > 0$ and $\|p^k\|_U > 0$, p^k is a descent direction for m_k in 0. \square

Remark 4.14. This lemma is essentially different from [118, Lem. 5.3] in the sense that we do not use any projection property of $\hat{P}_{U_{\text{ad}}}$. If the discrete projection is used, (4.30) is trivially fulfilled.

Lemma 4.15. *Let $\hat{P}_{U_{\text{ad}}}(w^k) \in U_{\text{ad}}$ be computed such that (4.24) holds with $w^k(t) = w^k$ and with ε_p fulfilling*

$$(\mathbf{c}_i c_p \varepsilon_p + \|2\mathbf{c}_i P_{U_{\text{ad}}}(w^k) - w^k + (1 - 2\mathbf{c}_i)u^k\|_U) c_p \varepsilon_p \leq (1 - \mathbf{c}_i) \|P_{U_{\text{ad}}}(w^k) - u^k\|_U^2. \quad (4.31)$$

Then, (4.30) holds.

Proof. With the choice of ε_p and writing $P = P_{U_{\text{ad}}}(w^k)$, $\hat{P} = \hat{P}_{U_{\text{ad}}}(w^k)$, $w = w^k$, $u = u^k$ we estimate

$$\begin{aligned} &(\hat{P}_{U_{\text{ad}}}(w^k) - w^k, \hat{P}_{U_{\text{ad}}}(w^k) - u^k)_U = (\hat{P} - w, \hat{P} - u)_U \\ &= (\hat{P} - P + P - w, \hat{P} - P + P - u)_U \\ &= \|\hat{P} - P\|_U^2 + (\hat{P} - P, 2P - w - u)_U + (P - w, P - u)_U \end{aligned}$$

$$\begin{aligned}
 &= (1 - \mathbf{c}_i) \left[\|\hat{\mathbf{P}} - P\|_U^2 + (\hat{\mathbf{P}} - P, 2P - 2u)_U + \|P - u\|_U^2 \right] + \mathbf{c}_i \|\hat{\mathbf{P}} - P\|_U^2 \\
 &\quad + (\hat{\mathbf{P}} - P, 2\mathbf{c}_i P - w + (1 - 2\mathbf{c}_i)u)_U + (\mathbf{c}_i P - w + (1 - \mathbf{c}_i)u, P - u)_U \\
 &= (1 - \mathbf{c}_i) (\hat{\mathbf{P}} - P + P - u, \hat{\mathbf{P}} - P + P - u)_U + \mathbf{c}_i \|\hat{\mathbf{P}} - P\|_U^2 \\
 &\quad + (\hat{\mathbf{P}} - P, 2\mathbf{c}_i P - w + (1 - 2\mathbf{c}_i)u)_U \\
 &\quad + (P - w, P - u)_U - (1 - \mathbf{c}_i) (P - u, P - u)_U \\
 &\leq (1 - \mathbf{c}_i) \|p^k\|_U^2 + \mathbf{c}_i \|\hat{\mathbf{P}} - P\|_U^2 + \|\hat{\mathbf{P}} - P\|_U \cdot \|2\mathbf{c}_i P - w + (1 - 2\mathbf{c}_i)u\|_U \\
 &\quad - (1 - \mathbf{c}_i) \|P - u\|_U^2 \\
 &\stackrel{(4.24)}{\leq} (1 - \mathbf{c}_i) \|p^k\|_U^2 + \mathbf{c}_i c_p^2 \varepsilon_p^2 + c_p \varepsilon_p \cdot \|2\mathbf{c}_i P - w + (1 - 2\mathbf{c}_i)u\|_U \\
 &\quad - (1 - \mathbf{c}_i) \|P - u\|_U^2 \\
 &\stackrel{(4.31)}{\leq} (1 - \mathbf{c}_i) \|p^k\|_U^2.
 \end{aligned}$$

For the first estimate we have used the variational inequality property of the exact projection and the Cauchy-Schwarz inequality. \square

Remark 4.16. If $\tilde{\chi}_k(0) > 0$, then $\|P_{U_{\text{ad}}}(w^k(t)) - u^k\|_U > 0$ holds for every $t > 0$: If there existed $t > 0$ such that $\|P_{U_{\text{ad}}}(w^k(t)) - u^k\|_U = 0$, then $0 \in U$ would be first-order stationary for problem (4.8), which would yield $\tilde{\chi}_k(0) = 0$. Therefore, the bound on the right hand side of (4.31) can be satisfied by taking ε_p small enough. If we use the exact projection to compute $\chi_k(0) = \tilde{\chi}_k(0)$, $\tilde{\chi}_k(0) > 0$ is ensured when performing the linesearch by the stopping criterion in Algorithm 1.

Definition 4.17 (Armijo condition). Choose the largest $t_k \in \{\mathbf{c}_{a,k} \cdot \mathbf{c}_f^j, j \in \mathbb{N}_0\}$ such that

$$m_k(p^k(t_k)) \leq m_k(0) - \frac{\mathbf{c}_f \mathbf{c}_e}{t_k} \|p^k(t_k)\|_U^2, \quad (4.32a)$$

$$\|p^k(t_k)\|_U \leq \Delta_k \quad (4.32b)$$

with two fixed parameters $\mathbf{c}_f, \mathbf{c}_e \in (0, 1)$ are satisfied. Here, $\mathbf{c}_{a,k} \in \mathbb{R}$ are constants satisfying $\mathbf{c}_a \leq \mathbf{c}_{a,k} \leq \tau$ for all $k \in \mathbb{N}_0$ with some $\mathbf{c}_a \in (0, \tau]$.

Remark 4.18. Note that

$$-\frac{\mathbf{c}_f \mathbf{c}_e}{t_k} \|p^k(t_k)\|_U^2 \geq \mathbf{c}_e (\nabla m_k(0), p^k(t_k))_U$$

holds by Lemma 4.13 if $\hat{P}_{U_{\text{ad}}}$ satisfies (4.30). Then, condition (4.32a) is implied by

$$m_k(p^k(t_k)) \leq m_k(0) + \mathbf{c}_e (\nabla m_k(0), p^k(t_k))_U.$$

Lemma 4.19. If Assumption 4.5 holds and if $p^k(t_k)$ is computed according to Lemma 4.13, condition (4.32a) is satisfied for all $t_k \in (0, \frac{2(1-\mathbf{c}_e)\varepsilon_1}{\mathbf{c}_{m_k}}]$, where $\mathbf{c}_{m_k} > 0$ is the Lipschitz constant of the model gradient ∇m_k .

Proof. We estimate using the fundamental theorem of calculus:

$$\begin{aligned}
 m_k(p^k(t_k)) - m_k(0) &= \int_0^1 (\nabla m_k(\sigma p^k(t_k)), p^k(t_k))_U \, d\sigma \\
 &\leq (\nabla m_k(0), p^k(t_k))_U + \int_0^1 \|\nabla m_k(\sigma p^k(t_k)) - \nabla m_k(0)\|_U \cdot \|p^k(t_k)\|_U \, d\sigma \\
 &\leq -\frac{\mathbf{c}_i}{t_k} \|p^k(t_k)\|_U^2 + \frac{\mathbf{c}_{m_k}}{2} \|p^k(t_k)\|_U^2 = \left(-\frac{\mathbf{c}_i}{t_k} + \frac{\mathbf{c}_{m_k}}{2}\right) \|p^k(t_k)\|_U^2.
 \end{aligned}$$

In the last estimate we have used Lemma 4.13 and the Lipschitz continuity of ∇m_k . With the given choice of t_k , (4.32a) follows. \square

Lemma 4.20. *Let Assumption 4.5 hold and let t_k be computed according to the Armijo condition (Definition 4.17), where the inexact projection $\hat{P}_{U_{\text{ad}}}$ satisfies (4.30) and additionally*

$$\|P_{U_{\text{ad}}}(w^k(t)) - \hat{P}_{U_{\text{ad}}}(w^k(t))\|_U \leq \frac{1-\mathbf{c}_{11}}{\mathbf{c}_{11}} \|\hat{P}_{U_{\text{ad}}}(w^k(t)) - u^k\|_U \quad (4.33)$$

as well as

$$\|P_{U_{\text{ad}}}(w^k(t)) - \hat{P}_{U_{\text{ad}}}(w^k(t))\|_U \leq \frac{1-\mathbf{c}_{12}}{\mathbf{c}_{12}} \|P_{U_{\text{ad}}}(w^k(t)) - u^k\|_U \quad (4.34)$$

for every $t = t_k$ tested during the Armijo linesearch with some constants $\mathbf{c}_{11}, \mathbf{c}_{12} \in (0, 1]$, cf. [118, Eq. (3.24)].

Then, the trial step $s_C^k := p^k(t_k)$ fulfills condition (4.13) with

$$\varrho_{t1}(t) := \frac{\mathbf{c}_i \mathbf{c}_{11}^2 \mathbf{c}_{12} \mathbf{c}_e \mathbf{c}_f}{\tau} t, \quad \varrho_{t2}(t) := \frac{1}{\mathbf{c}_{12} \tau} \cdot \min\left\{\frac{2(1-\mathbf{c}_e)\mathbf{c}_i}{\mathbf{c}_m}, \frac{\mathbf{c}_a}{\mathbf{c}_f}\right\} t,$$

and $\chi_k(0) = \tilde{\chi}_k(0)$.

Proof. From (4.33) we get that

$$\begin{aligned}
 &\mathbf{c}_{11} \|P_{U_{\text{ad}}}(w^k(t)) - u^k\|_U \\
 &\leq \mathbf{c}_{11} \|P_{U_{\text{ad}}}(w^k(t)) - \hat{P}_{U_{\text{ad}}}(w^k(t))\|_U + \mathbf{c}_{11} \|\hat{P}_{U_{\text{ad}}}(w^k(t)) - u^k\|_U \\
 &\leq (1 - \mathbf{c}_{11}) \|\hat{P}_{U_{\text{ad}}}(w^k(t)) - u^k\|_U + \mathbf{c}_{11} \|\hat{P}_{U_{\text{ad}}}(w^k(t)) - u^k\|_U \\
 &= \|\hat{P}_{U_{\text{ad}}}(w^k(t)) - u^k\|_U = \|p^k(t)\|_U
 \end{aligned} \quad (4.35)$$

for every tested $t = t_k$. Using

$$\frac{1}{t} \|P_{U_{\text{ad}}}(w^k(t)) - u^k\|_U \geq \frac{1}{\tau} \tilde{\chi}_k(0) \quad (4.36)$$

for $t \leq \tau$ by [60, Lem. 1.10 (e)], which states that the function $\phi(t) := \frac{1}{t} \|P_{U_{\text{ad}}}(u^k - t \nabla m_k(0)) - u^k\|_U$ ($t > 0$) is non-increasing, and (4.32a) as well as (4.35), we conclude that

$$\begin{aligned}
 \text{pred}_k &= m_k(0) - m_k(p^k(t_k)) \geq \frac{\mathbf{c}_i \mathbf{c}_e}{t_k} \|p^k(t_k)\|_U^2 \\
 &\geq \mathbf{c}_i \mathbf{c}_e \mathbf{c}_{11} \cdot \frac{1}{t_k} \|P_{U_{\text{ad}}}(w^k(t_k)) - u^k\|_U \cdot \|p^k(t_k)\|_U \\
 &\geq \mathbf{c}_i \mathbf{c}_e \mathbf{c}_{11} \cdot \frac{1}{\tau} \tilde{\chi}_k(0) \cdot \|p^k(t_k)\|_U
 \end{aligned} \quad (4.37)$$

holds by $t_k \leq \mathbf{c}_{a,k} \leq \tau$.

Now consider the case that t_k found by the standard projected Armijo linesearch for (4.32a) already satisfies (4.32b). Thus, Lemma 4.19 can be applied and $t_k \geq \min\left\{\frac{2(1-c_e)c_i c_f}{c_{m_k}}, c_{a,k}\right\}$ holds. From (4.37), (4.35), and (4.36) it now follows that

$$\begin{aligned} \text{pred}_k &\geq \frac{c_i c_e c_{11}}{\tau} \tilde{\chi}_k(0) \cdot \frac{c_{11} t_k}{\tau} \tilde{\chi}_k(0) \\ &\geq \frac{c_i c_e c_{11}}{\tau} \tilde{\chi}_k(0) \cdot \min\left\{\frac{2(1-c_e)c_i c_f}{c_{m_k}}, c_{a,k}\right\} \cdot \frac{c_{11}}{\tau} \tilde{\chi}_k(0) \\ &= \frac{c_i c_{11}^2 c_{12} c_e c_f}{\tau} \tilde{\chi}_k(0) \cdot \frac{1}{c_{12} \tau} \cdot \min\left\{\frac{2(1-c_e)c_i}{c_m}, \frac{c_a}{c_f}\right\} \tilde{\chi}_k(0) \\ &= \varrho_{t1}(\tilde{\chi}_k(0)) \cdot \varrho_{t2}(\tilde{\chi}_k(0)) \geq \varrho_{t1}(\tilde{\chi}_k(0)) \cdot \min\{\varrho_{t2}(\tilde{\chi}_k(0)), \Delta_k\}. \end{aligned}$$

In the case that the standard search for (4.32a) does not yield t_k satisfying (4.32b), t_k has to be decreased further. It follows that $\|p^k(\frac{t_k}{c_f})\|_U > \Delta_k$. In analogy to (4.35) we can conclude from (4.34) that

$$c_{12} \|p^k(t)\|_U = c_{12} \|\hat{P}_{U_{\text{ad}}}(w^k(t)) - u^k\|_U \leq \|P_{U_{\text{ad}}}(w^k(t)) - u^k\|_U \quad (4.38)$$

for every tested $t = t_k$. With (4.35), [60, Lem. 1.10 (e)], and (4.38) we obtain

$$\begin{aligned} \|p^k(t_k)\|_U &\geq c_{11} \|P_{U_{\text{ad}}}(w^k(t_k)) - u^k\|_U \\ &\geq c_{11} c_f \|P_{U_{\text{ad}}}(w^k(\frac{t_k}{c_f})) - u^k\|_U \\ &\geq c_{11} c_{12} c_f \|p^k(\frac{t_k}{c_f})\|_U > c_{11} c_{12} c_f \Delta_k. \end{aligned} \quad (4.39)$$

Hence, by (4.37), we get

$$\begin{aligned} \text{pred}_k &\geq \frac{c_i c_e c_{11}}{\tau} \tilde{\chi}_k(0) \cdot \|p^k(t_k)\|_U \\ &> \frac{c_i c_e c_{11}}{\tau} \tilde{\chi}_k(0) \cdot c_{11} c_{12} c_f \Delta_k \\ &= \varrho_{t1}(\tilde{\chi}_k(0)) \cdot \Delta_k \geq \varrho_{t1}(\tilde{\chi}_k(0)) \cdot \min\{\varrho_{t2}(\tilde{\chi}_k(0)), \Delta_k\}. \end{aligned}$$

Therefore, in both cases, (4.13) is satisfied with $\chi_k(0) = \tilde{\chi}_k(0)$. □

Remark 4.21.

- If (4.33) holds with $c_{11} \in (\frac{1}{2}, 1]$, (4.34) follows if we choose $c_{12} = \frac{2c_{11}-1}{c_{11}} \in (0, 1]$: From (4.33), we obtain

$$\begin{aligned} &\|P_{U_{\text{ad}}}(w^k(t)) - \hat{P}_{U_{\text{ad}}}(w^k(t))\|_U \\ &\leq \frac{1-c_{11}}{c_{11}} \|\hat{P}_{U_{\text{ad}}}(w^k(t)) - u^k\|_U \\ &\leq \frac{1-c_{11}}{c_{11}} \left(\|\hat{P}_{U_{\text{ad}}}(w^k(t)) - P_{U_{\text{ad}}}(w^k(t))\|_U + \|P_{U_{\text{ad}}}(w^k(t)) - u^k\|_U \right). \end{aligned}$$

This is equivalent to

$$\frac{2c_{11}-1}{c_{11}} \|P_{U_{\text{ad}}}(w^k(t)) - \hat{P}_{U_{\text{ad}}}(w^k(t))\|_U \leq \frac{1-c_{11}}{c_{11}} \|P_{U_{\text{ad}}}(w^k(t)) - u^k\|_U,$$

which yields (4.34) with the given value of c_{12} .

- The prerequisite (4.34) in Lemma 4.20 can be dropped if the discrete projection $\hat{P}_{U_{\text{ad}}}$ onto $U_{\text{ad}} \cap U$ is used (cf. [118]) since then

$$\|p^k(t_k)\|_U \geq \mathbf{c}_f \|p^k(\frac{t_k}{\mathbf{c}_f})\|_U$$

holds by the projection property of $\hat{P}_{U_{\text{ad}}}$, which we do not assume in general, and then replaces (4.39). The constants in ϱ_{t1} and ϱ_{t2} change accordingly.

- Instead of (4.33) and (4.34) one can directly assume (4.35) and (4.38) in Lemma 4.20, but these conditions can be more difficult to evaluate than (4.33) and (4.34). In particular, the verification of condition (4.35) only requires the estimation of the projection error and the evaluation of the inexact projection.
- Given $\mathbf{c}_a \in (0, \tau]$, it makes sense to choose

$$\mathbf{c}_{a,k} := \max\{\mathbf{c}_a, \min\{\tau, \frac{\mathbf{c}_{12}\Delta_k}{\|\nabla m_k(0)\|_U}\}\} \quad (4.40)$$

in the Armijo condition because with $t_k \leq \frac{\mathbf{c}_{12}\Delta_k}{\|\nabla m_k(0)\|_U}$, which is not always ensured by (4.40), and (4.38) we obtain

$$\begin{aligned} \|p^k(t_k)\|_U &\leq \frac{1}{\mathbf{c}_{12}} \|P_{U_{\text{ad}}}(u^k - t_k \nabla m_k(0)) - u^k\|_U \\ &= \frac{1}{\mathbf{c}_{12}} \|P_{U_{\text{ad}}}(u^k - t_k \nabla m_k(0)) - P_{U_{\text{ad}}}(u^k)\|_U \\ &\leq \frac{t_k}{\mathbf{c}_{12}} \|\nabla m_k(0)\|_U \leq \Delta_k, \end{aligned}$$

which is exactly (4.32b). Note that this choice of t_k is not necessary (only sufficient) for (4.32b). This observation yields together with Lemma 4.19 that the projected Armijo linesearch terminates after finitely many times decreasing t_k .

Lemma 4.22. *Let $t \leq \tau$ and let $\hat{P}_{U_{\text{ad}}}(w^k(t)) \in U_{\text{ad}}$ be computed such that (4.24) holds with ε_p fulfilling*

$$\varepsilon_p \leq \frac{(1-\mathbf{c}_{11})t}{\mathbf{c}_p\tau} \tilde{\chi}_k(0). \quad (4.41)$$

Then, (4.33) holds.

Proof. We write $\hat{P} = \hat{P}_{U_{\text{ad}}}(w^k(t))$, $P = P_{U_{\text{ad}}}(w^k(t))$, and $u = u^k$ and get

$$\begin{aligned} \frac{1-\mathbf{c}_{11}}{\mathbf{c}_{11}} \|\hat{P} - u\|_U &\geq \frac{1-\mathbf{c}_{11}}{\mathbf{c}_{11}} (\|P - u\|_U - \|P - \hat{P}\|_U) \\ &\geq \frac{1-\mathbf{c}_{11}}{\mathbf{c}_{11}} \left(\frac{t}{\tau} \tilde{\chi}_k(0) - \mathbf{c}_p \varepsilon_p \right) \stackrel{(4.41)}{\geq} \frac{\mathbf{c}_p}{\mathbf{c}_{11}} \varepsilon_p - \frac{1-\mathbf{c}_{11}}{\mathbf{c}_{11}} \mathbf{c}_p \varepsilon_p \\ &= \mathbf{c}_p \varepsilon_p \stackrel{(4.24)}{\geq} \|\hat{P} - P\|_U, \end{aligned}$$

which proves (4.33). In the second estimate we have used (4.24) and (4.36). \square

Remark 4.23. In the same way, we obtain that (4.34) is satisfied if

$$\varepsilon_p \leq \frac{(1-\mathbf{c}_{12})t}{\mathbf{c}_{12}\mathbf{c}_p\tau} \tilde{\chi}_k(0).$$

Having computed the generalized Cauchy point s_C^k yields a simple criterion for (4.13):

Lemma 4.24 (Fraction of generalized Cauchy decrease). *Let s_C^k be computed according to Lemma 4.20 and let $s^k \in U_{\text{ad}} - u^k$ satisfy*

$$m_k(0) - m_k(s^k) \geq \mathbf{c}_d(m_k(0) - m_k(s_C^k))$$

with some $\mathbf{c}_d \in (0, 1]$.

Then, s^k satisfies (4.13) with $\chi_k(0) = \tilde{\chi}_k(0)$ and

$$\varrho_{t1}(t) := \frac{\mathbf{c}_d \mathbf{c}_i \mathbf{c}_{11}^2 \mathbf{c}_{12} \mathbf{c}_e \mathbf{c}_f}{\tau} t, \quad \varrho_{t2}(t) := \frac{1}{\mathbf{c}_{12} \tau} \cdot \min\left\{\frac{2(1-\mathbf{c}_e)\mathbf{c}_i}{\mathbf{c}_m}, \frac{\mathbf{c}_a}{\mathbf{c}_f}\right\} t.$$

Proof. From Lemma 4.20 we know that s_C^k satisfies (4.13) with

$$\varrho_{t1}(t) := \frac{\mathbf{c}_i \mathbf{c}_{11}^2 \mathbf{c}_{12} \mathbf{c}_e \mathbf{c}_f}{\tau} t, \quad \varrho_{t2}(t) := \frac{1}{\mathbf{c}_{12} \tau} \cdot \min\left\{\frac{2(1-\mathbf{c}_e)\mathbf{c}_i}{\mathbf{c}_m}, \frac{\mathbf{c}_a}{\mathbf{c}_f}\right\} t,$$

and $\chi_k(0) = \tilde{\chi}_k(0)$. The stated result follows immediately. \square

Computed Reduction

The computed reduction is only evaluated as long as $\chi_k(0)$ is positive, which is ensured by the stopping criterion of Algorithm 1. Then the predicted reduction pred_k is positive by (4.13). The inequality (4.12) can be reduced to a bound on the inexact objective function evaluation:

Lemma 4.25. *Let*

$$|\hat{J}(u) - \hat{J}_k(u)| \leq \frac{1}{2} \varrho_r(\eta_3 \min\{\text{pred}_k, \mathbf{r}_k\})$$

hold for all $u \in \{u^k, u^k + s^k\} \subset U$. Then, condition (4.12) is fulfilled.

Proof. Using the definitions (4.6a) and (4.6b), we estimate

$$|\text{ared}_k - \text{cred}_k| \leq |\hat{J}(u^k) - \hat{J}_k(u^k)| + |\hat{J}_k(u^k + s^k) - \hat{J}(u^k + s^k)| \leq \varrho_r(\eta_3 \min\{\text{pred}_k, \mathbf{r}_k\}),$$

which is (4.12). \square

Choosing $\varepsilon_o \leq \tilde{\mathbf{c}}_o \cdot (\eta_3 \min\{\text{pred}_k, \mathbf{r}_k\})^{\mathbf{c}_o} > 0$ with some constant $\tilde{\mathbf{c}}_o > 0$ and $\varrho_r(t) = 2\mathbf{c}_o \tilde{\mathbf{c}}_o t^{\mathbf{c}_o}$, the condition given in Lemma 4.25 can be fulfilled by Assumption 4.9.

Update of the Trust-Region Radius

Due to $\text{pred}_k > 0$, the quantity $\frac{\text{cred}_k}{\text{pred}_k}$ is well-defined. Since we have $0 < \Delta_k \leq \Delta_{\max}$ and $0 < \nu_1 < 1 \leq \nu_2 < \nu_3$, the intervals from which the new radius Δ_{k+1} is chosen are always nonempty.

4.4. Solution of Subproblems by a Semismooth Newton Method

To obtain fast convergence of the algorithm, it is necessary to apply a more sophisticated solver to the trust-region subproblem (4.5). Since a semismooth Newton method [92, 110, 58, 111] gave very good results previously [46] and also in different contexts involvings PDEs [60, 111], we want to apply it here also. As it is difficult to handle, e.g., box constraints and an additional trust-region constraint at the same time, we replace the latter in (4.5) by quadratic regularization and obtain

$$\min_{s \in U} m_k(s) + \frac{c_{n,k}}{2\Delta_k} \|s\|_U^2 \quad \text{s. t.} \quad u^k + s \in U_{\text{ad}} \quad (4.42)$$

with some constant $c_{n,k} > 0$. A solution of this problem is not necessarily feasible for (4.5). Still, there is a connection between the two problems:

Lemma 4.26. *Let $\bar{s} \in U$ be a global solution of (4.42) with $\|\bar{s}\|_U = \Delta_k$. Then, \bar{s} solves (4.5).*

Proof. It holds that $m_k(\bar{s}) + \frac{c_{n,k}}{2\Delta_k} \|\bar{s}\|_U^2 \leq m_k(s) + \frac{c_{n,k}}{2\Delta_k} \|s\|_U^2$ for all $s \in U$ such that $u^k + s \in U_{\text{ad}}$. Thus, $m_k(\bar{s}) \leq m_k(s) + \frac{c_{n,k}}{2\Delta_k} (\|s\|_U^2 - \|\bar{s}\|_U^2) \leq m_k(s)$ for all $s \in U$ such that $u^k + s \in U_{\text{ad}}$ and $\|s\|_U^2 \leq \Delta_k^2$. \square

For the application of semismooth Newton we follow [60, Chap. 2] and [111]. A necessary optimality condition for (4.42) is given by

$$R(\bar{s}) := \bar{s} - P_{U_{\text{ad}} - u^k}(\bar{s} - \tau_{n,k}(\nabla m_k(\bar{s}) + \frac{c_{n,k}}{\Delta_k} \bar{s})) = 0 \quad (4.43)$$

with an arbitrary $\tau_{n,k} > 0$. We consider the case of a quadratic model

$$m_k(s) = m_k(0) + (\nabla m_k(0), s)_U + \frac{1}{2}(\nabla^2 m_k(0)s, s)_U$$

with gradient $\nabla m_k(s) = \nabla m_k(0) + \nabla^2 m_k(0)s$. In our application we have $\nabla m_k(0) = -B^* \tilde{T}(u^k) + \gamma u^k =: \hat{G}(u^k) + \gamma u^k$, where $\tilde{T}(u^k)$ is the inexact adjoint state. The Hessian of the model can be a positive multiple of the identity or a better approximation of the true Hessian. The Hessian is often of the form

$$\nabla^2 m_k(0)s = \hat{H}(u^k)s + \gamma s, \quad (4.44)$$

cf. (3.23), with $\hat{H}(u^k) \in \mathcal{L}(U, U)$. The operator $\hat{H}(u^k)$ need not necessarily correspond to the exact Hessian as derived in (3.23), but is only required to fulfill certain properties to establish semismoothness.

Choosing $\tau_{n,k} := (\gamma + \frac{c_{n,k}}{\Delta_k})^{-1}$, (4.43) becomes

$$\begin{aligned} R(\bar{s}) &= \bar{s} - P_{U_{\text{ad}} - u^k}(\bar{s} - \tau_{n,k}(\nabla m_k(\bar{s}) + \frac{c_{n,k}}{\Delta_k} \bar{s})) \\ &= \bar{s} - P_{U_{\text{ad}} - u^k}(\bar{s} - \tau_{n,k}(\nabla m_k(0) + \nabla^2 m_k(0)\bar{s} + \frac{c_{n,k}}{\Delta_k} \bar{s})) = \\ &= \bar{s} - P_{U_{\text{ad}} - u^k}(-\tau_{n,k}(\nabla m_k(0) + \hat{H}(u^k)\bar{s})) \\ &= \bar{s} + u^k - P_{U_{\text{ad}}}(u^k - \tau_{n,k}\nabla m_k(0) - \tau_{n,k}\hat{H}(u^k)\bar{s}) = 0. \end{aligned} \quad (4.45)$$

Now, assume that $u^k \in L^q(\Omega_u)$ and $\nabla m_k(0) \in L^q(\Omega_u)$ holds with some $q \geq 2$,⁸ which can be ensured by $u^k \in U_{\text{ad}} \subset L^q(\Omega_u)$ and $\hat{G}(u^k) \in L^q(\Omega_u)$ for example, as well as $\hat{H}(u^k) \in \mathcal{L}(L^2(\Omega_u), L^q(\Omega_u))$. Then,

$$f_1 : L^2(\Omega_u) \rightarrow L^q(\Omega_u), s \mapsto f_1(s) := u^k - \tau_{n,k} \nabla m_k(0) - \tau_{n,k} \hat{H}(u^k) s$$

is affine, bounded and thus $\{-\tau_{n,k} \hat{H}(u^k)\}$ -semismooth by [111, Prop. 3.4]. If additionally $P_{U_{\text{ad}}} : L^q(\Omega_u) \rightarrow L^2(\Omega_u)$ is $\partial P_{U_{\text{ad}}}$ -semismooth and bounded near $f_1(s)$, we can apply the chain rule [111, Prop. 3.8] to show that R is ∂R -semismooth at s with $\partial R(s) = \{I + \tau_{n,k} M_P \hat{H}(u^k) : M_P \in \partial P_{U_{\text{ad}}}(f_1(s))\}$.

Given a current iterate $s^{k,\ell}$ and $M_P \in \partial P_{U_{\text{ad}}}(f_1(s^{k,\ell}))$, the semismooth Newton equation for an update $\tilde{d}^{k,\ell}$ reads

$$\tilde{d}^{k,\ell} + \tau_{n,k} M_P \hat{H}(u^k) \tilde{d}^{k,\ell} = -R(s^{k,\ell}).$$

Having computed a solution of it (approximately), we set $s^{k,\ell+1} := s^{k,\ell} + \tilde{d}^{k,\ell}$.

Remark 4.27. To prove superlinear convergence of the semismooth Newton method, a regularity condition is required. Such a condition can be derived depending on the concrete choice of U_{ad} and the solution of (4.42), see [111, Sec. 9.1].

In contrast to the projected linesearch developed in Section 4.3, the (approximate) solution found by applying semismooth Newton to (4.42) may not be feasible for (4.5). Thus, we set $s_{\text{SSN}}^k := P_{\{s \in U : \|s\|_U \leq \Delta_k\}}(P_{U_{\text{ad}}-u^k}(\hat{s}^k))$, where \hat{s}^k is an (approximate) solution of (4.42).

Lemma 4.28. *Given $\hat{s} \in U$, $u \in U$ and $\Delta_k > 0$, $s := P_{\{s \in U : \|s\|_U \leq \Delta_k\}}(P_{U_{\text{ad}}-u}(\hat{s}))$ is feasible for (4.5), where*

$$P_{\{s \in U : \|s\|_U \leq \Delta_k\}}(s) = \frac{\Delta_k}{\max\{\|s\|_U, \Delta_k\}} s.$$

Proof. It is obvious that $\|\hat{s}\|_U \leq \Delta_k$ holds. Moreover, by $\frac{\Delta_k}{\max\{\|s\|_U, \Delta_k\}} \leq 1$ for all $s \in U$ and the convexity of U_{ad} , we have

$$\begin{aligned} u + s &= u + \frac{\Delta_k}{\max\{\|P_{U_{\text{ad}}-u}(\hat{s})\|_U, \Delta_k\}} P_{U_{\text{ad}}-u}(\hat{s}) \\ &= \frac{\Delta_k}{\max\{\|P_{U_{\text{ad}}-u}(\hat{s})\|_U, \Delta_k\}} \underbrace{(u + P_{U_{\text{ad}}-u}(\hat{s}))}_{\in U_{\text{ad}}} + \left(1 - \frac{\Delta_k}{\max\{\|P_{U_{\text{ad}}-u}(\hat{s})\|_U, \Delta_k\}}\right) u \in U_{\text{ad}}. \end{aligned}$$

□

⁸Often, depending on the concrete choice of U_{ad} , we need $q > 2$ to establish semismoothness.

5. Realization of the Required Error Estimates for the Model Problem

We want to apply Algorithm 1 to the example from Section 3.2. To ensure global convergence we have to control the following quantities, see Chapter 4:

- The inexactness of the model gradient:

$$\|\nabla m_k(0) - \nabla \hat{J}(u^k)\|_U \leq \varrho_g(\Delta_k), \quad (5.1)$$

where $\varrho_g : \mathbb{R}_{>0} \rightarrow \mathbb{R}_{\geq 0}$ is a function satisfying $\lim_{t \rightarrow 0^+} \varrho_g(t) = 0$, e. g., $\varrho_g(t) = \mathbf{c}_g t$, $\mathbf{c}_g > 0$.

- The inexactness of the approximate criticality measure:

$$|\chi_k(0) - \chi(u^k)| \leq \varrho_c(\chi_k(0)), \quad (5.2)$$

where $\varrho_c : \mathbb{R}_{\geq 0} \rightarrow \mathbb{R}_{\geq 0}$ is a function satisfying $\lim_{t \rightarrow 0^+} \varrho_c(t) = 0$ and $\varrho_c(0) = 0$, e. g., $\varrho_c(t) = \mathbf{c}_c t$, $\mathbf{c}_c > 0$.

- The quality of the computed reduction:

$$|\text{ared}_k - \text{cred}_k| \leq \varrho_r(\eta_3 \min\{\text{pred}_k, \mathbf{r}_k\}), \quad (5.3)$$

with $\eta_3 < \min\{\eta_1, 1 - \eta_2\}$, where $0 < \eta_1 < \eta_2 < 1$ are chosen a priori, and with a forcing sequence $(\mathbf{r}_k)_{k \in \mathbb{N}_0} \subset \mathbb{R}_{>0}$ fulfilling $\lim_{k \rightarrow \infty} \mathbf{r}_k = 0$ and a function $\varrho_r : \mathbb{R}_{>0} \rightarrow \mathbb{R}_{\geq 0}$ fulfilling $\varrho_r(t) \leq t$ for all $t \in (0, \bar{t}]$ with some fixed $\bar{t} > 0$, e. g., $\varrho_r(t) = \mathbf{c}_r t^{\mathbf{e}_r}$, $\mathbf{c}_r > 0$, $\mathbf{e}_r > 1$.

In this chapter, we discuss how the required error estimates (5.1), (5.2), and (5.3) can be ensured for the example from Section 3.2, but several estimates will hold in a more general setting.

5.1. Realization of the Error Estimates in the Deterministic Case

First, we focus on the deterministic case with fixed $\xi \in \Xi$, where the objective function

$$J[\xi](y, u) = \frac{1}{2} \|Q(\xi)y - \hat{q}(\xi)\|_H^2 + \frac{\gamma}{2} \|u\|_U^2 \quad (5.4)$$

is of tracking type with a real Hilbert space H , a desired state $\hat{q}(\xi) \in H$, $Q(\xi) \in \mathcal{L}(Y, H)$, and $\gamma > 0$ as in (3.4) and the state equation is

$$E[\xi](y, u) = A(\xi)y + N(y) - B(\xi)u - b(\xi) = 0 \quad (5.5)$$

with a strongly monotone (with constant $\underline{\kappa}$), linear, bounded operator $A(\xi) : Y \rightarrow Y^*$, a monotone and continuously differentiable operator $N : Y \rightarrow Y^*$, $B(\xi) \in \mathcal{L}(U, Y^*)$, and $b(\xi) \in Y^*$, cf. (3.10) and Assumption 3.3. In the following, we will skip the dependence on ξ because it is fixed. Still, all derived results are valid for almost every $\xi \in \Xi$. Furthermore, we skip the index k denoting the iteration number in the algorithm for readability purposes in this section as far as possible.

In the following, we describe the procedure which is employed to realize the error control in the deterministic case.

Model Gradient Error

The model gradient is computed by the (formal) adjoint approach with inexact solutions \tilde{y} and \tilde{z} of the state and (formal) adjoint equation, respectively. Let $y = S(u) \in Y$ be the exact state solving $E(S(u), u) = 0$ and let $\tilde{y} \in Y$ be an inexact solution. The perturbed adjoint equation reads

$$E_y(\tilde{y}, u)^* \hat{z} = -J_y(\tilde{y}, u). \quad (5.6)$$

Its exact solution is denoted by \hat{z} and its inexact solution by \tilde{z} , whereas z is the exact adjoint state solving

$$E_y(y, u)^* z = -J_y(y, u). \quad (5.7)$$

Theorem 5.1. *Let $J : Y \times U \rightarrow \mathbb{R}$ be defined as in (5.4), let $E : Y \times U \rightarrow Y^*$ be as in (5.5), and let $u \in U$ and $\tilde{y}, \tilde{z} \in Y$ be given. Moreover, let $y \in Y$ be the exact solution of $E(y, u) = 0$ and let $z \in Y$ and $\hat{z} \in Y$ be the exact solutions of (5.7) and (5.6), respectively. Choosing m such that $m'(0) = E_u(\tilde{y}, u)^* \tilde{z} + J_u(\tilde{y}, u)$, it then holds that*

$$\begin{aligned} \|m'(0) - \hat{J}'(u)\|_{U^*} &\leq \|B^*(\hat{z} - \tilde{z})\|_{U^*} + \frac{1}{\underline{\kappa}} \|B\|_{\mathcal{L}(U, Y^*)} \|Q\|_{\mathcal{L}(Y, H)} (\|Q(y - \tilde{y})\|_H \\ &\quad + \frac{1}{\underline{\kappa}} \|N'(\tilde{y}) - N'(y)\|_{\mathcal{L}(Y, Y^*)} (\|Q\tilde{y} - \hat{q}\|_H + \|Q(y - \tilde{y})\|_H)). \end{aligned} \quad (5.8)$$

Proof. From the choice of $m'(0)$ and (5.7) we obtain using that $E_u \equiv -B$ and that J_u is independent of y :

$$\begin{aligned} \|m'(0) - \hat{J}'(u)\|_{U^*} &= \|E_u(\tilde{y}, u)^* \tilde{z} - E_u(y, u)^* z + J_u(\tilde{y}, u) - J_u(y, u)\|_{U^*} \\ &= \|B^*(z - \tilde{z})\|_{U^*} \leq \|B^*(z - \hat{z})\|_{U^*} + \|B^*(\hat{z} - \tilde{z})\|_{U^*}. \end{aligned} \quad (5.9)$$

Since both equations (5.6) and (5.7) are uniquely solvable because of the strong monotonicity of $A^* + N'(y)^*$ for every $y \in Y$ by Proposition A.4, we can compute

$$\begin{aligned} \hat{z} - z &= -E_y(\tilde{y}, u)^{-*} J_y(\tilde{y}, u) + E_y(y, u)^{-*} J_y(y, u) \\ &= E_y(\tilde{y}, u)^{-*} (J_y(y, u) - J_y(\tilde{y}, u)) + (E_y(y, u)^{-*} - E_y(\tilde{y}, u)^{-*}) J_y(\tilde{y}, u) \\ &\quad + (E_y(y, u)^{-*} - E_y(\tilde{y}, u)^{-*}) (J_y(y, u) - J_y(\tilde{y}, u)). \end{aligned} \quad (5.10)$$

With (5.4) and (5.5), the perturbed adjoint equation (5.6) reads

$$E_y(\tilde{y}, u)^* z = A^* z + N'(\tilde{y})^* z = -Q^*(Q\tilde{y} - \hat{q}) = -J_y(\tilde{y}, u).$$

Compared to the exact adjoint equation (5.7), the error in the right-hand side is

$$J_y(y, u) - J_y(\tilde{y}, u) = Q^*Q(y - \tilde{y}). \quad (5.11)$$

For the estimation of the error caused by the approximate left-hand side operator, we introduce for $\tilde{b} \in Y^*$ the unique solutions $v, \tilde{v} \in Y$ of the equations

$$\begin{aligned} A^*v + N'(y)^*v &= \tilde{b}, \\ A^*\tilde{v} + N'(\tilde{y})^*\tilde{v} &= \tilde{b}, \end{aligned}$$

respectively. We have $v - \tilde{v} = (E_y(y, u)^{-*} - E_y(\tilde{y}, u)^{-*})\tilde{b}$. Using the monotonicity of $N'(y)$ (Proposition A.4) and the strong monotonicity of A with constant $\underline{\kappa}$, we estimate:

$$\begin{aligned} \underline{\kappa}\|v - \tilde{v}\|_Y^2 &\leq \langle v - \tilde{v}, A(v - \tilde{v}) \rangle_{Y, Y^*} \leq \langle v - \tilde{v}, (A + N'(y))(v - \tilde{v}) \rangle_{Y, Y^*} \\ &= \langle A^*v - A^*\tilde{v} + N'(y)^*v - N'(y)^*\tilde{v}, v - \tilde{v} \rangle_{Y^*, Y} \\ &= \langle \tilde{b} - A^*\tilde{v} - N'(\tilde{y})^*\tilde{v} + N'(\tilde{y})^*\tilde{v} - N'(y)^*\tilde{v}, v - \tilde{v} \rangle_{Y^*, Y} \\ &= \langle (N'(\tilde{y}) - N'(y))^*\tilde{v}, v - \tilde{v} \rangle_{Y^*, Y} \leq \|(N'(\tilde{y}) - N'(y))^*\tilde{v}\|_{Y^*} \|v - \tilde{v}\|_Y, \end{aligned}$$

cf. the proof of Theorem 3.12. This results in

$$\begin{aligned} \|(E_y(y, u)^{-*} - E_y(\tilde{y}, u)^{-*})\tilde{b}\|_Y &\leq \frac{1}{\underline{\kappa}} \|(N'(\tilde{y}) - N'(y))^*\tilde{v}\|_{Y^*} \\ &\leq \frac{1}{\underline{\kappa}^2} \|N'(\tilde{y}) - N'(y)\|_{\mathcal{L}(Y, Y^*)} \|\tilde{b}\|_{Y^*}, \end{aligned} \quad (5.12)$$

where the last estimate is due to the strong monotonicity of $A^* + N'(\tilde{y})^*$. Inserting (5.11) into (5.10) and using (5.12), we obtain (again using strong monotonicity):

$$\begin{aligned} \|\hat{z} - z\|_Y &\leq \frac{1}{\underline{\kappa}} \|Q^*Q(y - \tilde{y})\|_{Y^*} \\ &\quad + \frac{1}{\underline{\kappa}^2} \|N'(\tilde{y}) - N'(y)\|_{\mathcal{L}(Y, Y^*)} (\|Q^*(Q\tilde{y} - \hat{q})\|_{Y^*} + \|Q^*Q(y - \tilde{y})\|_{Y^*}). \end{aligned} \quad (5.13)$$

Combining this and (5.9) results in (5.8). \square

To bound the gradient error, we therefore have to control the error $\|B^*(\hat{z} - \tilde{z})\|_{U^*}$ caused by the inexact solution of the perturbed adjoint equation (5.6). If, e. g., $B \equiv \iota : L^2(\Omega) \hookrightarrow H^{-1}(\Omega)$, it is sufficient to control the adjoint state error in the $L^2(\Omega)$ -norm. This is no longer true if a boundary control problem is considered. We will therefore estimate $\|B^*(\hat{z} - \tilde{z})\|_{U^*} \leq \|B\|_{\mathcal{L}(U, Y^*)} \|\hat{z} - \tilde{z}\|_Y$ and control the error in the Y -norm. Another reason for this is that a posteriori error estimation techniques to estimate the $L^2(\Omega)$ -error require the PDE solution to have $H^2(\Omega)$ -regularity, see, e. g., [2, Sec. 2.4]. This cannot be guaranteed if the coefficient function $\kappa(\cdot, \xi)$ in the definition (3.10) of the operator $A(\xi)$ is only $L^\infty(\Omega)$ -regular, i. e., it can contain jumps along edges for example, or if the domain Ω is non-convex.

Moreover, we have to control the errors $\|Q(y - \tilde{y})\|_H$ or even $\|Q^*Q(y - \tilde{y})\|_{Y^*}$, see (5.13), and $\|N'(\tilde{y}) - N'(y)\|_{\mathcal{L}(Y, Y^*)}$ introduced by the inexact solution of the state equation. Again, if, e. g., $Q \equiv \iota : H_0^1(\Omega) \hookrightarrow L^2(\Omega)$, it is sufficient to control the state error in the $L^2(\Omega)$ -norm, which is no longer true if we have a problem with, e. g., boundary observation. Furthermore, we will see that an error estimate possibly in a stronger norm than the $L^2(\Omega)$ -norm is required to bound $\|N'(\tilde{y}) - N'(y)\|_{\mathcal{L}(Y, Y^*)}$.

The error in the computed state and adjoint state can be controlled by standard a posteriori techniques for elliptic PDEs. If N' is locally Lipschitz continuous w. r. t. y , bounding the error $\|N'(\tilde{y}) - N'(y)\|_{\mathcal{L}(Y, Y^*)}$ reduces to the estimation of the local Lipschitz constant $c_{N'}$ due to $\|N'(y, u) - N'(\tilde{y}, u)\|_{\mathcal{L}(Y, Y^*)} \leq c_{N'}\|\tilde{y} - y\|_Y$. But then, the state error in the possibly stronger Y -norm has to be estimated. For the example from Section 3.2, this local Lipschitz constant can be bounded as follows:

Lemma 5.2. *Let $N : Y \rightarrow Y^*$ be defined as in (3.10) (skipping the dependence on ξ) and let $\varphi : \mathbb{R} \rightarrow \mathbb{R}$ fulfill the respective conditions in Assumption 3.3. Then it holds that*

$$\|N'(\tilde{y}) - N'(y)\|_{\mathcal{L}(Y, Y^*)} \leq c_p^3 \left(a''_{\varphi''} \lambda(\Omega)^{(p-3)/p} + c''_{\varphi''} c_p^{p-3} (\|\tilde{y}\|_Y + \|y - \tilde{y}\|_Y)^{p-3} \right) \|y - \tilde{y}\|_Y, \quad (5.14)$$

where λ is the Lebesgue measure on Ω and $c_p > 0$ is the Sobolev constant such that $\|y\|_{L^p(\Omega)} \leq c_p \|y\|_{H_0^1(\Omega)}$ holds for every $y \in Y$.

Proof. We have that $\langle N'(y)v, \tilde{v} \rangle_{Y^*, Y} = \int_{\Omega} \varphi'(y)v\tilde{v} \, dx$ for $y, v, \tilde{v} \in Y$ and that $\varphi : \mathbb{R} \rightarrow \mathbb{R}$ is a twice continuously differentiable, increasing function, which fulfills the growth condition (3.7), i. e., $|\varphi''(t)| \leq a''_{\varphi''} + c''_{\varphi''}|t|^{p-3}$ with $a''_{\varphi''}, c''_{\varphi''} \geq 0$ and $p \in (3, \infty)$ for $n = 2$ and $p \in (3, 6]$ for $n = 3$. Thus, we can estimate with $r_i \in [1, \infty]$ ($i \in \{1, 2, 3, 4, 5\}$), $\frac{1}{r_1} + \frac{1}{r_2} + \frac{1}{r_3} = 1$, $\frac{1}{r_4} + \frac{1}{r_5} = \frac{1}{r_1}$, and $r_4(p-3) \geq 1$ (to be specified later):

$$\begin{aligned} & | \langle (N'(y) - N'(\tilde{y}))v, \tilde{v} \rangle_{Y^*, Y} | \\ & \leq \|\varphi'(y) - \varphi'(\tilde{y})\|_{L^{r_1}(\Omega)} \|v\|_{L^{r_2}(\Omega)} \|\tilde{v}\|_{L^{r_3}(\Omega)} \\ & \leq \int_0^1 \|\varphi''(\tilde{y} + \tau(y - \tilde{y}))\|_{L^{r_1}(\Omega)} \, d\tau \cdot \|v\|_{L^{r_2}(\Omega)} \|\tilde{v}\|_{L^{r_3}(\Omega)} \\ & \leq \left(a''_{\varphi''} \cdot \lambda(\Omega)^{1/r_4} + c''_{\varphi''} \cdot \sup_{\tau \in [0, 1]} \|\tilde{y} + \tau(y - \tilde{y})\|_{L^{r_4(p-3)}(\Omega)}^{p-3} \right) \\ & \quad \cdot \|y - \tilde{y}\|_{L^{r_5}(\Omega)} \|v\|_{L^{r_2}(\Omega)} \|\tilde{v}\|_{L^{r_3}(\Omega)} \\ & \leq c_{r_2} c_{r_3} \left(a''_{\varphi''} \cdot \lambda(\Omega)^{1/r_4} + c''_{\varphi''} \cdot (\max\{\|y\|_{L^{r_4(p-3)}(\Omega)}, \|\tilde{y}\|_{L^{r_4(p-3)}(\Omega)}\})^{p-3} \right) \\ & \quad \cdot \|y - \tilde{y}\|_{L^{r_5}(\Omega)} \|v\|_{H_0^1(\Omega)} \|\tilde{v}\|_{H_0^1(\Omega)}, \end{aligned} \quad (5.15)$$

where $c_{\hat{r}}$ is the constant from the Sobolev embedding $Y = H_0^1(\Omega) \hookrightarrow L^{\hat{r}}(\Omega)$ with adequately chosen $\hat{r} \in [1, \infty)$ or $\hat{r} \in [1, 6]$ dependent on n . Choosing $r_2 = r_3 = \hat{r}$ with $\hat{r} \in (2, \infty)$ for $n = 2$ and $\hat{r} \in (2, 6]$ for $n = 3$, $r_5 = \frac{\hat{r}(p-2)}{\hat{r}-2} = p-2 + \frac{2(p-2)}{\hat{r}-2} > p-2 > 1$, and $r_4 = \frac{r_5}{p-3} > 1 + \frac{1}{p-3}$,

this gives

$$\begin{aligned} & \|N'(y) - N'(\tilde{y})\|_{\mathcal{L}(Y, Y^*)} \\ & \leq c_{\tilde{r}}^2 \left(a''_{\varphi''} \cdot \lambda(\Omega)^{(p-3)/r_5} + c''_{\varphi''} \left(\max\{\|y\|_{L^{r_5}(\Omega)}, \|\tilde{y}\|_{L^{r_5}(\Omega)}\} \right)^{p-3} \right) \|y - \tilde{y}\|_{L^{r_5}(\Omega)}. \end{aligned} \quad (5.16)$$

The concrete choice $\tilde{r} = p$ in (5.16) (giving $r_5 = p$, $r_4 = \frac{p}{p-3}$) and the Sobolev embedding $Y \hookrightarrow L^p(\Omega)$ yield (5.14). \square

Remark 5.3. From (5.16) we see that it is enough to control the $L^{r_5}(\Omega)$ -error in the computed state. For $n = 2$ we can choose \tilde{r} arbitrarily large, but not $\tilde{r} = \infty$, giving that r_5 can be arbitrarily close to $p - 2$.

Criticality Measure Error

By Proposition 4.11, ensuring

$$\|\nabla m_k(0) - \nabla \hat{J}(u^k)\|_U \leq \min\left\{\frac{1-\mathbf{c}_s}{\tau} \varrho_c(\chi_k(0)), \varrho_g(\Delta_k)\right\}$$

as well as

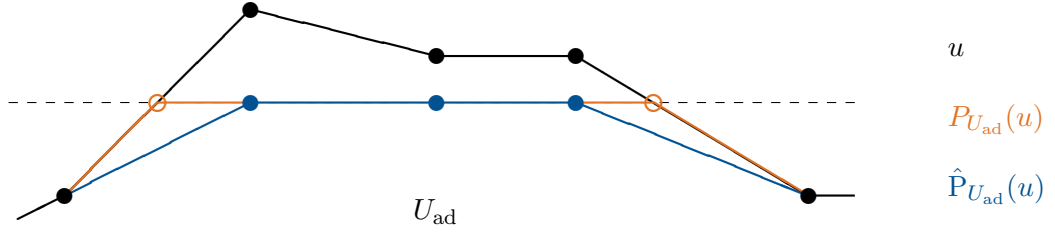
$$\|[P_{U_{\text{ad}}} - \hat{P}_{U_{\text{ad}}}] (u^k - \tau \nabla m_k(0))\|_U \leq \mathbf{c}_s \varrho_c(\chi_k(0))$$

for some constant $\mathbf{c}_s \in [0, 1]$ yields (4.10) and (4.11).

Thus, if we can control the error in the model gradient as discussed in Theorem 5.1, it remains to control the inexactness in the projection.

Remark 5.4. In certain cases, the difference between the exact and the approximate projection can be computed exactly:

- If $U = \mathbb{R}^{n_u}$ is finite dimensional and U_{ad} is a convex set for which the projection can be computed simply, e. g., a box or an ellipsoid, there is no need for introducing an approximate projection, i. e., $P_{U_{\text{ad}}} - \hat{P}_{U_{\text{ad}}} \equiv 0$.
- If $U = L^2(\Omega_u)$ with Ω_u being a measurable subset of Ω or $\partial\Omega$ and the discretization allows for the exact computation of norms, the projection onto a ball $U_{\text{ad}} := \{u \in U : \|u\|_U \leq \varrho\}$ with $\varrho > 0$ can also be computed exactly.
- The same holds true for $U = L^2(\Omega_u)$ and $U_{\text{ad}} := \{u \in U : u_l \leq u \leq u_u \text{ a. e.}\}$ with $u_l < u_u \in \mathbb{R}$ and a discretization by piecewise constant functions.
- Assume that $U = L^2(\Omega_u)$ with Ω_u being a subset of the domain Ω or its boundary $\partial\Omega$, and that continuous, linear finite elements with nodal bases are used for its discretization. Let U_{ad} be described by pointwise bound constraints with continuous, piecewise linear functions as bounds, which can be represented exactly in the current discretization. Then, the approximate projection is typically computed by pointwisely projecting the nodal function values onto the box. In [118] a method for computing the error exactly is presented. It is used that the L^2 -projection on the continuous level can be written in a pointwise fashion. The error is computed on each element of the FE


 Figure 5.1.: Exact and inexact L^2 -projection onto a box in 1D.

discretization separately and can be reduced by refining the elements with the largest error contribution. In Figure 5.1, the exact and the inexact projection of a linear FE function u onto a box with constant upper bound in 1D is depicted. It can be recognized that the error in the node-wise projection occurs exactly on two elements, namely the ones where the function u crosses the upper bound. This error can be computed exactly. The elements with the largest error contribution can be refined uniformly if a higher accuracy of the node-wise projection is required.

- If the projection is computed by solving

$$\bar{p} := P_{U_{\text{ad}}}(u) = \arg \min_{p \in U_{\text{ad}}} \frac{1}{2} \|u - p\|_U^2$$

approximately, the distance of an ε -solution $\hat{p} \in U_{\text{ad}}$ to $\bar{p} \in U_{\text{ad}}$ can be estimated as $\|\hat{p} - \bar{p}\|_U \leq \sqrt{2\varepsilon}$ due to

$$\begin{aligned} 2\varepsilon &\geq \|u - \hat{p}\|_U^2 - \|u - \bar{p}\|_U^2 = \|u - \bar{p} + \bar{p} - \hat{p}\|_U^2 - \|u - \bar{p}\|_U^2 \\ &= 2 \underbrace{(u - \bar{p}, \bar{p} - \hat{p})_U}_{\geq 0} + \|\hat{p} - \bar{p}\|_U^2. \end{aligned}$$

Note that in all cases where the projection can be computed exactly, we can set $c_P = 0$ in Proposition 4.11.

Objective Function Evaluation Error

We assume that the inexact reduced objective function \hat{J}_k is evaluated using an inexact solution $\tilde{y} \in Y$ of the state equation, i. e., $\hat{J}(u) = J(S(u), u)$ and $\hat{J}_k(u) = J(\tilde{y}, u)$ for some $\tilde{y} \in Y$. By Lemma 4.25, (5.3) holds if

$$|\hat{J}(u) - \hat{J}_k(u)| \leq \frac{1}{2} \varrho_r (\eta_3 \min\{\text{pred}_k, \mathfrak{r}_k\}) \text{ for all } u \in \{u^k, u^k + s^k\} \subset U.$$

Thus, it is enough to control the error $|\hat{J}(u) - \hat{J}_k(u)| = |J(y, u) - J(\tilde{y}, u)|$ for $u \in U$, where $y = S(u)$ is the exact solution of the state equation. If J is locally Lipschitz continuous w. r. t. y , it holds that $|J(y, u) - J(\tilde{y}, u)| \leq c_J \|y - \tilde{y}\|_Y$, and for error estimation we have to

estimate the error in the computed state and possibly the local Lipschitz constant c_J . For a tracking-type objective function, we have a more explicit estimate:

Proposition 5.5. *Let $J : Y \times U \rightarrow \mathbb{R}$ be of tracking type form (5.4) and let $y, \tilde{y} \in Y$, $u \in U$ be given. Then the following estimate holds true:*

$$|J(y, u) - J(\tilde{y}, u)| \leq \frac{1}{2} \|Q(y - \tilde{y})\|_H^2 + \|Q\tilde{y} - \hat{q}\|_H \|Q(y - \tilde{y})\|_H. \quad (5.17)$$

Proof. We compute

$$\begin{aligned} J(y, u) - J(\tilde{y}, u) &= \frac{1}{2} \|Qy - Q\tilde{y} + Q\tilde{y} - \hat{q}\|_H^2 - \frac{1}{2} \|Q\tilde{y} - \hat{q}\|_H^2 \\ &= \frac{1}{2} \|Q(y - \tilde{y})\|_H^2 + (Q(y - \tilde{y}), Q\tilde{y} - \hat{q})_H. \end{aligned}$$

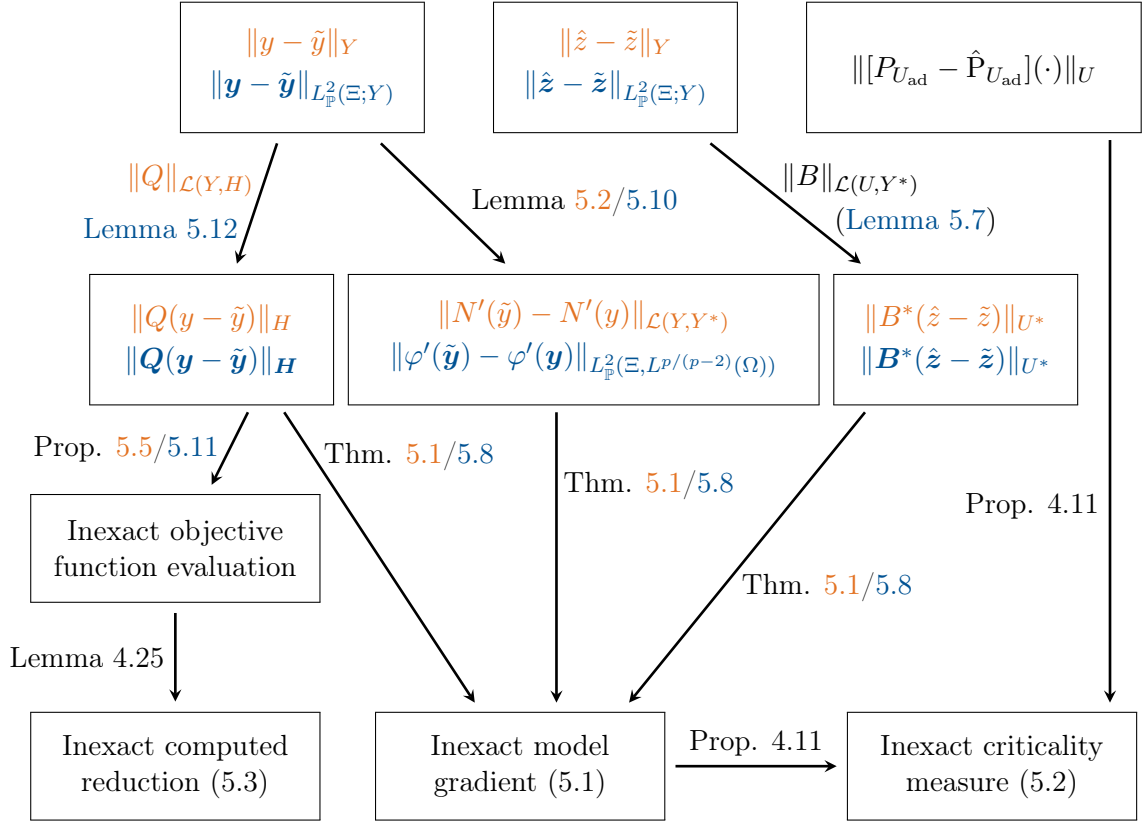
This shows (5.17). □

Combining Lemma 4.25 and Proposition 5.5, we see that we have to estimate the error $\|Q(y - \tilde{y})\|_H \leq \|Q\|_{\mathcal{L}(Y, H)} \|y - \tilde{y}\|_Y$, which again reduces to the question of error estimation for the computed state. An alternative to the estimate given in Proposition 5.5 would be the dual-weighted-residual method [18], which is well-suited for the estimation of the error in the objective function of an optimal control problem. We do not employ it here having our stochastic application in mind. We want to rely on already established error estimation techniques for PDEs with uncertain inputs, which can be implemented with low-rank tensors.

In conclusion, we only have to control a few errors to ensure (5.1), (5.2), and (5.3):

- The error $\|Q(y - \tilde{y})\|_H$ of the computed state observation. This error controls the accuracy of the objective function evaluation and the right-hand side in the perturbed adjoint equation. For the adjoint equation it would even be enough to control $\|Q^*Q(y - \tilde{y})\|_{Y^*}$.
- The error $\|B^*(\hat{z} - \tilde{z})\|_{U^*}$ introduced by the computed adjoint state. This influences the error in the computed gradient (together with the perturbation error of the adjoint equation). The accuracy of the computed gradient is also relevant for the accuracy of the computed criticality measure.
- The error $\|N'(y, u) - N'(\tilde{y}, u)\|_{\mathcal{L}(Y, Y^*)}$, which is relevant for the perturbed adjoint equation. A concrete estimator reducing to estimating $\|y - \tilde{y}\|_Y$ is given in Lemma 5.2.
- The error made by the discrete projection in the U -norm. In certain cases this error can be computed exactly and the U -grid can be refined appropriately.

Overall, everything can be reduced to estimating the errors in the inexact state and adjoint state as well as the error caused by the inexact projection. An overview of the error estimation procedure is given in Figure 5.2, where everything highlighted in orange corresponds to the deterministic case.



Additionally, the following quantities have to be computed or estimated:

$\|Q\tilde{\mathbf{y}} - \hat{q}\|_H$ (Theorem 5.1), $\|\tilde{\mathbf{y}}\|_Y$ (Lemma 5.2),
 $\|Q\tilde{\mathbf{y}} - \hat{q}\|_{L^2_{\mathbb{F}}(\Xi, H)}$ (Theorem 5.8), $\|\mathbf{y}\|_{L^\infty(\Xi; Y)}$ and $\|\tilde{\mathbf{y}}\|_{L^\infty(\Xi; Y)}$ (Lemma 5.10).

Figure 5.2.: Overview of the error estimation procedure for the deterministic problem and for the stochastic problem in the case $r_f = \infty$. Everything highlighted in orange corresponds to the deterministic case, whereas blue stands for the stochastic case.

5.2. Realization of the Error Estimates in the Stochastic Case

We now extend the considerations from Section 5.1 to the stochastic case, for the example from Section 3.2 with

$$\mathbf{E} : \mathbf{Y} \times U \rightarrow \mathbf{Y}^*, \quad \mathbf{E}(\mathbf{y}, u) = \mathbf{A}\mathbf{y} + \mathbf{N}(\mathbf{y}) - \mathbf{B}u - \mathbf{b}$$

defined in (3.11) and

$$\mathbf{J} : \mathbf{Y} \times U \rightarrow \mathbb{R}, \quad \mathbf{J}(\mathbf{y}, u) = \int_{\Xi} J[\xi](\mathbf{y}(\xi), u) \, d\mathbb{P} \quad (5.18)$$

from (3.4) and $J[\xi]$ from (5.4). In Section 3.4 we have already discussed that the adjoint state can be used to compute the gradient $\nabla \hat{\mathbf{J}}(u)$ of the reduced objective function, see (3.15).

Assumption 5.6. In this section we require the following ξ -regularities:

- $f \in L_{\mathbb{P}}^{r_f}(\Xi; L^2(\Omega))$, $\hat{\mathbf{q}} \in L_{\mathbb{P}}^{r_f}(\Xi; H)$ for some $r_f \in [p, \infty]$ with p from Assumption 3.3,
- $Q \in L_{\mathbb{P}}^{\infty}(\Xi; \mathcal{L}(Y, H))$, and
- \mathbf{B} is constant and defined as in (3.11).

Then we can ensure $\mathbf{y}, \mathbf{z} \in L_{\mathbb{P}}^{r_f}(\Xi; Y)$ by Corollary 3.15 and Lemma 3.23 and have the a priori bounds (3.9) and (3.20) on the state and the adjoint state. These estimates yield bounds for the respective $L_{\mathbb{P}}^{r_f}(\Xi; Y)$ -norm, which is at least as strong as the \mathbf{Y} -norm by Proposition A.2 because $r_f \geq p$.

Model Gradient Error

The model gradient is computed as described in Section 3.4, but with inexact solutions of the respective state and adjoint equations.

Quantities of the form $\|\mathbf{B}^*(\mathbf{z} - \tilde{\mathbf{z}})\|_{U^*}$ can be estimated as follows:

Lemma 5.7. Let $\mathbf{B} : U \rightarrow \mathbf{Y}^*$ be defined by

$$\langle \mathbf{B}u, \mathbf{v} \rangle_{\mathbf{Y}^*, \mathbf{Y}} = \int_{\Xi} \langle Bu, \mathbf{v}(\cdot, \xi) \rangle_{\mathbf{Y}^*, \mathbf{Y}} d\mathbb{P}$$

for some operator $B \in \mathcal{L}(U, Y^*)$ (cf. (3.11)) and let $\mathbf{z}, \tilde{\mathbf{z}} \in \mathbf{Y}$ be given. Then it holds that

$$\|\mathbf{B}^*(\mathbf{z} - \tilde{\mathbf{z}})\|_{U^*} \leq \|B\|_{\mathcal{L}(U, Y^*)} \cdot \|\tilde{\mathbf{z}} - \mathbf{z}\|_{L_{\mathbb{P}}^1(\Xi; Y)} \leq \|B\|_{\mathcal{L}(U, Y^*)} \cdot \|\tilde{\mathbf{z}} - \mathbf{z}\|_{L_{\mathbb{P}}^2(\Xi; Y)}.$$

Proof. For $\mathbf{v} \in \mathbf{Y}$ we have

$$\langle \mathbf{B}u, \mathbf{v} \rangle_{\mathbf{Y}^*, \mathbf{Y}} = \int_{\Xi} \langle Bu, \mathbf{v}(\cdot, \xi) \rangle_{\mathbf{Y}^*, \mathbf{Y}} d\mathbb{P} = \langle u, B^* \left(\int_{\Xi} \mathbf{v}(\cdot, \xi) d\mathbb{P} \right) \rangle_{U, U^*} = \langle u, \mathbf{B}^* \mathbf{v} \rangle_{U, U^*}.$$

We compute

$$\begin{aligned} \|\mathbf{B}^*(\mathbf{z} - \tilde{\mathbf{z}})\|_{U^*} &= \|B^* \left(\int_{\Xi} \tilde{\mathbf{z}}(\cdot, \xi) - \mathbf{z}(\cdot, \xi) d\mathbb{P} \right)\|_{U^*} \\ &\leq \|B\|_{\mathcal{L}(U, Y^*)} \cdot \int_{\Xi} \|(\tilde{\mathbf{z}}(\cdot, \xi) - \mathbf{z}(\cdot, \xi))\|_Y d\mathbb{P} \\ &\leq \|B\|_{\mathcal{L}(U, Y^*)} \cdot \|\tilde{\mathbf{z}} - \mathbf{z}\|_{L_{\mathbb{P}}^2(\Xi; Y)}, \end{aligned}$$

where the last inequality follows from Proposition A.2. \square

Theorem 5.8. Given Assumptions 3.3 and 5.6 and a control $u \in U$, let $\mathbf{y} = \mathbf{S}(u) \in L_{\mathbb{P}}^{r_f}(\Xi; Y)$ be the exact state (cf. Corollary 3.15) and let $\tilde{\mathbf{y}} \in L_{\mathbb{P}}^{r_f}(\Xi; Y)$. Furthermore, let $\mathbf{z} = \mathbf{T}(u) \in L_{\mathbb{P}}^{r_f}(\Xi; Y)$ be the exact adjoint state (cf. Lemma 3.23), let $\tilde{\mathbf{z}} \in L_{\mathbb{P}}^{r_f}(\Xi; Y)$ be the unique solution (cf. Proposition 3.24) of the perturbed adjoint equation

$$\mathbf{A}\mathbf{z} + \mathbf{N}'(\tilde{\mathbf{y}})\mathbf{z} = -\mathbf{Q}^*(\mathbf{Q}\tilde{\mathbf{y}} - \hat{\mathbf{q}})_{\mathbf{H}},$$

cf. (3.18) and (5.6), and let $\tilde{\mathbf{z}} \in L_{\mathbb{P}}^{r_f}(\Xi; Y)$.

Assume that $\mathbf{J} : \mathbf{Y} \times U \rightarrow \mathbb{R}$ is defined as in (5.18), and that $\hat{\mathbf{J}} : U \rightarrow \mathbb{R}$ is defined as in (3.15), and set $\mathbf{m}'(0) := -\mathbf{B}^* \tilde{\mathbf{z}} + \mathbf{J}_u(\tilde{\mathbf{y}}, u)$. Suppose that $\tilde{r} \in (2, \infty)$ (for $n = 2$) or even $\tilde{r} \in (2, 6]$ (for $n = 3$) and let $c_{\tilde{r}}$ be the Sobolev constant such that $\|y\|_{L^{\tilde{r}}(\Omega)} \leq c_{\tilde{r}} \|y\|_Y$ holds for every $y \in Y$. Let $p_1 \in [1, \frac{r_f}{p-2}]$, $p_2 \in [1, r_f]$ such that $\frac{1}{p_1} + \frac{1}{p_2} = 1$. Then the following estimate holds true:

$$\begin{aligned} \|\mathbf{m}'(0) - \hat{\mathbf{J}}'(u)\|_{U^*} &\leq \\ \|\mathbf{B}^*(\tilde{\mathbf{z}} - \hat{\mathbf{z}})\|_{U^*} + \frac{1}{\underline{\kappa}} \|B\|_{\mathcal{L}(U, Y^*)} \|Q\|_{L_{\mathbb{P}}^{\infty}(\Xi; \mathcal{L}(Y, H))} &\left(\|Q(\mathbf{y} - \tilde{\mathbf{y}})\|_{L_{\mathbb{P}}^1(\Xi, H)} \right. \\ &\left. + \frac{1}{\underline{\kappa}} c_{\tilde{r}}^2 \|\varphi'(\tilde{\mathbf{y}}) - \varphi'(\mathbf{y})\|_{L_{\mathbb{P}}^{p_1}(\Xi, L^{\tilde{r}/(\tilde{r}-2)}(\Omega))} (\|Q\tilde{\mathbf{y}} - \hat{\mathbf{q}}\|_{L_{\mathbb{P}}^{p_2}(\Xi, H)} + \|Q(\mathbf{y} - \tilde{\mathbf{y}})\|_{L_{\mathbb{P}}^{p_2}(\Xi, H)}) \right). \end{aligned} \quad (5.19)$$

Proof. Analogously to (5.9), it holds that

$$\|\mathbf{m}'(0) - \hat{\mathbf{J}}'(u)\|_{U^*} \leq \|\mathbf{B}^*(\tilde{\mathbf{z}} - \hat{\mathbf{z}})\|_{U^*} + \|\mathbf{B}^*(\mathbf{z} - \hat{\mathbf{z}})\|_{U^*}.$$

The second summand is estimated by Lemma 5.7: $\|\mathbf{B}^*(\mathbf{z} - \hat{\mathbf{z}})\|_{U^*} \leq \|B\|_{\mathcal{L}(U, Y^*)} \|\mathbf{z} - \hat{\mathbf{z}}\|_{L_{\mathbb{P}}^1(\Xi; Y)}$. Since \mathbf{z} and $\hat{\mathbf{z}}$ are defined pointwise (almost everywhere), the estimate (5.13) yields

$$\begin{aligned} \|\mathbf{z} - \hat{\mathbf{z}}\|_{L_{\mathbb{P}}^1(\Xi; Y)} &\leq \frac{1}{\underline{\kappa}} \|Q^* Q(\mathbf{y} - \tilde{\mathbf{y}})\|_{L_{\mathbb{P}}^1(\Xi; Y^*)} + \frac{1}{\underline{\kappa}^2} \|N'(\tilde{\mathbf{y}}) - N'(\mathbf{y})\|_{L_{\mathbb{P}}^{p_1}(\Xi; \mathcal{L}(Y, Y^*))} \\ &\cdot \left(\|Q^*(Q\tilde{\mathbf{y}} - \hat{\mathbf{q}})\|_{L_{\mathbb{P}}^{p_2}(\Xi; Y^*)} + \|Q^* Q(\mathbf{y} - \tilde{\mathbf{y}})\|_{L_{\mathbb{P}}^{p_2}(\Xi; Y^*)} \right). \end{aligned}$$

with $p_1, p_2 \in [1, \infty]$, $\frac{1}{p_1} + \frac{1}{p_2} = 1$. Using $\|N'(\tilde{\mathbf{y}}) - N'(\mathbf{y})\|_{\mathcal{L}(Y, Y^*)} \leq c_{\tilde{r}}^2 \|\varphi'(\tilde{\mathbf{y}}) - \varphi'(\mathbf{y})\|_{L^{\tilde{r}/(\tilde{r}-2)}(\Omega)}$ (cf. (5.15)) results in (5.19). The admissible values of p_1 and p_2 ensure together with the regularity of $\mathbf{y}, \tilde{\mathbf{y}}$ and the growth of φ' that every appearing quantity in (5.19) is finite. \square

Remark 5.9. The parameters p_1 and p_2 in Theorem 5.8 make different estimates involving $\|Q(\mathbf{y} - \tilde{\mathbf{y}})\|_{L_{\mathbb{P}}^{p_2}(\Xi; H)}$ with $2 \leq p_2 \leq r_f$ possible. For larger p_2 , a weaker norm w. r. t. ξ can be used to estimate the error in $N'(\tilde{\mathbf{y}})$.

We see that in general it is sufficient to control the adjoint state error in the $L_{\mathbb{P}}^2(\Xi; Y)$ -norm or even the $L_{\mathbb{P}}^1(\Xi; Y)$ -norm. Error control in a weaker norm can be sufficient depending on the example: For the concrete definition (3.11) of \mathbf{B} it holds that

$$\|\mathbf{B}^*(\mathbf{z} - \tilde{\mathbf{z}})\|_{U^*} = \|D^* \left(\int_{\Xi} \tilde{\mathbf{z}}(\cdot, \xi) - \mathbf{z}(\cdot, \xi) \, d\mathbb{P} \right)\|_{L^2(\Omega_u)},$$

identifying $L^2(\Omega_u)^* = L^2(\Omega_u)$. If, e. g., $D \equiv I : L^2(\Omega) \rightarrow L^2(\Omega)$, it would be sufficient to control the $L_{\mathbb{P}}^1(\Xi; L^2(\Omega))$ -error. This is no longer true if we consider, e. g., a boundary control problem. Thus, we will control the $L_{\mathbb{P}}^2(\Xi; Y)$ -error in the adjoint state to keep our algorithm flexible. Furthermore, this enables us to use the fact that the operator \mathbf{A} is strongly monotone with constant $\underline{\kappa}$ on $L_{\mathbb{P}}^2(\Xi; Y)$ (but not strongly monotone on $L_{\mathbb{P}}^p(\Xi; Y)$ for $p > 2$).

We will see that, in certain cases, it can be sufficient to control the $L_{\mathbb{P}}^2(\Xi; Y)$ -error in the computed state in order to bound the error caused by $\tilde{\mathbf{y}}$ entering the perturbed adjoint equation. Again, we need an estimate of the form

$$\|\varphi'(\tilde{\mathbf{y}}) - \varphi'(\mathbf{y})\|_{L_{\mathbb{P}}^{p_1}(\Xi; L^{\tilde{r}/(\tilde{r}-2)}(\Omega))} \leq C_{\varphi'} \|\tilde{\mathbf{y}} - \mathbf{y}\|_{L_{\mathbb{P}}^{p_3}(\Xi; Y)}$$

with a local Lipschitz constant $C_{\varphi'}$.

Lemma 5.10. *Let $\varphi : \mathbb{R} \rightarrow \mathbb{R}$ fulfill the respective conditions in Assumption 3.3, let $\mathbf{N} : \mathbf{Y} \rightarrow \mathbf{Y}^*$ be defined as in (3.11), and let $\mathbf{y}, \tilde{\mathbf{y}} \in \mathbf{Y}$ be given. Let p_1 be as in Theorem 5.8 and let $p_3, p_4 \in [1, \infty]$ such that $\frac{1}{p_3} + \frac{1}{p_4} = \frac{1}{p_1}$. Then it holds that*

$$\begin{aligned} & \|\varphi'(\tilde{\mathbf{y}}) - \varphi'(\mathbf{y})\|_{L_{\mathbb{P}}^{p_1}(\Xi; L^{p/(p-2)}(\Omega))} \leq \\ & c_p \|\tilde{\mathbf{y}} - \mathbf{y}\|_{L_{\mathbb{P}}^{p_3}(\Xi; Y)} \left(a''_{\varphi''} \lambda(\Omega)^{(p-3)/p} + c''_{\varphi''} c_p^{p-3} \max\{\|\mathbf{y}\|_{L_{\mathbb{P}}^{p_4(p-3)}(\Xi; Y)}, \|\tilde{\mathbf{y}}\|_{L_{\mathbb{P}}^{p_4(p-3)}(\Xi; Y)}\}^{p-3} \right). \end{aligned} \quad (5.20)$$

Proof. As in the proof of Lemma 5.2 we estimate

$$\begin{aligned} & \|\varphi'(\tilde{\mathbf{y}}) - \varphi'(\mathbf{y})\|_{L_{\mathbb{P}}^{p_1}(\Xi; L^{\tilde{r}/(\tilde{r}-2)}(\Omega))} \\ &= \left\| \int_0^1 \varphi''(\mathbf{y} + \tau(\tilde{\mathbf{y}} - \mathbf{y})) (\tilde{\mathbf{y}} - \mathbf{y}) \, d\tau \right\|_{L_{\mathbb{P}}^{p_1}(\Xi; L^{\tilde{r}/(\tilde{r}-2)}(\Omega))} \\ &\leq \|\tilde{\mathbf{y}} - \mathbf{y}\|_{L_{\mathbb{P}}^{p_3}(\Xi; L^{\tilde{r}}(\Omega))} \cdot \sup_{\tau \in [0,1]} \|\varphi''(\mathbf{y} + \tau(\tilde{\mathbf{y}} - \mathbf{y}))\|_{L_{\mathbb{P}}^{p_4}(\Xi; L^{\tilde{r}/(\tilde{r}-3)}(\Omega))} \end{aligned}$$

for some $\tilde{r} \in (3, \infty)$, using $\frac{1}{p_3} + \frac{1}{p_4} = \frac{1}{p_1}$. The second factor can be estimated as

$$\begin{aligned} & \sup_{\tau \in [0,1]} \|\varphi''(\mathbf{y} + \tau(\tilde{\mathbf{y}} - \mathbf{y}))\|_{L_{\mathbb{P}}^{p_4}(\Xi; L^{\tilde{r}/(\tilde{r}-3)}(\Omega))} \\ &\leq \sup_{\tau \in [0,1]} \|a''_{\varphi''} + c''_{\varphi''} |\mathbf{y} + \tau(\tilde{\mathbf{y}} - \mathbf{y})|^{p-3}\|_{L_{\mathbb{P}}^{p_4}(\Xi; L^{\tilde{r}/(\tilde{r}-3)}(\Omega))} \\ &= a''_{\varphi''} \cdot \lambda(\Omega)^{(\tilde{r}-3)/\tilde{r}} + c''_{\varphi''} \cdot \max\{\|\mathbf{y}\|_{L_{\mathbb{P}}^{p_4(p-3)}(\Xi; L^{r_5}(\Omega))}, \|\tilde{\mathbf{y}}\|_{L_{\mathbb{P}}^{p_4(p-3)}(\Xi; L^{r_5}(\Omega))}\}^{p-3} \end{aligned}$$

with $r_5 = \frac{\tilde{r}(p-3)}{\tilde{r}-3}$ by $|\varphi''(t)| \leq a''_{\varphi''} + c''_{\varphi''} |t|^{p-3}$ with $c''_{\varphi''} \geq 0$ and $p > 3$. If we choose $\tilde{r} = p$ and use $H_0^1(\Omega) \subset L^p(\Omega)$, we obtain (5.20). \square

Combining Theorem 5.8 and Lemma 5.10 and choosing $p_2 = p_3$, we see that we have to estimate the error $\|\tilde{\mathbf{y}} - \mathbf{y}\|_{L_{\mathbb{P}}^{p_2}(\Xi; Y)}$, compute the norm $\|\tilde{\mathbf{y}}\|_{L_{\mathbb{P}}^{p_4(p-3)}(\Xi; Y)}$ or bound it from above, and bound the norm $\|\mathbf{y}\|_{L_{\mathbb{P}}^{p_4(p-3)}(\Xi; Y)}$, e. g., by the a priori estimate (3.9).

For $r_f = \infty$ we can choose $p_2 = p_3 = 2$ and $p_4 = \infty$. Then it is enough to estimate the error in the inexact state in the $L_{\mathbb{P}}^2(\Xi; Y)$ -norm as long as we can compute or bound $\|\tilde{\mathbf{y}}\|_{L_{\mathbb{P}}^{\infty}(\Xi; Y)}$ and $\|\mathbf{y}\|_{L_{\mathbb{P}}^{\infty}(\Xi; Y)}$.

For $r_f < \infty$, we have to choose $p_4 \in [\frac{r_f}{r_f-2}, \frac{r_f}{p-3}]$ and $p_2 = p_3 = \frac{2p_4}{p_4-1} \in [\frac{2r_f}{r_f-p+3}, r_f]$.⁹ Observe that $p_1 = \frac{2p_4}{p_4+1} \leq \frac{2r_f}{r_f+p-3} \leq \frac{r_f}{p-\frac{3}{2}}$ follows then. On the one hand, for $r_f = p$ we can take $p_4 = \frac{p}{p-2}$. Then we have $p_2 = p_3 = p$ and the error in the computed state has to be estimated in the \mathbf{Y} -norm. On the other hand, for increasing r_f and using $p_4 = \frac{r_f}{p-3}$, the exponents p_2 and p_3 get close to 2. Then it is enough to estimate the error in the computed state in a weaker norm than the \mathbf{Y} -norm, but bounds on the exact and the inexact state have to be computed in a stronger norm.

⁹Note that $1 < \frac{r_f}{r_f-2} \leq \frac{r_f}{r_f-3} \leq \frac{r_f}{p-3}$ and $2 < \frac{2r_f}{r_f-p+3} \leq \frac{2}{3}r_f$ hold for $p \in (3, \infty)$, $r_f \in [p, \infty)$.

Criticality Measure Error

By Proposition 4.11, the control of the error in the approximate criticality measure reduces to controlling the inexactness of the approximate projection and the model gradient.

Objective Function Evaluation Error

Proposition 5.11. *Let $\mathbf{J} : \mathbf{Y} \times U \rightarrow \mathbb{R}$ be defined as in (3.4) and let $\mathbf{y}, \tilde{\mathbf{y}} \in Y$, $u \in U$ be given. Then it holds that*

$$|\mathbf{J}(\mathbf{y}, u) - \mathbf{J}(\tilde{\mathbf{y}}, u)| \leq \frac{1}{2} \|\mathbf{Q}(\mathbf{y} - \tilde{\mathbf{y}})\|_{L_{\mathbb{P}}^2(\Xi; H)}^2 + \|\mathbf{Q}\tilde{\mathbf{y}} - \hat{\mathbf{q}}\|_{L_{\mathbb{P}}^2(\Xi; H)} \|\mathbf{Q}(\mathbf{y} - \tilde{\mathbf{y}})\|_{L_{\mathbb{P}}^2(\Xi; H)}.$$

Proof. The estimate (5.17) (Proposition 5.5) and Hölder's inequality yield

$$\begin{aligned} |\mathbf{J}(\mathbf{y}, u) - \mathbf{J}(\tilde{\mathbf{y}}, u)| &= \left| \int_{\Xi} J[\xi](\mathbf{y}(\xi), u) - J[\xi](\tilde{\mathbf{y}}(\xi), u) \, d\mathbb{P} \right| \\ &\leq \int_{\Xi} \frac{1}{2} \|Q(\xi)(\mathbf{y}(\xi) - \tilde{\mathbf{y}}(\xi))\|_H^2 + \|Q(\xi)\tilde{\mathbf{y}}(\xi) - \hat{\mathbf{q}}(\xi)\|_H \|Q(\xi)(\mathbf{y}(\xi) - \tilde{\mathbf{y}}(\xi))\|_H \, d\mathbb{P} \\ &\leq \frac{1}{2} \|\mathbf{Q}(\mathbf{y} - \tilde{\mathbf{y}})\|_{L_{\mathbb{P}}^2(\Xi; H)}^2 + \|\mathbf{Q}\tilde{\mathbf{y}} - \hat{\mathbf{q}}\|_{L_{\mathbb{P}}^2(\Xi; H)} \|\mathbf{Q}(\mathbf{y} - \tilde{\mathbf{y}})\|_{L_{\mathbb{P}}^2(\Xi; H)}. \end{aligned}$$

□

We see that we have to compute or bound the $L_{\mathbb{P}}^2(\Xi; H)$ -norm of $\mathbf{Q}\tilde{\mathbf{y}} - \hat{\mathbf{q}}$ and have to estimate $\|\mathbf{Q}(\mathbf{y} - \tilde{\mathbf{y}})\|_{L_{\mathbb{P}}^2(\Xi; H)}$. If, e. g., $Q(\xi) \equiv \iota : H_0^1(\Omega) \hookrightarrow L^2(\Omega) = H$, it is enough to estimate the $L_{\mathbb{P}}^2(\Xi; L^2(\Omega))$ -norm of $\mathbf{y} - \tilde{\mathbf{y}}$. This is not true anymore if we consider a problem with boundary observation. Again, to have a flexible algorithm and to use strong monotonicity of \mathbf{A} , we estimate the $L_{\mathbb{P}}^2(\Xi; Y)$ -error and use:

Lemma 5.12. *Under Assumption 5.6,*

$$\|\mathbf{Q}(\mathbf{y} - \tilde{\mathbf{y}})\|_{L_{\mathbb{P}}^2(\Xi; H)} \leq \|Q\|_{L_{\mathbb{P}}^\infty(\Xi; \mathcal{L}(Y, H))} \|\mathbf{y} - \tilde{\mathbf{y}}\|_{L_{\mathbb{P}}^2(\Xi; Y)}$$

holds for all $\mathbf{y}, \tilde{\mathbf{y}} \in Y$.

Proof. This is a simple consequence of Hölder's inequality, cf. Proposition 3.1. It is important that Q has a higher regularity w. r. t. ξ than required for the operator $\mathbf{Q} : \mathbf{Y} \rightarrow \mathbf{H}$ to be well-defined. □

How the error in the computed state and adjoint state can be measured up to fixed, but possibly unknown constant factors, is discussed in the Chapter 7 for the example of a semilinear, elliptic PDE with stochastic coefficients from Section 3.2.

6. Discretization of the Model Problem

After discussing the functional analytic setting of the problem, a solution algorithm formulated in a Hilbert space, and the realization of the error estimates based on the error in the approximate solution of the state and the adjoint equation, it remains to discretize and adaptively solve the problem and the corresponding equations. The discretization is carried out almost exactly as in our paper [46], i. e., we use conforming finite element discretizations for the deterministic state and control spaces and polynomials for the spaces of random variables. Then the full tensor product of the respective finite-dimensional subspaces is built and used for a conforming stochastic Galerkin discretization. We use weighted Lagrange polynomials w. r. t. the Gaussian quadrature nodes as bases instead of the typical orthonormal polynomials with increasing degree. This has some useful consequences:

- These Lagrange polynomials are orthogonal and can be weighted such that they are orthonormal.
- A connection to stochastic collocation methods can be established, see [44, 46].
- Pointwise state constraints or certain nonsmooth risk measures can be handled by posing the constraint in every quadrature node in the discrete setting.
- Nonlinear dependence on the parameters and nonlinear operators can be approximated nicely. This was already done in [46, Example 3.6] and is used in this thesis for the discretization of the nonlinearity.

In [46], the discretization was fixed, but will now be adaptive. Hence, we will use sequences of nested discrete spaces with their respective bases and linear maps prolongating the coefficients in a coarser space to those in a finer space. For completeness, we review the discretization procedure described in [46] here, add some details about the discretization of all appearing nonlinearities, and point out necessary changes. The concrete adaptive approach is described in Chapter 7. Since constructing a finer space and prolongating the coefficients to it is a simple task for the considered FE and polynomial spaces, we describe only a fixed discretization in this chapter, having in mind that the mesh size and the degrees of the polynomials are parameters which will be adapted in the final implementation.

6.1. Space Discretization: Finite Elements

Assumption 6.1. For the discretization of the deterministic function spaces we assume that

- the open, bounded Lipschitz domain $\Omega \subset \mathbb{R}^2$ is polygonal and that
- the restriction $D|_U : U \rightarrow L^2(\Omega)$ of the given operator D to any used discrete subspace $U \subset U$ can be evaluated exactly.

Remark 6.2. We exclude the case $\Omega \subset \mathbb{R}^3$ here, which leads to a more compact presentation of a posteriori error estimates in Chapter 7. It is possible to generalize many results to the 3D case. The polygonal domain $\Omega \subset \mathbb{R}^2$ can be covered by a finite element triangulation so that we do not have to care about variational crimes caused by domain approximation.

The domain Ω is partitioned into a finite element mesh yielding a triangulation \mathcal{T} . Let $Y \subset Y$ denote the discrete subspace of piecewise linear, globally continuous finite element functions with zero boundary data. Generalizations of this would be possible, but are skipped here for the ease of presentation and implementation. The nodal FE basis of the discrete, deterministic state space is $\{\phi_{k_0}\}_{k_0=1}^{d_0} \subset Y \subset Y = H_0^1(\Omega)$ and the basis of discrete control space $U \subset U = L^2(\Omega_u)$ is denoted by $\{\psi_\ell\}_{\ell=1}^{d_u}$. This can also be a nodal FE basis if Ω_u is a subset of Ω with positive measure or—for finite-dimensional controls—the standard basis of $\mathbb{R}^{d_u} = L^2([d_u])$. In the latter case, no discretization is required. The basis functions shall sum to one, i. e., $\sum_{k_0=1}^{d_0} \phi_{k_0}(x) = 1$ and $\sum_{k_u=1}^{d_u} \psi_{k_u}(\tilde{x}) = 1$ for all $x \in \Omega$, $\tilde{x} \in \Omega_u$ to perform mass lumping in a meaningful way later.

In addition to assuming that $D|_U : L^2(\Omega_u) \rightarrow L^2(\Omega)$ can be evaluated exactly (Assumption 6.1), we match the discretizations of the state and the control space, i. e., we use the same grid on Ω and Ω_u if Ω_u is a subset of Ω with positive measure. This makes the computation of the gradient of the reduced objective function easier since typically $\hat{J}'(u) = B^*z + \gamma(u, \cdot)_U$. Then it is desirable that B^*z and u share the same grid. This is also important to be able to perform gradient-based updates of the control in an optimization method.

We define the following matrices:

- the stiffness matrix $K \in \mathbb{R}^{d_0 \times d_0}$ for Y : $K_{k_0 l_0} := (\nabla \phi_{k_0}, \nabla \phi_{l_0})_{L^2(\Omega)^n}$,
- the mass matrix $M \in \mathbb{R}^{d_0 \times d_0}$ for Y : $M_{k_0 l_0} := (\phi_{k_0}, \phi_{l_0})_{L^2(\Omega)}$,
- the lumped mass matrix $M_L \in \mathbb{R}^{d_0 \times d_0}$ for Y : $(M_L)_{k_0 k_0} := \sum_{l_0=1}^{d_0} M_{k_0 l_0} = \int_{\Omega} \phi_{k_0} dx$ and $(M_L)_{k_0 l_0} = 0$ for $k_0 \neq l_0$,
- the mass matrix $\tilde{M} \in \mathbb{R}^{d_u \times d_u}$ for U : $\tilde{M}_{k_u l_u} := (\psi_{k_u}, \psi_{l_u})_{L^2(\Omega_u)}$,
- the lumped mass matrix $\tilde{M}_L \in \mathbb{R}^{d_u \times d_u}$ for U : $(\tilde{M}_L)_{k_u k_u} := \sum_{l_u=1}^{d_u} \tilde{M}_{k_u l_u} = \int_{\Omega_u} \psi_{k_u} d\tilde{x}$ and $(\tilde{M}_L)_{k_u l_u} = 0$ for $k_u \neq l_u$.

Let $y \in \mathbb{R}^{d_0}$ and $u \in \mathbb{R}^{d_u}$ be the coefficients, representing the discrete state y and control u , respectively. Inserting $y(x) = \sum_{k_0=1}^{d_0} y_{k_0} \phi_{k_0}(x)$ and $u(x) = \sum_{k_u=1}^{d_u} u_{k_u} \psi_{k_u}(x)$ into (3.8), and testing with $v \equiv \phi_{k_0}$ for $k_0 \in [d_0]$, the discrete version of the deterministic state equation reads

$$A(\xi)y + N(y) = Bu + b(\xi) \quad (6.1)$$

with

$$\begin{aligned} A(\xi) &\in \mathbb{R}^{d_0 \times d_0}, & A_{k_0 l_0}(\xi) &= (\kappa(\cdot, \xi) \nabla \phi_{l_0}, \nabla \phi_{k_0})_{L^2(\Omega)^n}, \\ N : \mathbb{R}^{d_0} &\rightarrow \mathbb{R}^{d_0}, & N(y) &= M_L \varphi(y), \\ B &\in \mathbb{R}^{d_0 \times d_u}, & B_{k_0 l_u} &= (D\psi_{l_u}, \phi_{k_0})_{L^2(\Omega)}, \\ b(\xi) &\in \mathbb{R}^{d_0}, & b_{k_0}(\xi) &= (f(\xi), \phi_{k_0})_{L^2(\Omega)}. \end{aligned} \quad (6.2)$$

We will refer to $\mathbf{A}(\xi)$ as the *system matrix* in contrast to the *stiffness matrix*, which is induced by the $H_0^1(\Omega)$ -norm. Due to the ease of implementation and interpolation, a quadrature error is allowed to occur in the discretization of the nonlinearity which is connected to mass lumping: The integral $\int_{\Omega} \varphi(y) \phi_{k_0} dx$ is evaluated inexactly by a quadrature formula, the nodes of which are the finite element grid nodes and the weights of which are the respective entries of the lumped mass matrix. This quadrature formula is exact for integrals over a single linear finite element function. We obtain

$$\int_{\Omega} \varphi(y) \phi_{k_0} dx \approx \varphi(y_{k_0}) (\mathbf{M}_L)_{k_0 k_0}. \quad (6.3)$$

Since the trust-region algorithm (Algorithm 1) is formulated in the infinite-dimensional space U , it is desirable to not make additional errors by mass lumping in the objective function, but to evaluate the U - and H -inner product exactly. Let H be discretized such that $Q(\xi)|_Y$ can be evaluated exactly and such that $\hat{q}(\xi)$ can be represented exactly. The discrete subspace \mathbb{H} of H is isomorphic to \mathbb{R}^{d_H} ($d_H \in \mathbb{N}$) equipped with the inner product induced by the symmetric, positive definite matrix $\mathbf{M}_H \in \mathbb{R}^{d_H \times d_H}$. Let $\mathbf{Q}(\xi) \in \mathbb{R}^{d_H \times d_0}$ and $\hat{\mathbf{q}}(\xi) \in \mathbb{R}^{d_H}$ be the discrete versions of $Q(\xi)$ and $\hat{q}(\xi)$, respectively. Then, the discretized objective function from (3.4) reads

$$\mathbf{J}[\xi](y, \mathbf{u}) = \frac{1}{2} \|\mathbf{Q}(\xi)y - \hat{\mathbf{q}}(\xi)\|_{\mathbf{M}_H}^2 + \frac{\gamma}{2} \mathbf{u}^\top \tilde{\mathbf{M}} \mathbf{u}. \quad (6.4)$$

We note that under the stated assumptions, the evaluation of the objective function is exact so that the error in the reduced objective function depends only on the error in the discretized state and Proposition 5.5 can be applied. The discrete version of the deterministic adjoint equation (3.16) is

$$\mathbf{A}(\xi)z(\xi) + \mathbf{M}_L(\varphi'(y(\xi)) \odot z(\xi)) = -\mathbf{Q}(\xi)^\top \mathbf{M}_H(\mathbf{Q}(\xi)y(\xi) - \hat{\mathbf{q}}(\xi)), \quad (6.5)$$

where again the quadrature formula using the finite element nodes has been applied. In fact, $\mathbf{N}'(y)z = \mathbf{M}_L(\varphi'(y) \odot z)$ holds also in the discrete setting and we can identify $\mathbf{N}'(y) = \mathbf{M}_L \text{diag}(\varphi'(y))$. The gradient of the reduced, deterministic, discretized objective function is then given by

$$\nabla \hat{\mathbf{J}}[\xi](\mathbf{u}) = -\tilde{\mathbf{M}}^{-1} \mathbf{B}^\top z(\xi) + \gamma \mathbf{u}, \quad (6.6)$$

cf. (3.21). Note that in typical situations it is not necessary to invert the mass matrix $\tilde{\mathbf{M}}$ to compute the reduced gradient: If, e. g., $\Omega_u \subset \Omega$ is a subset of positive measure, $B^* : H_0^1(\Omega) \rightarrow L^2(\Omega_u)$, $z \mapsto z|_{\Omega_u}$ is the canonical embedding $\iota : H_0^1(\Omega) \hookrightarrow L^2(\Omega)$ combined with restriction of the function to Ω_u , and the grids on Ω and Ω_u match, the application of $\tilde{\mathbf{M}}^{-1} \mathbf{B}$ consists of a simple extraction of components of the vector $z(\xi)$ and/or adding zero components for the nodes on $\partial\Omega$. Such situations are favorable because the error in the discrete gradient then only depends on the error in the discrete adjoint state and it is sufficient to apply the error estimate from Theorem 5.1.

Once the equations (3.25) and (3.24) are discretized, the application of the Hessian operator to a direction can be computed via (3.23). For this purpose, it only remains to discretize the term in (3.24) involving the second derivative N_{yy} , which is again done by using the

FE nodes based quadrature. Then, if $\mathbf{s} \in \mathbb{R}^{d_u}$ represents a direction $\mathbf{s} \in \mathbf{U}$, and $\mathbf{y}(\xi)$, $\mathbf{z}(\xi)$ represent the current state and adjoint state, respectively, the application of the Hessian to this direction reads

$$\nabla^2 \hat{J}[\xi](\mathbf{u})\mathbf{s} = \tilde{\mathbf{M}}^{-1} \mathbf{B}^\top \mathbf{h}(\xi) + \gamma \mathbf{s}, \quad (6.7)$$

where $\mathbf{h}(\xi)$ solves

$$[\mathbf{A}(\xi) + \mathbf{N}'(\mathbf{y}(\xi))]\mathbf{h}(\xi) = \mathbf{Q}(\xi)^\top \mathbf{M}_H \mathbf{Q}(\xi) \mathbf{d}(\xi) + \mathbf{M}_L (\mathbf{z}(\xi) \odot \varphi''(\mathbf{y}(\xi)) \odot \mathbf{d}(\xi))$$

with $\mathbf{d}(\xi) = [\mathbf{A}(\xi) + \mathbf{N}'(\mathbf{y}(\xi))]^{-1} \mathbf{B} \mathbf{s}$.

6.2. Stochastic Discretization: Polynomial Chaos

We proceed with the discretization of the space $L_{\mathbb{P}}^p(\Xi)$ of vectors of independent random variables distributed on $\Xi := \times_{i=1}^m \Xi_i$ with the probability measure $\mathbb{P} := \otimes_{i=1}^m \mathbb{P}_i$.

Assumption 6.3. For the discretization of these spaces we assume the following:

- The sets $\Xi_i \subset \mathbb{R}$ are open and bounded¹⁰ intervals for all $i \in [m]$.
- Each probability measure \mathbb{P}_i does not consist of finitely many atoms such that discretization is really necessary.

For $i \in \{1, \dots, m\}$ and a fixed $d \in \mathbb{N}^m$, the spaces $L_{\mathbb{P}_i}^p(\Xi_i)$ are discretized by polynomials of degree $d_i - 1$. By Assumption 6.3, all polynomials of arbitrary degree defined on Ξ_i are \mathbb{P}_i -integrable and the space of polynomials is dense in $L_{\mathbb{P}_i}^p(\Xi_i)$, see [106, Chap. 8]. Furthermore, there exist sets $\{\beta_{k_i}^{(i)}\}_{k_i=1}^\infty \subset L_{\mathbb{P}_i}^2(\Xi_i)$ of orthonormal polynomials w. r. t. the $L_{\mathbb{P}_i}^2(\Xi_i)$ -inner product, where $\beta_{k_i}^{(i)}$ has degree $k_i - 1$ by [106, Thm. 8.5]. These sets are Hilbert bases of $L_{\mathbb{P}_i}^2(\Xi_i)$, respectively, and can be constructed by applying the Grad-Schmidt process to the monomial basis $\{1, \xi_i, \xi_i^2, \dots\}$ for example. We want to mention that some papers [40, 8] dealing with uncertainty quantification restrict the discussion to certain probability distributions, for which the “classical” orthonormal polynomials of increasing degree are well-known. Important examples are the Legendre polynomials for the uniform distribution, the Hermite polynomials for the normal distribution and the Jacobi polynomials for the beta distribution [106, Example 8.2]. In order to have a more general setting, we only assume that the orthonormal polynomials can be constructed and evaluated, e. g., by the three-term recurrence relation [106, Sec. 8.2]. In particular, we do not assume a purely continuous [44, 108, 28, 68] or a symmetric [38] distribution. Defining $\beta_k(\xi) := \prod_{i=1}^m \beta_{k_i}^{(i)}(\xi_i)$ we obtain a set $\{\beta_k\}_{k \in \mathbb{N}^m}$ of orthonormal polynomials which form a Hilbert basis of $L_{\mathbb{P}}^2(\Xi)$, see [116, Thm. 3.12(b)] and Section 2.2. Note that k is an index vector.

Let $\{\mathbf{a}_{l_i}^{(i)}\}_{l_i=1}^{d_i} \subset \Xi_i$ be the d_i pairwise distinct roots of the polynomial $\beta_{d_i+1}^{(i)}$ in ascending order, respectively. They exist and have the mentioned properties due to [106, Thm. 8.16] and are known as Gaussian quadrature nodes. Let $\{\mathbf{w}_{l_i}^{(i)}\}_{l_i=1}^{d_i}$ be the positive Gaussian quadrature

¹⁰Boundedness of the intervals is assumed because then the set of polynomials of arbitrary degree is a dense subset of $L_{\mathbb{P}_i}^p(\Xi_i)$. This condition can possibly be relaxed, cf. [106, Chap. 8].

weights associated to the nodes $\{\mathbf{a}_{l_i}^{(i)}\}_{l_i=1}^{d_i}$ for $i \in [m]$ and $l_i \in [d_i]$ defined by integration over the respective Lagrange polynomials, see also [106, Def. 9.2]. We use the Gaussian quadrature nodes to define weighted Lagrange polynomials $\{\theta_{k_i}^{(i)}\}_{k_i=1}^{d_i}$, which fulfill $\theta_{k_i}^{(i)}(\mathbf{a}_{l_i}^{(i)}) = \delta_{k_i l_i} \cdot \omega_{k_i}^{(i)}$ for some weights $\omega_{k_i}^{(i)} > 0$. If we choose $\omega_{k_i}^{(i)} = (\mathbf{w}_{k_i}^{(i)})^{-1/2}$ for $k_i \in \{1, \dots, d_i\}$, it follows from the exactness of Gaussian quadrature, that these weighted Lagrange polynomials are also orthonormal:

$$\int_{\Xi_i} \theta_{k_i}^{(i)} \theta_{l_i}^{(i)} d\mathbb{P}_i = \sum_{\hat{l}_i=1}^{d_i} \mathbf{w}_{\hat{l}_i}^{(i)} \theta_{k_i}^{(i)}(\mathbf{a}_{\hat{l}_i}^{(i)}) \theta_{l_i}^{(i)}(\mathbf{a}_{\hat{l}_i}^{(i)}) = \sum_{\hat{l}_i=1}^{d_i} \mathbf{w}_{\hat{l}_i}^{(i)} \delta_{k_i \hat{l}_i} (\mathbf{w}_{k_i}^{(i)})^{-1/2} \delta_{l_i \hat{l}_i} (\mathbf{w}_{l_i}^{(i)})^{-1/2} = \delta_{k_i l_i}.$$

The product of two Lagrange polynomials has degree $2d_i - 2$ and is integrated exactly by Gaussian quadrature, which is exact up to degree $2d_i - 1$ [106, Thm. 9.9]. Defining $\theta_{k_i}^{(i)} := \beta_{k_i}^{(i)}$ for $k_i \geq d_i + 1$, we get that $\{\theta_{k_i}^{(i)}\}_{k_i=1}^{\infty}$ is also a Hilbert basis of $L_{\mathbb{P}_i}^2(\Xi_i)$. As before, the Hilbert basis $\{\theta_k\}_{k \in \mathbb{N}^m}$ with $\theta_k(\xi) := \prod_{i=1}^m \theta_{k_i}^{(i)}(\xi_i)$ can be defined. Writing $\mathbf{a}_l := (\mathbf{a}_{l_1}^{(1)}, \dots, \mathbf{a}_{l_m}^{(m)})^\top$ and $\omega_l := \prod_{i=1}^m \omega_{l_i}^{(i)}$, it holds that $\theta_k(\mathbf{a}_l) = \delta_{kl} \cdot \omega_l$ for $k \leq d$ componentwise.

The state space $\mathbf{Y} = L_{\mathbb{P}}^p(\Xi; Y) \cong Y \otimes L_{\mathbb{P}_1}^p(\Xi_1) \otimes \dots \otimes L_{\mathbb{P}_m}^p(\Xi_m)$ is discretized by the full tensor product of the respective finite-dimensional subspaces with the basis

$$\{\phi_{k_0} \otimes \theta_{k_1}^{(1)} \otimes \dots \otimes \theta_{k_m}^{(m)}, k_i \in [d_i], i \in \{0, \dots, m\}\}.$$

This is a basis due to [54, Lem. 3.11], see Remark 2.7. A function $\mathbf{y} \in \mathbf{Y}$ belonging to the finite-dimensional space is represented by a coefficient tensor $\mathbf{y} \in \mathbb{R}^{d_0 \times \dots \times d_m}$ corresponding to weighted values of the function \mathbf{y} since nodal FE ansatz functions and Lagrange polynomials are used. This means that

$$\mathbf{y}(x, \xi_1, \dots, \xi_m) = \sum_{k_0=1}^{d_0} \sum_{k_1=1}^{d_1} \dots \sum_{k_m=1}^{d_m} \mathbf{y}(k_0, k_1, \dots, k_m) \phi_{k_0}(x) \theta_{k_1}^{(1)}(\xi_1) \dots \theta_{k_m}^{(m)}(\xi_m) \quad (6.8)$$

and in particular

$$\mathbf{y}(x, \mathbf{a}_{k_1}^{(1)}, \dots, \mathbf{a}_{k_m}^{(m)}) = \sum_{k_0=1}^{d_0} \mathbf{y}(\overset{\circ}{k}) \phi_{k_0}(x) \omega(k_1, \dots, k_m), \quad (6.9)$$

where we abbreviate $\overset{\circ}{k} = (k_0, k_1, \dots, k_m)$ here, as well as $\overset{\circ}{l} = (l_0, l_1, \dots, l_m)$ and $(\overset{\circ}{k}, \overset{\circ}{l}) = (k_0, k_1, \dots, k_m, l_0, l_1, \dots, l_m)$ etc. in the following in contrast to $l = (l_1, \dots, l_m)$ etc. The weight tensors $\omega, \mathbf{w} \in \mathbb{R}^{n_1 \times \dots \times n_d}$ are defined by $\omega(k) := \prod_{i=1}^m \omega_{k_i}^{(i)}$ and $\mathbf{w}(k) := \prod_{i=1}^m \mathbf{w}_{k_i}^{(i)}$ and obviously have rank 1.

The discretization of the space $L_{\mathbb{P}_1}^p(\Xi_1) \otimes \dots \otimes L_{\mathbb{P}_m}^p(\Xi_m)$ is obtained from the above considerations by setting $Y = \mathbb{R}$, $d_0 = 1$ and $\phi_1 = 1$. Let \mathbf{p} be a polynomial of *coordinate degree* $d - 1$ belonging to the finite-dimensional subspace $\mathcal{P}_{d-1}(\Xi) \subset L_{\mathbb{P}}^p(\Xi)$, i. e., \mathbf{p} is a polynomial of degree $d_i - 1$ in the variable ξ_i for all $i \in [d]$. We want to emphasize that this function is

represented by a tensor $\mathbf{p} \in \mathbb{R}^{d_1 \times \dots \times d_m}$ given by

$$\mathbf{p}(k_1, \dots, k_m) = \mathbf{p}(\mathbf{a}_{k_1}^{(1)}, \dots, \mathbf{a}_{k_m}^{(m)}) \boldsymbol{\omega}(k_1, \dots, k_m)^{-1} \quad (6.10)$$

due to (6.9). The standard $L_{\mathbb{P}}^2(\Xi)$ -inner product of two functions $\mathbf{p}, \tilde{\mathbf{p}}$ is discretized by applying Gaussian quadrature:

$$\begin{aligned} (\mathbf{p}, \tilde{\mathbf{p}})_{L_{\mathbb{P}}^2(\Xi)} &= \int_{\Xi} \mathbf{p} \tilde{\mathbf{p}} \, d\mathbb{P} = \sum_{k_1=1}^{d_1} \dots \sum_{k_m=1}^{d_m} \mathbf{w}(k_1, \dots, k_m) \mathbf{p}(\mathbf{a}_{k_1}^{(1)}, \dots, \mathbf{a}_{k_m}^{(m)}) \tilde{\mathbf{p}}(\mathbf{a}_{k_1}^{(1)}, \dots, \mathbf{a}_{k_m}^{(m)}) \\ &= \langle \mathbf{w} \odot \boldsymbol{\omega} \odot \mathbf{p}, \boldsymbol{\omega} \odot \tilde{\mathbf{p}} \rangle =: \langle \mathbf{p}, \tilde{\mathbf{p}} \rangle_{\mathbf{w} \odot \boldsymbol{\omega}^2}. \end{aligned} \quad (6.11)$$

We see that for the special choice of orthonormal polynomials with $\boldsymbol{\omega} = \mathbf{w}^{-1/2}$ the inner product of two functions is computed by a simple Frobenius inner product of tensors. A special case is the expectation of a function \mathbf{p} , computed as an inner product with the function $\mathbb{1}$, which is constant one and represented by the tensor $\boldsymbol{\omega}^{-1}$.

We now discretize the operators defined in (3.11) by testing with $\mathbf{v} \in \mathbf{Y}$ represented by $\mathbf{v} \in \mathbb{R}^{d_0 \times \dots \times d_m}$ and using the discretized deterministic operators (6.2), see [46]. This gives

$$\begin{aligned} \langle \mathbf{B}\mathbf{u}, \mathbf{v} \rangle_{\mathbf{Y}^*, \mathbf{Y}} &= \int_{\Xi} (D\mathbf{u}, \mathbf{v}(\cdot, \xi))_{L^2(\Omega)} \, d\mathbb{P} \\ &= \sum_{k_0, \dots, k_m=1}^{d_0, \dots, d_m} \mathbf{v}(k_0, \dots, k_m) (D\mathbf{u}, \phi_{k_0})_{L^2(\Omega)} \int_{\Xi} \prod_{i=1}^m \theta_{k_i}^{(i)}(\xi_i) \, d\mathbb{P} \\ &= \langle (\mathbf{B}\mathbf{u}) \otimes (\mathbf{w} \odot \boldsymbol{\omega}), \mathbf{v} \rangle =: \langle \mathbf{B}\mathbf{u}, \mathbf{v} \rangle_{\mathbb{1} \otimes (\mathbf{w} \odot \boldsymbol{\omega}^2)} \end{aligned}$$

with $\mathbf{B}\mathbf{u} = (\mathbf{B}\mathbf{u}) \otimes \boldsymbol{\omega}^{-1}$ and

$$\begin{aligned} \langle \mathbf{b}, \mathbf{v} \rangle_{\mathbf{Y}^*, \mathbf{Y}} &= \int_{\Xi} (f(\xi), \mathbf{v}(\cdot, \xi))_{L^2(\Omega)} \, d\mathbb{P} \\ &= \sum_{k_0, \dots, k_m=1}^{d_0, \dots, d_m} \mathbf{v}(k_0, \dots, k_m) \int_{\Xi} \mathbf{b}_{k_0}(\xi) \prod_{i=1}^m \theta_{k_i}^{(i)}(\xi_i) \, d\mathbb{P} =: \langle \mathbf{b}, \mathbf{v} \rangle_{\mathbb{1} \otimes (\mathbf{w} \odot \boldsymbol{\omega}^2)} \end{aligned}$$

with

$$\mathbf{b}(\mathring{k}) = \frac{1}{\prod_{i=1}^m \mathbf{w}_{k_i}^{(i)} (\boldsymbol{\omega}_{k_i}^{(i)})^2} \left(\int_{\Xi} \mathbf{b}_{k_0}(\xi) \prod_{i=1}^m \theta_{k_i}^{(i)}(\xi_i) \, d\mathbb{P} \right).$$

We assume that this tensor can be constructed and has sufficiently small rank. In particular, if $\mathbf{b}_{k_0}(\cdot)$ are polynomials of total degree at most d , Gaussian quadrature, which is exact up to coordinate degree $2d - 1$, can be applied to compute

$$\int_{\Xi} \mathbf{b}_{k_0}(\xi) \prod_{i=1}^m \theta_{k_i}^{(i)}(\xi_i) \, d\mathbb{P} = \sum_{l_1, \dots, l_m=1}^{d_1, \dots, d_m} \mathbf{w}(l) \mathbf{b}_{k_0}(\mathbf{a}_l) \prod_{i=1}^m \theta_{k_i}^{(i)}(\mathbf{a}_{l_i}^{(i)}) = \mathbf{w}(k) \boldsymbol{\omega}(k) \mathbf{b}_{k_0}(\mathbf{a}_k) \quad (6.12)$$

so that $\mathbf{b}(\cdot, k) = \boldsymbol{\omega}(k)^{-1} \mathbf{b}(\mathbf{a}_k)$ holds. Furthermore,

$$\begin{aligned} \langle \mathbf{A}\mathbf{y}, \mathbf{v} \rangle_{\mathbf{Y}^*, \mathbf{Y}} &= \int_{\Xi} (\kappa(\cdot, \xi) \nabla_x \mathbf{y}(\cdot, \xi), \nabla_x \mathbf{v}(\cdot, \xi))_{L^2(\Omega)^n} \, d\mathbb{P} \\ &= \sum_{k_0, \dots, k_m=1}^{d_0, \dots, d_m} \sum_{l_0, \dots, l_m=1}^{d_0, \dots, d_m} \mathbf{y}(\mathring{l}) \mathbf{v}(\mathring{k}) \int_{\Xi} (\mathbf{A}(\xi))_{k_0 l_0} \prod_{i=1}^m \theta_{l_i}^{(i)}(\xi_i) \theta_{k_i}^{(i)}(\xi_i) \, d\mathbb{P} \\ &= \langle \tilde{\mathbf{a}}, \mathbf{y} \rangle_{(m+2, \dots, 2m+2), [m+1]}, \mathbf{v} \rangle_{\mathbb{1} \otimes (\boldsymbol{\omega} \odot \boldsymbol{\omega}^2)} =: \langle \mathbf{A}\mathbf{y}, \mathbf{v} \rangle_{\mathbb{1} \otimes (\boldsymbol{\omega} \odot \boldsymbol{\omega}^2)} \end{aligned} \quad (6.13)$$

holds with the tensor $\tilde{\mathbf{a}} \in \mathbb{R}^{d_0 \times d_1 \times \dots \times d_m \times d_0 \times d_1 \times \dots \times d_m}$ defined by

$$\tilde{\mathbf{a}}(\mathring{k}, \mathring{l}) := \frac{1}{\prod_{i=1}^m \mathbf{w}_{k_i}^{(i)}(\boldsymbol{\omega}_{k_i}^{(i)})^2} \left(\int_{\Xi} (\mathbf{A}(\xi))_{k_0 l_0} \prod_{i=1}^m \theta_{l_i}^{(i)}(\xi_i) \theta_{k_i}^{(i)}(\xi_i) \, d\mathbb{P} \right). \quad (6.14)$$

The nonlinear part of the equation is approximated by Gaussian quadrature/interpolation to simplify the implementation. We will see later that this relates the whole approach to stochastic collocation.

$$\begin{aligned} \langle \mathbf{N}(\mathbf{y}), \mathbf{v} \rangle_{\mathbf{Y}^*, \mathbf{Y}} &= \int_{\Xi} \int_{\Omega} \varphi(\mathbf{y}(x, \xi)) \mathbf{v}(x, \xi) \, dx \, d\mathbb{P} \\ &= \sum_{k_0, \dots, k_m=1}^{d_0, \dots, d_m} \mathbf{v}(k_0, \dots, k_m) \int_{\Xi} \int_{\Omega} \varphi(\mathbf{y}(x, \xi)) \phi_{k_0}(x) \, dx \prod_{i=1}^m \theta_{k_i}^{(i)}(\xi_i) \, d\mathbb{P} \\ &\approx \sum_{k_0, \dots, k_m=1}^{d_0, \dots, d_m} \mathbf{v}(\mathring{k}) \sum_{l_1, \dots, l_m=1}^{d_1, \dots, d_m} \mathbf{w}(l) \int_{\Omega} \varphi(\mathbf{y}(x, \mathbf{a}_l)) \phi_{k_0}(x) \, dx \prod_{i=1}^m \theta_{k_i}^{(i)}(\mathbf{a}_{l_i}^{(i)}) \\ &= \sum_{k_0, \dots, k_m=1}^{d_0, \dots, d_m} \mathbf{v}(\mathring{k}) \mathbf{w}(k) \boldsymbol{\omega}(k) \int_{\Omega} \varphi \left(\sum_{l_0=1}^{d_0} \mathbf{y}(l_0, k_1, \dots, k_m) \phi_{l_0}(x) \boldsymbol{\omega}(k) \right) \phi_{k_0}(x) \, dx \\ &\approx \sum_{k_0, \dots, k_m=1}^{d_0, \dots, d_m} \mathbf{v}(\mathring{k}) \mathbf{w}(k) \boldsymbol{\omega}(k) \varphi \left(\mathbf{y}(\mathring{k}) \boldsymbol{\omega}(k) \right) (\mathbf{M}_L)_{k_0 k_0} \\ &= \langle (\mathbb{1} \otimes \boldsymbol{\omega}^{-1}) \odot (\mathbf{M}_L \circ_1 \varphi((\mathbb{1} \otimes \boldsymbol{\omega}) \odot \mathbf{y})), \mathbf{v} \rangle_{\mathbb{1} \otimes (\boldsymbol{\omega} \odot \boldsymbol{\omega}^2)} \\ &= \langle \hat{\boldsymbol{\omega}}^{-1} \odot (\mathbf{M}_L \circ_1 \varphi(\hat{\boldsymbol{\omega}} \odot \mathbf{y})), \mathbf{v} \rangle_{\mathbb{1} \otimes (\boldsymbol{\omega} \odot \boldsymbol{\omega}^2)} =: \langle \mathbf{N}(\mathbf{y}), \mathbf{v} \rangle_{\mathbb{1} \otimes (\boldsymbol{\omega} \odot \boldsymbol{\omega}^2)}, \end{aligned} \quad (6.15)$$

where the first approximate equality is due to Gaussian quadrature in the parameter space and the second one is due to the FE nodes based quadrature, which is related to mass lumping, see (6.3). We define the rank-1-tensor $\hat{\boldsymbol{\omega}} := \mathbb{1} \otimes \boldsymbol{\omega}$ for readability purposes. Here we see the importance of using mass lumping and weighted Lagrange polynomials: It yields that the nonlinear function $\varphi : \mathbb{R} \rightarrow \mathbb{R}$ can be applied componentwise to a tensor in the discrete setting (6.15).

The discrete state equation reads

$$\mathbf{A}\mathbf{y} + \mathbf{N}(\mathbf{y}) = \mathbf{B}\mathbf{u} + \mathbf{b}. \quad (6.16)$$

Assumption 6.4. In order to be able to evaluate the objective function (3.4) on the finite-dimensional subspace exactly in a simple way, we need the following additional assumptions:

- The operator $\mathbf{Q}(\xi) = \mathbf{Q}$ is constant.
- The desired state $\hat{\mathbf{q}}(\cdot)$ is a polynomial of coordinate degree at most $d - 1$.

The discrete version of the objective function is

$$\begin{aligned} \mathbf{J}(\mathbf{y}, \mathbf{u}) &:= \int_{\Xi} \mathbf{J}[\xi] \left(\sum_{k_1, \dots, k_m=1}^{d_1, \dots, d_m} \mathbf{y}(\cdot, k_1, \dots, k_m) \theta_{k_1}^{(1)}(\xi_1) \cdots \theta_{k_m}^{(m)}(\xi_m), \mathbf{u} \right) d\mathbb{P} \\ &= \int_{\Xi} \frac{1}{2} \left\| \sum_{k_1, \dots, k_m=1}^{d_1, \dots, d_m} \mathbf{Q}(\xi) \mathbf{y}(\cdot, k_1, \dots, k_m) \theta_{k_1}^{(1)}(\xi_1) \cdots \theta_{k_m}^{(m)}(\xi_m) - \hat{\mathbf{q}}(\xi) \right\|_{\mathbf{M}_H}^2 d\mathbb{P} + \frac{\gamma}{2} \mathbf{u}^\top \tilde{\mathbf{M}} \mathbf{u}, \end{aligned} \quad (6.17)$$

cf. (6.4). Since $\mathbf{Q}(\xi) = \mathbf{Q}$ is constant and $\hat{\mathbf{q}}(\cdot)$ is a polynomial of coordinate degree at most $d - 1$ (Assumption 6.4), the integrand in (6.17) has degree at most $2d - 2 \cdot 1$ and the integral can be evaluated exactly by Gaussian quadrature. This gives

$$\begin{aligned} \mathbf{J}(\mathbf{y}, \mathbf{u}) &= \sum_{l_1, \dots, l_m}^{d_1, \dots, d_m} \mathbf{w}(l) \frac{1}{2} \left\| \mathbf{Q} \mathbf{y}(\cdot, l_1, \dots, l_m) - \hat{\mathbf{q}}(\mathbf{a}_l) \right\|_{\mathbf{M}_H}^2 + \frac{\gamma}{2} \mathbf{u}^\top \tilde{\mathbf{M}} \mathbf{u} \\ &=: \langle \mathbf{M}_H (\mathbf{Q} \mathbf{y} - \hat{\mathbf{q}}), \mathbf{Q} \mathbf{y} - \hat{\mathbf{q}} \rangle + \frac{\gamma}{2} \mathbf{u}^\top \tilde{\mathbf{M}} \mathbf{u} \end{aligned} \quad (6.18)$$

with $\mathbf{Q} \mathbf{y} = \mathbf{Q} \circ_1 \mathbf{y}$, $\hat{\mathbf{q}} \in \mathbb{R}^{d_H \times d_1 \times \cdots \times d_m}$ defined by $\hat{\mathbf{q}}(\cdot, l) = \hat{\mathbf{q}}(\mathbf{a}_l) \boldsymbol{\omega}(l)^{-1}$, and

$$\mathbf{M}_H : \mathbb{R}^{d_H \times d_1 \times \cdots \times d_m} \rightarrow \mathbb{R}^{d_H \times d_1 \times \cdots \times d_m}, \quad \mathbf{M}_H \hat{\mathbf{q}} := (\mathbb{1} \otimes (\mathbf{w} \odot \boldsymbol{\omega}^2)) \odot (\mathbf{M}_H \circ_1 \hat{\mathbf{q}}). \quad (6.19)$$

Again, this evaluation is exact on the discrete subspace and the objective function evaluation error can be estimated using Proposition 5.11.

For the computation of the gradient of the reduced objective function, it remains to discretize the adjoint equation (3.18). The adjoint state \mathbf{z} is represented by the tensor \mathbf{z} analogously to the state. We have

$$\begin{aligned} \langle \mathbf{N}'(\mathbf{y}) \mathbf{z}, \mathbf{v} \rangle_{\mathbf{Y}^*, \mathbf{Y}} &= \int_{\Xi} \int_{\Omega} \varphi'(\mathbf{y}(x, \xi)) \mathbf{z}(x, \xi) \mathbf{v}(x, \xi) dx d\mathbb{P} \\ &\approx \sum_{k_0, \dots, k_m=1}^{d_0, \dots, d_m} \mathbf{v}(\overset{\circ}{k}) \sum_{l_1, \dots, l_m=1}^{d_1, \dots, d_m} \mathbf{w}(l) \int_{\Omega} \varphi'(\mathbf{y}(x, \mathbf{a}_l)) \mathbf{z}(x, \mathbf{a}_l) \phi_{k_0}(x) dx \prod_{i=1}^m \theta_{\overset{\circ}{k}_i}^{(i)}(\mathbf{a}_{l_i}^{(i)}) \\ &\approx \sum_{k_0, \dots, k_m=1}^{d_0, \dots, d_m} \mathbf{v}(\overset{\circ}{k}) \mathbf{w}(k) \boldsymbol{\omega}(k) \varphi'(\mathbf{y}(\overset{\circ}{k}) \boldsymbol{\omega}(k)) \mathbf{z}(\overset{\circ}{k}) \boldsymbol{\omega}(k) (\mathbf{M}_L)_{k_0 k_0} \\ &= \langle \hat{\boldsymbol{\omega}}^{-1} \odot (\mathbf{M}_L \circ_1 (\varphi'(\hat{\boldsymbol{\omega}} \odot \mathbf{y}) \odot (\hat{\boldsymbol{\omega}} \odot \mathbf{z}))), \mathbf{v} \rangle_{\mathbb{1} \otimes (\mathbf{w} \odot \boldsymbol{\omega}^2)} \\ &= \langle \mathbf{N}'(\mathbf{y}) \mathbf{z}, \mathbf{v} \rangle_{\mathbb{1} \otimes (\mathbf{w} \odot \boldsymbol{\omega}^2)}, \end{aligned}$$

cf. (6.15), where \mathbf{N}' is exactly the derivative of the discretized operator \mathbf{N} . Furthermore, the right-hand side is

$$\begin{aligned}
 & \langle -\mathbf{Q}^*(\mathbf{Q}\mathbf{y} - \hat{\mathbf{q}}), \mathbf{v} \rangle_{\mathbf{Y}^*, \mathbf{Y}} \\
 &= \sum_{k_1, \dots, k_m=1}^{d_1, \dots, d_m} \mathbf{v}(\cdot, k)^\top \int_{\Xi} \left(-\mathbf{Q}^\top \mathbf{M}_H \left(\sum_{l_1, \dots, l_m=1}^{d_1, \dots, d_m} \mathbf{Q}\mathbf{y}(\cdot, l) \prod_{i=1}^m \theta_{l_i}^{(i)}(\xi_i) - \hat{\mathbf{q}}(\xi) \right) \right) \prod_{i=1}^m \theta_{k_i}^{(i)}(\xi_i) \, d\mathbb{P} \\
 &= \sum_{k_1, \dots, k_m=1}^{d_1, \dots, d_m} \mathbf{v}(\cdot, k)^\top \mathbf{w}(k) \boldsymbol{\omega}(k) \left(-\mathbf{Q}^\top \mathbf{M}_H \left(\mathbf{Q}\mathbf{y}(\cdot, k) \boldsymbol{\omega}(k) - \hat{\mathbf{q}}(\mathbf{a}_k) \right) \right) \\
 &= \langle -(\mathbf{Q}^\top \mathbf{M}_H) \circ_1 (\mathbf{Q}\mathbf{y} - \hat{\mathbf{q}}), \mathbf{v} \rangle_{\mathbb{1} \otimes (\mathbf{w} \odot \boldsymbol{\omega}^2)}, \tag{6.20}
 \end{aligned}$$

cf. (6.5). Given the solution \mathbf{z} of the discretized adjoint equation

$$\mathbf{A}\mathbf{z} + \mathbf{N}'(\mathbf{y})\mathbf{z} = -(\mathbf{Q}^\top \mathbf{M}_H) \circ_1 (\mathbf{Q}\mathbf{y} - \hat{\mathbf{q}}), \tag{6.21}$$

the gradient of the reduced, discretized objective function reads

$$\nabla \hat{\mathbf{J}}(\mathbf{u}) = -(\tilde{\mathbf{M}}^{-1} \mathbf{B}^\top) \langle \mathbf{z}, \mathbf{w} \odot \boldsymbol{\omega} \rangle_{[m]+1, [m]} + \gamma \mathbf{u}, \tag{6.22}$$

cf. (6.6). Due to $\mathbf{w} \odot \boldsymbol{\omega} = (\mathbf{w} \odot \boldsymbol{\omega}^2) \odot \boldsymbol{\omega}^{-1}$, the term $\langle \mathbf{z}, \mathbf{w} \odot \boldsymbol{\omega} \rangle_{[m]+1, [m]}$ corresponds to computing the $L_{\mathbb{P}}^2(\Xi)$ -inner product (induced by $\mathbf{w} \odot \boldsymbol{\omega}^2$) of \mathbf{z} and the function $\mathbb{1}$ (represented by $\boldsymbol{\omega}^{-1}$), i. e., computing the expectation of \mathbf{z} , see also (6.11) and the paragraph below this equation. In [46] we used differently scaled formulations of (6.21) and (6.22), which coincide to the ones given here in the case $\mathbf{w} = \boldsymbol{\omega}^{-1/2}$ (orthonormal polynomials). Here the form (6.21) of the adjoint equation is closer to the continuous setting as derived in (6.20). As in (6.6), the matrix $(\tilde{\mathbf{M}}^{-1} \mathbf{B}^\top)$ can often be applied without inverting the mass matrix $\tilde{\mathbf{M}}$ so that Theorem 5.8 can be applied to estimate the gradient error.

In analogy to (6.7), we obtain the discrete Hessian

$$\nabla^2 \hat{\mathbf{J}}(\mathbf{u})\mathbf{s} = (\tilde{\mathbf{M}}^{-1} \mathbf{B}^\top) \langle \mathbf{h}, \mathbf{w} \odot \boldsymbol{\omega} \rangle_{[m]+1, [m]} + \gamma \mathbf{s}, \tag{6.23}$$

where $\mathbf{h} \in \mathbb{R}^{d_0 \times \dots \times d_m}$ solves

$$[\mathbf{A} + \mathbf{N}'(\mathbf{y})]\mathbf{h} = (\mathbf{Q}^\top \mathbf{M}_H) \circ_1 (\mathbf{Q}\mathbf{d}) + \hat{\boldsymbol{\omega}}^{-1} \odot (\mathbf{M}_L \circ_1 ((\hat{\boldsymbol{\omega}} \odot \mathbf{z}) \odot \varphi''(\hat{\boldsymbol{\omega}} \odot \mathbf{y}) \odot (\hat{\boldsymbol{\omega}} \odot \mathbf{d}))) \tag{6.24}$$

with $\mathbf{d} = [\mathbf{A} + \mathbf{N}'(\mathbf{y})]^{-1} \mathbf{B}\mathbf{s}$ and the current state \mathbf{y} and adjoint state \mathbf{z} . The second summand at the right hand side of (6.24) is exactly the derivative $\mathbf{N}''(\mathbf{y})$ applied to \mathbf{z} and \mathbf{d} . In (6.24) we could get rid of the componentwise multiplication by $\hat{\boldsymbol{\omega}}^{-1}$ by, e. g., dropping the multiplication of \mathbf{d} by $\hat{\boldsymbol{\omega}}$. But we want to emphasize that this approach is related to interpolation: The tensor $\hat{\boldsymbol{\omega}} \odot \mathbf{y}$ for example contains exactly the function values of the state at the FE nodes and Gaussian quadrature grid points. Therefore, pointwise operations, such as the application of φ'' or the multiplication of functions, carry over to componentwise operations on tensors in the discrete setting. Multiplying a tensor of functions values by $\hat{\boldsymbol{\omega}}^{-1}$ again transforms it back to the representation with weighted Lagrange polynomials.

Application of the Operators to Low-Rank Tensors

To make the solution of the state and the adjoint equation and further computations efficient, we represent all tensors, such as \mathbf{y} and \mathbf{z} , in a low-rank format (TT or HT, see Subsection 2.1.2). As already noted in [40], this is related to reduced basis methods [29, 28]. Let for example $\mathbf{y} \in \mathbb{R}^{d_0 \times \dots \times d_m}$ be given in the HT format. Then each leaf matrix $\mathbf{U}_i \in \mathbb{R}^{d_i \times r_i}$ (see Subsection 2.1.2) contains a generating system of a lower-dimensional subspace of \mathbb{R}^{d_i} , typically even a basis of an r_i -dimensional subspace. These bases can be interpreted as bases of subspaces of the discrete spaces \mathbf{Y} , $\mathcal{P}_{d_1-1}(\Xi_1)$, \dots , $\mathcal{P}_{d_m-1}(\Xi_m)$, respectively.

Equations (6.16) and (6.21) will be solved by a low-rank tensor solver such as AMEn, which chooses the tensor rank adaptively. In particular, the basis \mathbf{U}_0 for the subspace of the deterministic discrete state space \mathbf{Y} may be adapted by the low-rank tensor solver during the solution process and it may even change its size. This is in contrast to the standard reduced basis method, where the basis is chosen a priori based on solution snapshots. A second difference is that also reduced bases for the polynomial spaces $\mathcal{P}_{d_i-1}(\Xi_i)$ appear in the low-rank tensor.

In order to be able to apply such a low-rank tensor solver, it is important that all operators are efficiently applicable to low-rank tensors. This shall be discussed in the following. We see that the nonlinear function \mathbf{N} as well as its derivatives \mathbf{N}' , \mathbf{N}'' can be applied to low-rank tensors as long as the elementwise application of the functions φ , φ' , and φ'' can be implemented efficiently because the rest of the computations consists of multiplications with rank-1-tensors, i -mode matrix multiplications, and componentwise multiplications to apply $\mathbf{N}'(\mathbf{y})$ and $\mathbf{N}''(\mathbf{y})$. In particular, the componentwise multiplication by a rank-1-tensor can be written as i -mode multiplications with diagonal matrices. In our case, it holds that

$$\mathbf{N}(\mathbf{y}) = \mathbf{M}_L \circ_1 \text{diag}(\omega^{(1), -1}) \circ_2 \dots \text{diag}(\omega^{(m), -1}) \circ_{m+1} \varphi(\text{diag}(\omega^{(1)}) \circ_2 \dots \text{diag}(\omega^{(m)}) \circ_{m+1} \mathbf{y}). \quad (6.25)$$

Analogously, the operator \mathbf{M}_H defined in (6.19) can be implemented efficiently using i -mode matrix products only. More i -mode matrix products, outer and inner products appear in form of the operator \mathbf{Q} used in (6.18) and (6.21), the operator \mathbf{B} , and in (6.20), (6.22), and (6.23). There, the contraction with a rank-1-tensor can be carried out by

$$\langle \mathbf{z}, \mathbf{w} \odot \boldsymbol{\omega} \rangle_{[m]+1, [m]} = ((\mathbf{w}^{(1)} \odot \boldsymbol{\omega}^{(1)})^\top \circ_2 \dots (\mathbf{w}^{(m)} \odot \boldsymbol{\omega}^{(m)})^\top \circ_{m+1} \mathbf{z})(\bullet, 1, \dots, 1),$$

where the indexing at the end is done to remove the tensor modes of dimension one. It only remains to investigate the application of the operator \mathbf{A} to low-rank tensors. Then all operations in (6.16), (6.21), (6.22), and (6.23) can be realized in an efficient way.

By (6.13), the application of the operator \mathbf{A} can be performed as contraction with the tensor $\tilde{\mathbf{a}}$ defined in (6.14). If this tensor can be represented or approximated in the used low-rank format, this contraction can be computed easily. More concretely, it can also represent a TT matrix as described in Subsection 2.1.3. The contraction $\langle \tilde{\mathbf{a}}, \mathbf{y} \rangle_{(m+2, \dots, 2m+2), [m+1]}$ is exactly the application of a TT matrix to a TT tensor as depicted in Figure 2.6. Another favorable option is having the operator \mathbf{A} of small CP rank. We can achieve this by using additional structure.

Example 6.5 ([46, Example 3.3, Example 3.6]). In light of Example 3.17 and Lemma 3.18, we consider the specific form $\kappa(x, \xi) = \kappa_0(x)(1 + \sum_{i=1}^m \xi_i \eta_i(x))$ of the coefficient function. The definitions of the operators $A_0, A_i \in \mathcal{L}(Y, Y^*)$ from Lemma 3.18 carry over to the discrete setting via

$$\begin{aligned} \mathbf{A}_0 &\in \mathbb{R}^{d_0 \times d_0}, & (\mathbf{A}_0)_{k_0 l_0} &:= (\kappa_0 \nabla \phi_{l_0}, \nabla \phi_{k_0})_{L^2(\Omega)^n}, \\ \mathbf{A}_i &\in \mathbb{R}^{d_0 \times d_0}, & (\mathbf{A}_i)_{k_0 l_0} &:= (\eta_i \kappa_0 \nabla \phi_{l_0}, \nabla \phi_{k_0})_{L^2(\Omega)^n}, \end{aligned}$$

cf. (6.2), and it holds that $\mathbf{A}(\xi) = \mathbf{A}_0 + \sum_{i=1}^m \xi_i \mathbf{A}_i$. Writing $\theta_{l_j}^{(j)} \theta_{k_j}^{(j)} = \theta_{l_j}^{(j)}(\xi_j) \theta_{k_j}^{(j)}(\xi_j)$ for brevity, we obtain

$$\begin{aligned} &\int_{\Xi} (\mathbf{A}(\xi))_{k_0 l_0} \prod_{j=1}^m \theta_{l_j}^{(j)}(\xi_j) \theta_{k_j}^{(j)}(\xi_j) \, d\mathbb{P} \\ &= (\mathbf{A}_0)_{k_0 l_0} \prod_{j=1}^m \int_{\Xi_j} \theta_{l_j}^{(j)} \theta_{k_j}^{(j)} \, d\mathbb{P}_j + \sum_{i=1}^m (\mathbf{A}_i)_{k_0 l_0} \int_{\Xi_i} \xi_i \theta_{l_i}^{(i)} \theta_{k_i}^{(i)} \, d\mathbb{P}_i \prod_{j \neq i} \int_{\Xi_j} \theta_{l_j}^{(j)} \theta_{k_j}^{(j)} \, d\mathbb{P}_j \\ &= (\mathbf{A}_0)_{k_0 l_0} \prod_{j=1}^m \mathbf{w}_{l_j}^{(j)} (\omega_{l_j}^{(j)})^2 \delta_{l_j k_j} + \sum_{i=1}^m (\mathbf{A}_i)_{k_0 l_0} \mathbf{a}_{l_i}^{(i)} \mathbf{w}_{l_i}^{(i)} (\omega_{l_i}^{(i)})^2 \delta_{l_i k_i} \prod_{j \neq i} \mathbf{w}_{l_j}^{(j)} (\omega_{l_j}^{(j)})^2 \delta_{l_j k_j} \\ &= \prod_{j=1}^m \mathbf{w}_{l_j}^{(j)} (\omega_{l_j}^{(j)})^2 \left((\mathbf{A}_0)_{k_0 l_0} + \sum_{i=1}^m (\mathbf{A}_i)_{k_0 l_0} \mathbf{a}_{l_i}^{(i)} \right) \prod_{j=1}^m \delta_{l_j k_j} \end{aligned}$$

using that Gaussian quadrature is exact up to degree $2d_i - 1$. Therefore, (6.14) becomes

$$\tilde{\mathbf{a}}(\overset{\circ}{k}, \overset{\circ}{l}) = \left((\mathbf{A}_0)_{k_0 l_0} + \sum_{i=1}^m (\mathbf{A}_i)_{k_0 l_0} \mathbf{a}_{l_i}^{(i)} \right) \prod_{j=1}^m \delta_{l_j k_j}.$$

This gives

$$\begin{aligned} \mathbf{A} \mathbf{y} &= \langle \tilde{\mathbf{a}}, \mathbf{y} \rangle_{(m+2, \dots, 2m+2), [m+1]} = \sum_{l_0, \dots, l_m=1}^{d_0, \dots, d_m} \tilde{\mathbf{a}}(\overset{\circ}{k}, \overset{\circ}{l}) \mathbf{y}(\overset{\circ}{l}) \\ &= \sum_{l_0=1}^{d_0} \left((\mathbf{A}_0)_{k_0 l_0} + \sum_{i=1}^m (\mathbf{A}_i)_{k_0 l_0} \mathbf{a}_{l_i}^{(i)} \right) \mathbf{y}(l_0, k_1, \dots, k_m) \\ &= \mathbf{A}_0 \circ_1 \mathbf{y} + \sum_{i=1}^m \mathbf{A}_i \circ_1 \tilde{\mathbf{S}}_i \circ_{i+1} \mathbf{y} \end{aligned}$$

with $\tilde{\mathbf{S}}_i = \text{diag}(\mathbf{a}^{(i)})$ in analogy to Lemma 3.18. We see that the operator \mathbf{A} admits CP rank $m + 1$ and can be implemented using simple i -mode matrix multiplications and summation.

In [46] an additional example is given, where a similar structure is derived for an operator with additional domain parametrizations via interpolation. As mentioned at the end of Section 3.2, we skip this example in this thesis because the a posteriori error estimator derived in Chapter 7 will also rely on the structure of the coefficient function used in Example 6.5.

Relation to Stochastic Collocation

With the structure of the operator \mathbf{A} from Example 6.5 and the approximation of the non-linearity (6.15) and the right-hand side (6.12) by Gaussian quadrature, we obtain that the described stochastic Galerkin discretization of the state equation (6.16) yields in fact a completely decoupled system, one nonlinear equation for each combination $(\mathbf{a}_{k_1}^{(1)}, \dots, \mathbf{a}_{k_m}^{(m)})$ of parameter realizations:

$$\begin{aligned} (\mathbf{A}\mathbf{y})(\bullet, k) &= (\mathbf{A}_0 + \sum_{i=1}^m \mathbf{a}_{k_i}^{(i)} \mathbf{A}_i) \mathbf{y}(\bullet, k) = \boldsymbol{\omega}(k)^{-1} \mathbf{A}(\mathbf{a}_{k_i}^{(i)}) (\boldsymbol{\omega}(k) \mathbf{y}(\bullet, k)), \\ (\mathbf{N}(\mathbf{y}))(\bullet, k) &= (\hat{\boldsymbol{\omega}} \cdot^{-1} \odot (\mathbf{M}_L \circ_1 \varphi(\hat{\boldsymbol{\omega}} \odot \mathbf{y}))) (\bullet, k) = \boldsymbol{\omega}(k)^{-1} \mathbf{M}_L \varphi(\boldsymbol{\omega}(k) \mathbf{y}(\bullet, k)), \\ (\mathbf{B}\mathbf{u})(\bullet, k) &= ((\mathbf{B}\mathbf{u}) \otimes \boldsymbol{\omega} \cdot^{-1})(\bullet, k) = \boldsymbol{\omega}(k)^{-1} \mathbf{B}\mathbf{u}, \\ \mathbf{b}(\bullet, k) &= \boldsymbol{\omega}(k) \cdot^{-1} \mathbf{b}(\mathbf{a}_k). \end{aligned}$$

The vector $\boldsymbol{\omega}(k) \mathbf{y}(\bullet, k)$ represents the FE function $\mathbf{y}(\cdot, \mathbf{a}_{k_1}^{(1)}, \dots, \mathbf{a}_{k_m}^{(m)})$, see (6.9). This relates the approach to stochastic collocation, cf. [44], where one would have $\boldsymbol{\omega}(k) = 1$ for all k and then get the same weight 1 for all equations obtained by inserting the collocation points into (6.1). Taking orthonormal polynomials as basis, these single equations are weighted differently, namely by $\boldsymbol{\omega}(k)^{-1} = \mathbf{w}(k)^{1/2}$, in the tensor version.

6.3. Choice of the Trust-Region Model

We describe the choice of the trust-region model in the stochastic setting only because the deterministic setting works analogously. The main difference is that the stochastic setting requires further approximations of the Hessian to have efficient computations. Let $U \subset \mathcal{U}$ be the current discrete subspace and let $\mathbf{u}^k \in U$ be the current iterate, represented by the vector $\mathbf{u}^k \in \mathbb{R}^{d_u}$. The tensor $\tilde{\mathbf{y}}^k$ shall solve the state equation (6.16) with $\mathbf{u} = \mathbf{u}^k$ approximately and let $\tilde{\mathbf{z}}^k$ be an approximate solution of the adjoint equation (6.21) with $\mathbf{y} = \tilde{\mathbf{y}}^k$. We have to keep in mind that we apply iterative low-rank tensor solvers to these equations so that we cannot expect to obtain exact solutions.

In order to be able to apply the error estimates from Chapter 5, we choose a quadratic trust-region model m_k for any direction $\mathbf{s} \in U$ represented by the vector \mathbf{s} as follows:

$$m_k(\mathbf{s}) = \mathbf{m}_k(\mathbf{s}) := \nabla \mathbf{m}_k(0)^\top \tilde{\mathbf{M}} \mathbf{s} + \frac{1}{2} \mathbf{s}^\top \tilde{\mathbf{M}} \nabla^2 \mathbf{m}_k(0) \mathbf{s}$$

The space \mathbb{R}^{d_u} is equipped with the inner product induced by the mass matrix $\tilde{\mathbf{M}}$ and derivatives, such as $\nabla \mathbf{m}_k$, are computed w.r.t. this inner product. To approximate the true gradient sufficiently well, we choose

$$\nabla \mathbf{m}_k(0) = -(\tilde{\mathbf{M}}^{-1} \mathbf{B}^\top) \langle \tilde{\mathbf{z}}^k, \mathbf{w} \odot \boldsymbol{\omega} \rangle_{[m]+1, [m]} + \gamma \mathbf{u}^k$$

in light of (6.22) and Theorem 5.8. In principle, we could choose $\nabla^2 m_k(0)$ to be computed as in (6.23), but this evaluation requires two stochastic PDE solves, which is on the one hand costly and on the other hand inexact due to the low-rank tensor solver so that the Hessian

model would be noisy. To overcome this issue, we approximate the operator $\mathbf{A} + \mathbf{N}'(\tilde{\mathbf{y}}^k)$ by an operator of CP rank 1, which can be inverted easily. We do this by using a reference operator acting only on the space mode.

To approximate the operator \mathbf{A} , we can choose $\mathbf{A}_{\text{ref}} := A(\bar{\xi})$, where $\bar{\xi} = \int_{\Xi} \xi \, d\mathbb{P}$, or $\mathbf{A}_{\text{ref}} := \mathbb{E}[\mathbf{A}(\cdot)]$ for example. If $\mathbf{A} : \Xi \rightarrow \mathbb{R}^{d_0 \times d_0}$ is affine as in Example 6.5, the two choices coincide. We can approximate

$$\mathbf{A} \approx \mathbf{A}_{\text{ref}} := \mathbf{A}_{\text{ref}} \otimes \mathbf{I} \cdots \otimes \mathbf{I}$$

and have $\mathbf{A}_{\text{ref}}^{-1} = \mathbf{A}_{\text{ref}}^{-1} \otimes \mathbf{I} \cdots \otimes \mathbf{I}$, i. e., the inverse can be applied to low-rank tensors by computing a simple i -mode matrix product. In Chapter 7, it will be shown that it makes sense to use the operator \mathbf{A}_{ref} as a preconditioner for the solution of equations of the form (6.16) and (6.21). Then, the operator $\mathbf{A}_{\text{ref}}^{-1}$ has to be applied many times so that it is suitable to increase efficiency by, e. g., computing a sparse Cholesky decomposition of \mathbf{A}_{ref} once so that applying the inverse $\mathbf{A}_{\text{ref}}^{-1}$ to a vector consists only of permutations as well as forward and backward substitutions.

For the approximation of the operator $\mathbf{N}'(\mathbf{y})$ we introduce the parameter dependent vector

$$\tilde{\mathbf{y}}(\xi) := \sum_{k_1=1}^{d_1} \cdots \sum_{k_m=1}^{d_m} \tilde{\mathbf{y}}(\cdot, k_1, \dots, k_m) \theta_{k_1}^{(1)}(\xi_1) \cdots \theta_{k_m}^{(m)}(\xi_m)$$

cf. (6.8), and propose three options:

- Firstly, we can take $\mathbf{N}'_{\text{ref}}(\tilde{\mathbf{y}}) := \mathbf{N}'(\langle \tilde{\mathbf{y}}, \mathbf{w} \odot \boldsymbol{\omega} \rangle_{[m]+1, [m]})$, i. e., we evaluate \mathbf{N}' at the expected value $\mathbb{E}[\tilde{\mathbf{y}}(\cdot)]$ of the current state. This is the option we use in the implementation because after having computed the expected value, the code from the deterministic setting can be reused.
- Secondly, and similarly, $\mathbf{N}'_{\text{ref}}(\tilde{\mathbf{y}}) := \mathbf{N}'(\langle \tilde{\mathbf{y}}, \boldsymbol{\theta}(\bar{\xi}) \rangle_{[m]+1, [m]})$ can be taken, where $\boldsymbol{\theta}(\bar{\xi})(k) = \prod_{i=1}^m \theta_{k_i}^{(i)}(\bar{\xi}_i)$ meaning that \mathbf{N}' is evaluated at the “reference state” $\tilde{\mathbf{y}}(\bar{\xi})$.
- A third option would be to define $\mathbf{N}'_{\text{ref}}(\tilde{\mathbf{y}}) := \mathbf{M}_{\text{L}} \text{diag}(\langle \tilde{\boldsymbol{\varphi}}', \mathbf{w} \rangle_{[m]+1, [m]}) \approx \mathbb{E}[\mathbf{N}'(\tilde{\mathbf{y}}(\cdot))]$ with $\tilde{\boldsymbol{\varphi}}' \approx \varphi'(\hat{\boldsymbol{\omega}} \odot \tilde{\mathbf{y}})$. This option seems to be costly because the nonlinear function φ' has to be evaluated on the full tensor $\tilde{\mathbf{y}}$, but the quantity can be reused from the solution of the adjoint equation.

Using one of these options, the the final rank-1-approximation of the stochastic, linearized PDE operator is given by

$$\mathbf{A} + \mathbf{N}'(\tilde{\mathbf{y}}^k) \approx (\mathbf{A}_{\text{ref}} + \mathbf{N}'_{\text{ref}}(\tilde{\mathbf{y}}^k)) \otimes \mathbf{I} \otimes \cdots \otimes \mathbf{I} =: \mathbf{A}_{\text{ref}} + \mathbf{N}'_{\text{ref}}(\tilde{\mathbf{y}}^k).$$

With this approximation, the tensor \mathbf{d} from (6.24) is inexactly computed via

$$\mathbf{d} \approx \tilde{\mathbf{d}}^k := (\mathbf{A}_{\text{ref}} + \mathbf{N}'_{\text{ref}}(\tilde{\mathbf{y}}^k))^{-1} \mathbf{B} \mathbf{s} = ((\mathbf{A}_{\text{ref}} + \mathbf{N}'_{\text{ref}}(\tilde{\mathbf{y}}^k))^{-1} \mathbf{B} \mathbf{s}) \otimes \boldsymbol{\omega}^{-1}.$$

Given this rank-1-tensor, the tensor \mathbf{h} from (6.24) can be approximated via

$$\mathbf{h} \approx \tilde{\mathbf{h}}^k := (\mathbf{A}_{\text{ref}} + \mathbf{N}'_{\text{ref}}(\tilde{\mathbf{y}}^k))^{-1} \left[(\mathbf{Q}^{\top} \mathbf{M}_H \mathbf{Q}) \circ_1 \tilde{\mathbf{d}}^k + \mathbf{M}_{\text{L}} \circ_1 (\tilde{\boldsymbol{\zeta}}^k \odot \tilde{\mathbf{d}}^k) \right],$$

where $\tilde{\boldsymbol{\zeta}}^k \approx \hat{\boldsymbol{\omega}} \odot \tilde{\boldsymbol{z}}^k \odot \varphi''(\hat{\boldsymbol{\omega}} \odot \tilde{\boldsymbol{y}}^k)$. Using this approximation in (6.23) defines the model Hessian, applied to a direction \mathbf{s} :

$$\begin{aligned} \nabla^2 \mathbf{m}_k(0) \mathbf{s} &= (\tilde{\mathbf{M}}^{-1} \mathbf{B}^\top) \langle \tilde{\mathbf{h}}_k, \mathbf{w} \odot \boldsymbol{\omega} \rangle_{[m]+1, [m]} + \gamma \mathbf{s} \\ &= (\tilde{\mathbf{M}}^{-1} \mathbf{B}^\top) (\tilde{\mathbf{A}}_{\text{ref}}^{(k)})^{-1} \left[(\mathbf{Q}^\top \mathbf{M}_H \mathbf{Q}) + \mathbf{M}_L \text{diag}(\langle \tilde{\boldsymbol{\zeta}}^k, \mathbf{w} \rangle_{[m]+1, [m]}) \right] (\tilde{\mathbf{A}}_{\text{ref}}^{(k)})^{-1} \mathbf{B} \mathbf{s} + \gamma \mathbf{s} \end{aligned} \quad (6.26)$$

using $\langle \boldsymbol{\omega}^{-1}, \mathbf{w} \odot \boldsymbol{\omega} \rangle = 1$ (sum of quadrature weights) and with $\tilde{\mathbf{A}}_{\text{ref}}^{(k)} := \mathbf{A}_{\text{ref}} + \mathbf{N}'_{\text{ref}}(\tilde{\boldsymbol{y}}^k)$. This Hessian operator is symmetric on \mathbb{R}^{d_u} w. r. t. the inner product induced by $\tilde{\mathbf{M}}$. In contrast to [46], we use the concrete rank-1-property of all involved operators here to derive that after forming $\langle \tilde{\boldsymbol{\zeta}}^k, \mathbf{w} \rangle_{[m]+1, [m]}$, the application of the model Hessian can be done by pure matrix calculus. In our implementation we choose

$$\langle \tilde{\boldsymbol{\zeta}}^k, \mathbf{w} \rangle_{[m]+1, [m]} = \langle \tilde{\boldsymbol{z}}^k, \mathbf{w} \odot \boldsymbol{\omega} \rangle_{[m]+1, [m]} \odot \varphi''(\langle \tilde{\boldsymbol{y}}^k, \mathbf{w} \odot \boldsymbol{\omega} \rangle_{[m]+1, [m]})$$

for simplicity, i. e., we work with the expected state and adjoint state and reuse the deterministic Hessian computation.

6.4. Solution of the Discrete Semismooth Newton System

As pointed out in Section 4.4, a semismooth Newton method can be applied to improve the step found by a linesearch. In this section, we describe how this method is implemented in the discrete setting for the special choice $U_{\text{ad}} := \{u \in U : u_l(x) \leq u(x) \leq u_u(x) \text{ for a. e. } x \in \Omega_u\}$, where $u_l, u_u \in U \subset L^\infty(\Omega_u)$ are functions living in the finite-dimensional subspace and fulfilling $u_l(x) \leq u_u(x)$ for (almost) every $x \in \Omega_u$. This description follows our work [46]. The exact projection onto the feasible set is given by

$$P_{U_{\text{ad}}} : U \rightarrow U_{\text{ad}}, (P_{U_{\text{ad}}}(u))(x) = \min\{\max\{u_l(x), u(x)\}, u_u(x)\}. \quad (6.27)$$

It is known [111] that this superposition operator is semismooth as a mapping from $L^q(\Omega_u)$ to $L^2(\Omega_u)$ for any $q > 2$. Since $u^k \in U_{\text{ad}} \subset L^\infty(\Omega_u) \subset L^q(\Omega_u)$ and if $\iota \mathbf{B}^*$ maps \mathbf{Y} into $L^q(\Omega_u)$ continuously, where $\iota : L^2(\Omega_u)^* \rightarrow L^2(\Omega_u)$ denotes the Riesz representation operator, the residual R defined in (4.45) is semismooth w. r. t. the direction $\bar{\mathbf{s}}$ as noted in Section 4.4. In our concrete case with \mathbf{B} from (3.11), i. e., $\iota \mathbf{B}^* \mathbf{v} = D^* \int_{\Xi} \mathbf{v}(\cdot, \xi) d\mathbb{P}$, the requirement on $\iota \mathbf{B}^*$ is met if $D^* \in \mathcal{L}(L^2(\Omega), L^2(\Omega_u))$ (using Riesz representatives) maps $Y \subset L^2(\Omega)$ continuously into $L^q(\Omega_u)$. If, e. g., $\Omega_u = \Omega$ and $D \equiv I : L^2(\Omega_u) \rightarrow L^2(\Omega)$ one can use the Sobolev embedding $H_0^1(\Omega) \hookrightarrow L^p(\Omega)$ and $q \leq p$ to show this property.

In the flavor of (6.27), it is desirable to have a componentwise formulation of the projection also on the discrete subspace so that it can be computed efficiently. If we equip \mathbb{R}^{d_u} with the inner product induced by the *lumped* mass matrix $\tilde{\mathbf{M}}_L$, this holds true because this matrix is diagonal with positive entries. Let u_l, u_u be represented by the vectors $\mathbf{u}_l, \mathbf{u}_u \in \mathbb{R}^{d_u}$ which fulfill equivalently $u_l \leq u_u$ componentwise and define $\mathbf{U}_{\text{ad}} := \{\mathbf{u} \in \mathbb{R}^{d_u} : \mathbf{u}_l \leq \mathbf{u} \leq \mathbf{u}_u\}$. The projection onto this box w. r. t. the $\tilde{\mathbf{M}}_L$ -inner product then reads

$$P_{\mathbf{U}_{\text{ad}}} : \mathbb{R}^{d_u} \rightarrow \mathbf{U}_{\text{ad}}, (P_{\mathbf{U}_{\text{ad}}}(\mathbf{u}))_j = \min\{\max\{(\mathbf{u}_l)_j, \mathbf{u}_j\}, (\mathbf{u}_u)_j\} \text{ for every } j \in [d_u]. \quad (6.28)$$

Given discrete versions \mathbf{u}^k and $\nabla \mathbf{m}_k(0)$ of the current control and the model gradient, respectively, we formulate the discrete version of (4.45) w. r. t. the $\tilde{\mathbf{M}}_{\mathbf{L}}$ -inner product. This means that we use the discrete projection (6.28) and a discrete version $\tilde{\nabla}^2 \mathbf{m}_k(0)$ of the model Hessian w. r. t. the $\tilde{\mathbf{M}}_{\mathbf{L}}$ -inner product, which is obtained by replacing the mass matrix $\tilde{\mathbf{M}}$ by the lumped mass matrix $\tilde{\mathbf{M}}_{\mathbf{L}}$ in (6.26). This leads to a favorable structure in the semismooth Newton system as we will see later. In analogy to (4.44), we define

$$\hat{\mathbf{H}}(\mathbf{u}^k) := \mathbf{B}^\top (\tilde{\mathbf{A}}_{\text{ref}}^{(k)})^{-1} \left[(\mathbf{Q}^\top \mathbf{M}_H \mathbf{Q}) + \mathbf{M}_{\mathbf{L}} \text{diag}(\langle \tilde{\boldsymbol{\zeta}}^k, \mathbf{w} \rangle_{[m]_{+1}, [m]}) \right] (\tilde{\mathbf{A}}_{\text{ref}}^{(k)})^{-1} \mathbf{B},$$

where $\tilde{\mathbf{A}}_{\text{ref}}^{(k)}$ and $\tilde{\boldsymbol{\zeta}}^k$ are computed from \mathbf{u}^k as in Section 6.3. Then the alternative Hessian is given by $\tilde{\nabla}^2 \mathbf{m}_k(0) \mathbf{s} = \tilde{\mathbf{M}}_{\mathbf{L}}^{-1} \hat{\mathbf{H}}(\mathbf{u}^k) \mathbf{s} + \gamma \mathbf{s}$. Inverting the diagonal matrix $\tilde{\mathbf{M}}_{\mathbf{L}}$ is cheap compared to the inversion of $\tilde{\mathbf{M}}$. With these definitions, (4.45) becomes

$$\mathbf{R}(\bar{\mathbf{s}}) := \bar{\mathbf{s}} + \mathbf{u}^k - \mathbf{P}_{\text{U}_{\text{ad}}}(\mathbf{u}^k - \tau_{n,k} \nabla \mathbf{m}_k(0) - \tau_{n,k} \tilde{\mathbf{M}}_{\mathbf{L}}^{-1} \hat{\mathbf{H}}(\mathbf{u}^k) \bar{\mathbf{s}}) = 0. \quad (6.29)$$

with $\tau_{n,k} := (\gamma + \frac{\epsilon_{n,k}}{\Delta_k})^{-1} > 0$. The discrete residual \mathbf{R} is also semismooth w. r. t. the direction \mathbf{s} . We choose

$$\mathbf{DP}_{\text{U}_{\text{ad}}}(\mathbf{u}) \in \mathbb{R}^{d_u \times d_u} \text{ diagonal}, \quad (\mathbf{DP}_{\text{U}_{\text{ad}}}(\mathbf{u}))_{jj} = \begin{cases} 0 & \text{if } \mathbf{u}_j < (\mathbf{u}_l)_j, \\ 1 & \text{if } (\mathbf{u}_l)_j \leq \mathbf{u}_j \leq (\mathbf{u}_u)_j, \\ 0 & \text{if } (\mathbf{u}_u)_j < \mathbf{u}_j \end{cases}$$

as an element of the generalized Jacobian of $\mathbf{P}_{\text{U}_{\text{ad}}}$. Note that if $\mathbf{u}_j \in \{(\mathbf{u}_l)_j, (\mathbf{u}_u)_j\}$, we are in principle free to take $(\mathbf{DP}_{\text{U}_{\text{ad}}}(\mathbf{u}))_{jj} \in [0, 1]$, but choosing this value from $\{0, 1\}$ will make a block elimination possible. With this choice, the discrete semismooth Newton system for the iterative solution of (6.29) is

$$(\mathbf{s}^{\ell+1} - \mathbf{s}^\ell) + \tau_{n,k} \mathbf{DP}_{\text{U}_{\text{ad}}}(\mathbf{u}^k - \tau_{n,k} \nabla \mathbf{m}_k(0) - \tau_{n,k} \tilde{\mathbf{M}}_{\mathbf{L}}^{-1} \hat{\mathbf{H}}(\mathbf{u}^k) \mathbf{s}^\ell) \tilde{\mathbf{M}}_{\mathbf{L}}^{-1} \hat{\mathbf{H}}(\mathbf{u}^k) (\mathbf{s}^{\ell+1} - \mathbf{s}^\ell) = -\mathbf{R}(\mathbf{s}^\ell), \quad (6.30)$$

where \mathbf{s}^ℓ is the current approximation of the solution $\bar{\mathbf{s}}$ of (6.29) and $\mathbf{s}^{\ell+1} - \mathbf{s}^\ell$ is the semismooth Newton step. Since the typically large and dense matrix $\hat{\mathbf{H}}(\mathbf{u}^k)$ should never be formed explicitly, we have to apply an iterative method such as GMRES or CG working with the application of the matrix only to solve (6.30). This method should be formulated on the function space equivalent, i. e., work on \mathbb{R}^{d_u} equipped with the $\tilde{\mathbf{M}}_{\mathbf{L}}$ -inner product. Analogously and in favor of a simple implementation, we can multiply (6.30) by $\tilde{\mathbf{M}}_{\mathbf{L}}$ from the right and use the fact that diagonal matrices commute to obtain the equivalent, preconditioned system

$$\tilde{\mathbf{M}}_{\mathbf{L}} (\mathbf{s}^{\ell+1} - \mathbf{s}^\ell) + \tau_{n,k} \mathbf{DP}_{\text{U}_{\text{ad}}}(\mathbf{u}^k - \tau_{n,k} \nabla \mathbf{m}_k(0) - \tau_{n,k} \tilde{\mathbf{M}}_{\mathbf{L}}^{-1} \hat{\mathbf{H}}(\mathbf{u}^k) \mathbf{s}^\ell) \hat{\mathbf{H}}(\mathbf{u}^k) (\mathbf{s}^{\ell+1} - \mathbf{s}^\ell) = -\tilde{\mathbf{M}}_{\mathbf{L}} \mathbf{R}(\mathbf{s}^\ell), \quad (6.31)$$

which can be solved by iterative methods formulated on \mathbb{R}^{d_u} equipped with the standard Euclidean inner product. In particular, to apply the CG method successfully we symmetrize (6.31) by a block elimination. For this purpose, the index sets

$$\mathcal{I} = \mathcal{I}_k(\mathbf{s}^\ell) := \left\{ j \in [\tilde{N}] : (\mathbf{DP}_{\text{U}_{\text{ad}}}(\mathbf{u}^k - \tau_{n,k} \nabla \mathbf{m}_k(0) - \tau_{n,k} \tilde{\mathbf{M}}_{\mathbf{L}}^{-1} \hat{\mathbf{H}}(\mathbf{u}^k) \mathbf{s}^\ell))_{jj} = 0 \right\}$$

and $\mathcal{A} = \mathcal{A}_k(\mathbf{s}^\ell) := [d_u] \setminus \mathcal{I}_k(\mathbf{s}^\ell)$ are distinguished to define the partial solution $(\mathbf{s}^{\ell+1} - \mathbf{s}^\ell)_{\mathcal{I}} = -\mathbf{R}(\mathbf{s}^\ell)_{\mathcal{I}}$ of (6.30) or equivalently (6.31). We insert it into (6.31) to obtain the smaller system

$$((\tilde{\mathbf{M}}_{\mathbf{L}})_{\mathcal{A}\mathcal{A}} + \tau_{n,k} \hat{\mathbf{H}}(\mathbf{u}^k)_{\mathcal{A}\mathcal{A}})(\mathbf{s}^{\ell+1} - \mathbf{s}^\ell)_{\mathcal{A}} = \tau_{n,k} \hat{\mathbf{H}}(\mathbf{u}^k)_{\mathcal{A}\mathcal{I}} \mathbf{R}(\mathbf{s}^\ell)_{\mathcal{I}} - (\tilde{\mathbf{M}}_{\mathbf{L}})_{\mathcal{A}\mathcal{A}} \mathbf{R}(\mathbf{s}^\ell)_{\mathcal{A}}. \quad (6.32)$$

The matrix $(\tilde{\mathbf{M}}_{\mathbf{L}})_{\mathcal{A}\mathcal{A}} + \tau_{n,k} \hat{\mathbf{H}}(\mathbf{u}^k)_{\mathcal{A}\mathcal{A}}$ is always symmetric w. r. t. the Euclidean inner product and positive definite if e. g., the entries of the vector $\langle \tilde{\boldsymbol{\zeta}}^k, \mathbf{w} \rangle_{[m]+1, [m]}$ are large enough to make the matrix

$$(\mathbf{Q}^\top \mathbf{M}_H \mathbf{Q}) + \mathbf{M}_{\mathbf{L}} \text{diag}(\langle \tilde{\boldsymbol{\zeta}}^k, \mathbf{w} \rangle_{[m]+1, [m]})$$

(see (6.26)) positive definite. In practice, we do not experience a problem when applying CG to the system (6.32). In general, one would have to apply an adequate iterative solver such as MINRES if CG fails to compute an accurate enough solution. The matrix $\hat{\mathbf{H}}(\mathbf{u}^k)_{\mathcal{A}\mathcal{A}}$ should not be computed explicitly, but should be applied to a direction $\mathbf{s}_{\mathcal{A}}$ using the application of the full operator $\hat{\mathbf{H}}(\mathbf{u}^k)$ as follows:

$$\hat{\mathbf{H}}(\mathbf{u}^k)_{\mathcal{A}\mathcal{A}} \mathbf{s}_{\mathcal{A}} = \begin{pmatrix} \hat{\mathbf{H}}(\mathbf{u}^k)_{\mathcal{A}\mathcal{A}} & \hat{\mathbf{H}}(\mathbf{u}^k)_{\mathcal{A}\mathcal{I}} \\ & \end{pmatrix} \begin{pmatrix} \mathbf{s}_{\mathcal{A}} \\ \mathbf{0} \end{pmatrix} = \begin{pmatrix} \hat{\mathbf{H}}(\mathbf{u}^k) \\ \mathbf{0} \end{pmatrix}_{\mathcal{A}}.$$

7. A Posteriori Error Estimation and Adaptive Solution of a Class of Parametric PDEs Using Low-Rank Tensors

Following [22, 23, 38, 39, 40] we consider the parametric, elliptic, nonlinear operator equation

$$A(\xi)y(\xi) + N(y(\xi)) = b(\xi), \quad (7.1)$$

where ξ is a vector of independent random variables distributed on $\Xi := \times_{i=1}^m \Xi_i$ with the probability measure $\mathbb{P} := \otimes_{i=1}^m \mathbb{P}_i$ as in Section 6.2. Assumption 6.3 shall still be valid. The right-hand side $b(\xi) \in Y^*$ is given and strongly measurable w. r. t. ξ . Let $A(\xi) : Y \rightarrow Y^*$ be a linear, self-adjoint, and bounded operator between a Hilbert space and its dual for almost every ξ and let $N : Y \rightarrow Y^*$ be well-defined, continuous and monotone, but possibly nonlinear. We assume that $\xi \mapsto A(\xi)$ is also strongly measurable and that $A(\xi)$ is uniformly bounded and boundedly invertible:

$$\|A(\xi)\|_{\mathcal{L}(Y, Y^*)} \leq c_{\max}, \quad \|A(\xi)^{-1}\|_{\mathcal{L}(Y^*, Y)} \leq c_{\min}^{-1}.$$

This ensures, together with measurability, that $A \in L_{\mathbb{P}}^{\infty}(\Xi; \mathcal{L}(Y, Y^*))$ and thus $\mathbf{A} \in \mathcal{L}(\mathbf{Y}, \mathbf{Y}^*)$ by Proposition 3.1 for the operator

$$\langle \mathbf{A}y, v \rangle_{\mathbf{Y}^*, \mathbf{Y}} := \int_{\Xi} \langle A(\xi)y(\cdot, \xi), v(\cdot, \xi) \rangle_{Y^*, Y} d\mathbb{P} \quad (7.2)$$

and the spaces $\mathbf{Y} := L_{\mathbb{P}}^p(\Xi; Y)$, $\mathbf{Y}^* := L_{\mathbb{P}}^{p^*}(\Xi; Y^*)$ with adequately chosen $p > 3$. The operator $\hat{\mathbf{A}} \in \mathcal{L}(\hat{\mathbf{Y}}, \hat{\mathbf{Y}}^*)$ defined exactly as \mathbf{A} but with $\hat{\mathbf{Y}} := L_{\mathbb{P}}^2(\Xi; Y)$ and $\hat{\mathbf{Y}}^* := L_{\mathbb{P}}^2(\Xi; Y^*)$ is also well-defined and even boundedly invertible.

The ellipticity of the operator $A(\xi)$ gives rise to the inner product

$$(y, v)_{A(\xi)} := \langle A(\xi)y, v \rangle_{Y^*, Y}$$

and the corresponding energy norm $\|y\|_{A(\xi)} = \sqrt{(y, y)_{A(\xi)}}$.

For error estimation, we consider a linear, self-adjoint and elliptic “reference” operator $A_{\text{ref}} \in \mathcal{L}(Y, Y^*)$. This can be the nominal operator $A_{\text{ref}} = A(\bar{\xi})$, where $\bar{\xi} := \int_{\Xi} \xi d\mathbb{P}$, or any other norm-inducing operator, such as $A_{\text{ref}} = -\Delta$ for the $H_0^1(\Omega)$ -norm. Then, $(y, v)_{A_{\text{ref}}} := \langle A_{\text{ref}}y, v \rangle_{Y^*, Y}$ defines an alternative inner product on the space Y and the equivalence estimate

$$\lambda \langle A(\xi)v, v \rangle_{Y^*, Y} \leq \langle A_{\text{ref}}v, v \rangle_{Y^*, Y} \leq \Lambda \langle A(\xi)v, v \rangle_{Y^*, Y} \quad (7.3)$$

holds for a. e. $\xi \in \Xi$, every $v \in Y$ and some constants $\lambda, \Lambda \in (0, \infty)$, $\lambda \leq \Lambda$, due to the uniform ellipticity of $A(\xi)$. We define the alternative norm $\|y\|_{A_{\text{ref}}} := \sqrt{\langle A_{\text{ref}}y, y \rangle_{Y^*, Y}}$ on Y which is equivalent to the usual norm $\|\cdot\|_Y$. Analogously defining the inner products and norms on Y^* induced by the operators $A(\xi)^{-1}$ and A_{ref}^{-1} , we get the inverse estimates

$$\frac{1}{\Lambda} \|\hat{b}\|_{A(\xi)^{-1}}^2 \leq \langle \hat{b}, A_{\text{ref}}^{-1} \hat{b} \rangle_{Y^*, Y} = \|\hat{b}\|_{A_{\text{ref}}^{-1}}^2 \leq \frac{1}{\lambda} \|\hat{b}\|_{A(\xi)^{-1}}^2 \quad (7.4)$$

for every $\hat{b} \in Y^*$ by Proposition A.3.

Now let $y(\xi) \in Y$ be the unique and exact solution of (7.1), and let $\tilde{y}(\xi) \in Y \subset Y$ be an inexact solution living in a subspace Y of Y (e. g., a finite element subspace), fulfilling

$$A(\xi)\tilde{y}(\xi) + N(\tilde{y}(\xi)) - b(\xi) =: r(\xi), \quad (7.5)$$

where $r(\xi) \in Y^*$ is the residual. We assume that this residual can be evaluated exactly in the sense that $\langle r(\xi), v^+ \rangle_{Y^*, Y}$ can be evaluated for every given $v^+ \in Y^+ \supset Y$ in some finite-dimensional (e. g., FE) subspace $Y^+ \subset Y$.

Lemma 7.1. *Let $y(\xi) \in Y$ be the solution of (7.1) and let $\tilde{y}(\xi) \in Y$. Under the standing assumptions, it holds that*

$$\|\tilde{y}(\xi) - y(\xi)\|_{A_{\text{ref}}} \leq \Lambda \|r(\xi)\|_{A_{\text{ref}}^{-1}} \quad (7.6)$$

with $r(\xi)$ defined as in (7.5).

Proof. We can estimate using the monotonicity of N :

$$\begin{aligned} & \|\tilde{y}(\xi) - y(\xi)\|_{A(\xi)}^2 \\ & \leq \langle A(\xi)(\tilde{y}(\xi) - y(\xi)), \tilde{y}(\xi) - y(\xi) \rangle_{Y^*, Y} + \langle N(\tilde{y}(\xi)) - N(y(\xi)), \tilde{y}(\xi) - y(\xi) \rangle_{Y^*, Y} \\ & = \langle r(\xi), \tilde{y}(\xi) - y(\xi) \rangle_{Y^*, Y} = \langle A(\xi)A(\xi)^{-1}r(\xi), \tilde{y}(\xi) - y(\xi) \rangle_{Y^*, Y} \\ & = \langle A(\xi)^{-1}r(\xi), \tilde{y}(\xi) - y(\xi) \rangle_{A(\xi)} \leq \|A(\xi)^{-1}r(\xi)\|_{A(\xi)} \|\tilde{y}(\xi) - y(\xi)\|_{A(\xi)}. \end{aligned}$$

This gives $\|\tilde{y}(\xi) - y(\xi)\|_{A(\xi)} \leq \|r(\xi)\|_{A(\xi)^{-1}}$ and

$$\frac{1}{\sqrt{\Lambda}} \|\tilde{y}(\xi) - y(\xi)\|_{A_{\text{ref}}} \leq \|\tilde{y}(\xi) - y(\xi)\|_{A(\xi)} \leq \|r(\xi)\|_{A(\xi)^{-1}} \leq \sqrt{\Lambda} \|r(\xi)\|_{A_{\text{ref}}^{-1}}$$

by the norm estimates (7.3) and (7.4), which yields the result. \square

Remark 7.2. In the case $N \equiv 0$, we even have $\tilde{y}(\xi) - y(\xi) = A(\xi)^{-1}r(\xi)$ and thus $\|\tilde{y}(\xi) - y(\xi)\|_{A(\xi)} = \|r(\xi)\|_{A(\xi)^{-1}}$. Then the lower bound $\sqrt{\lambda} \|r(\xi)\|_{A_0^{-1}} \leq \|r(\xi)\|_{A(\xi)^{-1}}$ is useful.

For the computation of $\|r(\xi)\|_{A_{\text{ref}}^{-1}}$ we define $w(\xi) \in Y$ to be the unique solution of the equation $A_{\text{ref}}w(\xi) = r(\xi)$. Then it holds that

$$\|r(\xi)\|_{A_{\text{ref}}^{-1}}^2 = \langle r(\xi), A_{\text{ref}}^{-1}r(\xi) \rangle_{Y^*, Y} = \langle A_{\text{ref}}w(\xi), w(\xi) \rangle_{Y^*, Y} = \|w(\xi)\|_{A_{\text{ref}}}^2. \quad (7.7)$$

We compute a discrete solution $w(\xi) \in Y$ fulfilling

$$\langle A_{\text{ref}}w(\xi), v \rangle_{Y^*, Y} = \langle r(\xi), v \rangle_{Y^*, Y} \text{ for all } v \in Y, \quad (7.8)$$

i. e., $w(\xi) = A_{\text{ref}}^{-1}r(\xi)$, where $A_{\text{ref}} : Y \rightarrow Y^*$ is the restriction of the operator A_{ref} onto the space Y and its inverse is defined in the sense of (7.8). We assume that this equation is solved exactly having our application in mind. If this it not the case, the algebraic error caused by the inexact solution of the discrete equation has to be incorporated additionally.

Lemma 7.3. *Let $y(\xi), \tilde{y}(\xi), r(\xi)$ be as in Lemma 7.1, and let $w(\xi) = A_{\text{ref}}^{-1}r(\xi)$ and $w(\xi) \in Y$ defined by (7.8). Then, under the standing assumptions,*

$$\|\tilde{y}(\xi) - y(\xi)\|_{A_{\text{ref}}}^2 \leq \Lambda^2(\|w(\xi)\|_{A_{\text{ref}}}^2 + \|w(\xi) - w(\xi)\|_{A_{\text{ref}}}^2) \quad (7.9)$$

holds.

Proof. Combining (7.6) and (7.7) results in

$$\|\tilde{y}(\xi) - y(\xi)\|_{A_{\text{ref}}}^2 \leq \Lambda^2(\|w(\xi) + w(\xi) - w(\xi)\|_{A_{\text{ref}}}^2) = \Lambda^2(\|w(\xi)\|_{A_{\text{ref}}}^2 + \|w(\xi) - w(\xi)\|_{A_{\text{ref}}}^2),$$

where the last equality is due to (7.8) and $w(\xi) \in Y$, cf. the proof of [38, Thm. 5.1]. \square

The first summand in (7.9) turns out to be the purely algebraic error contribution caused by solving a discretized version of (7.1) inexactly:

$$\begin{aligned} \|w(\xi)\|_{A_{\text{ref}}}^2 &= \langle r(\xi), A_{\text{ref}}^{-1}r(\xi) \rangle_{Y^*, Y} \\ &= \langle A(\xi)\tilde{y}(\xi) + N(\tilde{y}(\xi)) - b(\xi), A_{\text{ref}}^{-1}(A(\xi)\tilde{y}(\xi) + N(\tilde{y}(\xi)) - b(\xi)) \rangle_{Y^*, Y}. \end{aligned}$$

The second summand will be estimated by standard a posteriori error estimates for (7.8).

7.1. Realization of the Deterministic a Posteriori Error Estimator

We discuss the realization of a deterministic a posteriori error estimator for the estimation of the term $\|w(\xi) - w(\xi)\|_{A_{\text{ref}}}$ for the example from Section 3.2 with the deterministic state equation (3.8) and the adjoint equation (3.16). This is an adaption of the ideas presented in [114, Chap. 1] and [2, Chap. 2] to our setting.

We define in analogy to (3.10), but with slight differences:

$$\begin{aligned} \langle A(\xi)y(\xi), v \rangle_{Y^*, Y} &:= \int_{\Omega} \kappa(\cdot, \xi) \nabla y(\xi) \cdot \nabla v + \chi(\cdot, \xi) y(\xi) v \, dx, \\ \langle N(y(\xi)), v \rangle_{Y^*, Y} &:= \int_{\Omega} \hat{\varphi}(y(\xi)) v \, dx, \\ \langle b(\xi), v \rangle_{Y^*, Y} &:= \int_{\Omega} \hat{f}(\xi) v \, dx, \\ \langle A_{\text{ref}}y(\xi), v \rangle_{Y^*, Y} &:= \int_{\Omega} \kappa_{\text{ref}} \nabla y(\xi) \cdot \nabla v + \chi_{\text{ref}} y(\xi) v \, dx \end{aligned} \quad (7.10)$$

with $\chi \in L_{\lambda \otimes \mathbb{P}}^{\infty}(\Omega \times \Xi)$, $\chi(x, \xi) \geq 0$ for a. e. $(x, \xi) \in \Omega \times \Xi$ and the deterministic reference coefficients κ_{ref} (uniformly positive) and χ_{ref} (nonnegative). The functions $\hat{\varphi} \in \mathcal{C}^2(\mathbb{R}, \mathbb{R})$ and $\hat{f}(\xi) \in L^2(\Omega)$ shall fulfill the same properties as φ and f in Assumption 3.3.

Remark 7.4. The setting covers all important equations and operators:

- For $\chi \equiv 0$, $\hat{\varphi} \equiv \varphi$, and $\hat{f}(\xi) = Du + f(\xi)$ we get the state equation (3.8).
- If, e. g., $Q(\xi) \equiv \iota : H_0^1(\Omega) \hookrightarrow L^2(\Omega)$, we get the adjoint equation (3.16) for $\chi(\cdot, \xi) = \varphi'(\tilde{y}(\xi))$, $\hat{\varphi} \equiv 0$ and $\hat{f}(\xi) = -(Q(\xi)\tilde{y}(\xi) - \hat{q}(\xi))$.
- For $\kappa_{\text{ref}}(x) = \kappa(x, \bar{\xi})$ and $\chi_{\text{ref}} \equiv 0$, we obtain the nominal operator $A(\bar{\xi})$. For $\kappa_{\text{ref}} \equiv 1$ and $\chi_{\text{ref}} \equiv 0$ we get the $H_0^1(\Omega)$ -norm inducing operator. Changing to $\chi_{\text{ref}} \equiv 1$, the $H^1(\Omega)$ -norm is induced.

Assumption 7.5.

- We consider a two-dimensional, polygonal domain $\Omega \subset \mathbb{R}^2$, cf. Assumption 6.1.
- As described in Section 6.1 we discretize the space $Y = H_0^1(\Omega)$ by linear finite elements (continuous on $\bar{\Omega}$) on a triangulation \mathcal{T} . The discrete subspace is denoted by $Y \subset Y$.
- The coefficient functions $\kappa(\cdot, \xi)$ and κ_{ref} are assumed to be piecewise constant on the triangles.

Remark 7.6. These assumptions are only made to simplify the explanations and results in this section. Quadrilateral elements, cf. [2], higher order FE functions, or not piecewise constant coefficient functions could be included quite simply, but would result in more complicated formulas or distinctions of cases.

With $v \in Y$ we get

$$\begin{aligned}
 (w(\xi) - w(\xi), v)_{A_{\text{ref}}} &= \langle r(\xi) - A_{\text{ref}}w(\xi), v \rangle_{Y^*, Y} = \\
 &= \sum_{T \in \mathcal{T}} \left(\int_T \kappa(\cdot, \xi) \nabla \tilde{y}(\xi) \cdot \nabla v + \chi(\cdot, \xi) \tilde{y}(\xi) v + \hat{\varphi}(\tilde{y}(\xi)) v - \hat{f}(\xi) v \right. \\
 &\quad \left. - \kappa_{\text{ref}} \nabla w(\xi) \cdot \nabla v - \chi_{\text{ref}} w(\xi) v \, dx \right) \\
 &= \sum_{T \in \mathcal{T}} \left(\int_T -\text{div}(\kappa(\cdot, \xi) \nabla \tilde{y}(\xi)) v + \text{div}(\kappa_{\text{ref}} \nabla w(\xi)) v \right. \\
 &\quad \left. + (\chi(\cdot, \xi) \tilde{y}(\xi) + \hat{\varphi}(\tilde{y}(\xi)) - \hat{f}(\xi) - \chi_{\text{ref}} w(\xi)) v \, dx \right. \\
 &\quad \left. + \int_{\partial T} (\kappa(\cdot, \xi) \nabla \tilde{y}(\xi) - \kappa_{\text{ref}} \nabla w(\xi)) \cdot \nu_T v \, dS \right),
 \end{aligned}$$

where ν_T is the outer unit normal of the triangle T . Since $\kappa(\cdot, \xi)$ and κ_{ref} are piecewise constant on the triangles and we use piecewise linear ansatz functions (Assumption 7.5), $\text{div}(\kappa_{\text{ref}} \nabla w(\xi)) \equiv 0 \equiv \text{div}(\kappa(\cdot, \xi) \nabla \tilde{y}(\xi))$ holds true on each element T . Therefore, using $(w(\xi) - w(\xi), v)_{A_{\text{ref}}} = 0$, we get

$$\begin{aligned}
 (w(\xi) - w(\xi), v)_{A_{\text{ref}}} &= (w(\xi) - w(\xi), v - v)_{A_{\text{ref}}} \\
 &= \sum_{T \in \mathcal{T}} \left(\int_T (\chi(\cdot, \xi) \tilde{y}(\xi) + \hat{\varphi}(\tilde{y}(\xi)) - \hat{f}(\xi) - \chi_{\text{ref}} w(\xi)) (v - v) \, dx \right. \\
 &\quad \left. + \int_{\partial T} (\kappa(\cdot, \xi) \frac{\partial}{\partial \nu_T} \tilde{y}(\xi) - \kappa_{\text{ref}} \frac{\partial}{\partial \nu_T} w(\xi)) (v - v) \, dS \right) \tag{7.11}
 \end{aligned}$$

for arbitrary $v \in Y$. Since \tilde{y} is bounded and continuous on $\overline{\Omega}$ and $\hat{\varphi} : \mathbb{R} \rightarrow \mathbb{R}$ is continuous, $\hat{\varphi}(\tilde{y}(\xi))$ belongs to $L^2(\Omega)$. Due to $\chi(\cdot, \xi), \chi_{\text{ref}} \in L^\infty(\Omega)$, we also have $\chi(\cdot, \xi)\tilde{y}(\xi), \chi_{\text{ref}}w(\xi) \in L^2(\Omega)$. The integrals over the triangle boundaries ∂T are considered edge-wise: If an edge E belongs to $\partial\Omega$, this part of the integral vanishes. Interior edges also appear in the integral over the boundary of a neighboring triangle. Thus, the sum over all triangle boundary integrals can be collected to integrals over all interior edges $E \in \mathcal{E}^0$, denoting by \mathcal{E} the set of all edges and by \mathcal{E}^0 the set of all interior edges. There, the normal jumps involving the discrete solutions \tilde{y} and w appear:

$$\llbracket \phi \rrbracket_E(x) := \lim_{t \rightarrow 0^+} \phi(x + t\nu_E) - \lim_{t \rightarrow 0^+} \phi(x - t\nu_E),$$

where ν_E is a unit normal vector corresponding to E , cf. [114, Sec. 1.1].

We estimate

$$\begin{aligned} (w(\xi) - w(\xi), v)_{A_{\text{ref}}} &\leq \sum_{T \in \mathcal{T}} \|\chi(\cdot, \xi)\tilde{y}(\xi) + \hat{\varphi}(\tilde{y}(\xi)) - \hat{f}(\xi) - \chi_{\text{ref}}w(\xi)\|_{L^2(T)} \|v - v\|_{L^2(T)} \\ &\quad + \sum_{E \in \mathcal{E}^0} \|\llbracket (\kappa(\cdot, \xi)\nabla\tilde{y} - \kappa_{\text{ref}}\nabla w(\xi)) \cdot \nu_E \rrbracket_E\|_{L^2(E)} \|v - v\|_{L^2(E)} \end{aligned} \quad (7.12)$$

If we insert the Clément interpolant $v \in Y$ of v and use that $(w - w, v)_{A_{\text{ref}}} \leq c \cdot \eta(w) \cdot \|v\|_{A_{\text{ref}}}$ for all $v \in Y$ yields $\|w - w\|_{A_{\text{ref}}} \leq c \cdot \eta(w)$, a similar estimation as in [114, Sec. 1.2] or [2, Sec. 2.2] can be performed, giving

$$\|w(\xi) - w(\xi)\|_{A_{\text{ref}}} \leq c_{\mathcal{T}} c_{A_{\text{ref}}} \left(\sum_{T \in \mathcal{T}} \eta_T(\tilde{y}(\xi))^2 + \sum_{E \in \mathcal{E}^0} \eta_E(\tilde{y}(\xi))^2 \right)^{1/2} \quad (7.13)$$

with

$$\eta_T(\tilde{y}(\xi)) := h_T \|\chi(\cdot, \xi)\tilde{y}(\xi) + \hat{\varphi}(\tilde{y}(\xi)) - \hat{f}(\xi) - \chi_{\text{ref}}w(\xi)\|_{L^2(T)} \quad (7.14)$$

and

$$\eta_E(\tilde{y}(\xi)) := h_E^{1/2} \|\llbracket (\kappa(\cdot, \xi)\nabla\tilde{y}(\xi) - \kappa_{\text{ref}}\nabla w(\xi)) \cdot \nu_E \rrbracket_E\|_{L^2(E)}. \quad (7.15)$$

The constant $c_{\mathcal{T}} > 0$ depends only on the smallest angle in the triangulation \mathcal{T} and the coercivity constant $c_{A_{\text{ref}}} > 0$ is chosen such that $\|y\|_{H^1(\Omega)} \leq c_{A_{\text{ref}}} \|y\|_{A_{\text{ref}}}$ holds for all $y \in Y$. The diameters of the triangles and edges are denoted by h_T and h_E , respectively.

Remark 7.7. This estimate can possibly be refined if the local properties of the interpolation operator are considered more carefully by including the coefficient function. Then the constant $c_{A_{\text{ref}}}$ could be improved.

The overall error estimator looks as follows:

Theorem 7.8. *Let $\tilde{y}(\xi) \in Y$ be given and let $y(\xi) \in Y$ be the exact solution of (7.1), where the respective operators are defined as in (7.10). Furthermore, let Assumption 7.5 hold and define $w(\xi)$ by (7.8) and (7.5).*

Then,

$$\|\tilde{y}(\xi) - y(\xi)\|_{A_{\text{ref}}}^2 \leq \Lambda^2 \|w(\xi)\|_{A_{\text{ref}}}^2 + \Lambda^2 c_{\mathcal{T}}^2 c_{A_{\text{ref}}}^2 \left(\sum_{T \in \mathcal{T}} \eta_T(\tilde{y}(\xi))^2 + \sum_{E \in \mathcal{E}^0} \eta_E(\tilde{y}(\xi))^2 \right) \quad (7.16)$$

holds with Λ from (7.3), $c_{\mathcal{T}}$ and $c_{A_{\text{ref}}}$ from (7.13), and $\eta_T(\tilde{y}(\xi))$ and $\eta_E(\tilde{y}(\xi))$ as in (7.14), (7.15).

Proof. Combining (7.9) and (7.13) yields the desired result. \square

7.2. A Posteriori Error Estimation in $L_{\mathbb{P}}^2(\Xi)$

In this subsection we discuss an a posteriori error estimator for $Y = \mathbb{R}$ and $\hat{\varphi} \equiv 0$, i. e., a linear equation on \mathbb{R} , discretized by a stochastic Galerkin approach. This is only a simple model problem. The essential ideas will be used later to combine both error estimators to one for the PDE with stochastic coefficients.

Let $A(\xi) : \mathbb{R} \rightarrow \mathbb{R}$ be a linear, uniformly coercive and bounded operator mapping from \mathbb{R} to \mathbb{R} , i. e., we identify $A(\xi) \in [\underline{\kappa}, \bar{\kappa}] \subset \mathbb{R}_{>0}$, and let $b \in L_{\mathbb{P}}^2(\Xi)$. We want to solve $A(\xi)y(\xi) = b(\xi)$ by a stochastic Galerkin method with polynomial chaos. Let $y \in L_{\mathbb{P}}^2(\Xi)$ be the unique weak solution fulfilling

$$(Ay, v)_{L_{\mathbb{P}}^2(\Xi)} = (b, v)_{L_{\mathbb{P}}^2(\Xi)} \quad \text{for all } v \in L_{\mathbb{P}}^2(\Xi),$$

i. e., $y(\xi) = \frac{b(\xi)}{A(\xi)}$ almost surely. We choose polynomials of degree at most $d_i - 1 \in \mathbb{N}_0$ as discretization subspace $\mathcal{P}_{d_i-1}(\Xi_i) \subset L_{\mathbb{P}}^2(\Xi_i)$ and take the full tensor product

$$\mathcal{P}_{d-1}(\Xi) := \bigotimes_{i=1}^m \mathcal{P}_{d_i-1}(\Xi_i) \subset \bigotimes_{i=1}^m L_{\mathbb{P}_i}^2(\Xi_i) = L_{\mathbb{P}}^2(\Xi),$$

where $d \in \mathbb{N}^m$ is a vector, cf. Section 6.2. Let $y \in \mathcal{P}_{d-1}(\Xi)$ be the Galerkin solution fulfilling

$$(Ay, v)_{L_{\mathbb{P}}^2(\Xi)} = (b, v)_{L_{\mathbb{P}}^2(\Xi)} \quad \text{for all } v \in \mathcal{P}_{d-1}(\Xi).$$

Let $\{\beta_{k_i}^{(i)}\}_{k_i=1}^{\infty} \subset L_{\mathbb{P}_i}^2(\Xi_i)$ be sets of orthonormal polynomials w. r. t. the $L_{\mathbb{P}_i}^2(\Xi_i)$ inner product, where $\beta_{k_i}^{(i)}$ has degree $k_i - 1$, and let $\{\beta_k\}_{k \in \mathbb{N}^m}$ be the corresponding orthonormal tensor product polynomials which form a Hilbert basis of $L_{\mathbb{P}}^2(\Xi)$, see Section 6.2. Then, for every $v \in L_{\mathbb{P}}^2(\Xi)$ there exist unique coefficients $(\mu_k)_k \in \ell^2(\mathbb{N}^m)$ such that $v(\xi) = \sum_{k \in \mathbb{N}^m} \mu_k \beta_k(\xi)$ holds almost surely and the series converges in the $L_{\mathbb{P}}^2(\Xi)$ sense.

We have

$$(A(y - y), v)_{L_{\mathbb{P}}^2(\Xi)} = (b - Ay, \sum_{k \in \mathbb{N}^m} \mu_k \beta_k)_{L_{\mathbb{P}}^2(\Xi)} = \sum_{k \not\leq d} \mu_k (b - Ay, \beta_k)_{L_{\mathbb{P}}^2(\Xi)}.$$

The summands for $k \leq d$ (componentwise) cancel out due to Galerkin orthogonality. This corresponds to testing with $v(\xi) := \sum_{k \leq d} \mu_k \beta_k(\xi)$. Note that we write $k \not\leq d$ and not $k > d$ since we compare index vectors here.

Now we use an approximation $\tilde{A} \in \mathcal{P}_{\tilde{d}}(\Xi)$ of A with $\tilde{d} \in \mathbb{N}_0^m$ and an approximation $\tilde{b} \in \mathcal{P}_{d+\tilde{d}-1}(\Xi)$ of b and compute

$$(b - Ay, v)_{L^2_{\mathbb{P}}(\Xi)} = (\tilde{b} - \tilde{A}y, v)_{L^2_{\mathbb{P}}(\Xi)} + ((\tilde{A} - A)y, v)_{L^2_{\mathbb{P}}(\Xi)} + (b - \tilde{b}, v)_{L^2_{\mathbb{P}}(\Xi)}.$$

Since $\tilde{b} - \tilde{A}y \in \mathcal{P}_{d+\tilde{d}-1}(\Xi)$, we have

$$\begin{aligned} \sum_{k \not\leq d} \mu_k (\tilde{b} - \tilde{A}y, \beta_k)_{L^2_{\mathbb{P}}(\Xi)} &= \sum_{\substack{k \leq d+\tilde{d} \\ k \not\leq d}} \mu_k (\tilde{b} - \tilde{A}y, \beta_k)_{L^2_{\mathbb{P}}(\Xi)} \\ &\leq \left(\sum_{\substack{k \leq d+\tilde{d} \\ k \not\leq d}} \mu_k^2 \right)^{1/2} \left(\sum_{\substack{k \leq d+\tilde{d} \\ k \not\leq d}} (\tilde{b} - \tilde{A}y, \beta_k)_{L^2_{\mathbb{P}}(\Xi)}^2 \right)^{1/2} \end{aligned}$$

due to orthogonality. The sum $\sum_{\substack{k \leq d+\tilde{d} \\ k \not\leq d}} (\tilde{b} - \tilde{A}y, \beta_k)_{L^2_{\mathbb{P}}(\Xi)}^2$ is finite and can be evaluated simply, e. g., by using a quadrature formula of high enough order or analytically. Inserting $v = y - y$ and using $\|(\mu_k)_k\|_{\ell^2(\mathbb{N})^m} = \|v\|_{L^2_{\mathbb{P}}(\Xi)}$ gives

$$\begin{aligned} \|y - y\|_A^2 &= (\tilde{b} - \tilde{A}y, y - y)_{L^2_{\mathbb{P}}(\Xi)} + ((\tilde{A} - A)y, y - y)_{L^2_{\mathbb{P}}(\Xi)} + (b - \tilde{b}, y - y)_{L^2_{\mathbb{P}}(\Xi)} \leq \\ &= \left(\sum_{\substack{k \leq d+\tilde{d} \\ k \not\leq d}} (\tilde{b} - \tilde{A}y, \beta_k)_{L^2_{\mathbb{P}}(\Xi)}^2 \right)^{1/2} + \|(\tilde{A} - A)y\|_{L^2_{\mathbb{P}}(\Xi)} + \|b - \tilde{b}\|_{L^2_{\mathbb{P}}(\Xi)} \|y - y\|_{L^2_{\mathbb{P}}(\Xi)}. \end{aligned}$$

By $\|v\|_{L^2_{\mathbb{P}}(\Xi)} \leq \frac{1}{\sqrt{\kappa}} \|v\|_A$ we get

$$\|y - y\|_A \leq \frac{1}{\sqrt{\kappa}} \left(\left(\sum_{\substack{k \leq d+\tilde{d} \\ k \not\leq d}} (\tilde{b} - \tilde{A}y, \beta_k)_{L^2_{\mathbb{P}}(\Xi)}^2 \right)^{1/2} + \|(\tilde{A} - A)y\|_{L^2_{\mathbb{P}}(\Xi)} + \|b - \tilde{b}\|_{L^2_{\mathbb{P}}(\Xi)} \right).$$

We see that even for $\tilde{d} = \mathbb{1}$ there are $2^m - 1$ different basis functions β_k to test with. Therefore, it is convenient to use further structure of the involved polynomials. If, e. g., $y \in \mathcal{P}_{d-1}(\Xi)$, $\tilde{b} \in \mathcal{P}_{d-1}(\Xi)$ and \tilde{A} is affine, i. e., $\tilde{A}(\xi) = \tilde{A}_0 + \sum_{i=1}^m \xi_i \tilde{A}_i$, we obtain for $k \not\leq d$, $k \leq d + \mathbb{1}$:

$$\begin{aligned} (\tilde{b} - \tilde{A}y, \beta_k)_{L^2_{\mathbb{P}}(\Xi)} &= - \left(\sum_{i=1}^m \xi_i \tilde{A}_i y(\xi), \beta_k \right)_{L^2_{\mathbb{P}}(\Xi)} \\ &= \begin{cases} - (\xi_i \tilde{A}_i y(\xi), \beta_k)_{L^2_{\mathbb{P}}(\Xi)} & \text{if } k = d + e_i, \\ 0 & \text{otherwise.} \end{cases} \end{aligned}$$

This gives

$$\sum_{\substack{k \leq d+\mathbb{1} \\ k \not\leq d}} (\tilde{b} - \tilde{A}y, \beta_k)_{L^2_{\mathbb{P}}(\Xi)}^2 = \sum_{i=1}^m (\xi_i \tilde{A}_i y(\xi), \beta_{d+e_i})_{L^2_{\mathbb{P}}(\Xi)}^2,$$

which can be interpreted as a decomposition of the error into contributions for each parameter.

7.3. Combination of Both Error Estimators

We combine the deterministic FE a posteriori error estimator (Section 7.1) with the one for the polynomial discretization of $L^2_{\mathbb{P}}(\Xi)$ (Section 7.2) to obtain an overall error estimator for the solution of the weak formulation of (7.1). For this purpose, we return to the notation from Section 7.1.

We consider the weak formulation of (7.1) w. r. t. the parameters:

$$\langle \mathbf{A}\mathbf{y} + \mathbf{N}(\mathbf{y}), \mathbf{v} \rangle_{\mathbf{Y}^*, \mathbf{Y}} = \langle \mathbf{b}, \mathbf{v} \rangle_{\mathbf{Y}^*, \mathbf{Y}} \quad \text{for all } \mathbf{v} \in \mathbf{Y}. \quad (7.17)$$

The operator \mathbf{A} is defined in (7.2). Analogously, we have

$$\begin{aligned} \langle \mathbf{N}(\mathbf{y}), \mathbf{v} \rangle_{\mathbf{Y}^*, \mathbf{Y}} &:= \int_{\Xi} \langle \mathbf{N}(\mathbf{y}(\cdot, \xi)), \mathbf{v}(\cdot, \xi) \rangle_{\mathbf{Y}^*, \mathbf{Y}} d\mathbb{P} \quad \text{and} \\ \langle \mathbf{b}, \mathbf{v} \rangle_{\mathbf{Y}^*, \mathbf{Y}} &:= \int_{\Xi} \langle b(\xi), \mathbf{v}(\cdot, \xi) \rangle_{\mathbf{Y}^*, \mathbf{Y}} d\mathbb{P} \quad \text{for all } \mathbf{v} \in \mathbf{Y}. \end{aligned}$$

The exact solution \mathbf{y} of (7.17) is contained in $\mathbf{Y} = L^p_{\mathbb{P}}(\Xi; Y)$ with $p \in (3, \infty)$, which we require for the weak formulation to be well-defined and for the nonlinear operator \mathbf{N} to be twice continuously differentiable, but we will provide only error estimation in $\hat{\mathbf{Y}} = L^2_{\mathbb{P}}(\Xi; Y)$, because we want to make use of (Galerkin) orthogonality.

Analogously to (7.2) we define the reference operator $\mathbf{A}_{\text{ref}} \in \mathcal{L}(\mathbf{Y}, \mathbf{Y}^*)$ by

$$\langle \mathbf{A}_{\text{ref}}\mathbf{y}, \mathbf{v} \rangle_{\mathbf{Y}^*, \mathbf{Y}} := \int_{\Xi} \langle \mathbf{A}_{\text{ref}}\mathbf{y}(\cdot, \xi), \mathbf{v}(\cdot, \xi) \rangle_{\mathbf{Y}^*, \mathbf{Y}} d\mathbb{P}$$

and its version $\hat{\mathbf{A}}_{\text{ref}} \in \mathcal{L}(\hat{\mathbf{Y}}, \hat{\mathbf{Y}}^*)$. The operators $\hat{\mathbf{A}}$, $\hat{\mathbf{A}}_{\text{ref}}$, and their inverses induce inner products and respective norms on $\hat{\mathbf{Y}}$ and $\hat{\mathbf{Y}}^*$, respectively. The estimates (7.3) and (7.4) carry over to this setting.

For every $\tilde{\mathbf{y}} \in \mathbf{Y} \subset \hat{\mathbf{Y}}$ we can conclude

$$\|\tilde{\mathbf{y}} - \mathbf{y}\|_{\hat{\mathbf{A}}_{\text{ref}}} \leq \Lambda \|\mathbf{r}\|_{\hat{\mathbf{A}}_{\text{ref}}^{-1}} \quad (7.18)$$

in analogy to Lemma 7.1 with $\mathbf{r} = \mathbf{A}\tilde{\mathbf{y}} + \mathbf{N}(\tilde{\mathbf{y}}) - \mathbf{b}$. Note that estimating the error in L^2 enables us to use the coercivity of the operators $\hat{\mathbf{A}}$ and $\hat{\mathbf{A}}_{\text{ref}}$. We need $\mathbf{r} \in L^2_{\mathbb{P}}(\Xi; Y^*)$, which can be concluded using integrability properties of $\tilde{\mathbf{y}}$ and the regularity of the data, see Sections 3.2 and 3.4, and thus $\mathbf{w} = \hat{\mathbf{A}}_{\text{ref}}^{-1}\mathbf{r} \in \hat{\mathbf{Y}}$. For the estimation of the total error we compute $\|\mathbf{w}\|_{\hat{\mathbf{A}}_{\text{ref}}}^2$ and have to estimate $\|\mathbf{w} - \mathbf{w}\|_{\hat{\mathbf{A}}_{\text{ref}}}^2$, cf. Lemma 7.3. Thus, it makes sense to minimize the error term $\|\mathbf{w}\|_{\hat{\mathbf{A}}_{\text{ref}}}^2$ by an iterative solver used for the discretized equation. We have that $\mathbf{w} \in Y \otimes \mathcal{P}_{d-1}(\Xi) \subset L^2_{\mathbb{P}}(\Xi; Y)$ solves

$$\langle \hat{\mathbf{A}}_{\text{ref}}\mathbf{w}, \mathbf{v} \rangle_{\hat{\mathbf{Y}}^*, \hat{\mathbf{Y}}} = \langle \mathbf{r}, \mathbf{v} \rangle_{\hat{\mathbf{Y}}^*, \hat{\mathbf{Y}}} \quad \text{for all } \mathbf{v} \in Y \otimes \mathcal{P}_{d-1}(\Xi). \quad (7.19)$$

We write an arbitrary test function in the form $\mathbf{v}(x, \xi) = \sum_{k \in \mathbb{N}^m} v_k^*(x) \vartheta_k(\xi)$ with some

Hilbert basis $\{\vartheta_k\}_{k \in \mathbb{N}^m} \subset L^2_{\mathbb{P}}(\Xi)$ and $v_k^* \in Y$ for all k . We have

$$\begin{aligned}
 & \langle \hat{\mathbf{A}}_{\text{ref}}(\mathbf{w} - \mathbf{w}), \mathbf{v} \rangle_{L^2_{\mathbb{P}}(\Xi; Y^*), L^2_{\mathbb{P}}(\Xi; Y)} \\
 &= \langle \mathbf{r} - \hat{\mathbf{A}}_{\text{ref}} \mathbf{w}, \sum_{k \in \mathbb{N}^m} v_k^* \vartheta_k \rangle_{\hat{Y}^*, \hat{Y}} \\
 &= \sum_{k \in \mathbb{N}^m} \langle \mathbf{r} - \hat{\mathbf{A}}_{\text{ref}} \mathbf{w}, v_k^* \vartheta_k \rangle_{\hat{Y}^*, \hat{Y}} \\
 &= \sum_{k \leq d} \langle \mathbf{r} - \hat{\mathbf{A}}_{\text{ref}} \mathbf{w}, (v_k^* - v_k^*) \vartheta_k \rangle_{\hat{Y}^*, \hat{Y}} + \sum_{k \not\leq d} \langle \mathbf{r} - \hat{\mathbf{A}}_{\text{ref}} \mathbf{w}, v_k^* \vartheta_k \rangle_{\hat{Y}^*, \hat{Y}}
 \end{aligned} \tag{7.20}$$

with $v_k^* \in Y$.

Assumption 7.9. We need the following prerequisites in addition to Assumption 7.5:

- The coefficient function is affine in ξ , i. e., $\kappa(x, \xi) = \kappa_0(x) + \sum_{i=1}^m \xi_i \kappa_i(x)$. This form can originate in a truncated Karhunen-Loève expansion. The functions κ_i ($i \in \{0, \dots, m\}$) are assumed to be piecewise constant on the triangles. We define $\langle A_i y, v \rangle_{Y^*, Y} := \int_{\Omega} \kappa_i \nabla y \cdot \nabla v \, dx$ for $i \in \{0, \dots, m\}$.
- Let $\tilde{\mathbf{y}} \in Y \otimes \mathcal{P}_{d-1}(\Xi)$. We assume that the nonlinear part $\hat{\varphi}(\tilde{\mathbf{y}})$ of the residual can be approximated sufficiently well by a function $\tilde{\varphi} \in L^2(\Omega) \otimes \mathcal{P}_{d-1}(\Xi)$, e. g., by interpolation. Moreover, we assume that the function $\chi \cdot \tilde{\mathbf{y}}$ can be approximated sufficiently well by a function $\tilde{\chi} \in L^2(\Omega) \otimes \mathcal{P}_{d-1}(\Xi)$.
- The right-hand side \mathbf{b} can be identified with a function $\hat{\mathbf{f}} \in L^2(\Omega) \otimes \mathcal{P}_{d-1}(\Xi)$.

Again, one can include the interpolation error as follows:

$$\begin{aligned}
 & \langle \mathbf{r} - \hat{\mathbf{A}}_{\text{ref}} \mathbf{w}, \mathbf{v} \rangle_{\hat{Y}^*, \hat{Y}} = \\
 & \int_{\Xi} \int_{\Omega} \kappa \nabla_x \tilde{\mathbf{y}} \cdot \nabla_x \mathbf{v} + \chi \cdot \tilde{\mathbf{y}} \cdot \mathbf{v} + \varphi(\tilde{\mathbf{y}}) \cdot \mathbf{v} - \hat{\mathbf{f}} \cdot \mathbf{v} - \kappa_{\text{ref}} \nabla_x \tilde{\mathbf{w}} \cdot \nabla_x \mathbf{v} \, dx \, d\mathbb{P} = \\
 & \int_{\Xi} \int_{\Omega} \kappa \nabla_x \tilde{\mathbf{y}} \cdot \nabla_x \mathbf{v} + \tilde{\chi} \cdot \mathbf{v} + \tilde{\varphi} \cdot \mathbf{v} - \hat{\mathbf{f}} \cdot \mathbf{v} - \kappa_{\text{ref}} \nabla_x \tilde{\mathbf{w}} \cdot \nabla_x \mathbf{v} \, dx \, d\mathbb{P} \\
 & + \int_{\Xi} \int_{\Omega} (\chi \cdot \tilde{\mathbf{y}} - \tilde{\chi} + \varphi(\tilde{\mathbf{y}}) - \tilde{\varphi}) \cdot \mathbf{v} \, dx \, d\mathbb{P}.
 \end{aligned} \tag{7.21}$$

In the following, we will neglect $\chi \cdot \tilde{\mathbf{y}} - \tilde{\chi}$ and $\varphi(\tilde{\mathbf{y}}) - \tilde{\varphi}$ for the ease of presentation and implementation. Moreover, we allow for different orthonormal polynomial bases in the derivation of the error estimates.

Conversion between Polynomial Basis Representations

To discretize the space $L^2_{\mathbb{P}_i}(\Xi_i)$, we choose orthonormal polynomials in tensor product form, as already described in Section 6.2 and used in Section 7.2. Typically, the orthonormal basis $\{\beta_{k_i}^{(i)}\}_{k_i=1}^{\infty}$, where $\beta_{k_i}^{(i)}$ has degree $k_i - 1$, and the tensor product basis $\{\beta_k\}_{k \in \mathbb{N}^m}$ of $L^2_{\mathbb{P}}(\Xi)$ is used. Alternatively, we have the Hilbert basis $\{\theta_k\}_{k \in \mathbb{N}^m}$, which results from replacing the

first d_i univariate polynomials in $\{\beta_{k_i}^{(i)}\}_{k_i=1}^\infty$ by weighted, orthonormal Lagrange polynomials w. r. t. the Gaussian quadrature nodes.

Now consider a function $\mathbf{w} \in Y \otimes \mathcal{P}_d(\Xi)$ written as

$$\mathbf{w}(x, \xi) = \sum_{k \leq d} \mathbf{w}_k^*(x) \beta_k(\xi) = \sum_{k \leq d} \mathbf{w}_k(x) \theta_k(\xi), \quad (7.22)$$

where $\mathbf{w}_k^*, \mathbf{w}_k \in Y$ are the coefficients. On the one hand, inserting $\xi = \mathbf{a}_l$ (Gaussian nodes) yields

$$\mathbf{w}_l(x) = \boldsymbol{\omega}_l^{-1} \cdot \sum_{k \leq d} \mathbf{w}_k^*(x) \beta_k(\mathbf{a}_l).$$

On the other hand, testing (7.22) with β_l gives

$$\mathbf{w}_l^*(x) = \sum_{k \leq d} \left(\mathbf{w}_k(x) (\theta_k, \beta_l)_{L^2_{\mathbb{P}^i}(\Xi)} \right) = \sum_{k \leq d} \left(\mathbf{w}_k(x) \prod_{i=1}^m (\theta_{k_i}^{(i)}, \beta_{l_i}^{(i)})_{L^2_{\mathbb{P}^i}(\Xi_i)} \right).$$

Defining the orthonormal matrices $\mathbf{Q}^{(i)} \in \mathbb{R}^{d_i+1 \times d_i+1}$ by $\mathbf{Q}_{k_i l_i}^{(i)} = (\theta_{k_i}^{(i)}, \beta_{l_i}^{(i)})_{L^2_{\mathbb{P}^i}(\Xi_i)}$, we see that the coefficients \mathbf{w}_l^* ($l \leq d$) of one orthonormal basis in tensor product form can be computed by applying an orthonormal rank-1-operator $\mathbf{Q}^{(1)} \otimes \dots \otimes \mathbf{Q}^{(m)}$ to the coefficients \mathbf{w}_k ($k \leq d$) of another orthonormal tensor product basis.

Discretization with Lagrange Polynomials

For a simpler approximation of the nonlinearity, it is beneficial to use the Lagrange basis $\{\theta_k\}_{k \in \mathbb{N}_0^m}$, as seen in Section 6.2. Still, we can convert the coefficients to representations in other bases if necessary. Furthermore, many results in this subsection still hold for arbitrary orthonormal polynomials. We will point out where we use special properties of the chosen Lagrange basis.

The polynomial basis representations of $\tilde{\mathbf{y}}$ and \mathbf{w} shall be

$$\tilde{\mathbf{y}}(x, \xi) = \sum_{l \leq d} \tilde{y}_l(x) \theta_l(\xi), \quad \mathbf{w}(x, \xi) = \sum_{l \leq d} \mathbf{w}_l(x) \theta_l(\xi)$$

with $\tilde{y}_l, \mathbf{w}_l \in Y$. Be aware of the fact that these sums are actually large since l is an index vector. Analogously, we write $\hat{\mathbf{f}}, \tilde{\boldsymbol{\varphi}}$, and $\tilde{\boldsymbol{\chi}}$ with the deterministic functions $f_l, \tilde{\varphi}_l, \tilde{\chi}_l \in L^2(\Omega)$ and an arbitrary test function $\mathbf{v} \in \hat{\mathbf{Y}}$ as $\mathbf{v}(x, \xi) = \sum_{k \in \mathbb{N}^m} v_k(x) \theta_k(\xi)$.

Moreover, we will use transformed representations $\tilde{\mathbf{y}}(x, \xi) = \sum_{l \leq d} \tilde{y}_l^*(x) \vartheta_l(\xi)$ as well as $\mathbf{v}(x, \xi) = \sum_{k \in \mathbb{N}^m} v_k^*(x) \vartheta_k(\xi)$, where $\{\vartheta_{k_i}^{(i)}\}_{k_i=1}^\infty$ are alternative, still to be specified orthonormal bases of $L^2_{\mathbb{P}^i}(\Xi_i)$ with $\vartheta_{k_i}^{(i)} \equiv \theta_{k_i}^{(i)} \equiv \beta_{k_i}^{(i)}$ for $k_i \geq d_i + 1$ and $\tilde{y}_l^* \in Y, v_k^* \in Y$. These coefficients can simply be computed from \tilde{y}_l and v_k by an orthonormal transformation. Especially for \mathbf{v} , we split up the representation into

$$\mathbf{v}(x, \xi) = \sum_{k \leq d} v_k(x) \theta_k(\xi) + \sum_{k \not\leq d} v_k^*(x) \vartheta_k(\xi).$$

For $k \leq d$, we obtain

$$\begin{aligned}
 & \langle \mathbf{r} - \hat{\mathbf{A}}_{\text{ref}} \mathbf{w}, v_k \theta_k \rangle_{L_{\mathbb{P}}^2(\Xi; Y^*), L_{\mathbb{P}}^2(\Xi; Y)} \\
 &= \int_{\Xi} \langle A(\xi) \tilde{\mathbf{y}}(\cdot, \xi) + N(\tilde{\mathbf{y}}(\cdot, \xi)) - b(\xi) - A_{\text{ref}} \mathbf{w}(\cdot, \xi), v_k \rangle_{Y^*, Y} \theta_k(\xi) \, d\mathbb{P} \\
 &\approx \sum_{l \leq d} \int_{\Xi} \left(\int_{\Omega} (\kappa_0 + \xi_1 \kappa_1 + \dots + \xi_m \kappa_m) \nabla \tilde{\mathbf{y}}_l \cdot \nabla v_k + \tilde{\chi}_l v_k + \tilde{\varphi}_l v_k \right. \\
 &\quad \left. - f_l v_k - \kappa_{\text{ref}} \nabla \mathbf{w}_l \cdot \nabla v_k - \chi_{\text{ref}} \mathbf{w}_l v_k \, dx \right) \theta_l(\xi) \theta_k(\xi) \, d\mathbb{P} \\
 &= \sum_{l \leq d} \left(\int_{\Omega} \kappa_0 \nabla \tilde{\mathbf{y}}_l \cdot \nabla v_k + \tilde{\chi}_l v_k + \tilde{\varphi}_l v_k \right. \\
 &\quad \left. - f_l v_k - \kappa_{\text{ref}} \nabla \mathbf{w}_l \cdot \nabla v_k - \chi_{\text{ref}} \mathbf{w}_l v_k \, dx \right) \left(\int_{\Xi} \theta_l(\xi) \theta_k(\xi) \, d\mathbb{P} \right) \\
 &\quad + \sum_{l \leq d} \sum_{i=1}^m \left(\int_{\Omega} \kappa_i \nabla \tilde{\mathbf{y}}_l \cdot \nabla v_k \, dx \right) \left(\int_{\Xi} \xi_i \theta_l(\xi) \theta_k(\xi) \, d\mathbb{P} \right) \\
 &= \int_{\Omega} (\kappa_0 + \mathbf{a}_{k_1}^{(1)} \kappa_1 + \dots + \mathbf{a}_{k_m}^{(m)} \kappa_m) \nabla \tilde{\mathbf{y}}_k \cdot \nabla v_k + \tilde{\chi}_k v_k \\
 &\quad + \tilde{\varphi}_k v_k - f_k v_k - \kappa_{\text{ref}} \nabla \mathbf{w}_k \cdot \nabla v_k - \chi_{\text{ref}} \mathbf{w}_k v_k \, dx. \tag{7.23}
 \end{aligned}$$

This is due to orthogonality of the polynomials $\{\theta_k\}_{k \in \mathbb{N}^m}$ and the fact

$$\begin{aligned}
 \int_{\Xi} \xi_i \theta_l(\xi) \theta_k(\xi) \, d\mathbb{P} &= \int_{\Xi} \xi_i \prod_{j=1}^m \theta_{l_j}^{(j)}(\xi_j) \theta_{k_j}^{(j)}(\xi_j) \, d\mathbb{P} \\
 &= \int_{\Xi_i} \xi_i \theta_{l_i}^{(i)}(\xi_i) \theta_{k_i}^{(i)}(\xi_i) \, d\mathbb{P}_i \cdot \prod_{j=1, j \neq i}^m \delta_{l_j k_j} \\
 &= \begin{cases} \mathbf{a}_{k_i}^{(i)} & \text{if } k = l, \\ 0 & \text{otherwise.} \end{cases}
 \end{aligned}$$

Recall that $\{\mathbf{a}_{k_i}^{(i)}\}_{k_i=1}^{d_i} \subset \Xi_i$ are the Gaussian quadrature nodes w.r.t. \mathbb{P}_i . For a different polynomial basis we would obtain different values here and possibly no decoupling of the deterministic equations which we have to solve.

For $k \not\leq d$ (yielding $k \neq l$) we get

$$\begin{aligned}
 & \langle \mathbf{r} - \hat{\mathbf{A}}_{\text{ref}} \mathbf{w}, v_k^* \vartheta_k \rangle_{L_{\mathbb{P}}^2(\Xi; Y^*), L_{\mathbb{P}}^2(\Xi; Y)} \\
 &\approx \sum_{l \leq d} \sum_{i=1}^m \left(\int_{\Omega} \kappa_i \nabla \tilde{\mathbf{y}}_l^* \cdot \nabla v_k^* \, dx \right) \left(\int_{\Xi} \xi_i \vartheta_l(\xi) \vartheta_k(\xi) \, d\mathbb{P} \right) \\
 &= \begin{cases} \sum_{l_i=1}^{d_i} \mathbf{c}_{l_i}^{(i)} \int_{\Omega} \kappa_i \nabla \tilde{\mathbf{y}}_{k_1, \dots, l_i, \dots, k_m}^* \cdot \nabla v_k^* \, dx & \text{if } \begin{cases} k_i = d_i + 1 \text{ for an } i \in [m], \\ k_j \leq d_j \text{ for } j \in [m] \setminus \{i\}, \end{cases} \\ 0 & \text{otherwise.} \end{cases} \tag{7.24}
 \end{aligned}$$

In this computation we used the following properties of any orthonormal polynomial basis $\{\vartheta_l\}_{l \in \mathbb{N}^m}$ in tensor product form, where $\vartheta_{l_i}^{(i)} \perp \mathcal{P}_{l_i-2}(\Xi_i)$ for all $l_i \geq d_i + 1$:

$$\begin{aligned} \int_{\Xi} \xi_i \vartheta_l(\xi) \vartheta_k(\xi) \, d\mathbb{P} &= \int_{\Xi_i} \xi_i \vartheta_{l_i}^{(i)}(\xi_i) \vartheta_{k_i}^{(i)}(\xi_i) \, d\mathbb{P}_i \cdot \prod_{j=1, j \neq i}^m \delta_{l_j k_j} \\ &= \begin{cases} c_{l_i}^{(i)} & \text{if } k_i = d_i + 1 \text{ and } k_j = l_j \text{ for all } j \in [m] \setminus \{i\} \\ 0 & \text{otherwise} \end{cases} \end{aligned}$$

with $c_{l_i}^{(i)} := \int_{\Xi_i} \xi_i \vartheta_{l_i}^{(i)}(\xi_i) \vartheta_{d_i+1}^{(i)}(\xi_i) \, d\mathbb{P}_i$ holds for any $l \leq d$ and $k \not\leq d$.

Example 7.10. For $\vartheta_{l_i}^{(i)} \equiv \beta_{l_i}^{(i)}$ it even holds $c_{l_i}^{(i)} = 0$ for all $l_i \leq d_i - 1$. Then we get

$$\sum_{l_i=1}^{d_i} c_{l_i}^{(i)} \int_{\Omega} \kappa_i \nabla \tilde{y}_{k_1, \dots, l_i, \dots, k_m}^* \cdot \nabla v_k^* \, dx = c_{d_i}^{(i)} \int_{\Omega} \kappa_i \nabla \tilde{y}_{k_1, \dots, d_i, \dots, k_m}^* \cdot \nabla v_k^* \, dx$$

in (7.24).

From (7.20), (7.23), and (7.24) we get

$$\begin{aligned} &\langle \hat{\mathbf{A}}_{\text{ref}}(\mathbf{w} - \mathbf{w}), \mathbf{v} \rangle_{\hat{\mathbf{Y}}^*, \hat{\mathbf{Y}}} \\ &= \sum_{k \leq d} \langle \mathbf{r} - \hat{\mathbf{A}}_{\text{ref}} \mathbf{w}, (v_k - v_k) \theta_k \rangle_{\hat{\mathbf{Y}}^*, \hat{\mathbf{Y}}} + \sum_{k \not\leq d} \langle \mathbf{r} - \hat{\mathbf{A}}_{\text{ref}} \mathbf{w}, v_k^* \vartheta_k \rangle_{\hat{\mathbf{Y}}^*, \hat{\mathbf{Y}}} \\ &\approx \sum_{k \leq d} \int_{\Omega} (\kappa(\cdot, \mathbf{a}_k) \nabla \tilde{y}_k - \kappa_{\text{ref}} \nabla \mathbf{w}_k) \cdot \nabla (v_k - v_k) + (\tilde{\chi}_k + \tilde{\varphi}_k - f_k - \chi_{\text{ref}} \mathbf{w}_k) (v_k - v_k) \, dx \\ &\quad + \sum_{i=1}^m \sum_{k \leq d} c_{k_i}^{(i)} \int_{\Omega} \kappa_i \nabla \tilde{y}_k^* \cdot \nabla v_{k_1, \dots, d_i+1, \dots, d_m}^* \, dx. \end{aligned} \tag{7.25}$$

Recall that we write $\mathbf{a}_k := (\mathbf{a}_{k_1}^{(1)}, \dots, \mathbf{a}_{k_m}^{(m)})^\top$.

We see that the first term at the end of (7.25) can be treated as in Section 7.1 and constitutes the deterministic part of the error. The second term can be viewed as the stochastic part of the error, where the i -th summand corresponds to the discretization error of $L_{\mathbb{P}_i}^2(\Xi_i)$.

The first term in (7.25) is estimated as in (7.11), (7.12) and is bounded as in (7.13) using the Clément interpolant v_k of v_k :

$$\begin{aligned} &\sum_{k \leq d} \int_{\Omega} (\kappa(\cdot, \mathbf{a}_k) \nabla \tilde{y}_k - \kappa_{\text{ref}} \nabla \mathbf{w}_k) \cdot \nabla (v_k - v_k) + (\tilde{\chi}_k + \tilde{\varphi}_k - f_k - \chi_{\text{ref}} \mathbf{w}_k) (v_k - v_k) \, dx \\ &= \sum_{k \leq d} \sum_{T \in \mathcal{T}} \left(\int_T (\tilde{\chi}_k + \tilde{\varphi}_k - f_k - \chi_{\text{ref}} \mathbf{w}_k) (v_k - v_k) \, dx \right. \\ &\quad \left. + \int_{\partial T} (\kappa(\cdot, \mathbf{a}_k) \frac{\partial}{\partial \nu_T} \tilde{y}_k - \kappa_{\text{ref}} \frac{\partial}{\partial \nu_T} \mathbf{w}_k) (v_k - v_k) \, dS \right) \end{aligned}$$

$$\begin{aligned}
 &\leq \sum_{T \in \mathcal{T}} \sum_{k \leq d} \|\tilde{\chi}_k + \tilde{\varphi}_k - f_k - \chi_{\text{ref}} \mathbf{w}_k\|_{L^2(T)} \|v_k - \mathbf{v}_k\|_{L^2(T)} \\
 &\quad + \sum_{E \in \mathcal{E}^0} \sum_{k \leq d} \|\llbracket (\kappa(\cdot, \mathbf{a}_k) \nabla \tilde{y}_k - \kappa_{\text{ref}} \nabla \mathbf{w}_k) \cdot \nu_E \rrbracket_E\|_{L^2(E)} \|v_k - \mathbf{v}_k\|_{L^2(E)} \\
 &\leq c_{\mathcal{T}} c_{A_{\text{ref}}} \left(\sum_{T \in \mathcal{T}} h_T^2 \sum_{k \leq d} \|\tilde{\chi}_k + \tilde{\varphi}_k - f_k - \chi_{\text{ref}} \mathbf{w}_k\|_{L^2(T)}^2 \right) \\
 &\quad + \sum_{E \in \mathcal{E}^0} h_E \sum_{k \leq d} \|\llbracket (\kappa(\cdot, \mathbf{a}_k) \nabla \tilde{y}_k - \kappa_{\text{ref}} \nabla \mathbf{w}_k) \cdot \nu_E \rrbracket_E\|_{L^2(E)}^2 \left(\sum_{k \leq d} \|v_k\|_{A_{\text{ref}}}^2 \right)^{\frac{1}{2}} \\
 &= c_{\mathcal{T}} c_{A_{\text{ref}}} \left(\sum_{T \in \mathcal{T}} \boldsymbol{\eta}_T(\tilde{\mathbf{y}})^2 + \sum_{E \in \mathcal{E}^0} \boldsymbol{\eta}_E(\tilde{\mathbf{y}})^2 \right)^{\frac{1}{2}} \left(\sum_{k \leq d} \|v_k\|_{A_{\text{ref}}}^2 \right)^{\frac{1}{2}}
 \end{aligned} \tag{7.26}$$

with

$$\boldsymbol{\eta}_T(\tilde{\mathbf{y}}) := h_T \left(\sum_{k \leq d} \|\tilde{\chi}_k + \tilde{\varphi}_k - f_k - \chi_{\text{ref}} \mathbf{w}_k\|_{L^2(T)}^2 \right)^{\frac{1}{2}} \tag{7.27}$$

and

$$\boldsymbol{\eta}_E(\tilde{\mathbf{y}}) := h_E^{1/2} \left(\sum_{k \leq d} \|\llbracket (\kappa(\cdot, \mathbf{a}_k) \nabla \tilde{y}_k - \kappa_{\text{ref}} \nabla \mathbf{w}_k) \cdot \nu_E \rrbracket_E\|_{L^2(E)}^2 \right)^{\frac{1}{2}}. \tag{7.28}$$

Note that $\tilde{\mathbf{y}}$ enters in these definitions via \mathbf{w}_k and \tilde{y}_k , and indirectly via $\tilde{\chi}_k$ and $\tilde{\varphi}_k$.

The second term in (7.25) can be treated as follows:

$$\begin{aligned}
 &\sum_{i=1}^m \sum_{k \leq d} c_{k_i}^{(i)} \int_{\Omega} \kappa_i \nabla \tilde{y}_k^* \cdot \nabla v_{k_1, \dots, d_i+1, \dots, k_m}^* \, dx \\
 &= \sum_{i=1}^m \sum_{k \leq d} c_{k_i}^{(i)} \langle A_i \tilde{y}_k^*, v_{k_1, \dots, d_i+1, \dots, k_m}^* \rangle_{Y^*, Y} \\
 &\leq \sum_{i=1}^m \sum_{k_j \leq d_j, j \neq i} \left\| \sum_{k_i \leq d_i} (c_{k_i}^{(i)} \cdot A_i \tilde{y}_k^*) \right\|_{A_{\text{ref}}^{-1}} \cdot \|v_{k_1, \dots, d_i+1, \dots, k_m}^*\|_{A_{\text{ref}}} \\
 &\leq \left(\sum_{i=1}^m \sum_{k_j \leq d_j, j \neq i} \left\| A_i \sum_{k_i \leq d_i} (c_{k_i}^{(i)} \tilde{y}_k^*) \right\|_{A_{\text{ref}}^{-1}}^2 \right)^{\frac{1}{2}} \left(\sum_{i=1}^m \sum_{k_j \leq d_j, j \neq i} \|v_{k_1, \dots, d_i+1, \dots, k_m}^*\|_{A_{\text{ref}}}^2 \right)^{\frac{1}{2}}.
 \end{aligned} \tag{7.29}$$

Example 7.11. For $\vartheta_k \equiv \beta_k$ we have $c_{k_i}^{(i)} = 0$ for all $k_i \leq d_i - 1$, and the estimate (7.29) becomes

$$\begin{aligned}
 &\sum_{i=1}^m \sum_{k \leq d} c_{k_i}^{(i)} \int_{\Omega} \kappa_i \nabla \tilde{y}_k^* \cdot \nabla v_{k_1, \dots, d_i+1, \dots, k_m}^* \, dx \leq \\
 &\left(\sum_{i=1}^m \sum_{k_j \leq d_j, j \neq i} (c_{d_i}^{(i)})^2 \|A_i \tilde{y}_{k_1, \dots, d_i, \dots, k_m}^*\|_{A_{\text{ref}}^{-1}}^2 \right)^{\frac{1}{2}} \left(\sum_{i=1}^m \sum_{k_j \leq d_j, j \neq i} \|v_{k_1, \dots, d_i+1, \dots, k_m}^*\|_{A_{\text{ref}}}^2 \right)^{\frac{1}{2}}.
 \end{aligned}$$

For general ϑ_k , we focus on the term $\sum_{k_i \leq d_i} (\mathbf{c}_{k_i}^{(i)} \tilde{y}_k^*) =: \tilde{y}_{k,i}^*$. It holds that

$$\sum_{k_i \leq d_i} (\mathbf{c}_{k_i}^{(i)} \tilde{y}_k^*) = \int_{\Xi_i} \xi_i \beta_{d_i+1}^{(i)}(\xi_i) \sum_{k_i \leq d_i} \vartheta_{k_i}^{(i)}(\xi_i) \tilde{y}_k^* \, d\mathbb{P}_i$$

due to $\vartheta_{d_i+1}^{(i)} \equiv \beta_{d_i+1}^{(i)}$, where the appearing integral is in the sense of Bochner. The term $\sum_{k_i \leq d_i} \vartheta_{k_i}^{(i)}(\xi_i) \tilde{y}_k^*$ is independent of the concrete choice of the basis $\{\vartheta_{k_i}^{(i)}\}_{k_i \in \mathbb{N}}$.

It still remains to compute

$$\|A_i \tilde{y}_{k,i}^*\|_{A_{\text{ref}}^{-1}}^2 = \langle A_{\text{ref}}^{-1} A_i \tilde{y}_{k,i}^*, A_i \tilde{y}_{k,i}^* \rangle_{Y, Y^*} = \langle \hat{y}_{k,i}^*, A_i \tilde{y}_{k,i}^* \rangle_{Y, Y^*} = \int_{\Omega} \kappa_i \nabla \tilde{y}_{k,i}^* \cdot \nabla \hat{y}_{k,i}^* \, dx$$

where we define $\hat{y}_{k,i}^* := A_{\text{ref}}^{-1} A_i \tilde{y}_{k,i}^*$, i. e.,

$$\int_{\Omega} \kappa_{\text{ref}} \nabla \hat{y}_{k,i}^* \cdot \nabla v + \chi_{\text{ref}} \hat{y}_{k,i}^* v \, dx = \int_{\Omega} \kappa_i \nabla \tilde{y}_{k,i}^* \cdot \nabla v \, dx$$

holds for all $v \in Y$. One could now compute a discrete version $\hat{y}_{k,i}^*$ of $\hat{y}_{k,i}^*$ and additionally estimate the error $\|\hat{y}_{k,i}^* - \hat{y}_{k,i}^*\|_{A_{\text{ref}}}$ with the techniques from Section 7.1. Alternatively, we can estimate

$$\begin{aligned} & \int_{\Omega} \kappa_i \nabla \tilde{y}_{k,i}^* \cdot \nabla \hat{y}_{k,i}^* \, dx \\ & \leq \left(\int_{\Omega} |\kappa_i| |\nabla \tilde{y}_{k,i}^* \cdot \nabla \tilde{y}_{k,i}^*| \, dx \right)^{\frac{1}{2}} \cdot \left(\int_{\Omega} |\kappa_i| |\nabla \hat{y}_{k,i}^* \cdot \nabla \hat{y}_{k,i}^*| \, dx \right)^{\frac{1}{2}} \\ & \leq \left(\int_{\Omega} |\kappa_i| |\nabla \tilde{y}_{k,i}^* \cdot \nabla \tilde{y}_{k,i}^*| \, dx \right)^{\frac{1}{2}} \cdot \left\| \frac{\kappa_i}{\kappa_{\text{ref}}} \right\|_{L^\infty(\Omega)}^{\frac{1}{2}} \cdot \left(\int_{\Omega} \kappa_{\text{ref}} \nabla \hat{y}_{k,i}^* \cdot \nabla \hat{y}_{k,i}^* \, dx \right)^{\frac{1}{2}} \\ & \leq \left(\int_{\Omega} |\kappa_i| |\nabla \tilde{y}_{k,i}^* \cdot \nabla \tilde{y}_{k,i}^*| \, dx \right)^{\frac{1}{2}} \cdot \left\| \frac{\kappa_i}{\kappa_{\text{ref}}} \right\|_{L^\infty(\Omega)}^{\frac{1}{2}} \cdot \left(\int_{\Omega} \kappa_{\text{ref}} \nabla \hat{y}_{k,i}^* \cdot \nabla \hat{y}_{k,i}^* + \chi_{\text{ref}} (\hat{y}_{k,i}^*)^2 \, dx \right)^{\frac{1}{2}} \\ & = \left(\int_{\Omega} |\kappa_i| |\nabla \tilde{y}_{k,i}^* \cdot \nabla \tilde{y}_{k,i}^*| \, dx \right)^{\frac{1}{2}} \cdot \left\| \frac{\kappa_i}{\kappa_{\text{ref}}} \right\|_{L^\infty(\Omega)}^{\frac{1}{2}} \cdot \left(\int_{\Omega} \kappa_i \nabla \tilde{y}_{k,i}^* \cdot \nabla \hat{y}_{k,i}^* \, dx \right)^{\frac{1}{2}}. \end{aligned}$$

It follows that

$$\|A_i \tilde{y}_{k,i}^*\|_{A_{\text{ref}}^{-1}}^2 = \int_{\Omega} \kappa_i \nabla \tilde{y}_{k,i}^* \cdot \nabla \hat{y}_{k,i}^* \, dx \leq \left\| \frac{\kappa_i}{\kappa_{\text{ref}}} \right\|_{L^\infty(\Omega)} \left(\int_{\Omega} |\kappa_i| |\nabla \tilde{y}_{k,i}^* \cdot \nabla \tilde{y}_{k,i}^*| \, dx \right). \quad (7.30)$$

Thus, we define

$$\zeta_i(\tilde{\mathbf{y}}) := \left\| \frac{\kappa_i}{\kappa_{\text{ref}}} \right\|_{L^\infty(\Omega)}^{1/2} \left(\sum_{k_j \leq d_j, j \neq i} \left(\int_{\Omega} |\kappa_i| \left(\sum_{k_i \leq d_i} \mathbf{c}_{k_i}^{(i)} \nabla \tilde{y}_k^* \right) \cdot \left(\sum_{k_i \leq d_i} \mathbf{c}_{k_i}^{(i)} \nabla \tilde{y}_k^* \right) \, dx \right) \right)^{\frac{1}{2}} \quad (7.31)$$

in order to formulate the following error estimation theorem.

Theorem 7.12. Let $\tilde{\mathbf{y}} \in Y \otimes \mathcal{P}_d(\Xi)$ be given such that $\mathbf{r} = \mathbf{A}\tilde{\mathbf{y}} + \mathbf{N}(\tilde{\mathbf{y}}) - \mathbf{b} \in L_{\mathbb{P}}^2(\Xi; Y^*)$. Suppose that Assumptions 7.5 and 7.9 are satisfied with the definitions (7.10).

Then,

$$\begin{aligned} \|\tilde{\mathbf{y}} - \mathbf{y}\|_{\hat{\mathbf{A}}_{\text{ref}}}^2 &\leq \Lambda^2 \left(c_{\mathcal{T}C_{A_{\text{ref}}}} \left(\sum_{T \in \mathcal{T}} \boldsymbol{\eta}_T(\tilde{\mathbf{y}})^2 + \sum_{E \in \mathcal{E}^0} \boldsymbol{\eta}_E(\tilde{\mathbf{y}})^2 \right)^{\frac{1}{2}} + \left(\sum_{i=1}^m \zeta_i(\tilde{\mathbf{y}})^2 \right)^{\frac{1}{2}} \right. \\ &\quad \left. + \left\| (\chi \cdot \tilde{\mathbf{y}} - \tilde{\boldsymbol{\chi}} + \varphi(\tilde{\mathbf{y}}) - \tilde{\boldsymbol{\varphi}}, \cdot) \right\|_{L_{\mathbb{P}}^2(\Xi; L^2(\Omega))} \right\|_{\hat{\mathbf{A}}_{\text{ref}}^{-1}}^2 + \Lambda^2 \|\mathbf{w}\|_{\hat{\mathbf{A}}_{\text{ref}}}^2, \end{aligned} \quad (7.32)$$

where \mathbf{w} is defined by (7.19), $\boldsymbol{\eta}_T$, $\boldsymbol{\eta}_E$ are defined by (7.27), (7.28), and ζ_i is defined in (7.31).

Proof. From (7.18) and with $\mathbf{w} = \hat{\mathbf{A}}_{\text{ref}}^{-1} \mathbf{r} \in \hat{\mathbf{Y}}$ and \mathbf{w} defined by (7.19) we conclude

$$\|\tilde{\mathbf{y}} - \mathbf{y}\|_{\hat{\mathbf{A}}_{\text{ref}}}^2 \leq \Lambda^2 (\|\mathbf{w} - \mathbf{w}\|_{\hat{\mathbf{A}}_{\text{ref}}}^2 + \|\mathbf{w}\|_{\hat{\mathbf{A}}_{\text{ref}}}^2) \quad (7.33)$$

due to Galerkin orthogonality.

For an arbitrary $\mathbf{v} \in \hat{\mathbf{Y}}$, $\mathbf{v}(x, \xi) = \sum_{k \in \mathbb{N}^m} v_k(x) \theta_k(\xi) = \sum_{k \in \mathbb{N}^m} v_k^*(x) \vartheta_k(\xi)$ with $v_k, v_k^* \in Y$ for all k , we have that

$$\|\mathbf{v}\|_{\hat{\mathbf{A}}_{\text{ref}}}^2 = \int_{\Xi} \langle A_{\text{ref}} \sum_{k \in \mathbb{N}^m} v_k(\cdot) \theta_k(\xi), \sum_{l \in \mathbb{N}^m} v_l(x) \theta_l(\xi) \rangle_{Y^*, Y} d\mathbb{P} = \sum_{k \in \mathbb{N}^m} \|v_k\|_{A_{\text{ref}}}^2$$

and $\|\mathbf{v}\|_{\hat{\mathbf{A}}_{\text{ref}}}^2 = \sum_{k \in \mathbb{N}^m} \|v_k^*\|_{A_{\text{ref}}}^2$ because the polynomials $\{\theta_k\}_k$, $\{\vartheta_k\}_k$ are orthonormal.

Combining (7.21), (7.25), (7.26), (7.29), (7.30) and the fact that

$$\langle \hat{\mathbf{A}}_{\text{ref}}(\mathbf{w} - \mathbf{w}), \mathbf{v} \rangle_{\hat{\mathbf{Y}}^*, \hat{\mathbf{Y}}} = (\mathbf{w} - \mathbf{w}, \mathbf{v})_{\hat{\mathbf{A}}_{\text{ref}}} \leq c \cdot \boldsymbol{\eta}(\mathbf{w}) \cdot \|\mathbf{v}\|_{\hat{\mathbf{A}}_{\text{ref}}} \quad \text{for all } \mathbf{v} \in \hat{\mathbf{Y}}$$

yields $\|\mathbf{w} - \mathbf{w}\|_{\hat{\mathbf{A}}_{\text{ref}}} \leq c \cdot \boldsymbol{\eta}(\mathbf{w})$, gives

$$\begin{aligned} \|\mathbf{w} - \mathbf{w}\|_{\hat{\mathbf{A}}_{\text{ref}}} &\leq c_{\mathcal{T}C_{A_{\text{ref}}}} \left(\sum_{T \in \mathcal{T}} \boldsymbol{\eta}_T(\tilde{\mathbf{y}})^2 + \sum_{E \in \mathcal{E}^0} \boldsymbol{\eta}_E(\tilde{\mathbf{y}})^2 \right)^{\frac{1}{2}} + \left(\sum_{i=1}^m \zeta_i(\tilde{\mathbf{y}})^2 \right)^{\frac{1}{2}} \\ &\quad + \left\| (\chi \cdot \tilde{\mathbf{y}} - \tilde{\boldsymbol{\chi}} + \varphi(\tilde{\mathbf{y}}) - \tilde{\boldsymbol{\varphi}}, \cdot) \right\|_{L_{\mathbb{P}}^2(\Xi; L^2(\Omega))} \Big\|_{\hat{\mathbf{A}}_{\text{ref}}^{-1}}, \end{aligned}$$

which yields the result together with (7.33). \square

Theorem 7.12 provides a nice split of the error into different contributions. The last term in (7.32) is exactly the algebraic error due to the inexact solution of the discrete system by, e. g., a low-rank tensor solver. The first term consists of

- a term related to the FE discretization error, which can itself be split into error contributions for each element,
- the error contribution due to the discretization of the stochastic space, which consists of terms for each stochastic parameter, and
- the interpolation error coming from the approximation of the nonlinear terms in the equation.

Remark 7.13. It remains to discuss the condition $\mathbf{r} \in L_{\mathbb{P}}^2(\Xi, Y^*)$ for the state and adjoint equation from Section 3.2, cf. Remark 7.4, provided $Q(\xi) \equiv \iota : H_0^1(\Omega) \hookrightarrow L^2(\Omega)$.

If $\tilde{\mathbf{y}} \in L_{\mathbb{P}}^{r_f}(\Xi; Y)$ ($r_f \geq p$) is an inexact solution of the state equation, the residual \mathbf{r} belongs to $L_{\mathbb{P}}^{r_f/(p-1)}(\Xi; Y^*)$ due to the growth condition (3.7) on the nonlinearity φ . Thus, we need $r_f \geq 2p-2$ to have $\mathbf{r} \in L_{\mathbb{P}}^2(\Xi, Y^*)$. Under the regularity assumption in the ‘‘Example’’ column of Table 3.1, i. e., $r_Q = \infty$ and $r_{\hat{q}} \geq r_f$, it is realistic that $\tilde{\mathbf{z}} \in L_{\mathbb{P}}^{r_f}(\Xi; Y)$ holds, because the exact adjoint state \mathbf{z} has exactly this regularity. The residual of the adjoint equation then also belongs to $L_{\mathbb{P}}^{r_f/(p-1)}(\Xi; Y^*)$ which is a subset of $L_{\mathbb{P}}^2(\Xi; Y^*)$ by $r_f \geq 2p-2$.

The error estimate given in Theorem 7.12 uses a deterministic reference operator as the one discussed in [22, 23]. The results in these papers rely on a saturation assumption, which we do not make here, but derive an error estimate similar to the one presented in [38] and used in [39, 40]. All mentioned papers present error estimates for linear equations whereas we consider a class of semilinear equations with a monotone nonlinearity, but have used a linear reference equation to derive the error estimate. In the linear case, it makes sense to discuss the efficiency of the estimator using a lower bound as indicated in Remark 7.2. Furthermore, the convergence of the adaptive solution technique should be investigated. Since it is already discussed in the linear case for a very similar estimator in [39], we skip this topic at this point.

For the case $\hat{\varphi} \equiv 0$ and $\chi \equiv 0$, i. e., we have a linear equation and no interpolation error, we briefly compare our Theorem 7.12 to [38, Thm. 6.2]: There, an additional discrete error term (analog of $\|\mathbf{w}\|_{\hat{\mathbf{A}}_{\text{ref}}}$) appears as an additional summand in the ‘‘large bracket’’ in (7.32). We do not get this term because we make use of the *exact* solution of the reference equation (7.19), which is realistic when we work with low-rank tensors, where the discrete rank-1-operator $\mathbf{A}_{\text{ref}}^{-1} = \mathbf{A}_{\text{ref}}^{-1} \otimes \mathbf{I} \cdots \otimes \mathbf{I}$ can be applied by an *i*-mode matrix product so that it can be used as a typically very good preconditioner, see below.

7.4. Realization with Low-Rank Tensors

In order to implement an adaptive solver for (7.17) based on the estimate (7.32) with low-rank tensors, we have to apply an iterative low-rank tensor solver such as truncated PCG [75], ALS [62], or AMEn [36] to the discretized system

$$\mathbf{A}\mathbf{y} + \mathbf{N}(\mathbf{y}) - \mathbf{b} = 0, \quad (7.34)$$

cf. (6.16) and (6.21). In addition, we have to implement the evaluation of the error indicators $\boldsymbol{\eta}_T(\tilde{\mathbf{y}})$, $\boldsymbol{\eta}_E(\tilde{\mathbf{y}})$, and $\boldsymbol{\zeta}_i(\tilde{\mathbf{y}})$ from the low-rank solution in an efficient manner. Based on the value of these error indicators, the discrete subspace with the largest error contribution is refined, provided the algebraic error $\|\mathbf{w}\|_{\hat{\mathbf{A}}_{\text{ref}}}$ does not dominate. In the latter case, additional iterations of the low-rank tensor solver have to be performed.

Solution of the Discrete System

Let us set $\mathbf{N} \equiv 0$ for a moment. As pointed out in Subsection 2.1.3, iterative solvers working with the low-rank tensor representation such as ALS or AMEn aim for solving the linear

system $\mathbf{A}\mathbf{y} = \mathbf{b}$ by minimizing the squared Frobenius norm $\|\mathbf{A}\mathbf{y} - \mathbf{b}\|_{\mathbb{F}}^2$ of the residual. This problem can be badly conditioned and as indicated by (7.32), we are in fact interested in minimizing

$$\|\mathbf{w}\|_{\hat{\mathbf{A}}_{\text{ref}}}^2 = \langle \mathbf{r}, \mathbf{w} \rangle_{\hat{\mathbf{Y}}^*, \hat{\mathbf{Y}}} = \langle \mathbf{r}, \hat{\mathbf{A}}_{\text{ref}}^{-1} \mathbf{r} \rangle =: \|\mathbf{r}\|_{\hat{\mathbf{A}}_{\text{ref}}^{-1}}^2,$$

where (7.19) translates to $\hat{\mathbf{A}}_{\text{ref}} \mathbf{w} = \mathbf{r} = \mathbf{A}\mathbf{y} - \mathbf{b}$ in the discrete tensor space. Let the symmetric, positive definite operator $\hat{\mathbf{A}}_{\text{ref}}$ now be decomposed as $\hat{\mathbf{A}}_{\text{ref}} = \hat{\mathbf{R}}_{\text{ref}}^* \hat{\mathbf{R}}_{\text{ref}}$, where $\hat{\mathbf{R}}_{\text{ref}}^*$ is the adjoint operator of $\hat{\mathbf{R}}_{\text{ref}} : \mathbb{R}^{n_0 \times \dots \times n_m} \rightarrow \mathbb{R}^{n_0 \times \dots \times n_m}$ w. r. t. the Frobenius inner product. Then,

$$\begin{aligned} \arg \min_{\mathbf{y}} \|\mathbf{A}\mathbf{y} - \mathbf{b}\|_{\hat{\mathbf{A}}_{\text{ref}}^{-1}}^2 &= \arg \min_{\mathbf{y}} \langle \mathbf{A}\mathbf{y} - \mathbf{b}, \hat{\mathbf{A}}_{\text{ref}}^{-1} (\mathbf{A}\mathbf{y} - \mathbf{b}) \rangle \\ &= \arg \min_{\mathbf{y}} \|\hat{\mathbf{R}}_{\text{ref}}^{-*} (\mathbf{A}\mathbf{y} - \mathbf{b})\|_{\mathbb{F}}^2 \\ &= \hat{\mathbf{R}}_{\text{ref}}^{-1} \arg \min_{\hat{\mathbf{y}}} \|\hat{\mathbf{R}}_{\text{ref}}^* \hat{\mathbf{A}}_{\text{ref}}^{-1} \hat{\mathbf{y}} - \hat{\mathbf{R}}_{\text{ref}}^* \mathbf{b}\|_{\mathbb{F}}^2. \end{aligned} \quad (7.35)$$

This means that the preconditioned, symmetric system $\hat{\mathbf{R}}_{\text{ref}}^{-*} \hat{\mathbf{A}}_{\text{ref}}^{-1} \hat{\mathbf{y}} = \hat{\mathbf{R}}_{\text{ref}}^{-*} \mathbf{b}$ can be solved by a standard implementation of, e. g., AMEn and its (approximate) solution $\hat{\mathbf{y}}$ can be transformed to the (approximate) solution $\mathbf{y} = \hat{\mathbf{R}}_{\text{ref}}^{-1} \hat{\mathbf{y}}$ of the original system $\mathbf{A}\mathbf{y} = \mathbf{b}$. In particular, we can use the tensor product form $\hat{\mathbf{A}}_{\text{ref}} = \hat{\mathbf{A}}_{\text{ref}} \otimes \mathbf{I} \otimes \dots \otimes \mathbf{I}$ with the positive definite, quadratic matrix $\hat{\mathbf{A}}_{\text{ref}}$, and compute, e. g., a sparse Cholesky decomposition $\hat{\mathbf{A}}_{\text{ref}} = \hat{\mathbf{R}}_{\text{ref}}^{\top} \hat{\mathbf{R}}_{\text{ref}}$, where $\hat{\mathbf{R}}_{\text{ref}}$ is the product of a triangular and a permutation matrix. Then, with $\hat{\mathbf{R}}_{\text{ref}} := \hat{\mathbf{R}}_{\text{ref}} \otimes \mathbf{I} \otimes \dots \otimes \mathbf{I}$ we obtain a decomposition of $\hat{\mathbf{A}}_{\text{ref}}$ as desired because $\hat{\mathbf{R}}_{\text{ref}}^* \hat{\mathbf{R}}_{\text{ref}} = (\hat{\mathbf{R}}_{\text{ref}}^{\top} \hat{\mathbf{R}}_{\text{ref}}) \otimes \mathbf{I} \otimes \dots \otimes \mathbf{I}$. It is then algorithmically relevant that $\hat{\mathbf{R}}_{\text{ref}}^{-1} = \hat{\mathbf{R}}_{\text{ref}}^{-1} \otimes \mathbf{I} \otimes \dots \otimes \mathbf{I}$, i. e., the inverse can be applied very efficiently to low-rank tensors. For our concrete implementation, we have extended the AMEn algorithm to be able to handle preconditioners of canonical rank 1 given in decomposed form. Since, e. g., the inverse Cholesky factor $\hat{\mathbf{R}}_{\text{ref}}^{-1}$ should never be formed explicitly, but can be applied to vectors efficiently by backward substitution and permutation, we have extended the {d, R}-format [35], which is essentially a collection of sparse matrices, each of which acts on one mode of the tensor, to function handles. We want to remark that it would be even more desirable to have a preconditioned version of AMEn which does only require the application of the preconditioner $\hat{\mathbf{A}}_{\text{ref}}^{-1}$ itself and not a decomposition of it (similar to the PCG method). But investigating and implementing such an algorithm is out of the scope of this thesis. Additionally, the operator $\hat{\mathbf{A}}_{\text{ref}}^{-1}$ is needed more frequently in the algorithm, e. g., to evaluate the norm of the residual of the nonlinear equation. Since $\hat{\mathbf{A}}_{\text{ref}}$ is symmetric, positive definite, and sparse, the application of its inverse is computed by the mentioned sparse Cholesky decomposition in MATLAB so that this can be precomputed and used in AMEn.

A similar idea can be used to truncate a given tensor \mathbf{y} to smaller rank w. r. t. the \mathbf{M} -norm, where $\mathbf{M} : \mathbb{R}^{n_0 \times \dots \times n_m} \rightarrow \mathbb{R}^{n_0 \times \dots \times n_m}$ induces some inner product and is decomposed as $\mathbf{M} = \tilde{\mathbf{R}}^* \tilde{\mathbf{R}}$ w. r. t. the Frobenius inner product. In fact, SVD-based truncation would quasi-minimize the function $\|\mathbf{z} - \mathbf{y}\|_{\mathbb{F}}^2$ w. r. t. \mathbf{z} , see Subsection 2.1.2. Instead, we can use the transformation (7.35) and approximate the tensor $\tilde{\mathbf{R}}\mathbf{y}$ by standard truncation w. r. t. the Frobenius norm to obtain the tensor $\hat{\mathbf{z}}$. Then, $\tilde{\mathbf{R}}^{-1} \hat{\mathbf{z}}$ is the respective approximation of \mathbf{y} w. r. t. the \mathbf{M} -norm.

If $\mathbf{N} \neq 0$ is nonlinear, we solve the system (7.34) by Newton's method, where each Newton step is computed approximately by the preconditioned AMEn solver. Especially in the last iterations of the trust-region algorithm, the state and the adjoint state do not change very much. Then it is often enough to take one inexact Newton step, which is computed by one AMEn sweep, to decrease the residual of the nonlinear equation sufficiently.

Evaluation of the Error Indicators

The discretization error indicators can be implemented as follows: If, e. g.,

$$\tilde{\mathbf{y}}(x, \xi) = \sum_{k \leq d} \tilde{\mathbf{y}}_k(x) \theta_k(\xi)$$

is represented by the tensor $\tilde{\mathbf{y}}$ analogously to (6.8),

$$\tilde{\mathbf{y}}_k(x) = \sum_{k_0=1}^{d_0} \tilde{\mathbf{y}}(k_0, k) \phi_{k_0}(x),$$

holds, i. e., the function $\tilde{\mathbf{y}}_k$ is represented by the vector $\tilde{\mathbf{y}}(\cdot, k)$.

Knowing that, we establish the evaluation of the triangle error contribution (7.27). Let $x^1, x^2, x^3 \in \Omega$ be the three vertices of the triangle T and let the function $\tilde{\mathbf{f}}$ be affine on T . Then,

$$\begin{aligned} \|\tilde{\mathbf{f}}\|_{L^2(T)}^2 &= \int_T \tilde{\mathbf{f}}(x)^2 dx = \frac{a_T}{6} (\tilde{\mathbf{f}}(x^1)^2 + \tilde{\mathbf{f}}(x^2)^2 + \tilde{\mathbf{f}}(x^3)^2 + \tilde{\mathbf{f}}(x^1)\tilde{\mathbf{f}}(x^2) + \tilde{\mathbf{f}}(x^1)\tilde{\mathbf{f}}(x^3) + \tilde{\mathbf{f}}(x^2)\tilde{\mathbf{f}}(x^3)) \\ &= \frac{a_T}{6} \sum_{j=1}^3 \sum_{l=j}^3 \tilde{\mathbf{f}}(x^j)\tilde{\mathbf{f}}(x^l), \end{aligned}$$

where a_T is the area of T . We approximate each $\tilde{\chi}_k + \tilde{\varphi}_k - f_k - \chi_{\text{ref}w_k} \approx \tilde{\mathbf{f}}_k$ by interpolation, where $\tilde{\mathbf{f}}_k$ are linear finite element functions for all k . These functions are represented altogether by a single tensor $\tilde{\mathbf{f}} \in \mathbb{R}^{\tilde{d} \times d_1 \times \dots \times d_m}$ of values at the FE nodes. Note that $\tilde{d} > d_0$ is the number of all FE nodes whereas d_0 counts only the interior nodes. Let $k_0^{(1)}, k_0^{(2)}, k_0^{(3)} \in [\tilde{d}]$ be the indices corresponding to the vertices x^1, x^2, x^3 , respectively, so that $\tilde{\mathbf{f}}(k_0^{(j)}, k) = \tilde{\mathbf{f}}_k(x^j)$ for all j and all k . We obtain

$$\begin{aligned} \eta_T(\tilde{\mathbf{y}})^2 &\approx h_T^2 \sum_{k \leq d} \|\tilde{\mathbf{f}}_k\|_{L^2(T)}^2 = h_T^2 \frac{a_T}{6} \sum_{k \leq d} \left(\sum_{j=1}^3 \sum_{l=j}^3 \tilde{\mathbf{f}}(k_0^{(j)}, k) \tilde{\mathbf{f}}(k_0^{(l)}, k) \right) \\ &= h_T^2 \frac{a_T}{6} \sum_{j=1}^3 \sum_{l=j}^3 \langle \mathbb{1}, \tilde{\mathbf{f}}(k_0^{(j)}, \cdot) \odot \tilde{\mathbf{f}}(k_0^{(l)}, \cdot) \rangle. \end{aligned}$$

The advantage of this formulation is that it can be vectorized using the tensors of evaluations at all first, second, and third triangle vertices, respectively, to evaluate the error indicator for all triangles T simultaneously. Then, only componentwise multiplication and contraction with the rank-1-tensor $\mathbb{1}$ are needed.

To compute the edge error contribution (7.28), one first computes the values of the partial derivatives $\partial_{x_1}\tilde{y}_k$ and $\partial_{x_2}\tilde{y}_k$ on the elements. These values are contained in the tensors $\mathbf{G}_1 \circ_1 \tilde{\mathbf{y}}$ and $\mathbf{G}_2 \circ_1 \tilde{\mathbf{y}}$, where $\mathbf{G}_1, \mathbf{G}_2 \in \mathbb{R}^{|\mathcal{T}| \times d_0}$ are sparse matrices mapping a vector of function values on the FE nodes to a vector of derivative values on the triangles numbered from 1 to $|\mathcal{T}|$. Analogously, we obtain the tensors $\mathbf{G}_1 \circ_1 \mathbf{w}$ and $\mathbf{G}_2 \circ_1 \mathbf{w}$ containing the values of $\partial_{x_1} w_k$ and $\partial_{x_2} w_k$. Furthermore, we create the tensors $\boldsymbol{\kappa}_{\text{ref}} = \boldsymbol{\kappa}_{\text{ref}} \otimes \mathbb{1} \in \mathbb{R}^{|\mathcal{T}| \times d_1 \times \dots \times d_m}$, where $\boldsymbol{\kappa}_{\text{ref}} \in \mathbb{R}^{|\mathcal{T}|}$ contains the values of κ_{ref} on the triangles, and $\boldsymbol{\kappa} \in \mathbb{R}^{|\mathcal{T}| \times d_1 \times \dots \times d_m}$ such that $\boldsymbol{\kappa}(\cdot, k)$ contains the respective values of $\kappa(\cdot, \mathbf{a}_k)$. Let now $\mathbf{v}^{j,1} \in \mathbb{R}^{|\mathcal{T}|}$ contain the first components of the outwards-pointing, normal vectors corresponding to the j -th edge of the triangles T , and let $\mathbf{v}^{j,2} \in \mathbb{R}^{|\mathcal{T}|}$ contain the second components for $j \in \{1, 2, 3\}$. This means that if T is the ℓ -th triangle, the outer normal vector $\nu_{T,j}$ corresponding to its j -th edge is given by $\nu_{T,j} = (\mathbf{v}_\ell^{j,1} \quad \mathbf{v}_\ell^{j,2})^\top$. Then, the values of

$$\kappa(\cdot, \mathbf{a}_k) \nabla \tilde{y}_k \cdot \nu_{T,j}$$

for all T are given by the single tensor

$$\boldsymbol{\kappa} \odot (\text{diag}(\mathbf{v}^{j,1}) \circ_1 \mathbf{G}_1 \circ_1 \tilde{\mathbf{y}} + \text{diag}(\mathbf{v}^{j,2}) \circ_1 \mathbf{G}_2 \circ_1 \tilde{\mathbf{y}}) = \boldsymbol{\kappa} \odot ((\text{diag}(\mathbf{v}^{j,1})\mathbf{G}_1 + \text{diag}(\mathbf{v}^{j,2})\mathbf{G}_2) \circ_1 \tilde{\mathbf{y}}).$$

With an analogous consideration we get that the values of

$$(\kappa(\cdot, \mathbf{a}_k) \nabla \tilde{y}_k - \kappa_{\text{ref}} \nabla w_k) \cdot \nu_{T,j}$$

are contained in the tensor

$$\tilde{\mathbf{g}}_j := \boldsymbol{\kappa} \odot ((\text{diag}(\mathbf{v}^{j,1})\mathbf{G}_1 + \text{diag}(\mathbf{v}^{j,2})\mathbf{G}_2) \circ_1 \tilde{\mathbf{y}}) - \boldsymbol{\kappa}_{\text{ref}} \odot ((\text{diag}(\mathbf{v}^{j,1})\mathbf{G}_1 + \text{diag}(\mathbf{v}^{j,2})\mathbf{G}_2) \circ_1 \mathbf{w}),$$

which can be computed using standard low-rank tensor arithmetics, namely i -mode matrix products as well as componentwise multiplication and subtraction. In a next step, we number the interior edges $E \in \mathcal{E}_0$ from 1 to $|\mathcal{E}_0|$ and create sparse matrices $\mathbf{H}_j \in \{0, 1\}^{|\mathcal{E}_0| \times |\mathcal{T}|}$ ($j \in \{1, 2, 3\}$) mapping vectors of values on the triangles to vectors of values on the respective j -th triangle edges in the correct order, provided they are interior edges. Then the tensor of jumps over the interior edges is given by

$$\tilde{\mathbf{h}} := \sum_{j=1}^3 \mathbf{H}_j \circ_1 \tilde{\mathbf{g}}_j \in \mathbb{R}^{|\mathcal{E}_0| \times d_1 \times \dots \times d_m}. \quad (7.36)$$

It is correct to sum here because the outer unit normals of two neighboring triangle point in the opposite direction. Using the presented definitions, the edge error indicator is given by

$$\boldsymbol{\eta}_E(\tilde{\mathbf{y}})^2 = h_E \sum_{k \leq d} \|[(\kappa(\cdot, \mathbf{a}_k) \nabla \tilde{y}_k - \kappa_{\text{ref}} \nabla w_k) \cdot \nu_E]_E\|_{L^2(E)}^2 = h_E^2 \langle \mathbb{1}, \tilde{\mathbf{h}}(l, \cdot)^2 \rangle,$$

where $l \in \{1, \dots, |\mathcal{E}_0|\}$ is the number of the edge E . As before, this procedure can be vectorized to compute the error contributions of all interior edges simultaneously. Finally, we assign half of the error to each of the two neighboring triangles. Based on that, we mark all

triangles with the largest error contributions which constitute a certain amount $\vartheta_\eta \in (0, 1)$ of the total error, see [38, Sec. 7.1], a so-called Dörfler strategy [37]. These triangles are refined regularly, i. e., divided into four triangles of the same shape. To avoid hanging nodes, additional triangles have to be divided into two new ones possibly.

To evaluate $\zeta_i(\tilde{\mathbf{y}})$ (see (7.31)) based on the tensor $\tilde{\mathbf{y}}$ representing the function $\tilde{\mathbf{y}}$, one first computes the values $\mathbf{c}_{l_i}^{(i)} = \int_{\Xi_i} \xi_i \vartheta_{l_i}^{(i)}(\xi_i) \vartheta_{d_i+1}^{(i)}(\xi_i) d\mathbb{P}_i$ by a quadrature formula of high enough order, e. g., by Gaussian quadrature with $d_i + 1$ nodes if $\vartheta_{l_i}^{(i)}$ has degree $d_i - 1$ and $\vartheta_{d_i+1}^{(i)}$ has degree d_i , or analytically and writes them as one vector $\mathbf{c}^{(i)} \in \mathbb{R}^{d_i}$. Furthermore, the matrix $\bar{\mathbf{A}}_i \in \mathbb{R}^{d_0 \times d_0}$ given by

$$(\bar{\mathbf{A}}_i)_{k_0 l_0} := (|\kappa_i| \nabla \phi_{l_0}, \nabla \phi_{k_0})_{L^2(\Omega)^n},$$

cf. Example 6.5, is assembled. Then we have

$$\begin{aligned} & \int_{\Omega} |\kappa_i| \left(\sum_{k_i \leq d_i} \mathbf{c}_{k_i}^{(i)} \nabla \tilde{\mathbf{y}}_{k_i}^* \right) \cdot \left(\sum_{k_i \leq d_i} \mathbf{c}_{k_i}^{(i)} \nabla \tilde{\mathbf{y}}_{k_i}^* \right) dx = \\ & \left((\mathbf{c}^{(i)\top} \circ_{i+1} \tilde{\mathbf{y}})(\bullet, k_1, \dots, k_{i-1}, k_{i+1}, \dots, k_m) \right)^\top \bar{\mathbf{A}}_i \left((\mathbf{c}^{(i)\top} \circ_{i+1} \tilde{\mathbf{y}})(\bullet, k_1, \dots, k_{i-1}, k_{i+1}, \dots, k_m) \right). \end{aligned}$$

Summing these values over all $k_j \in \{1, \dots, d_j\}$ ($j \neq i$) can be implemented as an inner product of tensors giving

$$\zeta_i(\tilde{\mathbf{y}})^2 = \left\| \frac{\kappa_i}{\kappa_{\text{ref}}} \right\|_{L^\infty(\Omega)} \langle \mathbf{c}^{(i)\top} \circ_{i+1} \tilde{\mathbf{y}}, \bar{\mathbf{A}}_i \circ_1 \mathbf{c}^{(i)\top} \circ_{i+1} \tilde{\mathbf{y}} \rangle.$$

If this is the largest contribution to the total error (7.32), we increase the respective polynomial degree $d_i - 1$ by 1.

Interpolation to a Finer Subspace

Lifting the current solution tensor $\tilde{\mathbf{y}}$ representing the function $\tilde{\mathbf{y}} \in \mathbf{Y}$ to a tensor representing the same function as an element of a finer subspace is easy: For the refinement of the linear FE space \mathbf{Y} , a sparse interpolation matrix can be constructed. Its upper block is the identity matrix because the coefficients belonging to the retained FE nodes are not changed. The coefficients of new nodes inserted on an edge are computed as a convex combination of the coefficients of the edge endpoints. Lifting the tensor $\tilde{\mathbf{y}}$ to a finer FE space then consists of applying the sparse interpolation matrix to the first tensor mode. The interpolation matrices of the polynomial spaces $\mathcal{P}_{d_i-1}(\Xi_i)$ with the weighted Lagrangian bases can be constructed by, e. g., evaluating the Lagrange polynomials of the coarse space at the Gaussian quadrature nodes of the finer space. Analogously, these dense but typically small interpolation matrices can be applied to the $(i + 1)$ -st mode of $\tilde{\mathbf{y}}$.

8. Implementation and Numerical Results

The inexact trust-region method (Algorithm 1) is implemented in MATLAB to solve different instances of the model problem presented in Chapter 3. In order to solve the arising PDEs with uncertain inputs, the adaptive solution technique described in Chapter 7 is used. The objective function evaluation and gradient error are estimated as derived in Chapter 5.

For this purpose, we consider the following concrete setup for the model problem (3.3) with state equation (3.8) and objective function (3.4):

- We choose uniformly distributed parameters, meaning that $\Xi_i = (-1, 1)$ and $\mathbb{P}_i = \frac{1}{2}\lambda$ for all $i \in [m]$, where λ is the Lebesgue measure on $(-1, 1)$. Assumption 6.3 is satisfied by this choice.
- The domain $\Omega := (-1, 1)^2 \setminus (-1, 0]^2 \subset \mathbb{R}^2$ is the polygonal (cf. Assumptions 6.1 and 7.5) L-shaped domain, frequently used to test adaptive FE codes. It is divided into $m = 6$ subdomains $\Omega_1, \dots, \Omega_6$ as depicted in Figure 8.1. The initial FE mesh, which contains 113 nodes, and thus all refined meshes respect this partition of the domain, meaning that each element $T \in \mathcal{T}$ is contained in exactly one subdomain Ω_i , see also Figure 8.1.
- The coefficient function is chosen to be $\kappa(x, \xi) = \kappa_0(x)(1 + \sum_{i=1}^m \xi_i \eta_i(x))$ (cf. Example 3.17) with $\kappa_0 \equiv 1$ and $\eta_i(x) = \sigma_i 1_{\Omega_i(x)}$, where the vector $\sigma \in [0, 1]^m$ describing the amounts of influence of the uncertain parameters is chosen differently for different tests. The reference coefficient is $\kappa_{\text{ref}} \equiv 1$. Thus, Assumption 7.5 and the first part of Assumption 7.9 with $\kappa_i = \kappa_0 \eta_i$ hold true. Furthermore, $\|\cdot\|_{A_{\text{ref}}} \equiv \|\cdot\|_{H_0^1(\Omega)}$ and the constant $c_{A_{\text{ref}}}$ introduced in (7.13) comes from the Poincaré inequality $\|\cdot\|_{H^1(\Omega)} \leq c_{A_{\text{ref}}} \|\cdot\|_{H_0^1(\Omega)}$.
- We use $U = L^2(\Omega)$, i. e., $\Omega_u = \Omega$, $D \equiv I : L^2(\Omega) \rightarrow L^2(\Omega)$, giving that the second part of Assumption 6.1 is satisfied, $f(\cdot) \equiv 0$, and $\varphi(t) := t^3$. The reason for the latter choice is that this nonlinearity and its derivative can be evaluated simply with low-rank tensors because only componentwise multiplication is needed.

Overall, Assumption 3.3 is fulfilled with $\underline{\kappa} = 1 - \max_{i \in [m]} \sigma_i$, $\bar{\kappa} = 1 + \max_{i \in [m]} \sigma_i$, $a''_{\varphi''} = 0$, $c''_{\varphi''} = 6$, $p = 4$, $r_f = \infty$. Additionally, (7.3) holds with $\lambda = \underline{\kappa}$ and $\Lambda = \bar{\kappa}$.

- The observation space is $H = L^2(\Omega)$ and $Q(\cdot) \equiv \iota : H_0^1(\Omega) \hookrightarrow L^2(\Omega)$ is constant.
- The desired state $\hat{q}(\cdot) \equiv \hat{q} \in L^2(\Omega)$ is constant and chosen differently for different tests.

Hence, we have the integrability exponents $r_Q = \infty$ and $r_{\hat{q}} = \infty$ and thus $r_{Qy} = \infty$, $\hat{r} = \infty$, and $r_z = \infty$ (see Table 3.1). Assumption 5.6 holds true and the error estimation scheme

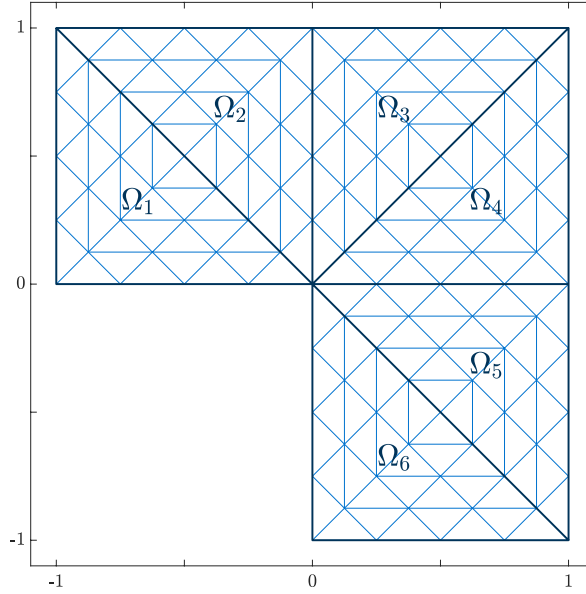


Figure 8.1.: The used physical domain Ω , its partition into subdomains Ω_i , and the initial finite element mesh.

given in Figure 5.2 can be applied. Assumption 6.1 is satisfied for arbitrary degree vectors $d \in \mathbb{N}^m$. Additionally, the right-hand sides of the state and adjoint equation are polynomials of coordinate degree at most $d - 1$ if the current approximate state $\tilde{\mathbf{y}}$ is such a polynomial, see Assumption 7.9 and also Remark 7.4.

- We take $\gamma = 10^{-3}$ and the set of admissible controls shall be $U_{\text{ad}} := \{u \in L^2(\Omega) : u(x) \leq 14 \text{ for a. e. } x \in \Omega\}$, where the upper bound is chosen such that it becomes active, but the optimal state is in the order of magnitude of the desired state. Therefore, the prerequisites of Theorem 3.19 as well as Assumption 4.2 are fulfilled and we can apply the exact and the inexact projection formula given in Section 6.4. The constant upper bound makes it possible to evaluate the error due to the “node-wise” projection for any FE mesh (Remark 5.4). In all tests, we initialize the algorithm with $u^0 = 0$.

The constants required for error estimation are set as follows: In Theorem 5.8 and in Lemma 5.10, we choose the Sobolev constant $c_{\tilde{r}} = c_p = c_4 = 0.5$ because we estimated $c_2 \approx 0.3$ numerically. The norm $\|\mathbf{y}\|_{L_{\mathbb{F}}^{\infty}(\Xi; Y)}$ in (5.20) is estimated by (3.9) with the Poincaré constant set to $C_{\Omega} = 1$. Additionally, we assume that $\|\tilde{\mathbf{y}}\|_{L_{\mathbb{F}}^{\infty}(\Xi; Y)}$ is approximately of the same size as $\|\mathbf{y}\|_{L_{\mathbb{F}}^{\infty}(\Xi; Y)}$ and neglect this term in (5.20). For a more rigorous implementation it could be estimated by a method for tensor maximum estimation as described in Subsection 2.1.3, but we skip this step for efficiency reasons. We also ignore the interpolation error in (7.12) and set $c_{\mathcal{T}} c_{A_{\text{ref}}} = 10^{-3}$ in this estimate to have all error contributions in the same order of magnitude. Even if the real, unknown constants are underestimated by the described choices, the algorithm still works because we need to know the error only up to a fixed, but possibly unknown multiplicative constant. In fact, an unrealistic choice of constants can be compensated by the choice of the error functions in the trust-region algorithm.

8.1. Implementation Details

The implementation is done in MATLAB R2017b because the use of various toolboxes facilitates it. The Partial Differential Equation Toolbox in MATLAB is used for basic FE tasks, such as mesh generation and refinement, or assembling matrices. Some finite element related code, such as the FE error estimator or the exact L^2 -projection onto a box, is not available within this toolbox and has to be written additionally. Since the FE error estimator has to be applied to low-rank tensors, we need to create suitable sparse matrices, see (7.36) for example. The exact projection with the suitable grid refinement has to be vectorized in order to have an efficient MATLAB code. This is not a trivial task because many decisions such as which triangles have to be divided into how many new ones and how these new triangles are chosen have to be encoded as vectors. We use the `htucker` toolbox 1.2 [76] for computations with hierarchical Tucker tensors and choose a “TT-like” dimension tree, the `TT-Toolbox` 2.2 [89] for Tensor Train tensors and an own implementation for the conversion between the two formats, see Subsection 2.1.3. This has the advantage that we can benefit from algorithms in both toolboxes. AMEn [36] including operators in the $\{d, R\}$ -format from the `tamen`-package [35] is used for the rank-adaptive solution of linear systems. Componentwise multiplication operators and a preconditioner are implemented as described in Subsection 2.1.3. These extensions of AMEn are needed due to the nonlinear term in the PDE and for minimizing the $\hat{\mathbf{A}}_{\text{ref}}^{-1}$ -norm of the residual, see Theorem 7.12.

Algorithm 1 is implemented such that it can work with “Hilbert space objects”, i. e., MATLAB objects for which the usual vector space operations are defined and an inner product is given. Additionally, we pay attention to the fact that the algorithm shall be used for optimal control applications. Therefore, the objective function and gradient evaluations accept initializations for the state and the adjoint state, a refined version of which is returned and used for further computations. All relevant functions, such as the control, the state, and the adjoint state are kept on matched grids to facilitate the computation. We keep all grid refinements stemming from the gradient computation and the projected linesearch for further iterations. In the concrete application it is found that sometimes it is necessary to evaluate the objective function up to high accuracy so that the state needs to be computed on a very fine FE mesh and/or with a high polynomial degree. Keeping this fine FE and polynomial spaces for further iterations would make the algorithm slow so that it turns out that it is better to discard them. We only reuse grids and polynomial degrees in the stochastic case for the objective function evaluation at the possible next iterate. To find a possibly better step than the generalized Cauchy step found by the projected Armijo linesearch with suitably refined projection, our implementation of Algorithm 1 requires an “advanced solver” for the trust-region subproblem. Here we apply a semismooth Newton method as outlined in Section 6.4 with the generalized Cauchy point as initial iterate. We perform at most 5 semismooth Newton steps and stop the iteration earlier if the discrete L^2 -norm of the residual (6.29) is smaller than 10^{-6} . For computing the semismooth Newton direction, at most 20 CG iterations are performed until the relative CG-residual is smaller than 0.1. If CG performs 20 iterations without reaching this tolerance and the relative CG-residual is in fact larger than 0.5, the algorithm stops and returns the iterate with the smallest semismooth Newton residual. Otherwise it checks whether the step computed by CG is capable of reducing the semismooth Newton residual. If this is not the case, the semismooth Newton iteration stops.

Note that for the performance of the overall algorithm, the semismooth Newton solver is not a bottleneck because it works on a fixed, discrete subspace of U . Compared to it, grid refinement and the solution of the stochastic PDE with low-rank tensors is costly. That is the reason why we use more than one semismooth Newton iteration. With this approach we can expect to obtain a larger model decrease pred_k . This allows then to evaluate the objective function less exactly, see Lemma 4.25, so that typically more computing time can be saved than by taking less semismooth Newton iterations.

To track the mesh refinement, we write MATLAB classes for FE functions carrying, e. g., the mesh, the respective assembled matrices, the coefficients, and functions for grid refinement. Vector space operations are overloaded such that the trust-region algorithm can work with these objects. For the state and the adjoint state we need an additional MATLAB class of functions of the form (6.8) with one FE space mode and $m \in \mathbb{N}$ polynomial parameter modes. The coefficients are represented by an HT tensor and we implement functions for FE and polynomial refinement, pointwise evaluations and computing the expectation etc.

To make sure that all required error bounds in Algorithm 1 can be satisfied, we use the error functions proposed in Section 4.3. The parameter settings for the trust-region algorithm as well as the inexact projection and the projected linesearch are listed in Table 8.1.

Algorithm 1, see Sections 4.1 and 4.3:

τ	$\tau = \frac{1}{\gamma} = 10^3$
ϱ_c	$\varrho_c(t) = \mathbf{c}_c t$ with $\mathbf{c}_c = c_g \tau \tilde{\mathbf{c}}_c$, where $\tilde{\mathbf{c}}_c > 0$ is chosen problem-specific
ϱ_g	$\varrho_g(t) = \mathbf{c}_g t$ with $\mathbf{c}_g = c_g \tilde{\mathbf{c}}_g$, where $\tilde{\mathbf{c}}_g > 0$ is chosen problem-specific
ϱ_{t1}	$\varrho_{t1}(t) = \mathbf{c}_{t1} t$ with $\mathbf{c}_{t1} = \frac{c_1 c_{11}^2 c_{12} c_e c_f}{\tau}$, see Lemma 4.20
ϱ_{t2}	$\varrho_{t1}(t) = \mathbf{c}_{t2} t$ with $\mathbf{c}_{t2} = \frac{1}{c_{12} \tau} \cdot \min\left\{\frac{2(1-c_e)c_1}{c_m}, \frac{c_a}{c_f}\right\}$, see Lemma 4.20
ϱ_r	$\varrho_r(t) = 2c_o \tilde{\mathbf{c}}_o t^{c_o}$ with $\tilde{\mathbf{c}}_o > 0$ chosen problem-specific and $c_o = 1.1$
$(\mathbf{r}_k)_{k \in \mathbb{N}_0}$	$\mathbf{r}_k = \frac{1000}{k+1}$
Δ	$\Delta_{\max} = 10^4$, $\Delta_0 = 1$
η_i	$\eta_1 = 0.3$, $\eta_2 = 0.7$, $\eta_3 = 0.2$
ν_i	$\nu_1 = 0.5$, $\nu_2 = 1.0$, $\nu_3 = 2.0$
$c_\chi \varepsilon_{\text{tol}}$	$c_\chi \varepsilon_{\text{tol}} = 10^{-4}$ (deterministic problems), $c_\chi \varepsilon_{\text{tol}} = 10^{-3}$ (stochastic problems)

Inexact projection and projected linesearch, see Section 4.3:

\mathbf{c}_s	$\mathbf{c}_s = 0$ (exact projection used for computing the inexact criticality measure)
\mathbf{c}_i	$\mathbf{c}_i = 0.5$
$\mathbf{c}_f, \mathbf{c}_e, \mathbf{c}_a$	$\mathbf{c}_f = 0.5$, $\mathbf{c}_e = 10^{-2}$, $\mathbf{c}_a = 10^{-3}$
$\mathbf{c}_{a,k}$	$\mathbf{c}_{a,k} = \max\left\{\mathbf{c}_a, \min\left\{\tau, \frac{c_{12} \Delta_k}{\ \nabla_{m_k}(0)\ _U}\right\}\right\}$, see (4.40)
\mathbf{c}_{11}	$\mathbf{c}_{11} = 0.7$
\mathbf{c}_{12}	$\mathbf{c}_{12} = \frac{2c_{11}-1}{c_{11}} = \frac{4}{7}$, see Remark 4.21
\mathbf{c}_d	$\mathbf{c}_d = 10^{-2}$

Semismooth Newton, see Sections 4.4 and 6.4:

$\mathbf{c}_{n,k}$	$\mathbf{c}_{n,k} = \mathbf{c}_n$ (constant over all iterations), problem-specific
--------------------	--

Table 8.1.: Parameters used in Algorithm 1

In particular, we stop the algorithm if $\chi_k(0) < 10^{-4}$ or $\chi_k(0) < 10^{-3}$ holds in the deterministic or in the stochastic case, respectively. Note that in the unconstrained case this would correspond to having $\|\nabla m_k(0)\| < 10^{-7}$ or $\|\nabla m_k(0)\| < 10^{-6}$ because $\tau = 10^3$ is chosen, cf. (4.9). Almost all parameters have the same value over all tests, but the error bound parameters $\tilde{\mathbf{c}}_c > 0$, $\tilde{\mathbf{c}}_g > 0$, and $\tilde{\mathbf{c}}_o > 0$ for the criticality measure, the gradient, and the objective function evaluation, respectively, are chosen problem-specific because they reflect how good the applied error estimators are or compensate unknown constants. This makes some experimentation necessary. On the one hand, one should not choose the parameters $\tilde{\mathbf{c}}_c$ and $\tilde{\mathbf{c}}_g$ too small because grid refinement happens early then and makes the algorithm slow. On the other hand, if $\tilde{\mathbf{c}}_c$ and $\tilde{\mathbf{c}}_g$ are too large, the trust-region model is not accurate enough to compute a direction of objective function decrease. Then, unsuccessful iterations happen and the trust-region radius is decreased until the gradient is accurate enough. This can take some time and should be avoided. A similar situation appears when choosing $\tilde{\mathbf{c}}_o$. If it is too large, a very inexact objective function evaluation can cause an unsuccessful iteration. Conversely, if it is too small, too much time is spent on the objective function evaluations. Additionally, we sometimes adapt the regularization parameter $\mathbf{c}_{n,k}$ appearing in the semismooth Newton problem (4.42), but choose it constant over all iterations. For example, it can be increased to have a better condition number and typically better convergence of the semismooth Newton iteration, but should not be too large because then the computed steps are too small.

The adaptive solution of the deterministic PDE (Section 7.1) is implemented as follows: The discretized PDE is always solved by Newton's method until the residual norm is below 10^{-6} . Often, we only need one Newton step to fulfill this criterion. This tolerance is chosen such that the FE error dominates in Theorem 7.8 during the whole computation. The triangles contributing 30% of the error are refined. The procedure is repeated until the error estimate given in Theorem 7.8 is small enough. To avoid an infeasible number of FE nodes and unknowns, the refinement is stopped if $2 \cdot 10^5$ FE nodes would be exceeded.

The stochastic PDE is solved adaptively based on Theorem 7.12. We start with the coarse FE mesh shown in Figure 8.1 and with the polynomial degrees $d = \mathbb{1}$, and refine iteratively. Depending on the highest error contribution in (7.32), one of the following tasks is pursued:

- If the algebraic error dominates, one inexact Newton step is computed by performing one AMEn sweep on the Newton equation. To avoid rank growth, the updated tensor is truncated to rank at most 100 w. r. t. to the $\hat{\mathbf{A}}_{\text{ref}}$ -norm as described in Section 7.4.
- If the polynomial discretization error is the largest one, the polynomial degree of one parameter, namely the one with the largest error contribution, is increased by 1.
- If the FE error dominates, the triangles contributing 30% of the error are refined. To make the evaluation of the FE error efficient, we have to round the componentwisely multiplied tensors to avoid too large tensor ranks.

The refinement is programmed such that the last iteration is always an AMEn sweep to have a meaningful solution on the current discrete subspace. Again, the number of FE nodes is bounded by $2 \cdot 10^5$ in order to avoid complexity issues so that the algorithm adapts to the semi-discrete solution using the finest FE mesh at the end. In some cases, AMEn stagnates and does not manage to compute a solution with small enough algebraic error in the last iterations of Algorithm 1. Then the adaptive solution procedure is stopped.

The projection onto the box U_{ad} is computed adaptively by computing the exact projection and the corresponding refined grid first. Then the exact error made by the node-wise projection can be evaluated on each element. If the node-wise projection does not meet one of the exactness requirements, the triangles contributing 30% of the projection error are refined. We use uniform refinement to not destroy the quality of the mesh, which is important to uniformly bound the constants $c_{\mathcal{T}}$ in (7.16) and (7.32), which depend on the smallest angle in the triangulations. The refinement is repeated until the inexact projection is accurate enough.

All computations are run on a Linux cluster with 1 TB RAM and 4 Intel Xeon E7-8857 processors, each of which has 12 kernels and a base frequency of 3.00 GHz. However, only serial but vectorized implementations are used in MATLAB because current toolboxes do not offer, e. g., parallel computations with low-rank tensors or parallel FE codes.

8.2. Results for the Deterministic Problem

First, we implement the deterministic problem, i. e., the optimal control problem with the objective function $J[\xi]$ (3.4) and the state equation (3.8), where the reference coefficient $\xi = \bar{\xi} = 0$ is inserted. Equivalently, we could solve the stochastic problem with $\sigma = 0$. We do this in order to test the error estimates from Sections 5.1 and 7.1, to visualize the adaptivity of the FE grid and to have reference solutions for the problem under uncertainty.

Since the discretized PDE is often solved to good accuracy by only one Newton step, there is no need to balance the errors in Theorem 7.8. Therefore, we choose $c_{\mathcal{T}}c_{A_{\text{ref}}} = 1$ in (7.13) and adapt the error estimation parameters adequately.

The first setup uses the desired state $\hat{q} \equiv 1$. Here we choose the parameters $\tilde{c}_c = 200$, $\tilde{c}_g = 200$, $\tilde{c}_o = 10^8$, and $c_n = 0.1$. The high value for \tilde{c}_o is needed to avoid too severe grid refinement due to inexact objective function evaluation. This value comes also from the fact that we overestimate the error given in Proposition 5.5, because we measure the state error in the $H_0^1(\Omega)$ -norm although the $L^2(\Omega)$ -error would be sufficient. The computed optimal state and the control are depicted in Figure 8.2. It can be observed that the optimal state is in the order of magnitude of the desired state. Furthermore, the three active sets of the control and the symmetry of the problem can be recognized. The FE mesh obtained in the final iteration of the algorithm is shown in Figure 8.3. Around the corner $x = 0$, the mesh is refined locally due to the adaptive solution of the PDE. Furthermore, the boundary of the active set of the control is resolved by the mesh due to the refined projection.

The convergence is shown in Figure 8.4. This is the typical convergence behavior so that we do not show convergence plots for the other deterministic setups. The number of FE nodes increases until it reaches the allowed upper bound. The criticality measure decreases in general until the desired tolerance is reached. But by refining the mesh the inexactness of the computed criticality measure is recognized sometimes. Then the computed criticality measure on the refined mesh is larger than the one on the coarser grid. We have to mention that it is not worth counting iterations in this setting. Typically, the first iterations run very fast within some seconds. For instance, the first 16 iterations needed to decrease the computed optimality measure from 23.2 to $3.09 \cdot 10^{-2}$ take about 19 seconds in total. Only the last iterations required to obtain the desired high accuracy last increasingly longer so that

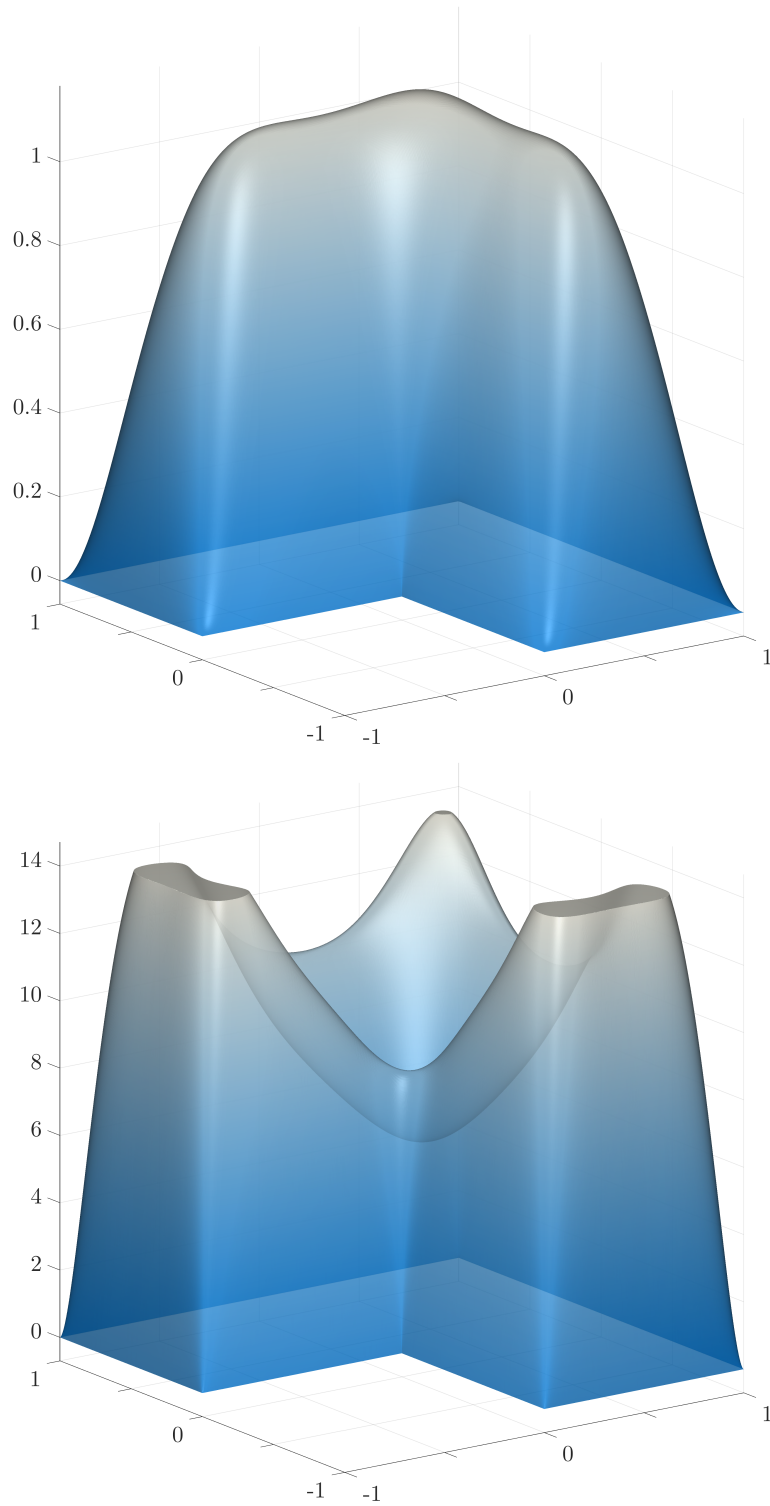


Figure 8.2.: Optimal state (top) and optimal control (bottom) for the deterministic problem with $\hat{q} \equiv 1$.

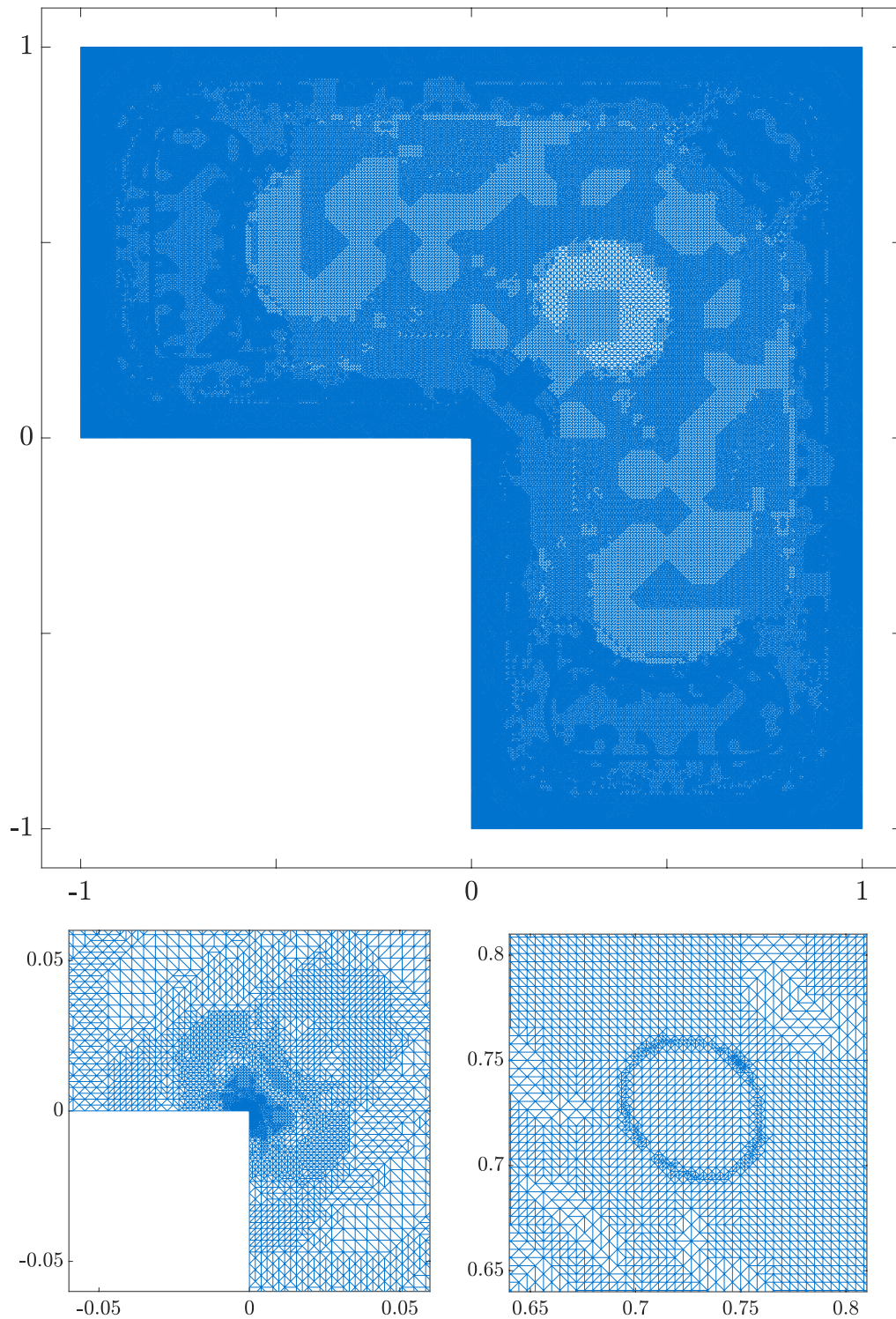


Figure 8.3.: Final mesh (with details at the bottom) for the deterministic problem with $\hat{q} \equiv 1$.

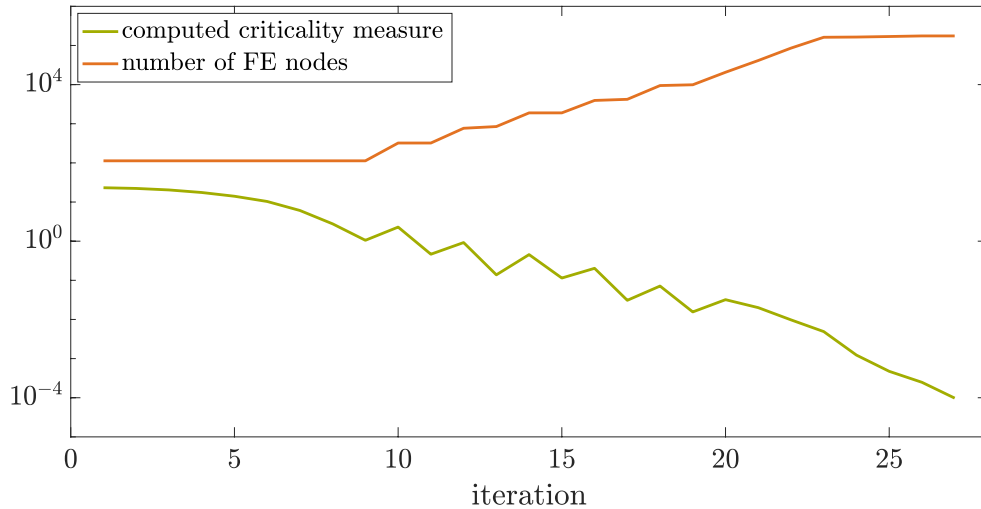


Figure 8.4.: Convergence and refinement plot for the deterministic problem with $\hat{q} \equiv 1$.

the total computing time is around 22 minutes. The last three iterations, during which the computed criticality measure is decreased from $1.23 \cdot 10^{-3}$ to $0.98 \cdot 10^{-4}$, take approximately 13 minutes. Often, the grid refinement due to the projection errors is done to more extent during the last iterations. Figure 8.5 shows the grid 5 iterations before the last one. It has approximately half the number of FE nodes compared to the one shown in Figure 8.3 and the resolution of the active sets is still missing.

As a second setup, we change the desired state to the polynomial $\hat{q}(x) = 9(x_1^3 - x_1)(x_2^3 - x_2)$. This is a smooth function fulfilling the zero boundary conditions on Ω . Additionally, it is very smooth around the corner point $x = 0$. We do not change the parameters compared to the first setup. The computed optimal state and control are depicted in Figure 8.6. The active set of the control is also nicely resolved in the final mesh (Figure 8.7). Around the corner $x = 0$, the mesh is not locally refined as before because of the smoothness of the optimal state around this point.

In the third setup, the desired state is $\hat{q} \equiv 1_{\{x \in \Omega: x_2 > 0.5\}}$. Here, it is necessary to choose the parameters $\tilde{\mathbf{c}}_c = 20$, $\tilde{\mathbf{c}}_g = 20$, $\tilde{\mathbf{c}}_o = 10^8$, and $\mathbf{c}_n = 1.0$. The gradient is computed more exactly by this choice so that that the algorithm does not run into unsuccessful iterations. The larger regularization parameter for semismooth Newton yields a more robust iteration. In this case, the optimal state is almost zero for $x_2 < 0$, see Figure 8.8, and also the optimal control has very small values in this area. By the first-order optimality condition, the adjoint state is also small there. This yields that the final mesh is quite coarse for $x_2 < 0$ except for the area around the corner $x = 0$, see Figure 8.9. This can be explained by taking a look at the error estimates (7.14) and (7.15). They become small on the triangles and their respective edges if the solution \tilde{y} and the jump in its gradient, the residual w and the jump in its gradient, and the right-hand side \hat{f} are small. This is true for both the state and the adjoint equation. Again, the active set of the control is resolved nicely, cf. Figures 8.8 and 8.9.

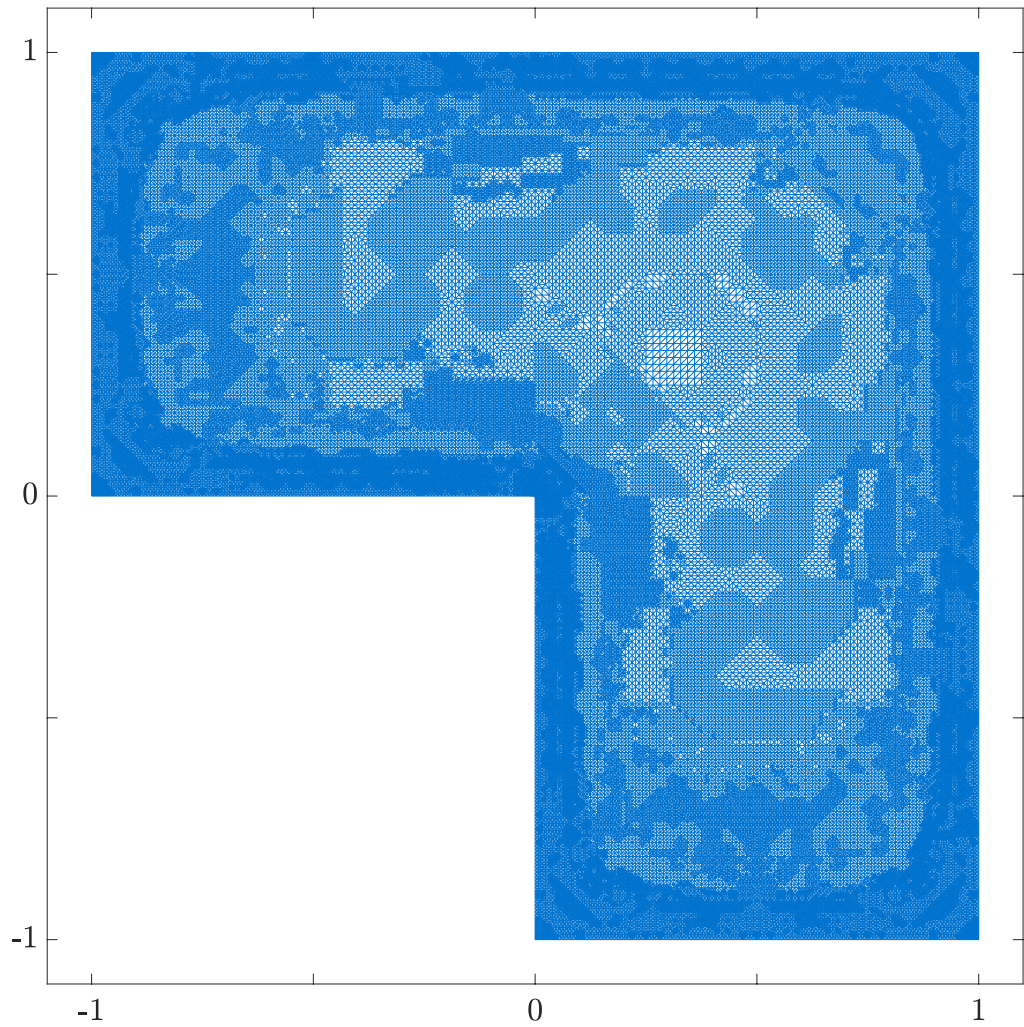


Figure 8.5.: Mesh from iteration 22 (of 27) for the deterministic problem with $\hat{q} \equiv 1$.

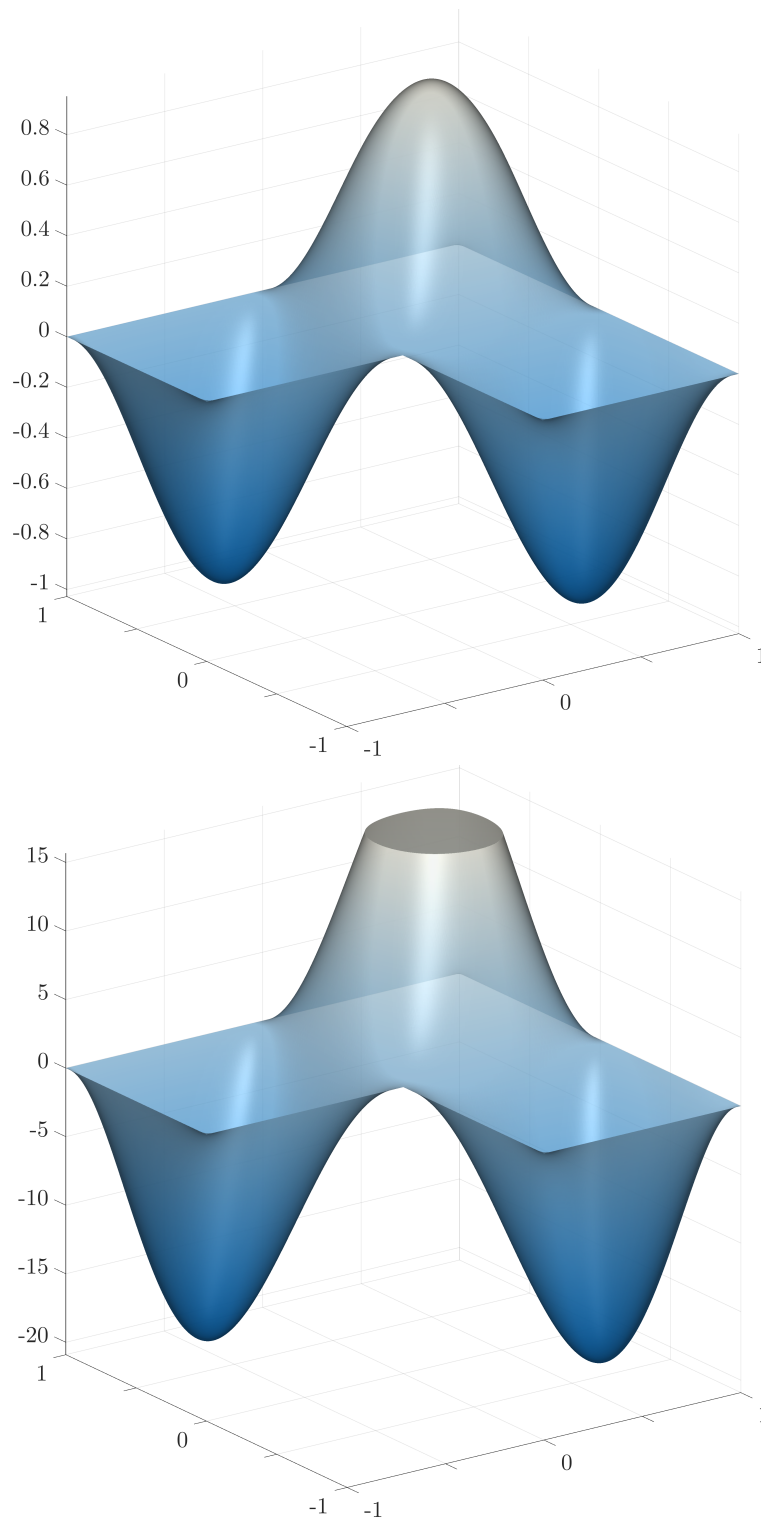


Figure 8.6.: Optimal state (top) and optimal control (bottom) for the deterministic problem with $\hat{q}(x) = 9(x_1^3 - x_1)(x_2^3 - x_2)$.

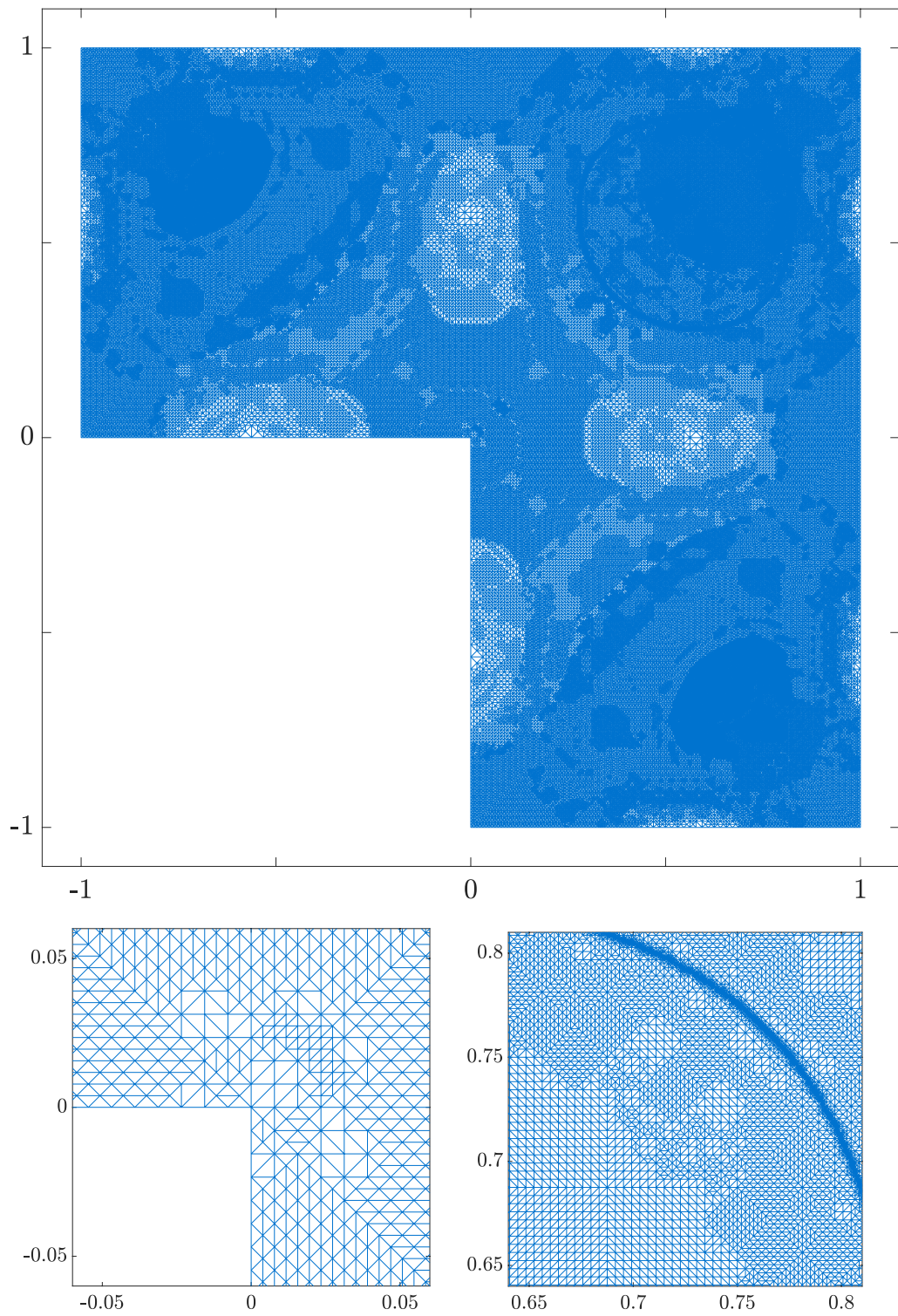


Figure 8.7.: Final mesh (with details at the bottom) for the deterministic problem with $\hat{q}(x) = 9(x_1^3 - x_1)(x_2^3 - x_2)$.

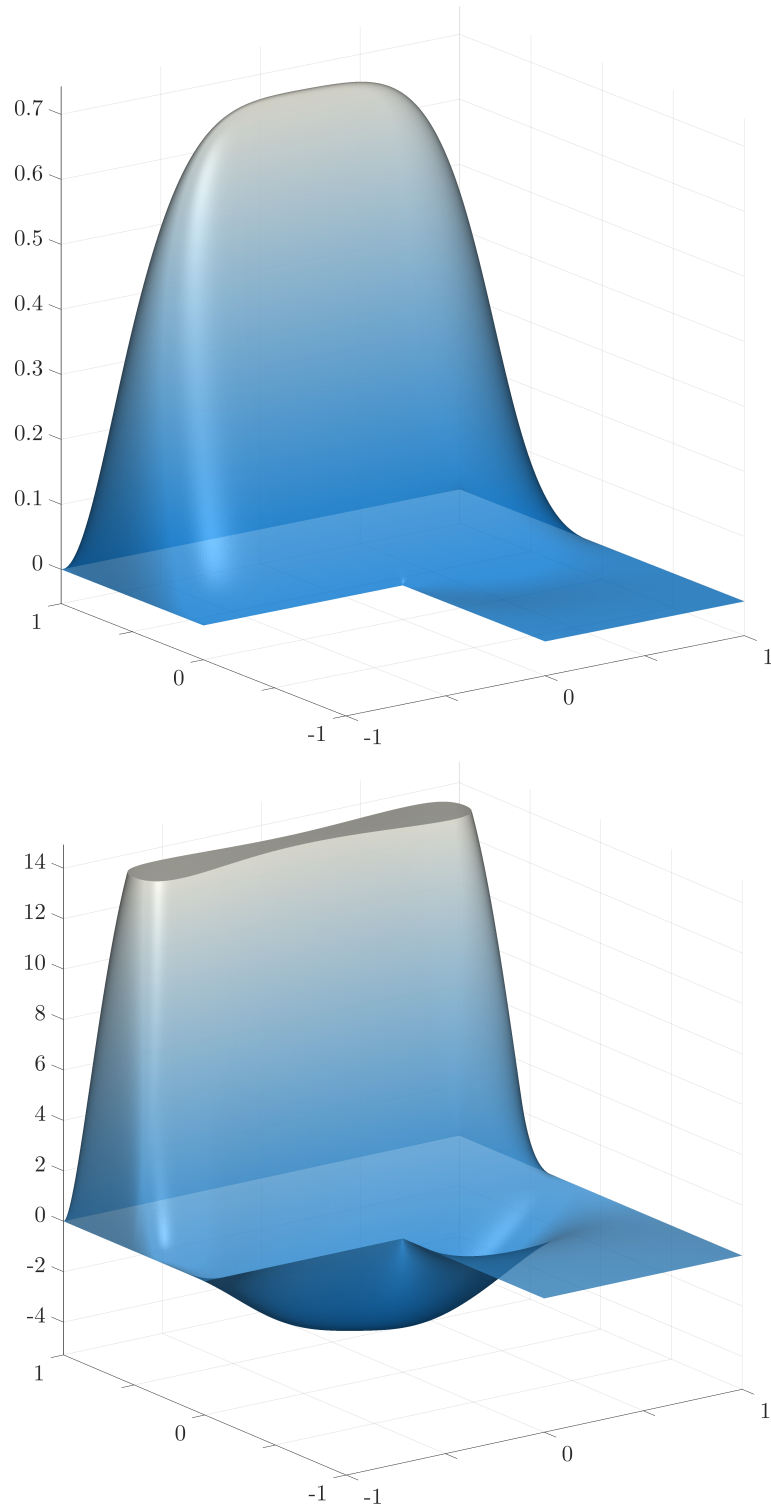


Figure 8.8.: Optimal state (top) and optimal control (bottom) for the deterministic problem with $\hat{q} \equiv 1_{\{x \in \Omega: x_2 > 0.5\}}$.

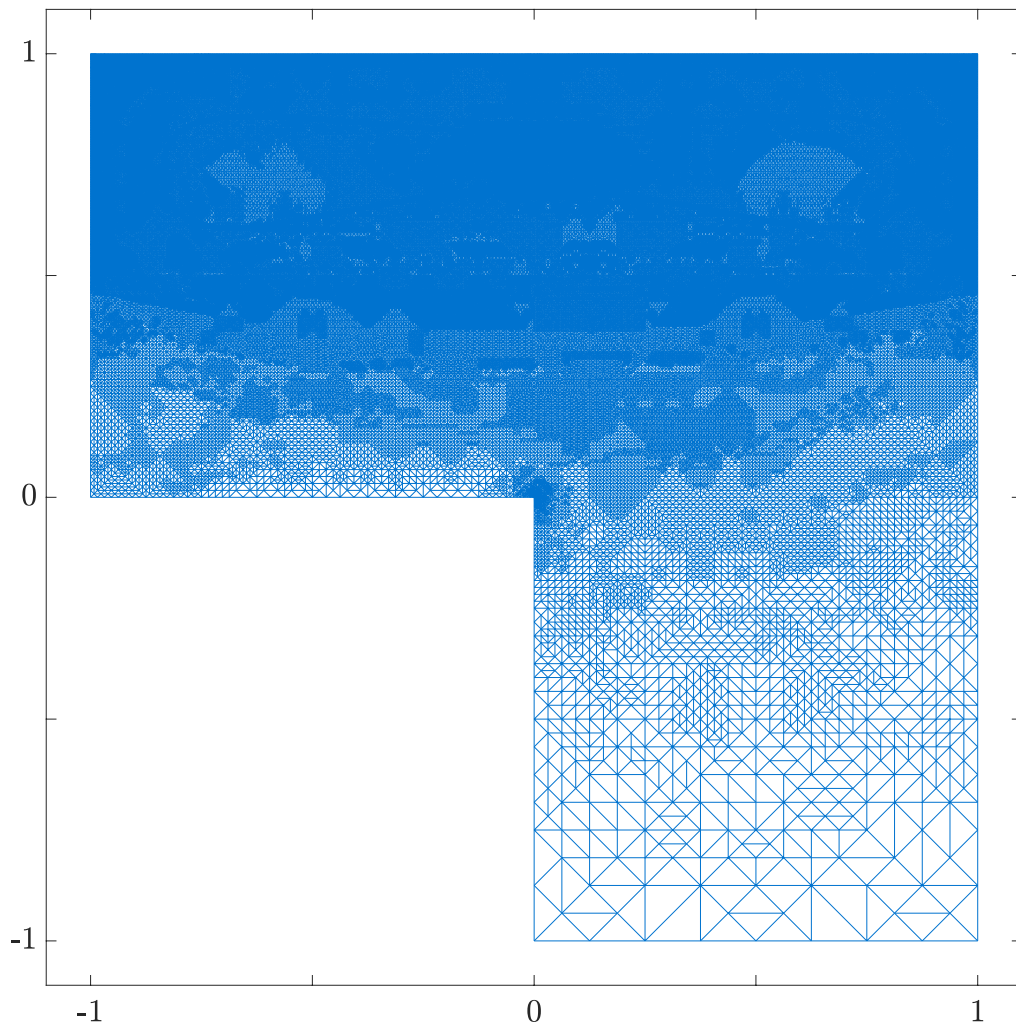


Figure 8.9.: Final mesh for the deterministic problem with $\hat{q} \equiv 1_{\{x \in \Omega: x_2 > 0.5\}}$.

8.3. Results for the Problem with Uncertainties

Now we consider the stochastic problem (3.3) with the state equation defined in (3.11) and vary the problem data to see how the algorithm adapts to the situation. As already noted, we choose $c_{\mathcal{T}}c_{A_{\text{ref}}} = 10^{-3}$ to balance the error contributions coming from the FE and the polynomial chaos discretization. This makes different choices of the error control constants necessary.

Again, we start with the desired state $\hat{q} \equiv 1$ and with the vector of influence amounts $\sigma = (0.25, 0.25, 0.25, 0.25, 0.25, 0.25)^\top$, i. e., the coefficient function varies $\pm 25\%$ equally on each subdomain independently and each parameter has the same influence on it. We choose $\tilde{c}_c = 0.05$, $\tilde{c}_g = 0.05$, $\tilde{c}_o = 10^3$, and $\mathbf{c}_n = 0.01$. The obtained optimal control and the expected optimal state are quite similar to the ones from the deterministic setup, see Figure 8.10, but the active sets are a bit smaller so that the smallest one almost disappears. The final mesh (Figure 8.11) is very similar to the one obtained before. Grid refinements due to the active sets can be recognized, but are not as pronounced as in the deterministic case because the iteration is already stopped when $\chi_k(0) < 10^{-3}$ instead of $\chi_k(0) < 10^{-4}$. The reason for this is that the high fidelity PDE solves in the last iterations become quite expensive. The algorithm takes about 42 hours to compute this solution, where 37 hours are spent for the last two iterations decreasing the computed criticality measure from $1.39 \cdot 10^{-2}$ to $5.82 \cdot 10^{-4}$. In practice, one would probably refrain from performing the last two iterations because less accurate solutions are sufficient if, e. g., the problem data or the parameter distribution is not known exactly. The convergence and refinement behavior is shown in Figure 8.12. We see that the computed criticality measure decreases monotonically and the algorithm stops in iteration 13, whereas a criticality measure smaller than 10^{-3} is obtained in iteration 25 in the deterministic case, see Figure 8.4. This is because the error control constants \tilde{c}_c and $\tilde{c}_g = 0.05$ are chosen a bit smaller even if we take account of the smaller constant $c_{\mathcal{T}}c_{A_{\text{ref}}}$ so that better steps are computed. The number of FE nodes grows until it reaches its upper bound. Furthermore, the polynomial degrees for each parameter increase almost equally during all iterations because the parameters have roughly the same influence on the optimal uncertain state. Figure 8.14 shows the difference between the deterministic and the robust control. We see the symmetry of the problem and can recognize the boundaries of the parameter influence sets Ω_i . The difference between the controls is zero where their active sets coincide.

The situation changes if we keep the desired state $\hat{q} \equiv 1$, but choose the influence amounts $\sigma = (0.05, 0.10, 0.20, 0.30, 0.40, 0.45)^\top$. The average influence of each uncertain parameter is still 25%, but now it is increasing clockwise, i. e., parameter 1 has the smallest influence of 5% on the coefficient function, whereas parameter 6 has the largest influence of 45%. The convergence plot (Figure 8.13) is changed by this choice in the sense that higher polynomial degrees are needed for the parameters with more influence on the system. This becomes clear if we review the stochastic error indicator (7.31), which becomes smaller if the coefficient function κ_i is smaller. Additionally, it can be seen in Figure 8.15 that the difference between the deterministic and the robust control is larger in areas where the variation of the coefficient function is larger. Clearly, the problem is not symmetric anymore and one can recognize the size change of the active sets, especially the one on the bottom right.

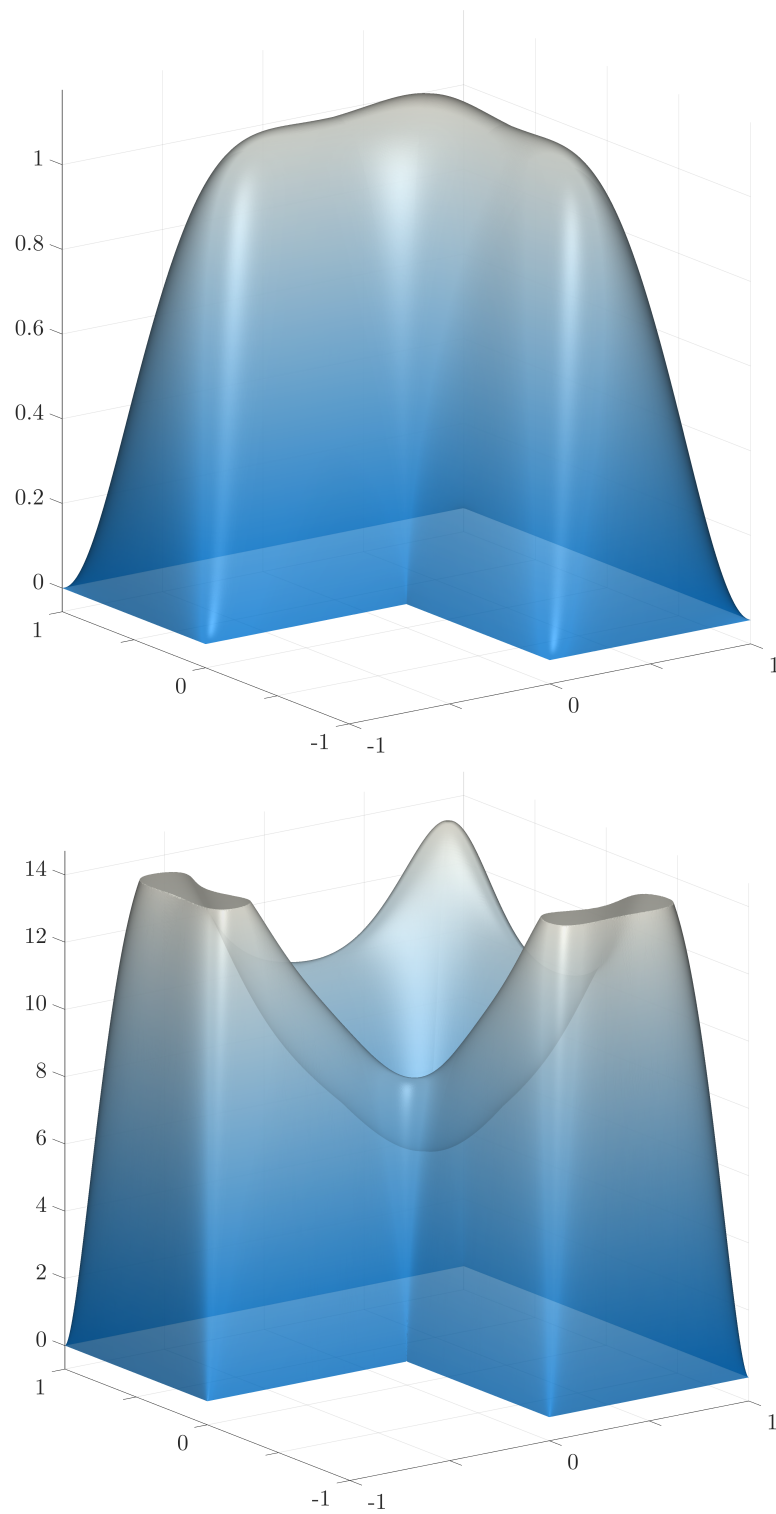


Figure 8.10.: Optimal control (bottom) and expected optimal state (top) for the stochastic problem with $\hat{q} \equiv 1$ and $\sigma = (0.25, 0.25, 0.25, 0.25, 0.25, 0.25)^\top$.

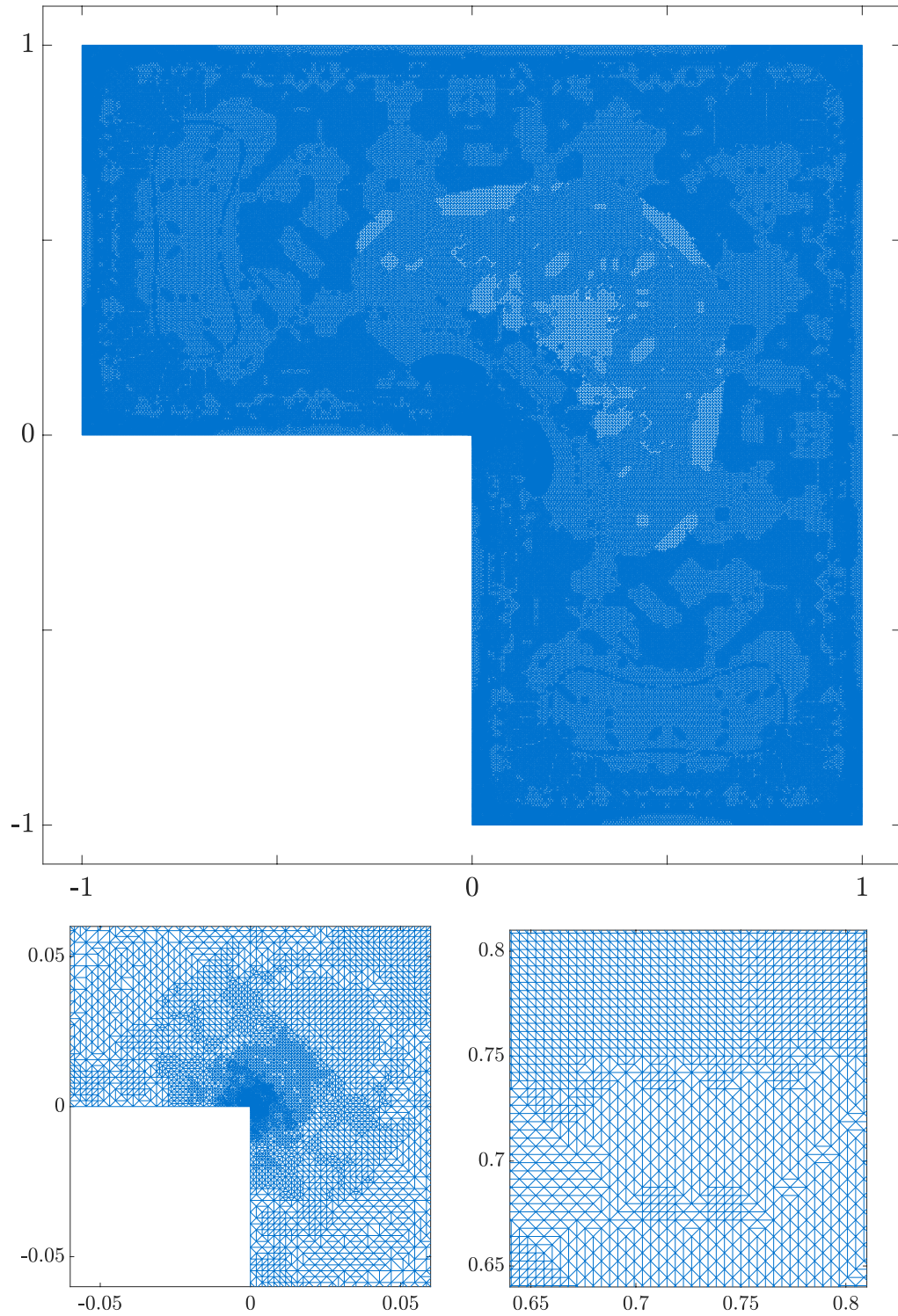


Figure 8.11.: Final mesh (with details at the bottom) for the stochastic problem with $\hat{q} \equiv 1$ and $\sigma = (0.25, 0.25, 0.25, 0.25, 0.25, 0.25)^\top$.

8. Implementation and Numerical Results

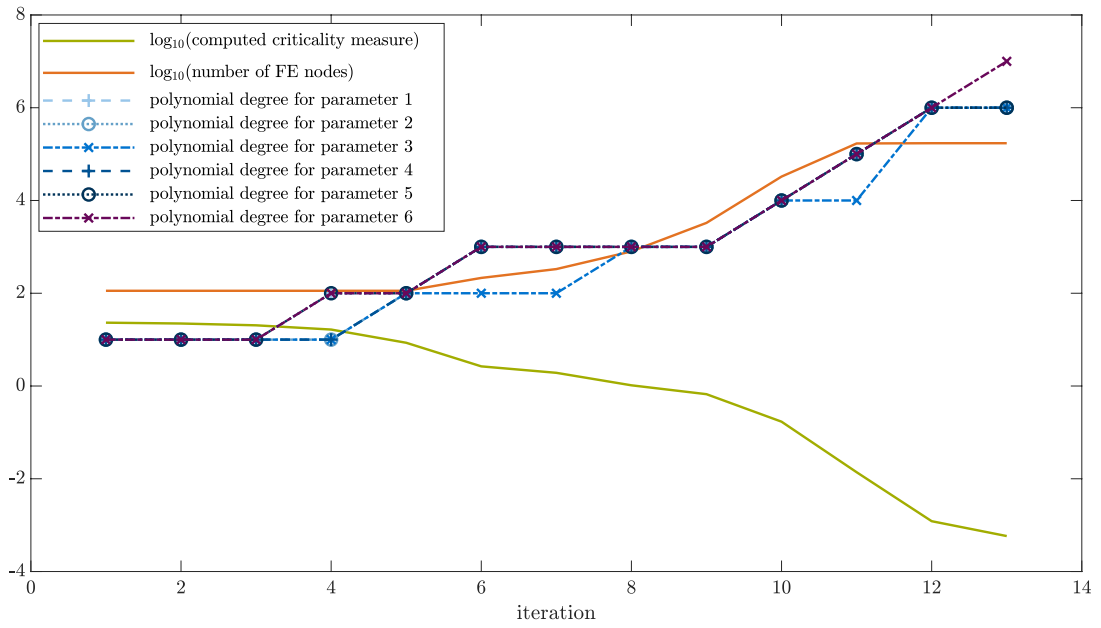


Figure 8.12.: Convergence and refinement plot for the stochastic problem with $\hat{q} \equiv 1$ and $\sigma = (0.25, 0.25, 0.25, 0.25, 0.25, 0.25)^\top$.

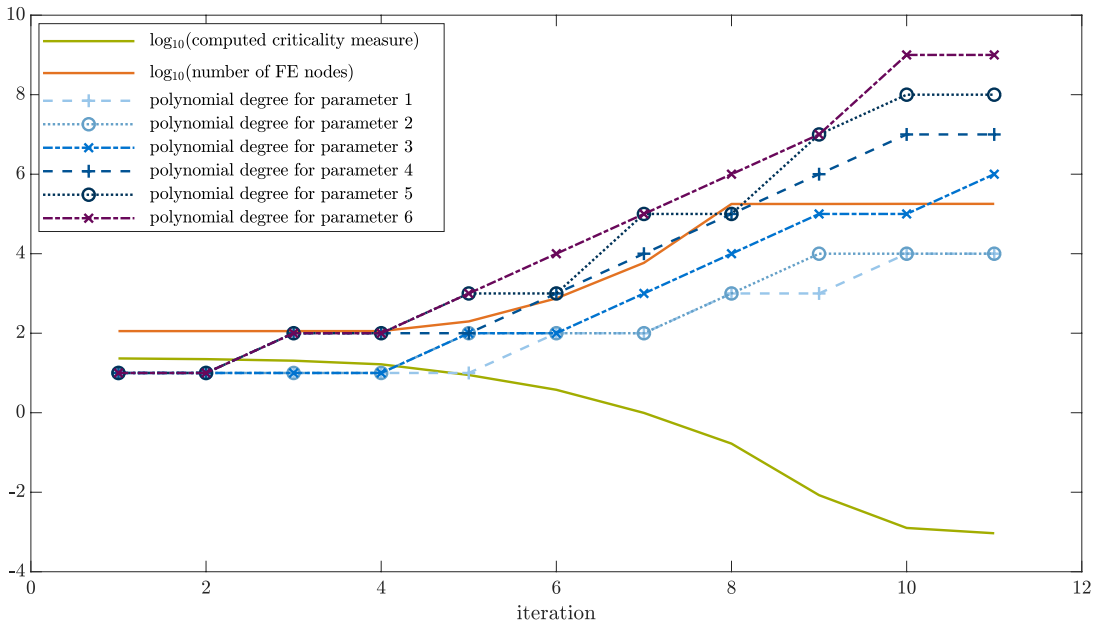


Figure 8.13.: Convergence plot and refinement for the stochastic problem with $\hat{q} \equiv 1$ and $\sigma = (0.05, 0.10, 0.20, 0.30, 0.40, 0.45)^\top$.

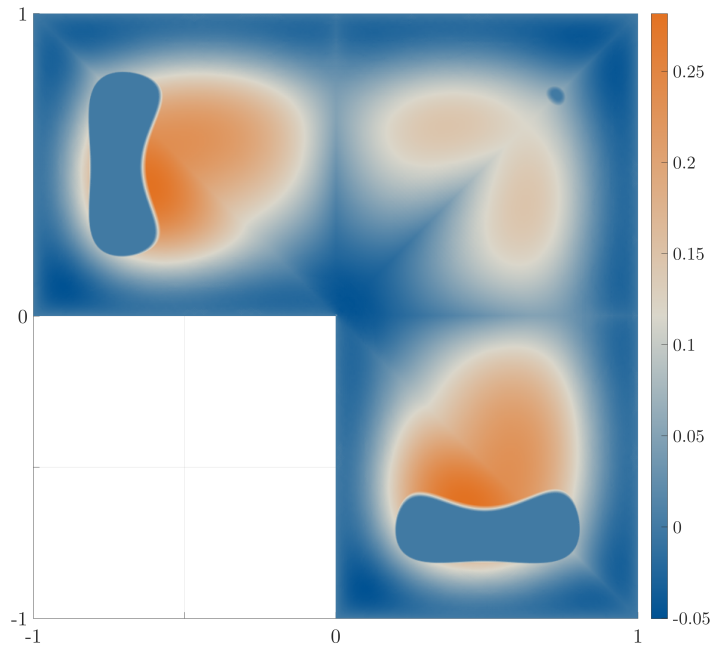


Figure 8.14.: Difference between the deterministic and the robust control for $\hat{q} \equiv 1$ and $\sigma = (0.25, 0.25, 0.25, 0.25, 0.25, 0.25)^\top$.

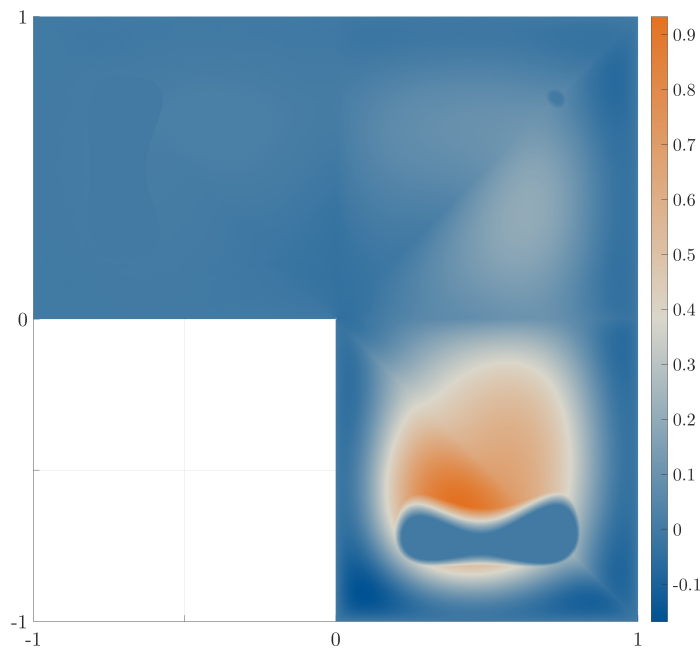


Figure 8.15.: Difference between the deterministic and the robust control for $\hat{q} \equiv 1$ and $\sigma = (0.05, 0.10, 0.20, 0.30, 0.40, 0.45)^\top$.

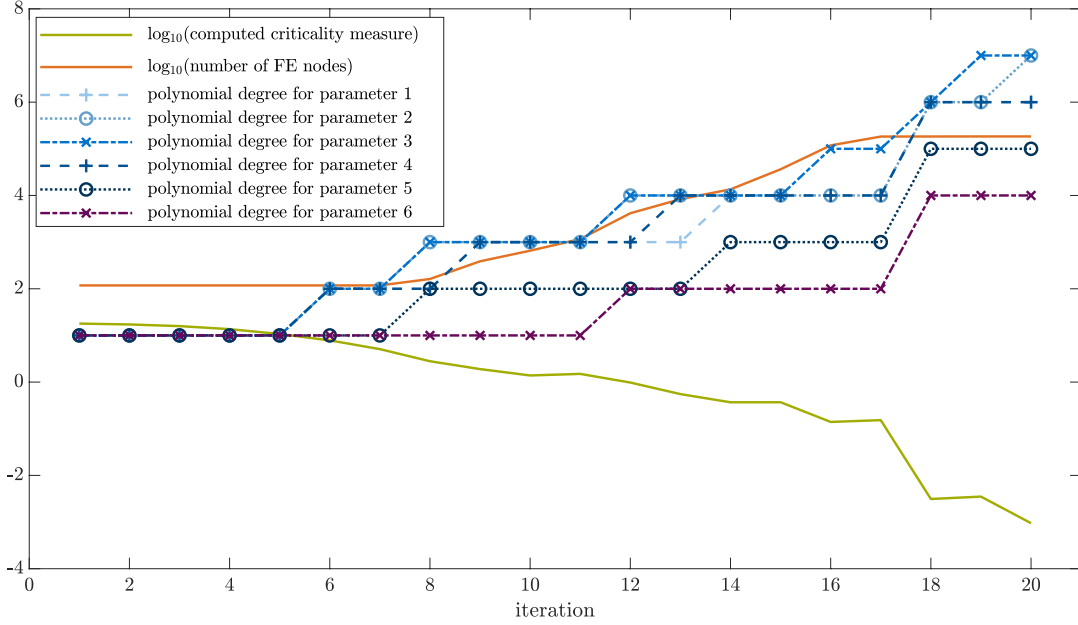


Figure 8.16.: Convergence and refinement plot for the stochastic problem with $\hat{q} \equiv 1_{\{x \in \Omega: x_2 > 0.5\}}$ and $\sigma = (0.25, 0.25, 0.25, 0.25, 0.25, 0.25)^\top$.

As a last setup, we consider the desired state $\hat{q} \equiv 1_{\{x \in \Omega: x_2 > 0.5\}}$ with equal influence $\sigma = (0.25, 0.25, 0.25, 0.25, 0.25, 0.25)^\top$. We choose the parameters $\tilde{\mathbf{c}}_c = 0.01$, $\tilde{\mathbf{c}}_g = 0.01$, $\tilde{\mathbf{c}}_o = 10^3$, and $\mathbf{c}_n = 0.1$ in analogy to the choice in the deterministic setup. Again, we obtain a similar optimal control and expected state. The convergence and refinement plot in Figure 8.16 reveals that smaller polynomial degrees are chosen especially for parameters 5 and 6, which act in the area where the optimal state and its gradient are almost zero. Because of this fact, the stochastic error indicator (7.31) becomes smaller for these parameters since it depends on the size of the state gradient in the area where the coefficient κ_i is not small. The difference to the deterministic control is larger in the subdomains $\Omega_1, \dots, \Omega_4$, see Figure 8.17, which is another indicator for the fact that parameters 5 and 6 do not influence the solution as much as parameters 1 to 4.

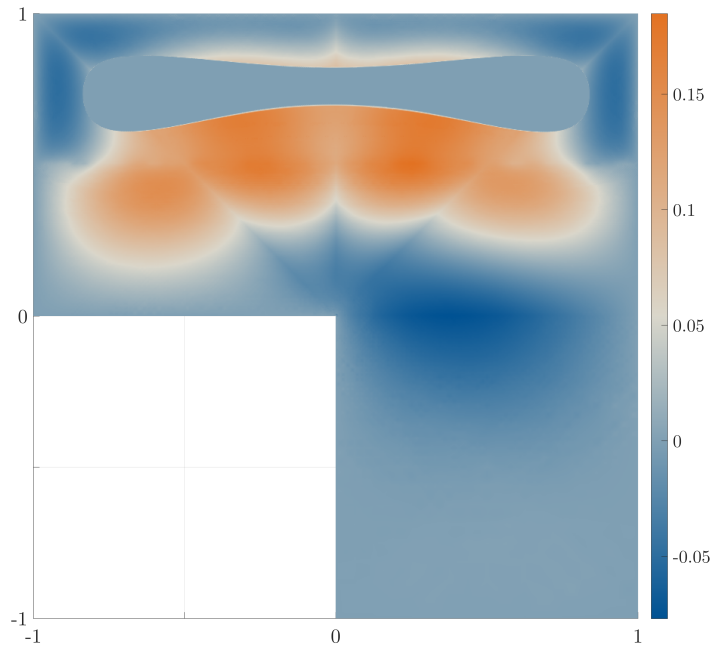


Figure 8.17.: Difference between the deterministic and the robust control for $\hat{q} \equiv 1_{\{x \in \Omega: x_2 > 0.5\}}$ and $\sigma = (0.25, 0.25, 0.25, 0.25, 0.25, 0.25)^\top$.

9. Alternative Risk Measures

Until now, all considerations in this thesis have referred to the expectation as risk measure, i. e., $\mathcal{R} \equiv \mathbb{E}$ in (1.2). This means that we want the controlled system to perform well *on average*. On the one hand, this is a viable approach if, e. g., the objective function stands for a cost one has to pay and this cost shall be small averaged over a long time horizon during which the uncertain parameters may change their value according to their distribution. On the other hand, in certain engineering applications, such as aeronautics or nuclear reactors, a high cost function value can correspond to a system fail, meaning that the plane crashes or the reactor core melts for example. Such situations should clearly be avoided with high probability. This is not taken into account if $\mathcal{R} \equiv \mathbb{E}$ is used as a risk measure because those scenarios are typically unlikely to occur so that they do not influence the expected value very much. Analogously, the expectation might not be the risk measure of choice if the randomness is only observed once and the probability of having a high cost shall be small.

In this chapter, we give an introduction to risk measures, which are functionals rating the probability distribution of a random variable, and their properties. Furthermore, we provide two examples which we discuss in the context of optimal control. Both lead to *risk-averse* choices of the optimal control instead of the *risk-neutral* controls computed with $\mathcal{R} \equiv \mathbb{E}$. The mean-variance risk measure $\mathcal{R} \equiv \mathbb{E} + \lambda \text{Var}$ with $\lambda > 0$ (Section 9.2) is smooth so that the trust-region algorithm (Algorithm 1) can be applied to solve optimal control problems involving it. We discuss this in theory and show how the corresponding required error estimates (cf. Chapter 5) change in this case. The conditional value-at-risk (CVaR, Section 9.3) is a risk measure with better theoretical properties than the mean-variance risk measure and is used frequently in financial mathematics [94, 95, 77, 103]. For instance, in the example control problem from the previous chapters the mean-variance risk measure would penalize upward and downward deviations of the cost term from its mean equally although downwards deviations lead to a probably desired smaller cost. In contrast, the CVaR is computed based on large deviations from the desired state only since it is the expectation of the upper tail of the considered random variable. Since CVaR is nonsmooth, different solution techniques have to be applied for solving optimal control problems involving it. We discuss the application of a primal interior-point method and its implementation with low-rank tensors.

9.1. Definition and Properties of Risk Measures

Definition 9.1. Let $(\Xi, \mathcal{F}, \mathbb{P})$ be a complete probability space. A *risk measure* \mathcal{R} is a functional $\mathcal{R} : L_{\mathbb{P}}^p(\Xi) \rightarrow \mathbb{R} \cup \{\infty\}$, where $p \in [1, \infty]$.

Remark 9.2. Instead of $L_{\mathbb{P}}^p(\Xi)$, we could also consider more general spaces or sets of random variables. Furthermore, the value $-\infty$ could be incorporated, see [103, Sec. 6.3].

Example 9.3. Prominent examples for risk measures are (see also [96, Sec. 2])

- the expectation $\mathcal{R} \equiv \mathbb{E}$, which is well-defined and always finite on $L_{\mathbb{P}}^1(\Xi)$,
- the *mean-variance risk measure* $\mathcal{R}[X] = \mathbb{E}[X] + \lambda \text{Var}[X]$ for some $\lambda > 0$, which is well-defined on $L_{\mathbb{P}}^2(\Xi)$,
- the risk measure $\mathcal{R}[X] = \mathbb{E}[X] + \lambda \sigma[X]$ with $\sigma[X] = \sqrt{\text{Var}[X]}$, which is well-defined on $L_{\mathbb{P}}^2(\Xi)$, and
- the β -quantile or *value-at-risk* $\mathcal{R}[X] = \text{VaR}_{\beta}[X]$ for some confidence level $\beta \in (0, 1)$, which is defined by

$$\text{VaR}_{\beta}[X] := \min\{t \in \mathbb{R} : F_X(t) \geq \beta\}$$

[95, Def. 1] with the distribution function $F_X(t) := \mathbb{P}(X \leq t)$. Note that VaR_{β} is well-defined since F_X is non-decreasing and right-continuous. Lastly, we have

- the *conditional value-at-risk* $\mathcal{R}[X] = \text{CVaR}_{\beta}[X]$ defined by

$$\text{CVaR}_{\beta}[X] := \text{“mean of the } \beta\text{-tail distribution of } X\text{”},$$

i. e., the mean of a random variable with distribution function $F_X^{\beta}(t) := \frac{1}{1-\beta}(F_X(t) - \beta) 1_{\{t \geq \text{VaR}_{\beta}[X]\}}(t)$ [95, Def. 3]. We will see how the CVaR can be computed.

Definition 9.4. Let \mathcal{R} be a risk measure. It may have the following properties (see [96]):

- P1. *Convexity:* $\mathcal{R}[(1-\tau)X + \tau Y] \leq (1-\tau)\mathcal{R}[X] + \tau\mathcal{R}[Y]$ for all $\tau \in [0, 1]$ and all $X, Y \in L_{\mathbb{P}}^p(\Xi)$.
- P2. *Positive homogeneity:* $\mathcal{R}[0] = 0$ and $\mathcal{R}[\lambda X] = \lambda\mathcal{R}[X]$ for all $\lambda \in (0, \infty)$ and all $X \in L_{\mathbb{P}}^p(\Xi)$.
- P3. *Subadditivity:* $\mathcal{R}[X + Y] \leq \mathcal{R}[X] + \mathcal{R}[Y]$ for all $X, Y \in L_{\mathbb{P}}^p(\Xi)$.
- P4. *Monotonicity:* $\mathcal{R}[X] \leq \mathcal{R}[Y]$ for all $X, Y \in L_{\mathbb{P}}^p(\Xi)$ such that $X \leq Y$ almost surely (a. s.)
- P5. *Translation equivariance:* $\mathcal{R}[X + c] = \mathcal{R}[X] + c$ for all $X \in L_{\mathbb{P}}^p(\Xi)$ and all $c \in \mathbb{R}$.
- P6. *Closedness:* The set $\{X \in L_{\mathbb{P}}^p(\Xi) : \mathcal{R}[X] \leq c\} \subset L_{\mathbb{P}}^p(\Xi)$ is closed for all $c \in \mathbb{R}$.
- P7. *Aversity to risk:* $\mathcal{R}[X] > \mathbb{E}[X]$ for all non-constant $X \in L_{\mathbb{P}}^p(\Xi)$.

Some of these properties are used in the definition of other classes of risk measures:

Definition 9.5. A proper, convex, positively homogeneous, monotonic, and translation-equivariant risk measure $\mathcal{R} : L_{\mathbb{P}}^p(\Xi) \rightarrow \mathbb{R} \cup \{\infty\}$ is called *coherent* [103, Sec. 6.3].

Remark 9.6. By [17, Prop. 10.3], a coherent risk measure is automatically *subadditive* because it is convex and positively homogeneous. Moreover, in the original definition [6, Def. 2.4], it is not explicitly mentioned that the risk measure has to be proper because only *finite* risk measures are considered. As noted in [96], a finite, convex and monotonic risk measure is automatically continuous and subdifferentiable on $L_{\mathbb{P}}^p(\Xi)$ by [99, Prop. 3.1].

Definition 9.7. A risk measure $\mathcal{R} : L_{\mathbb{P}}^p(\Xi) \rightarrow \mathbb{R} \cup \{\infty\}$ is called *regular* [96] if it is closed convex, risk averse, and fulfills $\mathcal{R}[C] = C$ for all $C \in \mathbb{R}$.

Example 9.8.

- $\mathcal{R} \equiv \mathbb{E} : L_{\mathbb{P}}^p(\Xi) \rightarrow \mathbb{R}$ ($p \geq 1$) fulfills properties P1, P2, P3, P4, P5, and P6, but not P7.
- $\mathcal{R} \equiv \mathbb{E} + \lambda \text{Var} : L_{\mathbb{P}}^p(\Xi) \rightarrow \mathbb{R}$ ($\lambda > 0, p \geq 2$) is regular and fulfills property P5, but not P2, P3, and P4.
- $\mathcal{R} \equiv \mathbb{E} + \lambda \sigma : L_{\mathbb{P}}^p(\Xi) \rightarrow \mathbb{R}$ ($\lambda > 0, p \geq 2$) is regular and fulfills properties P2, P3, and P5, but not P4.

For the last two risk measures, convexity and closedness can be observed by $\text{Var}[X] = \|X - \mathbb{E}[X]\|_{L_{\mathbb{P}}^2(\Xi)}^2$ and $\sigma[X] = \|X - \mathbb{E}[X]\|_{L_{\mathbb{P}}^2(\Xi)}$, i. e., they both are a composition of a closed, convex and a linear, bounded function. In addition, $\mathbb{E} + \lambda \sigma$ is subadditive, because it is convex and positively homogeneous, see Remark 9.6. But $\mathbb{E} + \lambda \text{Var}$ is not subadditive: Consider a random variable X with positive variance. Then, $\mathbb{E}[X + X] + \lambda \text{Var}[X + X] = 2\mathbb{E}[X] + 4\lambda \text{Var}[X] > 2(\mathbb{E}[X] + \lambda \text{Var}[X])$.

The lacking monotonicity for both $\mathbb{E} + \lambda \text{Var}$ and $\mathbb{E} + \lambda \sigma$ is one of the less obvious properties: Consider, e. g., a random variable X such that $X = d < 0$ with probability $p \in (0, 1)$ and $X = 0$ with probability $1 - p \in (0, 1)$. Furthermore, consider $Y \equiv 0$. Then, $X \leq Y$ a. s.,

$$\mathbb{E}[X] = pd, \quad \text{Var}[X] = pd^2 - (pd)^2 = p(1 - p)d^2,$$

$\mathbb{E}[Y] = 0$, and $\text{Var}[Y] = 0$. Thus,

$$0 = \mathbb{E}[Y] + \lambda \text{Var}[Y] < \mathbb{E}[X] + \lambda \text{Var}[X] = pd + \lambda p(1 - p)d^2 = pd(1 + \lambda(1 - p)d)$$

if and only if $d \in (-\frac{1}{(1-p)\lambda}, 0)$. Analogously, $0 < \mathbb{E}[X] + \lambda \sigma[X] = pd - \lambda \sqrt{p(1-p)}d$ if and only if $p \in (0, \frac{\lambda^2}{1+\lambda^2})$.

Furthermore, VaR_{β} is *not* subadditive in general [6], but CVaR_{β} fulfills all properties given in Definition 9.4 as shown later.

9.2. Mean-Variance Risk Measure

We consider the reduced form (1.3) of an optimal control problem with the risk measure $\mathbb{E} + \lambda \text{Var}$ for some $\lambda > 0$ as done in, e. g., [3]. This gives the problem

$$\min_{u \in U} \hat{\mathcal{J}}(u) := \mathbb{E}[\hat{J}_1(u, \cdot)] + \lambda \text{Var}[\hat{J}_1(u, \cdot)] + J_2(u) \quad \text{s. t.} \quad u \in U_{\text{ad}}. \quad (9.1)$$

Compared to the risk-neutral case, this approach shall additionally reduce the variability (measured by the variance) of the tracking term $\hat{J}_1(u, \cdot)$. But it should be mentioned that the objective function $\hat{\mathcal{J}}$ may be non-convex even though the functions $\hat{J}_1(\cdot, \xi)$ and $\hat{J}_2(\cdot)$ are convex for almost every $\xi \in \Xi$ and the variance itself is a convex functional. Hence, it is possible that only a local minimizer is computed by a local optimization algorithm, such as

the proposed trust-region method (Algorithm 1). The basic assumptions from Chapter 1 are not enough to discuss this problem. We require:

Assumption 9.9.

- $(\Xi, \mathcal{F}, \mathbb{P})$ is a complete probability space.
- U is a Hilbert space.
- $U_{\text{ad}} \subset U$ is nonempty, closed and convex.
- $J_2 : U \rightarrow \mathbb{R}$ is twice continuously differentiable.
- $\hat{J}_1 : U \times \Xi \rightarrow \mathbb{R}$ is such that $\hat{J}_1(u, \cdot) =: \hat{J}_1[\cdot](u) \in L_{\mathbb{P}}^{p_{\text{fun}}}(\Xi)$ for all $u \in U$ with some $p_{\text{fun}} \in [1, \infty]$.
- The function $\hat{J}_1[\cdot] : U \rightarrow L_{\mathbb{P}}^{p_{\text{fun}}}(\Xi)$ is continuously differentiable. Its first derivative $\hat{J}_1[\cdot]'(u) \in \mathcal{L}(U, L_{\mathbb{P}}^{p_{\text{fun}}}(\Xi))$ can be identified with the partial gradient $\nabla_u \hat{J}_1(u, \cdot) \in L_{\mathbb{P}}^{p_{\text{fun}}}(\Xi; U)$ via

$$\hat{J}_1[\cdot]'(u)w = (\nabla_u \hat{J}_1(u, \cdot), w)_U.$$

Here, the U -inner product is taken separately for almost every $\xi \in \Xi$.

- The function $\hat{J}_1[\cdot] : U \rightarrow L_{\mathbb{P}}^{p_{\text{fun}}}(\Xi)$ is twice continuously differentiable. Its second derivative $\hat{J}_1[\cdot]''(u) \in \mathcal{L}(U, \mathcal{L}(U, L_{\mathbb{P}}^{p_{\text{fun}}}(\Xi)))$ can be identified with the partial Hessian $\nabla_{uu}^2 \hat{J}_1(u, \cdot) \in L_{\mathbb{P}}^{p_{\text{fun}}}(\Xi; \mathcal{L}(U, U))$ via

$$[\hat{J}_1[\cdot]''(u)s]w = (\nabla_{uu}^2 \hat{J}_1(u, \cdot)s, w)_U.$$

In particular, we need $p_{\text{fun}} \in [2, \infty]$ for the mean-variance risk measure to be finite. Assumption 9.9 allows to apply the chain rule to compute the derivative of the objective function \hat{J} . We note that due to

$$\mathbb{E} : L_{\mathbb{P}}^{p_{\text{fun}}}(\Xi) \rightarrow \mathbb{R}, \quad \mathbb{E}[X] = \int_{\Xi} \mathbb{1} \cdot X \, d\mathbb{P}$$

we can identify the derivative $\mathbb{E}'[X] = \mathbb{1} \in L_{\mathbb{P}}^{p_{\text{fun}}^*}(\Xi)$ for all $X \in L_{\mathbb{P}}^{p_{\text{fun}}}(\Xi)$, where $p_{\text{fun}} \in [1, \infty]$ and $\frac{1}{p_{\text{fun}}} + \frac{1}{p_{\text{fun}}^*} = 1$. Furthermore, it holds that

$$\text{Var} : L_{\mathbb{P}}^{p_{\text{fun}}}(\Xi) \rightarrow \mathbb{R}, \quad \text{Var}[X] = \|X - \mathbb{E}[X] \cdot \mathbb{1}\|_{L_{\mathbb{P}}^2(\Xi)}^2 = \|X\|_{L_{\mathbb{P}}^2(\Xi)}^2 - \mathbb{E}[X]^2$$

and therefore $\text{Var}'[X] = 2X - 2\mathbb{E}[X] \cdot \mathbb{1} \in L_{\mathbb{P}}^{p_{\text{fun}}^*}(\Xi)$ with $p_{\text{fun}} \in [2, \infty]$. Note that we embed $L_{\mathbb{P}}^{p_{\text{fun}}}(\Xi) \hookrightarrow L_{\mathbb{P}}^2(\Xi)$ and use also the adjoint embedding $L_{\mathbb{P}}^2(\Xi) \hookrightarrow L_{\mathbb{P}}^{p_{\text{fun}}^*}(\Xi)$ to compute this Fréchet derivative involving the derivative of the squared $L_{\mathbb{P}}^2(\Xi)$ -norm. In general, $L_{\mathbb{P}}^{\infty}(\Xi)^*$ contains also functionals which cannot be identified with $L_{\mathbb{P}}^1(\Xi)$ functions, but all respective objects are indeed in $L_{\mathbb{P}}^1(\Xi)$ for $p_{\text{fun}} = \infty$. The second derivative of Var is $\text{Var}''[X]Y = 2Y - 2\mathbb{E}[Y] \cdot \mathbb{1} \in L_{\mathbb{P}}^{p_{\text{fun}}^*}(\Xi)$ for $Y \in L_{\mathbb{P}}^{p_{\text{fun}}}(\Xi)$.

By the chain rule, we obtain the gradient of $\hat{\mathbf{J}}$ as

$$\begin{aligned}\nabla \hat{\mathbf{J}}(u) &= \mathbb{E}[\nabla_u \hat{J}_1(u, \cdot)] + 2\lambda \int_{\Xi} (\hat{J}_1(u, \cdot) - \mathbb{E}[\hat{J}_1(u, \cdot)] \cdot \mathbb{1}) \nabla_u \hat{J}_1(u, \cdot) \, d\mathbb{P} + \nabla J_2(u) \\ &= (1 - 2\lambda \mathbb{E}[\hat{J}_1(u, \cdot)]) \mathbb{E}[\nabla_u \hat{J}_1(u, \cdot)] + 2\lambda \int_{\Xi} \hat{J}_1(u, \cdot) \nabla_u \hat{J}_1(u, \cdot) \, d\mathbb{P} + \nabla J_2(u) \quad (9.2) \\ &= \mathbb{E}[\nabla_u \hat{J}_1(u, \cdot)] + 2\lambda \text{Cov}[\hat{J}_1(u, \cdot), \nabla_u \hat{J}_1(u, \cdot)] + \nabla J_2(u)\end{aligned}$$

and the Hessian

$$\begin{aligned}\nabla^2 \hat{\mathbf{J}}(u)s &= \mathbb{E}[\nabla_{uu}^2 \hat{J}_1(u, \cdot)s] + \lambda \langle \text{Var}''[\hat{J}_1(u, \cdot)](\nabla_u J_1(u, \cdot), s)_U, \nabla_u J_1(u, \cdot) \rangle_{L_{\mathbb{P}}^{p_{\text{fun}}^*}(\Xi), L_{\mathbb{P}}^{p_{\text{fun}}}(\Xi)} \\ &\quad + \lambda \langle \text{Var}'[\hat{J}_1(u, \cdot)], \nabla_{uu}^2 J_1(u, \cdot)s \rangle_{L_{\mathbb{P}}^{p_{\text{fun}}^*}(\Xi), L_{\mathbb{P}}^{p_{\text{fun}}}(\Xi)} + \nabla^2 J_2(u)s \\ &= \mathbb{E}[\nabla_{uu}^2 \hat{J}_1(u, \cdot)s] + 2\lambda \int_{\Xi} ((\nabla_u J_1(u, \cdot), s)_U - \mathbb{E}[(\nabla_u J_1(u, \cdot), s)_U] \cdot \mathbb{1}) \nabla_u \hat{J}_1(u, \cdot) \, d\mathbb{P} \\ &\quad + 2\lambda \int_{\Xi} (\hat{J}_1(u, \cdot) - \mathbb{E}[\hat{J}_1(u, \cdot)] \cdot \mathbb{1}) \nabla_{uu}^2 \hat{J}_1(u, \cdot)s \, d\mathbb{P} + \nabla^2 J_2(u)s \\ &= \mathbb{E}[\nabla_{uu}^2 \hat{J}_1(u, \cdot)s] + 2\lambda \text{Cov}[(\nabla_u J_1(u, \cdot), s)_U, \nabla_u \hat{J}_1(u, \cdot)] \\ &\quad + 2\lambda \text{Cov}[(\hat{J}_1(u, \cdot), \nabla_{uu}^2 \hat{J}_1(u, \cdot)s) + \nabla^2 J_2(u)s.\end{aligned}$$

The advantage of these derivative formulations is that they use the partial derivatives w. r. t. u of the parameter-dependent reduced objective function. Often, the computation of these derivatives is already implemented by the adjoint approach and can be reused for an implementation of the mean-variance risk measure. Overall, we do not have to define an alternative adjoint state or adjoint equation for the objective function involving the variance term.

9.2.1. Discussion of the Example

We discuss under which conditions Assumption 9.9 is satisfied for the example from Chapter 3, i. e., we have

$$\hat{J}[\xi](u) = J_1[\xi](S[\xi](u), u, \xi) + J_2(u) = \frac{1}{2} \|Q(\xi)S[\xi](u) - \hat{q}(\xi)\|_H^2 + \frac{\gamma}{2} \|u\|_U^2,$$

see (3.5). According to (3.17), the derivatives are given by

$$\hat{J}_1[\xi]'(u) = -B(\xi)^* T[\xi](u) \quad \text{and} \quad \nabla J_2(u) = \gamma u,$$

where $\hat{J}_1[\xi](u) = J_1[\xi](S[\xi](u), u, \xi)$. Due to Corollary 3.15 it holds that $S[\cdot](u) \in L^{r_f}(\Xi; Y)$ and Lemma 3.23 yields $T[\cdot](u) \in L_{\mathbb{P}}^{r_z}(\Xi; Y)$ with r_z depending on the integrability exponents r_f , r_Q , and $r_{\hat{q}}$ of the state, the state-to-observation map $Q(\cdot)$, and the desired state $\hat{q}(\cdot)$, respectively. Especially, the functions $\xi \mapsto \hat{J}_1[\xi](u)$ and $\xi \mapsto \hat{J}_1[\xi]'(u)$ are measurable and one obtains the integrability exponents

$$\hat{J}_1[\cdot](u) \in L_{\mathbb{P}}^{\hat{r}/2}(\Xi) \quad \text{and} \quad \hat{J}_1[\cdot]'(u) \in L_{\mathbb{P}}^{r_z}(\Xi; U^*)$$

with \hat{r} from Table 3.1 and because $B(\cdot) \in L^\infty(\Xi; \mathcal{L}(U, Y^*))$. If $\hat{r} \geq 2p_{\text{fun}}$ and $r_z \geq p_{\text{fun}}$, the objective function and the gradient are both p_{fun} -integrable with $p_{\text{fun}} \geq 2$. In the ‘‘Example’’ column of Table 3.1 this would correspond to having $r_f \geq \max\{p, 4\}$, $r_Q = \infty$, and $r_{\hat{q}} \geq r_f$.

The differentiability assumption on $\hat{J}_1[\cdot]$ can be proven using a priori estimates on the adjoint state and the same techniques as in the proof of Theorem 3.26. For this purpose, let the regularity exponents $r_Q, r_f, r_{\hat{q}} \in [1, \infty]$ be such that $\frac{2}{r_Q} + \frac{1}{r_f} \leq \frac{1}{p_{\text{fun}}}$ and $\frac{1}{r_Q} + \frac{1}{r_{\hat{q}}} \leq \frac{1}{p_{\text{fun}}}$ hold. One computes for $w \in U \setminus \{0\}$:

$$\begin{aligned} 0 &\leq \lim_{\|w\|_U \rightarrow 0} \frac{\|\hat{J}_1(u+w, \cdot) - \hat{J}_1(u, \cdot) - (\nabla_u \hat{J}_1(u, \cdot), w)_U\|_{L^\infty(\Xi)}}{\|w\|_U} \\ &= \lim_{\|w\|_U \rightarrow 0} \frac{\|\int_0^1 (\nabla_u \hat{J}_1(u + \tau w, \cdot) - \nabla_u \hat{J}_1(u, \cdot), w)_U d\tau\|_{L^\infty(\Xi)}}{\|w\|_U} \\ &\leq \lim_{\|w\|_U \rightarrow 0} \left\| \int_0^1 \|\nabla_u \hat{J}_1(u + \tau w, \cdot) - \nabla_u \hat{J}_1(u, \cdot)\|_U d\tau \right\|_{L^\infty(\Xi)} \\ &\leq \lim_{\|w\|_U \rightarrow 0} \left\| \sup_{\tau \in [0,1]} \|\nabla_u \hat{J}_1(u + \tau w, \cdot) - \nabla_u \hat{J}_1(u, \cdot)\|_U \right\|_{L^\infty(\Xi)} = 0. \end{aligned}$$

In the following we argue why this limit is indeed zero.

For $p_{\text{fun}} \in [1, \infty)$, this follows from the dominated convergence theorem because

$$\sup_{\tau \in [0,1]} \|\nabla_u \hat{J}_1(u + \tau w, \xi) - \nabla_u \hat{J}_1(u, \xi)\|_U \longrightarrow 0 \quad \text{as} \quad \|w\|_U \rightarrow 0$$

for almost every $\xi \in \Xi$ due to the continuity of $u \mapsto \nabla_u \hat{J}_1(u, \xi)$. Additionally, we have

$$\begin{aligned} &\sup_{\tau \in [0,1]} \|\nabla_u \hat{J}_1(u + \tau w, \xi) - \nabla_u \hat{J}_1(u, \xi)\|_U = \sup_{\tau \in [0,1]} \| -B(\xi)^*(T[\xi](u + \tau w) - T[\xi](u)) \|_{U^*} \\ &\leq \|B(\xi)\|_{\mathcal{L}(U, Y^*)} \left(\sup_{\tau \in [0,1]} \|T[\xi](u + \tau w)\|_Y + \|T[\xi](u)\|_Y \right) \\ &\leq \frac{C_\Omega}{\kappa^2} \|B(\xi)\|_{\mathcal{L}(U, Y^*)} \|Q(\xi)\|_{\mathcal{L}(Y, H)}^2 (\|D\|_{\mathcal{L}(U, L^2(\Omega))} (2\|u\|_U + 1) + \|f(\xi) - \varphi(0)\|_{L^2(\Omega)}) \\ &\quad + \frac{1}{\kappa} \|B(\xi)\|_{\mathcal{L}(U, Y^*)} \|Q(\xi)\|_{\mathcal{L}(Y, H)} \|\hat{q}(\xi)\|_H \end{aligned}$$

for $\|w\|_U \leq 1$. The upper bound is due to the a priori estimate

$$\begin{aligned} \|T[\xi](u)\|_Y &\leq \frac{1}{\kappa} \|Q(\xi)\|_{\mathcal{L}(Y, H)} \|(Q(\xi)S[\xi](u) - \hat{q}(\xi))\|_H \\ &\leq \frac{1}{\kappa} \|Q(\xi)\|_{\mathcal{L}(Y, H)} (\|Q(\xi)\|_{\mathcal{L}(Y, H)} \|S[\xi](u)\|_Y + \|\hat{q}(\xi)\|_H) \\ &\leq \frac{C_\Omega}{\kappa^2} \|Q(\xi)\|_{\mathcal{L}(Y, H)}^2 (\|D\|_{\mathcal{L}(U, L^2(\Omega))} \|u\|_U + \|f(\xi) - \varphi(0)\|_{L^2(\Omega)}) \\ &\quad + \frac{1}{\kappa} \|Q(\xi)\|_{\mathcal{L}(Y, H)} \|\hat{q}(\xi)\|_H \end{aligned}$$

on the adjoint state $T[\xi](u)$, which is obtained by combining (3.20) and (3.9). It is an $L^\infty(\Xi)$ -function w. r. t. ξ because $\frac{2}{r_Q} + \frac{1}{r_f} \leq \frac{1}{p_{\text{fun}}}$, $\frac{1}{r_Q} + \frac{1}{r_{\hat{q}}} \leq \frac{1}{p_{\text{fun}}}$, and $B(\cdot)$ is essentially bounded w. r. t. ξ .

For $p_{\text{fun}} = \infty$, the uniform Lipschitz continuity of the gradient is used to show the result. In analogy to (3.9), one can estimate

$$\|S[\xi](u+w) - S[\xi](u)\|_Y \leq \frac{C_\Omega}{\kappa} \|Dw\|_{L^2(\Omega)}$$

independently of ξ . Using the derivations in the proofs of Theorem 5.1 and Lemma 5.2 with $y = S[\xi](u)$, $\tilde{y} = S[\xi](u+w)$, $z = T[\xi](u)$, $\hat{z} = T[\xi](u+w)$ as well as (3.9), this yields

$$\begin{aligned} & \|T[\xi](u+w) - T[\xi](u)\|_Y \\ & \leq \frac{1}{\kappa} \|Q(\xi)\|_{\mathcal{L}(Y,H)}^2 \|S[\xi](u) - S[\xi](u+w)\|_Y \\ & \quad + \frac{1}{\kappa^2} \|N'(S[\xi](u+w)) - N'(S[\xi](u))\|_{\mathcal{L}(Y,Y^*)} \|Q(\xi)^*(Q(\xi)S[\xi](u) - \hat{q}(\xi))\|_{Y^*} \\ & \leq \frac{C_\Omega}{\kappa^2} \|Q(\xi)\|_{\mathcal{L}(Y,H)}^2 \|Dw\|_{L^2(\Omega)} \\ & \quad + \frac{1}{\kappa^2} c_p^3 \left(a''_{\varphi''} \cdot \lambda(\Omega)^{(p-3)/p} + c''_{\varphi''} c_p^{p-3} \cdot \left(\frac{C_\Omega}{\kappa} (\|Du + f(\xi) - \varphi(0)\|_{L^2(\Omega)} + \|Dw\|_{L^2(\Omega)}) \right)^{p-3} \right) \\ & \quad \cdot \frac{C_\Omega}{\kappa} \|Dw\|_{L^2(\Omega)} \|Q(\xi)\|_{\mathcal{L}(Y,H)} (\|Q(\xi)\|_{\mathcal{L}(Y,H)} \|Du + f(\xi) - \varphi(0)\|_{L^2(\Omega)} + \|\hat{q}(\xi)\|_H), \end{aligned}$$

cf. (3.20). Since $r_Q = r_{\hat{q}} = r_f = \infty$, this bound converges to zero in $L_{\mathbb{P}}^\infty(\Xi)$ as $\|w\|_U \rightarrow 0$. Hence,

$$\begin{aligned} 0 & \leq \left\| \sup_{\tau \in [0,1]} \|\nabla_u \hat{J}_1(u + \tau w, \cdot) - \nabla_u \hat{J}_1(u, \cdot)\|_U \right\|_{L_{\mathbb{P}}^\infty(\Xi)} \\ & \leq \left\| \|B(\cdot)\|_{\mathcal{L}(U,Y^*)} \right\|_{L_{\mathbb{P}}^\infty(\Xi)} \left\| \sup_{\tau \in [0,1]} \|T[\xi](u + \tau w) - T[\xi](u)\|_Y \right\|_{L_{\mathbb{P}}^\infty(\Xi)} \longrightarrow 0 \quad \text{as } \|w\|_U \rightarrow 0 \end{aligned}$$

can be concluded since $B(\cdot)$ is essentially bounded w. r. t. ξ .

The considerations show the Fréchet approximation condition for the first derivative. For the second derivative, it can be proven analogously using the a priori bound (3.26). Continuity of the second derivative can also be shown with similar arguments, i. e., using pointwise convergence and the dominated convergence theorem for $p_{\text{fun}} \in [1, \infty)$ or refined bounds in the case $p_{\text{fun}} = \infty$.

9.2.2. Possible Error Estimates for Reduced Objective Function and Gradient

In the following, we investigate how the required error estimates from Section 5.2 change in the mean-variance case.

Objective Function Evaluation Error

We assume that the function $\hat{J}_1(u, \cdot)$ in (9.1) is evaluated inexactly due to an inexact state computation, giving the inexact version $\tilde{J}_1(u, \cdot)$, whereas J_2 is evaluated exactly. Altogether, we obtain the inexact version $\tilde{\mathbf{J}}$ of the objective function $\hat{\mathbf{J}}$. In the proof of Proposition 5.11, the error in $\tilde{\mathbf{J}}$ is derived from the error in \tilde{J}_1 by means of

$$|\hat{\mathbf{J}}(u) - \tilde{\mathbf{J}}(u)| = |\mathbb{E}[\hat{J}_1(u, \cdot) - \tilde{J}_1(u, \cdot)]| \leq \|\hat{J}_1(u, \cdot) - \tilde{J}_1(u, \cdot)\|_{L_{\mathbb{P}}^1(\Xi)},$$

which leads to an error bound based on the $L_{\mathbb{P}}^2(\Xi)$ -error $\|\mathbf{Q}(\tilde{\mathbf{y}} - \mathbf{y})\|_{L_{\mathbb{P}}^2(\Xi;H)}$ in the inexact state. In the mean-variance case, we have

$$\begin{aligned}
 |\hat{\mathbf{J}}(u) - \tilde{\mathbf{J}}(u)| &= \left| \mathbb{E}[\hat{\mathbf{J}}_1(u, \cdot) - \tilde{\mathbf{J}}_1(u, \cdot)] + \lambda(\text{Var}[\hat{\mathbf{J}}_1(u, \cdot)] - \text{Var}[\tilde{\mathbf{J}}_1(u, \cdot)]) \right| \\
 &= \left| \mathbb{E}[\hat{\mathbf{J}}_1(u, \cdot) - \tilde{\mathbf{J}}_1(u, \cdot)] + \lambda \mathbb{E}[(\hat{\mathbf{J}}_1(u, \cdot) - \tilde{\mathbf{J}}_1(u, \cdot))(\hat{\mathbf{J}}_1(u, \cdot) + \tilde{\mathbf{J}}_1(u, \cdot))] \right. \\
 &\quad \left. + \lambda \mathbb{E}[\tilde{\mathbf{J}}_1(u, \cdot) - \hat{\mathbf{J}}_1(u, \cdot)] \mathbb{E}[\hat{\mathbf{J}}_1(u, \cdot) + \tilde{\mathbf{J}}_1(u, \cdot)] \right| \\
 &= \left| \mathbb{E}[\mathbb{1} - \lambda \hat{\mathbf{J}}_1(u, \cdot) - \lambda \tilde{\mathbf{J}}_1(u, \cdot)] \mathbb{E}[\hat{\mathbf{J}}_1(u, \cdot) - \tilde{\mathbf{J}}_1(u, \cdot)] \right. \\
 &\quad \left. + \lambda \mathbb{E}[(\hat{\mathbf{J}}_1(u, \cdot) - \tilde{\mathbf{J}}_1(u, \cdot))(\hat{\mathbf{J}}_1(u, \cdot) + \tilde{\mathbf{J}}_1(u, \cdot))] \right| \\
 &\leq \|\mathbb{1} - \lambda \hat{\mathbf{J}}_1(u, \cdot) - \lambda \tilde{\mathbf{J}}_1(u, \cdot)\|_{L_{\mathbb{P}}^1(\Xi)} \|\hat{\mathbf{J}}_1(u, \cdot) - \tilde{\mathbf{J}}_1(u, \cdot)\|_{L_{\mathbb{P}}^1(\Xi)} \\
 &\quad + \lambda \|\hat{\mathbf{J}}_1(u, \cdot) - \tilde{\mathbf{J}}_1(u, \cdot)\|_{L_{\mathbb{P}}^1(\Xi)} \|\hat{\mathbf{J}}_1(u, \cdot) + \tilde{\mathbf{J}}_1(u, \cdot)\|_{L_{\mathbb{P}}^\infty(\Xi)}.
 \end{aligned}$$

The estimate requires $p_{\text{fun}} = \infty$ for both the function $\hat{\mathbf{J}}_1$ and its inexact version $\tilde{\mathbf{J}}_1$. In Table 3.1, we would therefore need $r_f = \infty$, $r_Q = \infty$ and $r_{\hat{q}} = \infty$ to be able to apply it. But this is the only way to obtain an estimate which depends on the $L_{\mathbb{P}}^1(\Xi)$ -error in $\hat{\mathbf{J}}_1$ so that an $L_{\mathbb{P}}^2(\Xi)$ -error estimate for the inexact state is sufficient. For the more natural choice that $\|\hat{\mathbf{J}}_1(u, \cdot) - \tilde{\mathbf{J}}_1(u, \cdot)\|_{L_{\mathbb{P}}^2(\Xi)}$ has to be controlled, the $L_{\mathbb{P}}^4(\Xi)$ -error in the state would have to be estimated, see the proof of Proposition 5.11.

Model gradient error

The gradient (9.2) is computed inexactly due to an inexact objective function evaluation $\tilde{\mathbf{J}}(u, \cdot)$ and an inexact gradient, denoted by $\tilde{\nabla}_u \hat{\mathbf{J}}(u, \cdot) \approx \nabla_u \hat{\mathbf{J}}(u, \cdot)$. This yields the inexact gradient $\tilde{\nabla} \hat{\mathbf{J}}(u)$, the error in which shall be estimated in the following.

$$\begin{aligned}
 &\|\nabla \hat{\mathbf{J}}(u) - \tilde{\nabla} \hat{\mathbf{J}}(u)\|_U \\
 &= \left\| \mathbb{E}[\nabla_u \hat{\mathbf{J}}_1(u, \cdot)] + 2\lambda \int_{\Xi} (\hat{\mathbf{J}}_1(u, \cdot) - \mathbb{E}[\hat{\mathbf{J}}_1(u, \cdot)] \cdot \mathbb{1}) \nabla_u \hat{\mathbf{J}}_1(u, \cdot) \, d\mathbb{P} \right. \\
 &\quad \left. - \mathbb{E}[\nabla_u \tilde{\mathbf{J}}_1(u, \cdot)] - 2\lambda \int_{\Xi} (\tilde{\mathbf{J}}_1(u, \cdot) - \mathbb{E}[\tilde{\mathbf{J}}_1(u, \cdot)] \cdot \mathbb{1}) \tilde{\nabla}_u \hat{\mathbf{J}}_1(u, \cdot) \, d\mathbb{P} \right\|_U \\
 &\leq \mathbb{E}[\|\nabla_u \hat{\mathbf{J}}_1(u, \cdot) - \nabla_u \tilde{\mathbf{J}}_1(u, \cdot)\|_U] \\
 &\quad + 2\lambda \left\| \int_{\Xi} (\hat{\mathbf{J}}_1(u, \cdot) - \mathbb{E}[\hat{\mathbf{J}}_1(u, \cdot)] \cdot \mathbb{1}) \nabla_u \hat{\mathbf{J}}_1(u, \cdot) - (\tilde{\mathbf{J}}_1(u, \cdot) - \mathbb{E}[\tilde{\mathbf{J}}_1(u, \cdot)] \cdot \mathbb{1}) \tilde{\nabla}_u \hat{\mathbf{J}}_1(u, \cdot) \, d\mathbb{P} \right\|_U \\
 &\leq \|\nabla_u \hat{\mathbf{J}}_1(u, \cdot) - \nabla_u \tilde{\mathbf{J}}_1(u, \cdot)\|_{L_{\mathbb{P}}^1(\Xi;U)} \\
 &\quad + 2\lambda \|\hat{\mathbf{J}}_1(u, \cdot) - \tilde{\mathbf{J}}_1(u, \cdot) + \mathbb{E}[\tilde{\mathbf{J}}_1(u, \cdot) - \hat{\mathbf{J}}_1(u, \cdot)] \cdot \mathbb{1}\|_{L_{\mathbb{P}}^{\tilde{r}}(\Xi)} \|\nabla_u \hat{\mathbf{J}}_1(u, \cdot) - \tilde{\nabla}_u \hat{\mathbf{J}}_1(u, \cdot)\|_{L_{\mathbb{P}}^{\tilde{r}^*}(\Xi;U)} \\
 &\quad + 2\lambda \|\tilde{\mathbf{J}}_1(u, \cdot) - \mathbb{E}[\tilde{\mathbf{J}}_1(u, \cdot)] \cdot \mathbb{1}\|_{L_{\mathbb{P}}^{\tilde{r}}(\Xi)} \|\nabla_u \hat{\mathbf{J}}_1(u, \cdot) - \tilde{\nabla}_u \hat{\mathbf{J}}_1(u, \cdot)\|_{L_{\mathbb{P}}^{\tilde{r}^*}(\Xi;U)} \\
 &\quad + 2\lambda \|\hat{\mathbf{J}}_1(u, \cdot) - \tilde{\mathbf{J}}_1(u, \cdot) + \mathbb{E}[\tilde{\mathbf{J}}_1(u, \cdot) - \hat{\mathbf{J}}_1(u, \cdot)] \cdot \mathbb{1}\|_{L_{\mathbb{P}}^{\tilde{r}}(\Xi)} \|\tilde{\nabla}_u \hat{\mathbf{J}}_1(u, \cdot)\|_{L_{\mathbb{P}}^{\tilde{r}^*}(\Xi;U)}
 \end{aligned}$$

with $\tilde{r}, \tilde{r}^* \in [1, \infty]$, $\frac{1}{\tilde{r}} + \frac{1}{\tilde{r}^*} = 1$ and provided the respective quantities enjoy the required integrability properties. Controlling the objective function evaluation in the $L_{\mathbb{P}}^{\tilde{r}}(\Xi)$ -norm requires to estimate the $L_{\mathbb{P}}^{2\tilde{r}}(\Xi)$ -error in the state. Even if the state equation is linear, i. e., φ'

is constant, the $L_{\mathbb{P}}^{\tilde{r}^*}(\Xi)$ -error in the state and adjoint state is needed to estimate the gradient error in the $L_{\mathbb{P}}^{\tilde{r}^*}(\Xi; U)$ -norm, cf. Theorem 5.8. The “best” choice of \tilde{r} , which demands the weakest integrability properties of the state, is therefore $\tilde{r} = \frac{3}{2}$, which yields that the $L_{\mathbb{P}}^3(\Xi)$ -error in the state has to be controlled because $2\tilde{r} = 3 = \tilde{r}^*$.

Overall, we see that the risk-neutral implementation can be extended to the mean-variance case by only changing the evaluation of the reduced objective function, its gradient and its Hessian based on the already derived (linearized) state and adjoint state computation. In addition, the error estimation would have to be changed accordingly. Here, the issue arises that it becomes necessary to control the error in the state at least in the $L_{\mathbb{P}}^3(\Xi)$ -norm even for $r_f = \infty$ and a linear state equation. This is not a simple task and would require a new error estimation procedure. We leave this for future research. Additionally, we want to mention again that the mean-variance risk measure is *not* monotonic. This means that we could have a control for which the tracking term is larger than the one for a different control almost surely although the corresponding risk measured by $\mathbb{E} + \lambda \text{Var}$ is smaller. We therefore refrain from presenting a numerical implementation and results here and consider a class of risk measures with better theoretical properties.

9.3. Convex Combination of Mean and Conditional Value-at-Risk

In this section, we introduce a class of risk measures, which are convex combinations of the expected value and the conditional value-at-risk CVaR_{β} , and derive the respective properties for them. These risk measures are nonsmooth, but minimizing them can be reformulated introducing a pointwisely constrained auxiliary variable. We propose solving this problem by an interior-point method. The introduced barrier term leads to a smoothed version of the original risk measure. It is investigated which properties of it are retained by the smoothing procedure. Optimal control problems with this smoothed risk measure are investigated in theory and numerical results of an implementation with low-rank tensors are presented.

9.3.1. Definition and Derivation of the Properties

In [95] it is shown that the conditional value-at-risk of an $L_{\mathbb{P}}^1(\Xi)$ random variable X can be computed as follows:

$$\text{CVaR}_{\beta}[X] = \inf_{t \in \mathbb{R}} t + \frac{1}{1-\beta} \mathbb{E}[(X - t)^+]$$

with $(s)^+ := \max\{0, s\}$. This gives rise to considering risk measures of the form

$$\mathcal{R}_v : L_{\mathbb{P}}^p(\Xi) \rightarrow \mathbb{R} \cup \{\infty\}, \quad \mathcal{R}_v[X] = \inf_{t \in \mathbb{R}} t + \mathbb{E}[v(X - t)], \quad (9.3)$$

where $v : \mathbb{R} \rightarrow \mathbb{R}$ is a given function. Risk measures of this kind correspond to the expectation quadrangle presented in [96]. In the following, we say that the risk measure \mathcal{R}_v is *induced* by the function v . Two examples are the expectation \mathbb{E} for $v(s) = s$ and the conditional value-at-risk CVaR_{β} for $v(s) = \frac{1}{1-\beta}(s)^+$.

The risk measure \mathcal{R}_v inherits several properties from the function v as shown in the following.

Lemma 9.10. *The risk measure \mathcal{R}_v is proper, i. e., $\mathcal{R}_v[X] > -\infty$ for all X and $\mathcal{R}_v[\tilde{X}] < \infty$ for some \tilde{X} if and only if there exists $c \in \mathbb{R}$ such that $v(s) \geq s + c$ holds for all $s \in \mathbb{R}$.*

Proof. For the “if” part observe that from $v(s) \geq s + c$ it follows that $t + v(X - t) \geq t + (X - t + c) = X + c$. Therefore, $\mathcal{R}_v[X] = \inf_{t \in \mathbb{R}} \mathbb{E}[t + v(X - t)] \geq \inf_{t \in \mathbb{R}} \mathbb{E}[X + c] = \mathbb{E}[X] + c > -\infty$ holds because the expectation is finite on $L_{\mathbb{P}}^p(\Xi)$ for all $p \in [1, \infty]$. Furthermore, for any constant random variable $\tilde{X} \equiv d \in \mathbb{R}$ and some given $\hat{t} \in \mathbb{R}$ we have that $\mathcal{R}_v[\tilde{X}] = \inf_{t \in \mathbb{R}} t + \mathbb{E}[v(\tilde{X} - t)] = \inf_{t \in \mathbb{R}} t + v(d - t) \leq \hat{t} + v(d - \hat{t}) < \infty$.

Now we consider the “only if” part: Assume that there does not exist any $c \in \mathbb{R}$ such that $v(s) \geq s + c$ holds for all $s \in \mathbb{R}$, i. e., for every $c \in \mathbb{R}$ there exists $\tilde{s} \in \mathbb{R}$ with $v(\tilde{s}) < \tilde{s} + c$. Now let $(c_k) \subset \mathbb{R}$ be a sequence diverging to $-\infty$, e. g., $c_k = -k$, and let $X \equiv d \in \mathbb{R}$. Then there exists $(t_k) \subset \mathbb{R}$ such that $t_k + v(X - t_k) < d + c_k$ holds for all $k \in \mathbb{N}$. It follows that $\inf_{t \in \mathbb{R}} t + \mathbb{E}[v(X - t)] = -\infty$ for all constant X . \square

Corollary 9.11. *If v is convex and there exists $\tilde{s} \in \mathbb{R}$ such that $1 \in \partial v(\tilde{s})$ (convex subdifferential), the risk measure \mathcal{R}_v is proper.*

Proof. Lemma 9.10 can be applied because $v(s) \geq v(\tilde{s}) + 1 \cdot (s - \tilde{s}) = s + c$ holds with $c = v(\tilde{s}) - \tilde{s}$. \square

Lemma 9.12. *If there exist $c, d \in \mathbb{R}$ such that $|v(s)| \leq c|s|^p + d$ holds for all $s \in \mathbb{R}$, we have $\mathcal{R}_v[X] < \infty$ for all $X \in L_{\mathbb{P}}^p(\Xi)$.*

Proof. We have $\mathbb{E}[v(X)] \leq \mathbb{E}[|v(X)|] \leq \mathbb{E}[c|X|^p] + d = c\|X\|_{L_{\mathbb{P}}^p(\Xi)}^p + d < \infty$. Therefore, the argument in the infimum defining \mathcal{R}_v is finite for $t = 0$. \square

Lemma 9.13. *The risk measure \mathcal{R}_v is always translation-equivariant.*

Proof. The statement follows from

$$\mathcal{R}_v[X + c] = \inf_{t \in \mathbb{R}} t + \mathbb{E}[v(X + c - t)] \stackrel{\tilde{t} = t - c}{=} \inf_{\tilde{t} \in \mathbb{R}} \tilde{t} + c + \mathbb{E}[v(X - \tilde{t})] = \mathcal{R}_v[X] + c.$$

\square

Lemma 9.14. *If the function v is convex, the risk measure \mathcal{R}_v is convex.*

For the proof of this statement, we use the following general result:

Proposition 9.15. *Let U and W be convex subsets of \mathbb{R} -vector spaces and let $f : U \times W \rightarrow (-\infty, \infty]$ be a convex function. Then the function $g : U \rightarrow [-\infty, \infty]$, $g(u) := \inf_{w \in W} f(u, w)$ is convex.*

Proof. This claim follows as in the proof of [17, Prop. 8.26]. \square

Proof of Lemma 9.14. The result follows by applying Proposition 9.15 with $W = \mathbb{R}$, $U = L_{\mathbb{P}}^p(\Xi)$ and $f(u, w) := w + \mathbb{E}[v(u - w)]$. \square

Lemma 9.16. *If v is positively homogeneous and $\inf_{t \in \mathbb{R}} v(t) - t = 0$, the risk measure \mathcal{R}_v is positively homogeneous.*

Proof. We use that $v(\lambda s) = \lambda v(s)$ holds for all $x \in \mathbb{R}$ and all $\lambda > 0$. Let $X \in L_{\mathbb{P}}^p(\Xi)$ be a random variable and $\lambda > 0$. Then:

$$\begin{aligned} \mathcal{R}_v[\lambda X] &= \inf_{\tau \in \mathbb{R}} \tau + \mathbb{E}[v(\lambda X - \tau)] \stackrel{v \text{ pos. hom.}}{=} \inf_{\tau \in \mathbb{R}} \tau + \mathbb{E}[\lambda v(X - \frac{\tau}{\lambda})] \\ &\stackrel{\tau \equiv \lambda t}{=} \inf_{t \in \mathbb{R}} \lambda t + \lambda \mathbb{E}[v(X - t)] \stackrel{\lambda \geq 0}{=} \lambda \mathcal{R}_v[X]. \end{aligned}$$

Furthermore, $\mathcal{R}_v[0] = \inf_{\tau \in \mathbb{R}} \tau + \mathbb{E}[v(0 - \tau)] \stackrel{t = -\tau}{=} \inf_{t \in \mathbb{R}} v(t) - t = 0$. □

Lemma 9.17. *If v is subadditive, the risk measure \mathcal{R}_v is subadditive.*

Proof. We use that $v(s + \tilde{s}) \leq v(s) + v(\tilde{s})$ holds for all $s, \tilde{s} \in \mathbb{R}$. Let $X, \tilde{X} \in L_{\mathbb{P}}^p(\Xi)$ be random variables. Then:

$$\begin{aligned} \mathcal{R}_v[X + \tilde{X}] &= \inf_{\tau \in \mathbb{R}} \tau + \mathbb{E}[v(X + \tilde{X} - \tau)] \stackrel{\tau = t + \tilde{t}}{=} \inf_{t, \tilde{t} \in \mathbb{R}} t + \tilde{t} + \mathbb{E}[v(X + \tilde{X} - t - \tilde{t})] \\ &\stackrel{v \text{ subadditive}}{\leq} \inf_{t, \tilde{t} \in \mathbb{R}} t + \tilde{t} + \mathbb{E}[v(X - t) + v(\tilde{X} - \tilde{t})] \\ &= \inf_{t \in \mathbb{R}} t + \mathbb{E}[v(X - t)] + \inf_{\tilde{t} \in \mathbb{R}} \tilde{t} + \mathbb{E}[v(\tilde{X} - \tilde{t})] = \mathcal{R}_v[X] + \mathcal{R}_v[\tilde{X}]. \end{aligned}$$
□

Lemma 9.18. *If v is increasing, the risk measure \mathcal{R}_v is monotonic.*

Proof. Let $X \leq Y$ a. s. Then, $v(X - t) \leq v(Y - t)$ a. s. for all $t \in \mathbb{R}$ because v is increasing. Using the monotonicity of the expectation it follows that

$$\mathcal{R}_v[X] = \inf_{t \in \mathbb{R}} t + \mathbb{E}[v(X - t)] \leq \inf_{t \in \mathbb{R}} t + \mathbb{E}[v(Y - t)] = \mathcal{R}_v[Y].$$
□

Lemma 9.19. *Let some constant $d \in \mathbb{R}$ be given. Then the function $\hat{v} : \mathbb{R} \rightarrow \mathbb{R}$, $\hat{v}(s) := v(s + d) - d$ induces the risk measure $\mathcal{R}_{\hat{v}}$, i. e., $\mathcal{R}_{\hat{v}} \equiv \mathcal{R}_v$.*

Proof. For any random variable X we have

$$\mathcal{R}_{\hat{v}}[X] = \inf_{t \in \mathbb{R}} t + \mathbb{E}[v(X - t + d) - d] \stackrel{\tilde{t} = t - d}{=} \inf_{\tilde{t} \in \mathbb{R}} \tilde{t} + \mathbb{E}[v(X - \tilde{t})] = \mathcal{R}_v[X].$$
□

Theorem 9.20. *If $v : \mathbb{R} \rightarrow \mathbb{R}$ is closed convex with $v(0) = 0$ and $v(s) > s$ for all $s \neq 0$, we obtain a regular risk measure \mathcal{R}_v by the construction (9.3). It is coherent if and only if it holds $v(s) = \max\{a_1 s, a_2 s\}$ with (w. l. o. g.) $a_1 \in [0, 1)$ and $a_2 \in (1, \infty)$, i. e., $v(s) = a_1 s$ for $s < 0$, $v(0) = 0$ and $v(s) = a_2 s$ for $s > 0$.*

Proof. This result is an excerpt of the ‘‘Expectation Theorem’’ in [96]. Note that the function v is convex, monotonically increasing, and positively homogeneous. □

It follows that CVaR_{β} with $\beta \in (0, 1)$ is a coherent and regular risk measure since it belongs to the mentioned class with $a_1 = 0$ and $a_2 = \frac{1}{1-\beta}$.

9.3.2. Smoothing by a Log-Barrier Approach

Motivated by Theorem 9.20, we consider regular, coherent risk measures by choosing $a_1 \in [0, 1)$ and $a_2 \in (1, \infty)$ and defining

$$\mathcal{R} : L_{\mathbb{P}}^p(\Xi) \rightarrow \mathbb{R}, \quad \mathcal{R}[X] := \inf_{t \in \mathbb{R}} t + \mathbb{E}[\max\{a_1(X - t), a_2(X - t)\}], \quad (9.4)$$

i. e., we have $\mathcal{R} \equiv \mathcal{R}_v$ with $v(s) = \max\{a_1 s, a_2 s\}$. In fact, as noted in [72, Sec. 2.4.1], this class of risk measures consists exactly of convex combinations of the expectation and the CVaR_{β} for every quantile parameter $\beta \in (0, 1)$ and every combination parameter $\lambda \in (0, 1]$:

$$\begin{aligned} \lambda \text{CVaR}_{\beta}[X] + (1 - \lambda) \mathbb{E}[X] &= \inf_{t \in \mathbb{R}} \lambda t + \lambda \mathbb{E}\left[\frac{1}{1-\beta}(X - t)^+\right] + (1 - \lambda) \mathbb{E}[X] \\ &= \inf_{t \in \mathbb{R}} t + \mathbb{E}\left[\frac{\lambda}{1-\beta}(X - t)^+ + (1 - \lambda)(X - t)\right] \\ &= \inf_{t \in \mathbb{R}} t + \mathbb{E}[\max\{(1 - \lambda)(X - t), (1 - \lambda + \frac{\lambda}{1-\beta})(X - t)\}], \end{aligned} \quad (9.5)$$

which fits into the setting (9.4) with $a_1 = 1 - \lambda \in [0, 1)$ and $a_2 = 1 + \frac{\beta}{1-\beta}\lambda \in (1, \infty)$. Conversely, $\lambda = 1 - a_1$ and $\beta = \frac{a_2 - 1}{a_2 - a_1}$ can be computed. Due to the max-term in (9.4) one has to solve a nonsmooth optimization problem over $t \in \mathbb{R}$ to compute the risk measure \mathcal{R} . Furthermore, \mathcal{R} itself is nonsmooth, see, e. g., [71]. Thus, we typically need a smoothing procedure that will allow us to develop derivative-based optimization methods. In [71] it is suggested to smooth CVaR using a suitable, smooth approximation of the $(\cdot)^+$ -function so that many properties of the risk measure, such as convexity and monotonicity, can be preserved. Another advantage of this approach is the improved accuracy of quadrature formulas for the evaluation of the expectation. The convergence of sparse grid quadrature, for example, depends strongly on the smoothness of the integrand, see, e. g., [86, 47].

Alternatively, we can reformulate problem (9.4) and solve

$$\inf_{W \in L_{\mathbb{P}}^p(\Xi), t \in \mathbb{R}} t + \mathbb{E}[W] \quad \text{s. t.} \quad W \geq a_1(X - t) \text{ a. s.}, \quad W \geq a_2(X - t) \text{ a. s.}, \quad (9.6)$$

which is a linear optimization problem. Solving such a problem also requires the discretization of the space of random variables $L_{\mathbb{P}}^p(\Xi)$ to be able to optimize over W .

In [46], a problem with similar constraints as in (9.6), namely a stochastic obstacle problem, was solved successfully using a stochastic Galerkin discretization with a low-rank tensor representation of the coefficients, cf. Chapter 6, and a primal interior-point method implemented with low-rank tensors. Thus, we want to follow a similar approach here and replace the inequality constraints in (9.6) by a log-barrier term with a barrier parameter $\mu > 0$, see also [85, Sec. 19.6]. We obtain

$$\inf_{W \in L_{\mathbb{P}}^p(\Xi), t \in \mathbb{R}} t + \mathbb{E}[W] - \mu \mathbb{E}[\ln(W - a_1(X - t))] - \mu \mathbb{E}[\ln(W - a_2(X - t))] + \zeta(\mu) \quad (9.7)$$

or equivalently $\inf_{W \in L_{\mathbb{P}}^p(\Xi), t \in \mathbb{R}} F_{\mu}(X, t, W)$ with

$$F_{\mu}(X, t, W) := \mathbb{E}[t + W - \mu \ln(W - a_1(X - t)) - \mu \ln(W - a_2(X - t))] + \zeta(\mu),$$

where $\zeta(\mu) := \mu(\ln(\frac{a_2-a_1}{a_2-1}\mu) + \ln(\frac{a_2-a_1}{1-a_1}\mu) - 2) \in \mathbb{R}$ is a constant shift. This shift is chosen such that problem (9.7) has exactly the minimal value $C \in \mathbb{R}$ whenever $X \equiv C$.

In the numerical experiments in [46] it was observed that the pointwise reciprocals of the functions $W - a_i(X - t)$, which appear in the barrier Newton system, cannot be represented sufficiently well by a low-rank tensor if the functions are too close to zero. Therefore, it would be beneficial to not push the barrier parameter μ arbitrarily close towards zero during the algorithm but to keep it on a fixed, positive level instead. Then it can be expected that the functions $W - a_i(X - t)$ are far enough away from zero during the iteration. We give a theoretical verification of this claim in (9.17).

Keeping the parameter μ bounded away from zero means that we solve a perturbed problem and obtain an approximation of the original risk measure. In the following, we want to investigate the properties of this approximate risk measure. Note that the goal of this procedure is not to construct new risk measures, but to provide an efficient solution algorithm based on low-rank tensor methods for, e.g., optimal control problems under uncertainty where risk-averse solutions are desired. Since this algorithm changes the properties of the underlying risk measure, which could have a great effect on the resulting solution, we analyze which properties are preserved by the log-barrier smoothing.

Proposition 9.21. *It holds that*

$$\inf_{W \in L_{\mathbb{P}}^p(\Xi), t \in \mathbb{R}} F_{\mu}(X, t, W) = \inf_{t \in \mathbb{R}} \mathbb{E}[\inf_{w \in \mathbb{R}} f_{\mu}(X(\cdot), t, w)]$$

with $f_{\mu} : \mathbb{R} \times \mathbb{R} \times \mathbb{R} \rightarrow (-\infty, +\infty]$,

$$f_{\mu}(x, t, w) = t + w - \mu \ln(w - a_1(x - t)) - \mu \ln(w - a_2(x - t)) + \zeta(\mu), \quad (9.8)$$

where we set $-\mu \ln(s) = +\infty$ for every $\mu > 0$ and every $s \leq 0$.

*Proof.*¹¹ Observe that

$$\inf_{W \in L_{\mathbb{P}}^p(\Xi), t \in \mathbb{R}} F_{\mu}(X, t, W) = \inf_{t \in \mathbb{R}} \inf_{W \in L_{\mathbb{P}}^p(\Xi)} \mathbb{E}[f_{\mu}(X(\cdot), t, W(\cdot))] \quad (9.9)$$

holds. The space $L_{\mathbb{P}}^p(\Xi)$ is decomposable in the sense of [97, Def. 14.59], the measure \mathbb{P} is σ -finite, and the integrand $\hat{f}_{\mu} : \Xi \times \mathbb{R} \rightarrow \bar{\mathbb{R}}$, $\hat{f}_{\mu}(\xi, w) := f_{\mu}(X(\xi), t, w)$ with fixed $t \in \mathbb{R}$ and fixed X is a normal integrand in the sense of [97, Def. 14.27], which we show in the following: Firstly, the function $(w, x) \mapsto f_{\mu}(x, t, w)$ is lower semicontinuous and convex for all $\xi \in \Xi$ because it is independent of ξ . Secondly, the interior of its domain is nonempty for all ξ . Thirdly, the function $\xi \mapsto f_{\mu}(x, t, w)$ is measurable for all $w \in \mathbb{R}$ since it is constant w. r. t. ξ . These properties yield that the function $(\xi, w, x) \mapsto f_{\mu}(x, t, w)$ is a normal integrand by [97, Prop. 14.39]. The composition rule [97, Prop. 14.45(c)] now shows that the function \hat{f}_{μ} is a normal integrand. Using this and having that there exists $\tilde{W} \in L_{\mathbb{P}}^p(\Xi)$ such that $\int_{\Xi} \hat{f}_{\mu}(\xi, \tilde{W}(\xi)) d\mathbb{P} < \infty$, we can apply the “interchangeability theorem” [97, Thm. 14.60] to

¹¹The essential content of this proof with the reference to [97] was provided by Prof. Dr. Thomas M. Surowiec.

derive

$$\begin{aligned} \inf_{W \in L_{\mathbb{P}}^p(\Xi)} \mathbb{E}[f_{\mu}(X(\cdot), t, W(\cdot))] &= \inf_{W \in L_{\mathbb{P}}^p(\Xi)} \int_{\Xi} \hat{f}_{\mu}(\xi, W(\xi)) \, d\mathbb{P} \\ &= \int_{\Xi} \inf_{w \in \mathbb{R}} \hat{f}_{\mu}(\xi, w) \, d\mathbb{P} = \mathbb{E}[\inf_{w \in \mathbb{R}} f_{\mu}(X(\cdot), t, w)], \end{aligned}$$

which shows the desired result together with (9.9). \square

In light of Proposition 9.21, we consider for each $t \in \mathbb{R}$, $\xi \in \Xi$, and $x = X(\xi)$ the one-dimensional problem

$$\min_{w \in \mathbb{R}} f_{\mu}(x, t, w),$$

where the unknown w stands for $W(\xi)$.

Proposition 9.22. *The function $w \mapsto f_{\mu}(x, t, w)$ with f_{μ} defined in (9.8) has the unique minimizer*

$$\bar{w} = w_{\mu}(x - t) := \mu + \frac{(a_1 + a_2)(x - t) + \sqrt{(a_1 - a_2)^2(x - t)^2 + 4\mu^2}}{2} \quad (9.10)$$

for every $x, t \in \mathbb{R}$.

Proof. The function f_{μ} is finite for $w > a_i(x - t)$ ($i \in \{1, 2\}$) and convex w. r. t. w . Furthermore, it is continuously differentiable w. r. t. w on the set $\{(w, x, t) : w > a_i(x - t) \text{ for } i \in \{1, 2\}\}$ and has the partial derivative

$$\frac{\partial}{\partial w} f_{\mu}(x, t, w) = 1 - \frac{\mu}{w - a_1(x - t)} - \frac{\mu}{w - a_2(x - t)},$$

which is zero only at the given stationary point \bar{w} . Therefore, this is the unique minimizer. The concrete computation of \bar{w} can be found in Section A.6 in the appendix. \square

By Propositions 9.21 and 9.22 we can define the partial solution of problem (9.7) given X and t as

$$\bar{W}_{X,t}(\xi) := \mu + \frac{1}{2} \left((a_1 + a_2)(X(\xi) - t) + \sqrt{(a_1 - a_2)^2(X(\xi) - t)^2 + 4\mu^2} \right).$$

Inserting this into F_{μ} gives the reduced problem

$$\inf_{t \in \mathbb{R}} F_{\mu}(X, t, \bar{W}_{X,t}) = \inf_{t \in \mathbb{R}} t + \mathbb{E}[v_{\mu}(X - t)] \quad (9.11)$$

which is similar to the initial problem and therefore defines a risk measure $\mathcal{R}_{\mu} \equiv \mathcal{R}_{v_{\mu}}$ induced by the function

$$v_{\mu}(s) := w_{\mu}(s) - \mu \ln(w_{\mu}(s) - a_1 s) - \mu \ln(w_{\mu}(s) - a_2 s) + \zeta(\mu)$$

having w_{μ} defined in (9.10). It remains to discuss and prove interesting properties of the corresponding risk measure $\mathcal{R}_{\mu}[X] := \inf_{t \in \mathbb{R}} t + \mathbb{E}[v_{\mu}(X - t)]$, which we call *the log-barrier risk measure*.

Theorem 9.23. *The log-barrier risk measure $\mathcal{R}_\mu : L_{\mathbb{P}}^1(\Xi) \rightarrow \mathbb{R}$ is well-defined, translation-equivariant, monotonic, and regular.*

Proof. The function v_μ is twice continuously differentiable with derivatives

$$v'_\mu(s) = w'_\mu(s) - \mu \frac{w'_\mu(s) - a_1}{w_\mu(s) - a_1 s} - \mu \frac{w'_\mu(s) - a_2}{w_\mu(s) - a_2 s},$$

$$v''_\mu(s) = \frac{\mu(a_2 - a_1)^2}{2\mu\sqrt{(a_2 - a_1)^2 s^2 + 4\mu^2} + (a_2 - a_1)^2 s^2 + 4\mu^2},$$

where $w_\mu(s) = \mu + \frac{(a_1+a_2)s + \sqrt{(a_1-a_2)^2 s^2 + 4\mu^2}}{2}$ and $w'_\mu(s) = \frac{a_1+a_2}{2} + \frac{(a_1-a_2)^2 s}{2\sqrt{(a_1-a_2)^2 s^2 + 4\mu^2}}$. For the computation of the derivatives see Section A.6.

- Translation equivariance follows directly from Lemma 9.13.
- We have that $\lim_{s \rightarrow -\infty} v'_\mu(s) = a_1$ and $\lim_{s \rightarrow +\infty} v'_\mu(s) = a_2$, see Section A.6. Since v'_μ is strictly increasing ($v''_\mu(s) > 0$), $v'_\mu(\mathbb{R}) = (a_1, a_2) \subset (0, \infty)$ follows. Therefore, the function v_μ itself is strictly increasing. Lemma 9.18 yields that \mathcal{R}_μ is monotonic.
- Since $v''_\mu(s) > 0$ holds for all $s \in \mathbb{R}$, the function v_μ is strictly convex. Now, we choose $d := \frac{2-a_1-a_2}{(1-a_1)(a_2-1)}\mu \in \mathbb{R}$. By Lemma 9.19, it holds that \mathcal{R}_μ is also induced by the function $\hat{v}_\mu(s) = v_\mu(s + d) - d$, which is also strictly convex. One can compute

$$\hat{v}_\mu(0) = v_\mu(d) - d = \mu \left(2 - \ln\left(\frac{a_2-a_1}{a_2-1}\mu\right) - \ln\left(\frac{a_2-a_1}{1-a_1}\mu\right) \right) + \zeta(\mu) = 0 \quad (9.12)$$

and $\hat{v}'_\mu(0) = v'_\mu(d) = 1$, see Section A.6. This yields—together with strict convexity—that $\hat{v}_\mu(s) > s$ for all $s \in \mathbb{R} \setminus \{0\}$. From Theorem 9.20 we get that \mathcal{R}_μ is regular. Equation (9.12) again motivates the choice of the shift $\zeta(\mu)$.

- Corollary 9.11 yields that the log-barrier risk measure is proper. Due to Lemma 9.12 it is always finite on $L_{\mathbb{P}}^1(\Xi)$:

$$|v_\mu(s)| = \left| \int_0^s v'_\mu(s) \, ds + v_\mu(0) \right| \leq a_2 |s| + |v_\mu(0)|$$

holds for all $s \in \mathbb{R}$ due to the computed range of v'_μ . □

The function \hat{v}_μ defined in the proof of Theorem 9.23 is plotted in Figure 9.1 two different values of $a_1 = 1 - \lambda$, $a_2 = 1 + \frac{\beta}{1-\beta}\lambda$, and μ . The properties $\hat{v}_\mu(0) = 0$ and $\hat{v}'_\mu(0) = 1$ are clearly visible and the approximation quality depending on the log-barrier parameter μ is depicted.

Remark 9.24. For $\mu > 0$ and any choice of $a_1 \in [0, 1)$ and $a_2 \in (1, \infty)$, the log-barrier risk measure is not a coherent measure of risk by Theorem 9.20 since v_μ is clearly *not* the maximum of two linear functions.

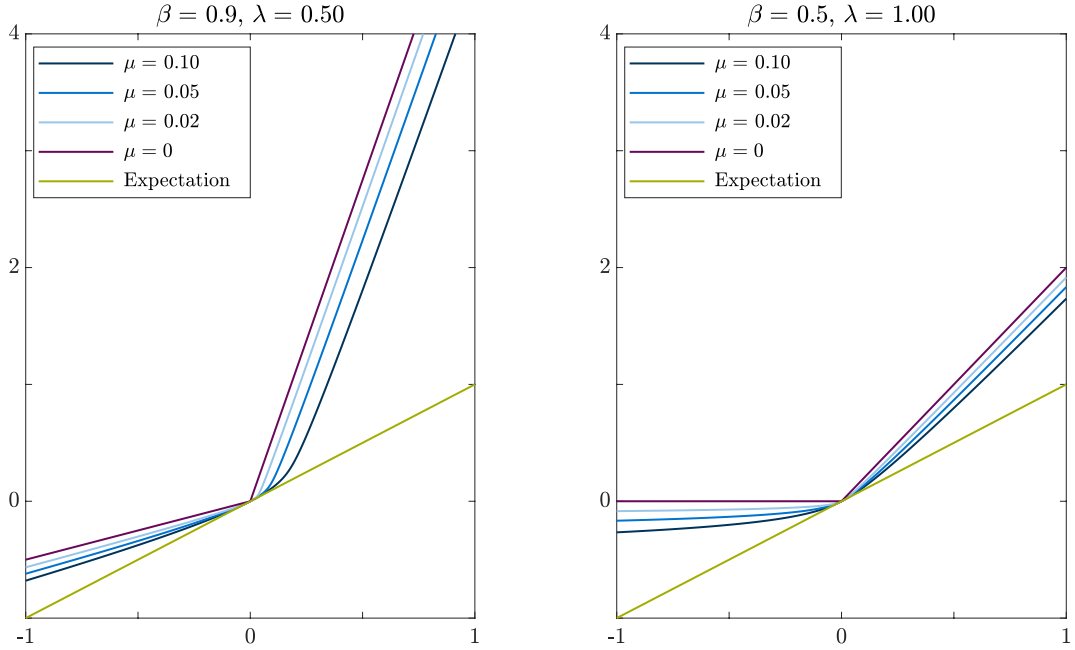


Figure 9.1.: Plots of the function \hat{v}_μ inducing the risk measure \mathcal{R}_μ for different values of $a_1 = 1 - \lambda$, $a_2 = 1 + \frac{\beta}{1-\beta}\lambda$, and μ . The green line corresponds to $v(t) = t$, which induces the risk measure $\mathcal{R}_v \equiv \mathbb{E}$.

9.3.3. Application to Optimal Control under Uncertainty

Consider now the setting from Section 9.2 and let Assumption 9.9 hold with $p_{\text{fun}} \in [1, \infty]$. Furthermore, we define $\hat{J} : U \times \Xi \rightarrow \mathbb{R}$, $\hat{J}(u, \xi) := \hat{J}_1(u, \xi) + J_2(u)$ and set $U_{\text{ad}} = U$ for simplicity. We have that $\hat{J}(u, \cdot) \in L_{\mathbb{P}}^{p_{\text{fun}}}(\Xi)$ for all $u \in U$. We consider the optimization problem under uncertainty

$$\min_{u \in U} \mathcal{R}_\mu[\hat{J}(u, \cdot)] \quad (9.13)$$

with the derived log-barrier risk measure. Since \mathcal{R}_μ is translation-equivariant and J_2 does not depend on ξ , this corresponds to the general setting (1.3) with $U_{\text{ad}} \equiv U$. Clearly, (9.13) is an approximation of an original problem with the nonsmooth risk measure \mathcal{R} , e.g., an optimal control problem with the conditional value-at-risk.

As derived before, we can write (9.13) as

$$\min_{u \in U, t \in \mathbb{R}, W \in L_{\mathbb{P}}^p(\Xi)} \hat{F}_\mu(u, t, W) \quad \text{s. t.} \quad W - a_i(\hat{J}(u, \cdot) - t) \geq 0 \quad \text{a. s. for } i \in \{1, 2\}, \quad (9.14)$$

with

$$\hat{F}_\mu(u, t, W) := \mathbb{E}[t + W - \mu \sum_{i=1}^2 \ln(W - a_i(\hat{J}(u, \cdot) - t))] + \zeta(\mu), \quad (9.15)$$

$a_1 \in [0, 1)$, $a_2 \in (1, \infty)$, $\mu > 0$, and $\zeta(\mu) = \mu(\ln(\frac{a_2 - a_1}{a_2 - 1}\mu) + \ln(\frac{a_2 - a_1}{1 - a_1}\mu) - 2)$.

Lemma 9.25. *The function \hat{F}_μ defined in (9.15) is convex if the function $u \mapsto \hat{J}(u, \xi)$ is convex for a. e. $\xi \in \Xi$.*

Proof. Let $\lambda \in [0, 1]$, $u^1, u^2 \in U$, $t_1, t_2 \in \mathbb{R}$, $W_1, W_2 \in L_{\mathbb{P}}^p(\Xi)$. Clearly, the affine part of the function \hat{F}_μ is convex. Since $\mu > 0$, it remains to show that

$$(u, t, W) \mapsto \mathbb{E}[-\ln(W - a_i \hat{J}(u, \cdot) + a_i t)]$$

is convex for $i \in \{1, 2\}$. Due to convexity, $\hat{J}(\lambda u^1 + (1 - \lambda)u^2, \xi) \leq \lambda \hat{J}(u^1, \xi) + (1 - \lambda) \hat{J}(u^2, \xi)$ holds for a. e. $\xi \in \Xi$. It follows that

$$\begin{aligned} & \lambda W_1(\xi) + (1 - \lambda)W_2(\xi) - a_i \hat{J}(\lambda u^1 + (1 - \lambda)u^2, \xi) + a_i \lambda t_1 + a_i(1 - \lambda)t_2 \\ & \geq \lambda(W_1(\xi) - a_i \hat{J}(u^1, \xi) + a_i t_1) + (1 - \lambda)(W_2(\xi) - a_i \hat{J}(u^2, \xi) + a_i t_2) \end{aligned}$$

for a. e. $\xi \in \Xi$. Applying the negative logarithm, which is decreasing and convex, on both sides yields

$$\begin{aligned} & -\ln(\lambda W_1(\xi) + (1 - \lambda)W_2(\xi) - a_i \hat{J}(\lambda u^1 + (1 - \lambda)u^2, \xi) + a_i \lambda t_1 + a_i(1 - \lambda)t_2) \\ & \leq -\lambda \ln(W_1(\xi) - a_i \hat{J}(u^1, \xi) + a_i t_1) + (1 - \lambda) \ln(W_2(\xi) - a_i \hat{J}(u^2, \xi) + a_i t_2) \end{aligned}$$

for a. e. $\xi \in \Xi$. The monotonicity of the expectation establishes the convexity of \hat{F}_μ . \square

Example 9.26. If $\varphi' \equiv 0$ in the example from Chapter 3, we have a linear state equation and thus a convex reduced objective function $u \mapsto \hat{J}(u, \xi)$ for a. e. $\xi \in \Xi$.

Given $u \in U$, $t \in \mathbb{R}$, we consider the partial problem

$$\min_{W \in L_{\mathbb{P}}^p(\Xi)} \hat{F}_\mu(u, t, W) \quad \text{s. t.} \quad W \geq a_i(\hat{J}(u, \cdot) - t) \quad \text{a. s. for } i \in \{1, 2\}.$$

We have already proven that this problem has the unique global solution $\bar{W}_{u,t} \in L_{\mathbb{P}}^p(\Xi)$ given by

$$\bar{W}_{u,t}(\xi) := \mu + \frac{a_1 + a_2}{2}(\hat{J}(u, \xi) - t) + \frac{1}{2} \sqrt{(a_1 - a_2)^2(\hat{J}(u, \xi) - t)^2 + 4\mu^2}. \quad (9.16)$$

This partial solution fulfills

$$\bar{W}_{u,t} - a_i(\hat{J}(u, \cdot) - t) = \mu \pm \frac{a_2 - a_1}{2}(\hat{J}(u, \xi) - t) + \frac{1}{2} \sqrt{(a_1 - a_2)^2(\hat{J}(u, \xi) - t)^2 + 4\mu^2} \geq \mu, \quad (9.17)$$

i. e., if (9.14) has a solution $(\bar{u}, \bar{t}, \bar{W})$, the constraints are uniformly inactive in an L^∞ sense there. Therefore, we can restrict (9.14) to

$$\min_{u \in U, t \in \mathbb{R}, W \in L_{\mathbb{P}}^p(\Xi)} \hat{F}_\mu(u, t, W) \quad \text{s. t.} \quad W - a_i(\hat{J}(u, \cdot) - t) \geq \mu \quad \text{a. s. for } i \in \{1, 2\} \quad (9.18)$$

without changing its global solution, provided it exists. We see that the objective function is finite on the feasible set.

Existence of a Solution

We derive conditions for the existence of a solution of (9.14) or, equivalently, (9.18). For this purpose, we first bound the auxiliary function f_μ :

Lemma 9.27. *For any $d \in (a_1, a_2)$ and all $x, w, t \in \mathbb{R}$ it holds that*

$$f_\mu(x, t, w) \geq dx + (1 - d)t + \tilde{\zeta}_d(\mu) \quad (9.19)$$

with the function f_μ defined in (9.8), and with $\tilde{\zeta}_d(\mu) := \mu(\ln(\frac{a_2-d}{a_2-1}) + \ln(\frac{d-a_1}{1-a_1}))$.

Proof. We use the fact that $-\mu \ln(s) \geq -\mu \ln(b) - \mu \ln'(b)(s - b) = -\frac{\mu}{b}s + \mu - \mu \ln(b)$ holds for all $s, b \in (0, \infty)$ due to convexity of the negative logarithm and $\mu > 0$. Now with $b_1 = c_1\mu$ and $b_2 = c_2\mu$, $c_1, c_2 > 0$, we obtain

$$\begin{aligned} f_\mu(x, t, w) &= t + w + \sum_{i=1}^2 (-\mu \ln(w - a_i(x - t))) + \zeta(\mu) \geq \\ &t + w + \sum_{i=1}^2 \left(-\frac{1}{c_i}(w - a_i(x - t)) + \mu - \mu \ln(c_i\mu)\right) + \zeta(\mu) = \\ &\left(\frac{a_1}{c_1} + \frac{a_2}{c_2}\right)x + \left(1 - \frac{a_1}{c_1} - \frac{a_2}{c_2}\right)t + \left(1 - \frac{1}{c_1} - \frac{1}{c_2}\right)w + \mu(2 - \ln(c_1\mu) - \ln(c_2\mu)) + \zeta(\mu). \end{aligned} \quad (9.20)$$

Choosing $c_1 = \frac{a_2-a_1}{a_2-d} > 0$ and $c_2 = \frac{a_2-a_1}{d-a_1} > 0$ (giving $\frac{1}{c_1} + \frac{1}{c_2} = 1$ and $\frac{a_1}{c_1} + \frac{a_2}{c_2} = d$) in (9.20) yields

$$\begin{aligned} f_\mu(x, t, w) &\geq dx + (1 - d)t + \mu(2 - \ln(\frac{a_2-a_1}{a_2-d}\mu) - \ln(\frac{a_2-a_1}{d-a_1}\mu)) + \zeta(\mu) \\ &= dx + (1 - d)t + \tilde{\zeta}_d(\mu). \end{aligned}$$

□

Corollary 9.28. *It holds that*

$$\hat{F}_\mu(u, t, W) = \mathbb{E}[f_\mu(\hat{J}(u, \cdot), t, W(\cdot))] \geq \mathbb{E}[\hat{J}(u, \cdot)]$$

for all $(u, t, W) \in U \times \mathbb{R} \times L_{\mathbb{P}}^p(\Xi)$ with \hat{F}_μ defined in (9.15).

Proof. Choosing $d = 1 \in (a_1, a_2)$ in (9.19) results in $f_\mu(x, t, w) \geq x$ for all $x, t, w \in \mathbb{R}$, which shows the result. □

Thus, if the function $u \mapsto \mathbb{E}[\hat{J}(u, \cdot)]$ is bounded from below on U , the function \hat{F}_μ is bounded from below on $U \times \mathbb{R} \times L_{\mathbb{P}}^p(\Xi)$. To prove existence of a solution, we restrict problem (9.18) to a bounded, closed, convex feasible set. Since the partial solution $\bar{W}_{u,t}$ is given for any $u \in U$, $t \in \mathbb{R}$, we bound the feasible set for u and t .

Lemma 9.29. *Let $(\tilde{u}, \tilde{t}, \tilde{W})$ and (u^*, t^*, W^*) be feasible points¹² of (9.18) such that*

$$\hat{F}_\mu(u^*, t^*, W^*) \leq \hat{F}_\mu(\tilde{u}, \tilde{t}, \tilde{W}),$$

and let $u \mapsto \mathbb{E}[\hat{J}(u, \cdot)]$ be bounded from below on U .

¹²A feasible point can be constructed simply: Given some $\tilde{u} \in U$, $\tilde{t} = 0$, we choose, e.g., $\tilde{W}(\xi) = \max_{i \in \{1,2\}} a_i(\hat{J}(\tilde{u}, \xi)) + \mu$.

Then it holds that

$$\mathbb{E}[\hat{J}(u^*, \cdot)] \leq \hat{F}_\mu(\tilde{u}, \tilde{t}, \tilde{W}), \text{ and}$$

$$t^* \in [\tilde{t}_2, \tilde{t}_1], \quad \tilde{t}_i := \frac{1}{1-d_i} (\hat{F}_\mu(\tilde{u}, \tilde{t}, \tilde{W}) - d_i \inf_{\hat{u} \in U} \mathbb{E}[\hat{J}(\hat{u}, \cdot)] - \tilde{\zeta}_{d_i}(\mu)) \quad (i \in \{1, 2\})$$

with any fixed $d_1 \in (a_1, 1)$, $d_2 \in (1, a_2)$, and $\tilde{\zeta}_d(\mu) = \mu(\ln(\frac{a_2-d}{a_2-1}) + \ln(\frac{d-a_1}{1-a_1}))$.

Proof. With the feasible point $(\tilde{u}, \tilde{t}, \tilde{W})$ and with $u \mapsto \mathbb{E}[\hat{J}(u, \cdot)]$ being bounded from below on U , we can bound u^* and t^* as follows: Observe that $\hat{F}_\mu(u^*, t^*, W^*) \leq \hat{F}_\mu(\tilde{u}, \tilde{t}, \tilde{W})$ implies

$$d \mathbb{E}[\hat{J}(u^*, \cdot)] + (1-d)t^* + \tilde{\zeta}_d(\mu) \leq \hat{F}_\mu(\tilde{u}, \tilde{t}, \tilde{W}) \text{ for all } d \in (a_1, a_2)$$

by Lemma 9.27. Choosing $d_0 = 1$ as in the proof of Corollary 9.28, $d_1 \in (a_1, 1)$ and $d_2 \in (1, a_2)$, this gives

$$\begin{aligned} \mathbb{E}[\hat{J}(u^*, \cdot)] &\leq \hat{F}_\mu(\tilde{u}, \tilde{t}, \tilde{W}), \\ t^* &\leq \frac{1}{1-d_1} (\hat{F}_\mu(\tilde{u}, \tilde{t}, \tilde{W}) - d_1 \mathbb{E}[\hat{J}(u^*, \cdot)] - \tilde{\zeta}_{d_1}(\mu)), \\ t^* &\geq \frac{1}{1-d_2} (\hat{F}_\mu(\tilde{u}, \tilde{t}, \tilde{W}) - d_2 \mathbb{E}[\hat{J}(u^*, \cdot)] - \tilde{\zeta}_{d_2}(\mu)) \end{aligned}$$

and thus

$$\begin{aligned} \mathbb{E}[\hat{J}(u^*, \cdot)] &\leq \hat{F}_\mu(\tilde{u}, \tilde{t}, \tilde{W}), \\ t^* &\leq \frac{1}{1-d_1} (\hat{F}_\mu(\tilde{u}, \tilde{t}, \tilde{W}) - d_1 \inf_{\hat{u} \in U} \mathbb{E}[\hat{J}(\hat{u}, \cdot)] - \tilde{\zeta}_{d_1}(\mu)), \\ t^* &\geq \frac{1}{1-d_2} (\hat{F}_\mu(\tilde{u}, \tilde{t}, \tilde{W}) - d_2 \inf_{\hat{u} \in U} \mathbb{E}[\hat{J}(\hat{u}, \cdot)] - \tilde{\zeta}_{d_2}(\mu)), \end{aligned}$$

which shows the result. Note that $-\frac{d_1}{1-d_1} < 0$ whereas $-\frac{d_2}{1-d_2} > 0$. \square

We see that, if $u \mapsto \mathbb{E}[\hat{J}(u, \cdot)]$ is coercive, i. e., $\lim_{\|u\|_U \rightarrow \infty} \mathbb{E}[\hat{J}(u, \cdot)] = \infty$, we can restrict the variable u in problem (9.14) to a nonempty, convex, closed, bounded set $U_{\text{ad}} \subset U$ with

$$\tilde{U}_{\text{ad}} := \{u : \mathbb{E}[\hat{J}(u, \cdot)] \leq \hat{F}_\mu(\tilde{u}, \tilde{t}, \tilde{W})\} \subset U_{\text{ad}}$$

and the variable t to the compact interval given in Lemma 9.29 without changing the global solution of problem (9.18), provided it exists. Note that the set \tilde{U}_{ad} is bounded.

Remark 9.30. Alternatively, the theory would also work if the constraint $u \in U_{\text{ad}}$ with a nonempty, bounded, closed, convex set $U_{\text{ad}} \subset U$ was already posed in (9.14) and if $u \mapsto \mathbb{E}[\hat{J}(u, \cdot)]$ was bounded from below on this set, but not necessarily coercive on U .

Moreover, the optimal partial solution $\bar{W}_{u,t}$ given (u, t) can be inserted. This yields the problem

$$\begin{aligned} &\min_{u \in U, t \in \mathbb{R}, W \in L^p_{\mathbb{F}}(\Xi)} \hat{F}_\mu(u, t, W) \\ \text{s. t.} &\quad \mathbb{E}[\hat{J}(u, \cdot)] \leq \hat{F}_\mu(\tilde{u}, \tilde{t}, \tilde{W}), \quad \tilde{t}_2 \leq t \leq \tilde{t}_1, \\ &\quad W(\cdot) = \mu + \frac{a_1+a_2}{2} (\hat{J}(u, \cdot) - t) + \frac{1}{2} \sqrt{(a_1 - a_2)^2 (\hat{J}(u, \cdot) - t)^2 + 4\mu^2} \text{ a. s.} \end{aligned} \tag{9.21}$$

Based on (9.21), we derive an existence result, which can be applied to convex problems. For a more general framework for risk-averse, PDE-constrained optimization including existence and optimality conditions we refer to [72].

Proposition 9.31. *Let the function $(u, t) \mapsto \hat{F}_\mu(u, t, \bar{W}_{u,t}) \in \mathbb{R}$ with $\bar{W}_{u,t}$ defined in (9.16) be sequentially weakly lower semicontinuous. Furthermore, let the function $u \mapsto \mathbb{E}[\hat{J}(u, \cdot)]$ be coercive on U . Then, problem (9.14) has a solution.*

Proof. Since $u \mapsto \mathbb{E}[\hat{J}(u, \cdot)]$ is coercive on U , it is bounded from below and problem (9.14) is equivalent to problem (9.21). As derived, this is equivalent to

$$\min_{u \in U, t \in \mathbb{R}} \hat{F}_\mu(u, t, \bar{W}_{u,t}) \quad \text{s. t.} \quad u \in U_{\text{ad}}, \quad t \in [\tilde{t}_2, \tilde{t}_1].$$

Now, a feasible, infimizing sequence for this problem is bounded and, therefore, it has a weakly convergent subsequence with limit $(\bar{u}, \bar{t}) \in U \times \mathbb{R}$ because U is reflexive and U_{ad} is closed, convex, bounded. Sequential lower semicontinuity of the objective function w. r. t. the weak topology yields that (\bar{u}, \bar{t}) solves the problem. Moreover, $(\bar{u}, \bar{t}, \bar{W}_{\bar{u}, \bar{t}})$ solves (9.14). \square

Lemma 9.32. *Let the function $u \mapsto \hat{J}(u, \xi) \in \mathbb{R}$ be convex and strongly continuous for a. e. $\xi \in \Xi$ and let the function $u \mapsto \hat{J}(u, \cdot) \in L^1_{\mathbb{P}}(\Xi)$ be such that for every $u \in U$ there exist $\varepsilon_u > 0$ and a function $g_u \in L^1_{\mathbb{P}}(\Xi)$ such that $|\hat{J}(\tilde{u}, \xi)| \leq g_u(\xi)$ holds for every \tilde{u} with $\|\tilde{u} - u\|_U < \varepsilon_u$ and a. e. $\xi \in \Xi$. Then, the function $(u, t) \mapsto \hat{F}_\mu(u, t, \bar{W}_{u,t}) \in \mathbb{R}$ is weakly lower semicontinuous.*

Proof. We observe that $\hat{F}_\mu(u, t, \bar{W}_{u,t}) = t + \mathbb{E}[v_\mu(\hat{J}(u, \cdot) - t)]$ holds, see (9.11). The convexity and the monotonicity of the function v_μ as well as the monotonicity of the expectation yield that $(u, t) \mapsto \hat{F}_\mu(u, t, \bar{W}_{u,t})$ is convex, cf. Lemma 9.25. Now let $(u^k, t_k)_{k \in \mathbb{N}} \subset U \times \mathbb{R}$ be a sequence converging strongly to some $(u, t) \in U \times \mathbb{R}$. W.l.o.g. we assume that $\|u^k - u\|_U < \varepsilon_u$ holds for every k . Then we have that $\hat{J}(u^k, \xi) - t_k \rightarrow \hat{J}(u, \xi) - t$ and thus $v_\mu(\hat{J}(u^k, \xi) - t_k) \rightarrow v_\mu(\hat{J}(u, \xi) - t)$ for a. e. $\xi \in \Xi$. Furthermore, $|v_\mu(\hat{J}(u^k, \xi) - t_k)| \leq a_2(g_u(\xi) + \sup_{k \in \mathbb{N}} |t_k|) + |v_\mu(0)|$ holds for every k and a. e. $\xi \in \Xi$, cf. the proof of Theorem 9.23. As $\sup_{k \in \mathbb{N}} |t_k| < \infty$, this bound is an $L^1_{\mathbb{P}}(\Xi)$ -function and independent of k . Hence, the dominated convergence theorem yields that $t_k + \mathbb{E}[v_\mu(\hat{J}(u^k, \xi) - t_k)] \rightarrow t + \mathbb{E}[v_\mu(\hat{J}(u, \cdot) - t)]$ as $k \rightarrow \infty$. This proves continuity of $(u, t) \mapsto \hat{F}_\mu(u, t, \bar{W}_{u,t})$ w. r. t. strong convergence. Together with convexity, this yields lower semicontinuity w. r. t. the weak topology. \square

Example 9.33. In the example from Chapter 3, the function $u \mapsto \mathbb{E}[\hat{J}(u, \cdot)]$ is coercive because it consists of a non-negative tracking term and the coercive regularization term $\frac{\gamma}{2}\|u\|_U^2$. The function $u \mapsto \hat{J}(u, \xi)$ is strongly continuous for a. e. $\xi \in \Xi$ and we have

$$\begin{aligned} |\hat{J}(\tilde{u}, \xi)| &\leq \frac{1}{2}(\|Q(\xi)\|_{\mathcal{L}(Y,H)}\|S[\xi](\tilde{u})\|_Y + \|\hat{q}(\xi)\|_H)^2 + \frac{\gamma}{2}\|\tilde{u}\|_U^2 \leq \\ &\frac{1}{2}\left(\frac{C_\Omega}{\kappa}\|Q(\xi)\|_{\mathcal{L}(Y,H)}(\|D\|_{\mathcal{L}(U,L^2(\Omega))}\|\tilde{u}\|_U + \|f(\xi) - \varphi(0)\|_{L^2(\Omega)}) + \|\hat{q}(\xi)\|_H\right)^2 + \frac{\gamma}{2}\|\tilde{u}\|_U^2 \leq \\ &\frac{1}{2}\left(\frac{C_\Omega}{\kappa}\|Q(\xi)\|_{\mathcal{L}(Y,H)}(\|D\|_{\mathcal{L}(U,L^2(\Omega))}(\|u\|_U + \varepsilon_u) + \|f(\xi) - \varphi(0)\|_{L^2(\Omega)}) + \|\hat{q}(\xi)\|_H\right)^2 \\ &\quad + \frac{\gamma}{2}(\|u\|_U + \varepsilon_u)^2, \end{aligned}$$

for all $u, \tilde{u} \in U$ and $\varepsilon_u > 0$ with $\|\tilde{u} - u\|_U < \varepsilon_u$, cf. the discussion in Subsection 9.2.1. The derived bound is an $L_{\mathbb{P}}^{\hat{r}}(\Xi)$ -function w. r. t. ξ with \hat{r} from Table 3.1 and hence an $L_{\mathbb{P}}^1(\Xi)$ -function. Furthermore, if $\varphi' \equiv 0$, the function $u \mapsto \hat{J}(u, \xi)$ is convex for a. e. $\xi \in \Xi$, cf. Example 9.26. We therefore can apply Lemma 9.32 and Proposition 9.31 to derive the existence of a solution.

Differentiability Properties of the Function \hat{F}_{μ}

For deriving first-order optimality conditions, we restrict the discussion to the case $p = \infty$. Then, the feasible set of (9.18) has interior points and Fréchet differentiability can be proven. In contrast to that, the case $p < \infty$ requires more sophisticated concepts, see [112].

We consider the function

$$\begin{aligned} \tilde{F}_{\mu} : \{(X, t, W) \in L_{\mathbb{P}}^{\infty}(\Xi) \times \mathbb{R} \times L_{\mathbb{P}}^{\infty}(\Xi) : W \geq a_i(X - t) + c_{\mu}\mu \text{ a. s. for } i \in \{1, 2\}\} &\rightarrow \mathbb{R}, \\ \tilde{F}_{\mu}(X, t, W) &:= \mathbb{E}[t + W - \mu \ln(W - a_1(X - t)) - \mu \ln(W - a_2(X - t)) + \zeta(\mu)] \end{aligned} \quad (9.22)$$

with some $c_{\mu} \in (0, 1)$. This function is well-defined on the given feasible set. We have enlarged the latter compared to the feasible set of (9.18) to prove differentiability properties on an L^{∞} -neighborhood.

Proposition 9.34. *The function \tilde{F}_{μ} defined in (9.22) is twice continuously differentiable on the given feasible set with the F-derivatives*

$$\begin{aligned} \nabla_X \tilde{F}_{\mu}(X, t, W) &= \mu \sum_{i=1}^2 \frac{a_i}{W - a_i(X - t)} \in L_{\mathbb{P}}^1(\Xi), \\ \nabla_t \tilde{F}_{\mu}(X, t, W) &= 1 - \mu \sum_{i=1}^2 \mathbb{E}\left[\frac{a_i}{W - a_i(X - t)}\right] \in \mathbb{R}, \\ \nabla_W \tilde{F}_{\mu}(X, t, W) &= \mathbb{1} - \mu \sum_{i=1}^2 \frac{1}{W - a_i(X - t)} \in L_{\mathbb{P}}^1(\Xi), \\ \nabla_{XX}^2 \tilde{F}_{\mu}(X, t, W)S &= \mu \sum_{i=1}^2 \frac{a_i^2 S}{(W - a_i(X - t))^2} \in L_{\mathbb{P}}^1(\Xi), \\ \nabla_{tX}^2 \tilde{F}_{\mu}(X, t, W)\tau &= -\mu \sum_{i=1}^2 \frac{a_i^2}{(W - a_i(X - t))^2} \tau \in L_{\mathbb{P}}^1(\Xi), \\ \nabla_{WX}^2 \tilde{F}_{\mu}(X, t, W)S &= -\mu \sum_{i=1}^2 \frac{a_i S}{(W - a_i(X - t))^2} \in L_{\mathbb{P}}^1(\Xi), \\ \nabla_{Xt}^2 \tilde{F}_{\mu}(X, t, W)S &= -\mu \sum_{i=1}^2 \frac{a_i^2 S}{(W - a_i(X - t))^2} \in L_{\mathbb{P}}^1(\Xi), \\ \nabla_{tt}^2 \tilde{F}_{\mu}(X, t, W)\tau &= \mu \sum_{i=1}^2 \mathbb{E}\left[\frac{a_i^2}{(W - a_i(X - t))^2}\right] \tau \in \mathbb{R}, \end{aligned}$$

$$\begin{aligned}\nabla_{Wt}^2 \tilde{F}_\mu(X, t, W)S &= \mu \sum_{i=1}^2 \frac{a_i S}{(W - a_i(X-t))^2} \in L_{\mathbb{P}}^1(\Xi), \\ \nabla_{XW}^2 \tilde{F}_\mu(X, t, W)S &= -\mu \sum_{i=1}^2 \frac{a_i S}{(W - a_i(X-t))^2} \in L_{\mathbb{P}}^1(\Xi), \\ \nabla_{tW}^2 \tilde{F}_\mu(X, t, W)\tau &= \mu \sum_{i=1}^2 \frac{a_i}{(W - a_i(X-t))^2} \tau \in L_{\mathbb{P}}^1(\Xi), \\ \nabla_{WW}^2 \tilde{F}_\mu(X, t, W)S &= \mu \sum_{i=1}^2 \frac{S}{(W - a_i(X-t))^2} \in L_{\mathbb{P}}^1(\Xi),\end{aligned}$$

where $S \in L_{\mathbb{P}}^\infty(\Xi)$ and $\tau \in \mathbb{R}$. We write $\nabla_X \tilde{F}_\mu(X, t, W) \in L_{\mathbb{P}}^1(\Xi)$ for the representative of the partial derivative $(\tilde{F}_\mu)_X(X, t, W) \in L_{\mathbb{P}}^\infty(\Xi)^*$ etc.

Proof. Clearly, the linear part $(X, t, W) \mapsto t + \mathbb{E}[W] + \zeta(\mu)$ is bounded and twice continuously differentiable. Therefore, we consider only $(X, t, W) \mapsto \mathbb{E}[\ln(W - a_i(X-t))]$. Let $(X, t, W) \in L_{\mathbb{P}}^\infty(\Xi) \times \mathbb{R} \times L_{\mathbb{P}}^\infty(\Xi)$ be such that $W \geq a_i(X-t) + c_\mu \mu$ holds a.s. Consider a sequence $(S_k)_{k \in \mathbb{N}} \subset L_{\mathbb{P}}^\infty(\Xi)$ such that $0 < \|S_k\|_{L_{\mathbb{P}}^\infty(\Xi)} \leq \frac{c_\mu \mu}{2}$ for all $k \in \mathbb{N}$ and $\|S_k\|_{L_{\mathbb{P}}^\infty(\Xi)} \rightarrow 0$ as $k \rightarrow \infty$. We have

$$\begin{aligned}& \frac{1}{\|S_k\|_{L_{\mathbb{P}}^\infty(\Xi)}} \left| \mathbb{E} \left[\ln(W + S_k - a_i(X-t)) - \ln(W - a_i(X-t)) - \frac{S_k}{W - a_i(X-t)} \right] \right| \\ & \leq \mathbb{E} \left[\int_0^1 \left| \frac{1}{W + \sigma S_k - a_i(X-t)} - \frac{1}{W - a_i(X-t)} \right| d\sigma \right] \\ & = \mathbb{E} \left[\int_0^1 \frac{\sigma |S_k|}{|(W + \sigma S_k - a_i(X-t))(W - a_i(X-t))|} d\sigma \right] \\ & \leq \mathbb{E} \left[\int_0^1 \frac{\sigma}{|(W + \sigma S_k - a_i(X-t))(W - a_i(X-t))|} d\sigma \right] \|S_k\|_{L_{\mathbb{P}}^\infty(\Xi)} \leq \frac{2}{c_\mu^2 \mu^2} \|S_k\|_{L_{\mathbb{P}}^\infty(\Xi)}.\end{aligned}$$

Note that all appearing functions are L^∞ -functions due to

$$W + \sigma S_k - a_i(X-t) \in \left[\frac{(2-\sigma)c_\mu \mu}{2}, \|W\|_{L_{\mathbb{P}}^\infty(\Xi)} + \frac{\sigma c_\mu \mu}{2} + a_i \|X-t\|_{L_{\mathbb{P}}^\infty(\Xi)} \right] \text{ a.s.}$$

for all $\sigma \in [0, 1]$ etc. The estimation proves F-differentiability w.r.t. W with the given derivative. Analogously, the function \tilde{F}_μ is F-differentiable w.r.t. X and t .

Similarly to the estimation above, we have

$$\begin{aligned}& \frac{1}{\|S_k\|_{L_{\mathbb{P}}^\infty(\Xi)}} \left\| \frac{1}{W + S_k - a_i(X-t)} - \frac{1}{W - a_i(X-t)} + \frac{S_k}{(W - a_i(X-t))^2} \right\|_{L_{\mathbb{P}}^1(\Xi)} \\ & \leq \left\| \int_0^1 \left| -\frac{1}{(W + \sigma S_k - a_i(X-t))^2} + \frac{1}{(W - a_i(X-t))^2} \right| d\sigma \right\|_{L_{\mathbb{P}}^1(\Xi)} \\ & \leq \left\| \int_0^1 \frac{|2\sigma W + \sigma^2 S_k - 2a_i \sigma(X-t)|}{(W + \sigma S_k - a_i(X-t))^2 (W - a_i(X-t))^2} d\sigma \right\|_{L_{\mathbb{P}}^1(\Xi)} \|S_k\|_{L_{\mathbb{P}}^\infty(\Xi)} \\ & \leq \frac{4}{c_\mu^4 \mu^4} \left(2\|W\|_{L_{\mathbb{P}}^\infty(\Xi)} + \frac{c_\mu \mu}{2} + 2a_i \|X-t\|_{L_{\mathbb{P}}^\infty(\Xi)} \right) \|S_k\|_{L_{\mathbb{P}}^\infty(\Xi)}\end{aligned}$$

with $0 < \|S_k\|_{L_{\mathbb{P}}^{\infty}(\Xi)} \leq \frac{c_{\mu}\mu}{2}$, which shows that the first derivative $\nabla_W \hat{F}_{\mu}$ is F-differentiable and thus continuous w. r. t. W at (X, t, W) . Again, the remaining second derivatives are shown to be F-derivatives in an analogous way. Continuity of them follows very similarly. \square

First-Order Necessary Optimality Conditions

If the function $U \ni u \mapsto \hat{J}(u, \cdot) \in L_{\mathbb{P}}^{\infty}(\Xi)$ is F-differentiable, cf. Assumption 9.9, we can apply the chain rule to compute the first F-derivative of the function \hat{F}_{μ} given in (9.15) and obtain

$$\begin{aligned} \nabla_u \hat{F}_{\mu}(u, t, W) &= \mu \sum_{i=1}^2 a_i \mathbb{E} \left[\frac{\nabla_u \hat{J}(u, \cdot)}{W - a_i(\hat{J}(u, \cdot) - t)} \right], \\ \nabla_t \hat{F}_{\mu}(u, t, W) &= 1 - \mu \sum_{i=1}^2 a_i \mathbb{E} \left[\frac{1}{W - a_i(\hat{J}(u, \cdot) - t)} \right], \\ \nabla_W \hat{F}_{\mu}(u, t, W) &= \mathbb{1} - \mu \sum_{i=1}^2 \frac{1}{W - a_i(\hat{J}(u, \cdot) - t)} \in L_{\mathbb{P}}^1(\Xi). \end{aligned} \quad (9.23)$$

Again, we write $\nabla_W \hat{F}_{\mu}(u, t, W) \in L_{\mathbb{P}}^1(\Xi)$ for the representative of $(\hat{F}_{\mu})_W(u, t, W) \in L_{\mathbb{P}}^{\infty}(\Xi)^*$. Having computed this function as well as $\nabla_u \hat{J}(u, \cdot)$, the rest of the evaluation of $\nabla \hat{F}_{\mu}$ consists of pointwise multiplications, computing expectations and vector space operations.

The first order necessary conditions for problem (9.18) are given by

$$\nabla_u \hat{F}_{\mu}(\bar{u}, \bar{t}, \bar{W}) = 0, \quad \nabla_t \hat{F}_{\mu}(\bar{u}, \bar{t}, \bar{W}) = 0, \quad \nabla_W \hat{F}_{\mu}(\bar{u}, \bar{t}, \bar{W}) = 0. \quad (9.24)$$

By Lemma 9.25 they are even sufficient if the function $u \mapsto \hat{J}(u, \xi)$ is convex for a. e. ξ .

Barrier-Newton System

To solve (9.24), we apply Newton's method. This is a suitable approach if the function $U \ni u \mapsto \hat{J}(u, \cdot) \in L_{\mathbb{P}}^{\infty}(\Xi)$ is twice continuously differentiable, cf. Assumption 9.9, because this property carries over to the function \hat{F}_{μ} by Proposition 9.34. Let $(u, t, W) \in U \times \mathbb{R} \times L_{\mathbb{P}}^{\infty}(\Xi)$ be the current iterates fulfilling

$$W - a_i(\hat{J}(u, \cdot) - t) \geq c_{\mu}\tilde{\mu} \quad \text{a. s. for } i \in \{1, 2\} \quad (9.25)$$

with $c_{\mu} \in (0, 1)$ as in (9.22) and $\tilde{\mu} = \mu$. We say that a triple (u, t, W) satisfying (9.25) is *approximately feasible* for (9.18) w. r. t. c_{μ} and $\tilde{\mu} \in (0, \infty)$. Then, the Newton directions $(s, \tau, S) \in U \times \mathbb{R} \times L_{\mathbb{P}}^{\infty}(\Xi)$ are given as a solution of the following barrier-Newton system:

$$\begin{pmatrix} \nabla_{uu}^2 \hat{F}_{\mu} & \nabla_{iu}^2 \hat{F}_{\mu} & \nabla_{Wu}^2 \hat{F}_{\mu} \\ \nabla_{ut}^2 \hat{F}_{\mu} & \nabla_{tt}^2 \hat{F}_{\mu} & \nabla_{Wt}^2 \hat{F}_{\mu} \\ \nabla_{uW}^2 \hat{F}_{\mu} & \nabla_{tW}^2 \hat{F}_{\mu} & \nabla_{WW}^2 \hat{F}_{\mu} \end{pmatrix} \begin{pmatrix} s \\ \tau \\ S \end{pmatrix} = \begin{pmatrix} -\nabla_u \hat{F}_{\mu} \\ -\nabla_t \hat{F}_{\mu} \\ -\nabla_W \hat{F}_{\mu} \end{pmatrix}, \quad (9.26)$$

where we skip the arguments (u, t, W) for the sake of brevity.

We compute the derivatives by the chain rule:

$$\begin{aligned}
 \nabla_{uu}^2 \hat{F}_\mu(u, t, W) s &= \mu \sum_{i=1}^2 \left(a_i \mathbb{E} \left[\frac{\nabla_{uu}^2 \hat{J}(u, \cdot) s}{(W - a_i(\hat{J}(u, \cdot) - t))} \right] + a_i^2 \mathbb{E} \left[\frac{(\nabla_u \hat{J}(u, \cdot), s)_U \nabla_u \hat{J}(u, \cdot)}{(W - a_i(\hat{J}(u, \cdot) - t))^2} \right] \right) \\
 \nabla_{tu}^2 \hat{F}_\mu(u, t, W) \tau &= -\mu \sum_{i=1}^2 a_i^2 \mathbb{E} \left[\frac{\nabla_u \hat{J}(u, \cdot)}{(W - a_i(\hat{J}(u, \cdot) - t))^2} \right] \tau \\
 \nabla_{Wu}^2 \hat{F}_\mu(u, t, W) S &= -\mu \sum_{i=1}^2 a_i \mathbb{E} \left[\frac{\nabla_u \hat{J}(u, \cdot) S}{(W - a_i(\hat{J}(u, \cdot) - t))^2} \right] \\
 \nabla_{ut}^2 \hat{F}_\mu(u, t, W) s &= -\mu \sum_{i=1}^2 a_i^2 \mathbb{E} \left[\frac{(\nabla_u \hat{J}(u, \cdot), s)_U}{(W - a_i(\hat{J}(u, \cdot) - t))^2} \right] \\
 \nabla_{tt}^2 \hat{F}_\mu(u, t, W) \tau &= \mu \sum_{i=1}^2 a_i^2 \mathbb{E} \left[\frac{1}{(W - a_i(\hat{J}(u, \cdot) - t))^2} \right] \tau \\
 \nabla_{Wt}^2 \hat{F}_\mu(u, t, W) S &= \mu \sum_{i=1}^2 a_i \mathbb{E} \left[\frac{S}{(W - a_i(\hat{J}(u, \cdot) - t))^2} \right] \\
 \nabla_{uW}^2 \hat{F}_\mu(u, t, W) s &= -\mu \sum_{i=1}^2 a_i \frac{(\nabla_u \hat{J}(u, \cdot), s)_U}{(W - a_i(\hat{J}(u, \cdot) - t))^2} \\
 \nabla_{tW}^2 \hat{F}_\mu(u, t, W) \tau &= \mu \sum_{i=1}^2 a_i \frac{1}{(W - a_i(\hat{J}(u, \cdot) - t))^2} \tau \\
 \nabla_{WW}^2 \hat{F}_\mu(u, t, W) S &= \mu \sum_{i=1}^2 \frac{1}{(W - a_i(\hat{J}(u, \cdot) - t))^2} S.
 \end{aligned}$$

Note that $\frac{1}{(W - a_i(\hat{J}(u, \cdot) - t))^2} \in L_{\mathbb{P}}^\infty(\Xi)$.

We solve (9.26) inexactly by applying a forward-backward block Gauss-Seidel iteration. This has the advantage that we can apply a standard low-rank tensor solver to compute the direction S approximately. We select $(s^0, \tau_0, S_0) = (0, 0, 0)$ as initial guess for the solution of the barrier-Newton system. Let $(s^\ell, \tau_\ell, S_\ell)$ with $\ell \in \mathbb{N}_0$ be the current iterate. The next iterate $(s^{\ell+1}, \tau_{\ell+1}, S_{\ell+1})$ is computed as follows:

In the forward solve, a solution update $(\tilde{s}^{\ell+1}, \tilde{\tau}_{\ell+1}, \tilde{S}_{\ell+1})$ for (9.26) is computed based on the current residual and by replacing the exact Hessian by its lower block triangle. This means that we solve

$$\nabla_{uu}^2 \hat{F}_\mu(\tilde{s}^{\ell+1} - s^\ell) = -\nabla_u \hat{F}_\mu - \nabla_{uu}^2 \hat{F}_\mu s^\ell - \nabla_{tu}^2 \hat{F}_\mu \tau_\ell - \nabla_{Wu}^2 \hat{F}_\mu S_\ell. \quad (9.27)$$

to compute $\tilde{s}^{\ell+1}$. Again, we neglect the argument (u, t, W) here and in the following. Equation (9.27) is rewritten as

$$\nabla_{uu}^2 \hat{F}_\mu \tilde{s}^{\ell+1} = -\nabla_u \hat{F}_\mu - \nabla_{tu}^2 \hat{F}_\mu \tau_\ell - \nabla_{Wu}^2 \hat{F}_\mu S_\ell \quad (9.28)$$

and an iterative solver for it, e. g., a PCG method, is initialized with s^ℓ . Then we compute $\tilde{\tau}_{\ell+1}$ by solving

$$\nabla_{tt}^2 \hat{F}_\mu (\tilde{\tau}_{\ell+1} - \tau_\ell) = -\nabla_t \hat{F}_\mu - \nabla_{ut}^2 \hat{F}_\mu s^\ell - \nabla_{it}^2 \hat{F}_\mu \tau_\ell - \nabla_{Wt}^2 \hat{F}_\mu S_\ell - \nabla_{ut}^2 \hat{F}_\mu (\tilde{s}^{\ell+1} - s^\ell),$$

which is

$$\nabla_{tt}^2 \hat{F}_\mu \tilde{\tau}_{\ell+1} = -\nabla_t \hat{F}_\mu - \nabla_{ut}^2 \hat{F}_\mu \tilde{s}^{\ell+1} - \nabla_{Wt}^2 \hat{F}_\mu S_\ell. \quad (9.29)$$

Since $\tilde{\tau}_{\ell+1}$ is a scalar, this equation can be solved directly. Analogously, $\tilde{S}_{\ell+1}$ is computed by solving

$$\nabla_{WW}^2 \hat{F}_\mu \tilde{S}_{\ell+1} = -\nabla_W \hat{F}_\mu - \nabla_{uW}^2 \hat{F}_\mu \tilde{s}^{\ell+1} - \nabla_{iW}^2 \hat{F}_\mu \tilde{\tau}_{\ell+1}. \quad (9.30)$$

An iterative tensor solver for this equation is initialized with S_ℓ .

In the backward solve, the solution update and next iterate $(s^{\ell+1}, \tau_{\ell+1}, S_{\ell+1})$ is computed by solving the system with the upper block triangle and the new residual on the right-hand side. In fact, $S_{\ell+1}$ would be computed via

$$\nabla_{WW}^2 \hat{F}_\mu (S_{\ell+1} - \tilde{S}_{\ell+1}) = -\nabla_W \hat{F}_\mu - \nabla_{uW}^2 \hat{F}_\mu \tilde{s}^{\ell+1} - \nabla_{iW}^2 \hat{F}_\mu \tilde{\tau}_{\ell+1} - \nabla_{WW}^2 \hat{F}_\mu \tilde{S}_{\ell+1}.$$

This system has the solution $S_{\ell+1} = \tilde{S}_{\ell+1}$ so that it does not have to be solved numerically. Next, $\tau_{\ell+1}$ is computed by

$$\nabla_{tt}^2 \hat{F}_\mu \tau_{\ell+1} = -\nabla_t \hat{F}_\mu - \nabla_{ut}^2 \hat{F}_\mu \tilde{s}^{\ell+1} - \nabla_{Wt}^2 \hat{F}_\mu S_{\ell+1}. \quad (9.31)$$

and $s^{\ell+1}$ is given as the solution of

$$\nabla_{uu}^2 \hat{F}_\mu s^{\ell+1} = -\nabla_u \hat{F}_\mu - \nabla_{iu}^2 \hat{F}_\mu \tau_{\ell+1} - \nabla_{Wu}^2 \hat{F}_\mu S_{\ell+1}. \quad (9.32)$$

The iterative solver for this equation is initialized with $\tilde{s}^{\ell+1}$. In the next iteration, we have $\tilde{s}^{\ell+2} = s^{\ell+1}$ so that (9.28) does not have to be solved numerically.

9.3.4. Implementation and Numerical Results

Using the above considerations and ideas from our paper [46], we implement a Newton-type method for the log-barrier problem (9.14) or equivalently (9.18) with the reduced objective function taken from the example from Chapter 3, where we choose $\varphi \equiv 0$ to have convexity. Then the Hessian operator in (9.26) is at least positive semidefinite since we are solving a convex problem due to Lemma 9.25. We work with a fixed discretization using low-rank tensors as described in Chapter 6, i. e., we use linear FE functions to discretize $H_0^1(\Omega)$ and $L^2(\Omega)$ and represent each random variable by its values at the Gaussian quadrature nodes or—equivalently—a multivariate polynomial. In particular, the FE function $u \in U \subset U$ is represented by a vector $\mathbf{u} \in \mathbb{R}^{d_u}$ and the polynomial $W \in \mathcal{P}_d(\Xi)$ is represented by a tensor $\mathbf{W} \in \mathbb{R}^{d_1 \times \dots \times d_m}$ of weighted function values. Since the weights are positive, pointwise constraints such as (9.25) translate to componentwise constraints on the tensor \mathbf{W} . The setup is the same as in Chapter 8 with the following exceptions: The domain is chosen to be the square $\Omega = (-1, 1)^2$. It is divided into $m = 4$ strips Ω_i as shown in Figure 9.2, where also the used FE mesh consisting of 16641 nodes is drawn. We use the coefficient

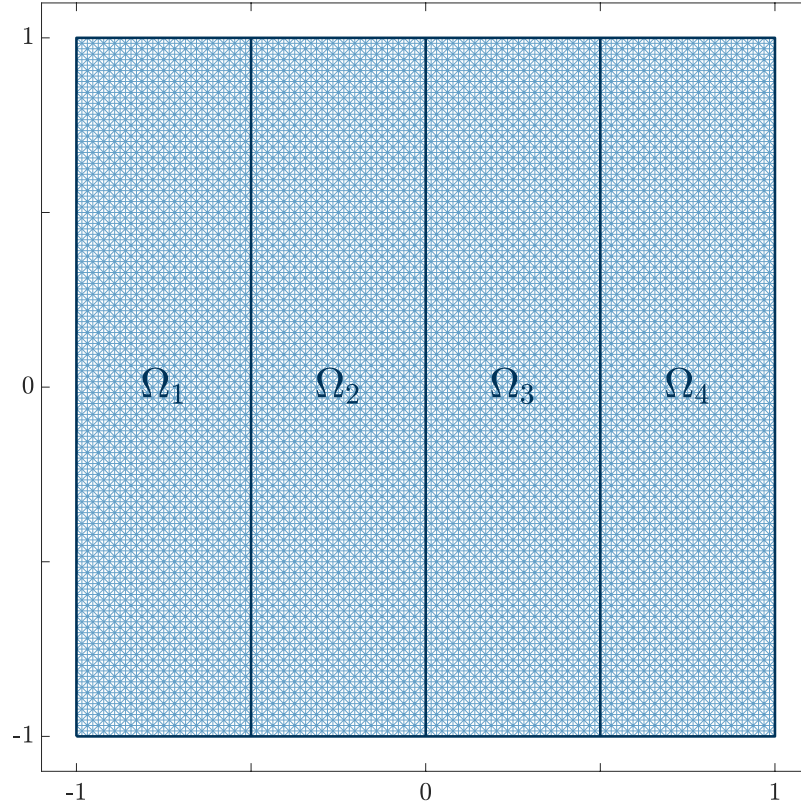


Figure 9.2.: The physical domain Ω , its partition into subdomains Ω_i , and the used finite element mesh

deviations $\sigma = (0.3, 0.4, 0.5, 0.6)^\top$ and discretize each parameter with polynomials of degree $d_i - 1 = 4i + 12$. A high resolution of the stochastic space is chosen because we want to be able to compute quantiles with high accuracy and generate meaningful plots of distribution functions. As mentioned, the “nonlinearity” is set to $\varphi \equiv 0$. The desired state is $\hat{q} \equiv 1$ and we have the set of admissible controls $U_{\text{ad}} = U = L^2(\Omega)$.

Algorithm 2 describes the implemented procedure. It is a Newton-type log-barrier method with stopping criteria stemming from practical experience with low-rank tensor implementations. As observed in [46], the approximate componentwise reciprocals computed by an iterative low-rank tensor method can become too inexact, especially during the last iterations of the log-barrier method. Then, the error in, e. g., the computed gradient of \hat{F}_μ is too large so that a stopping criterion based on its norm is not reliable anymore. Therefore, we also stop the algorithm if the error in the computed reciprocals is too large or if the stepsizes for retaining approximate feasibility become too small, which is a sign for not accurately enough computed Newton steps. Instead of working with a fixed value $\mu_k = \mu$, we use continuation and start with $\mu_0 \geq \mu$ and decrease this value in each iteration by a factor $\mu_{\text{fac}} \in (0, 1)$ until $\mu_k = \mu$ is reached. The algorithm is formulated such that the iterates (u^k, t_k, W_k) are approximately feasible for (9.18) w. r. t. c_μ and μ_k . We give some more details about the substeps in the following.

Algorithm 2: Log-Barrier Method for Solving Problem (9.14)

Parameters : $\mu_0 \geq \mu$, $\mu_{\text{fac}} \in (0, 1)$, $c_\mu \in (0, 1)$, $\varepsilon_{\text{grad}} > 0$

Input: Initial iterates $\mathbf{u}^0 \in U$, $t_0 \in \mathbb{R}$, $\mathbf{W}_0 \in \mathcal{P}_d(\Xi)$ s. t. $(\mathbf{u}^0, t_0, \mathbf{W}_0)$ is approximately feasible for (9.18) w. r. t. c_μ and μ_0 .

Output: $\bar{\mathbf{u}} \in U$, $\bar{t} \in \mathbb{R}$, $\bar{\mathbf{W}} \in \mathcal{P}_d(\Xi)$

for $k := 0, 1, 2, \dots$ **do**

 Compute the tensors corresponding to $\hat{J}(\mathbf{u}^k, \cdot)$ and $\nabla_{\mathbf{u}} \hat{J}(\mathbf{u}^k, \cdot)$.

 Compute the tensors representing $Y_{k,i} := \mathbf{W}_k - a_i(\hat{J}(\mathbf{u}^k, \cdot) - t_k)$ and $Z_{k,i} = Y_{k,i}^{-1}$, $i \in \{1, 2\}$.

if the tensors for $Z_{k,i}$ cannot be computed exactly enough, **then**

 | STOP and return $(\bar{\mathbf{u}}, \bar{t}, \bar{\mathbf{W}}) = (\mathbf{u}^k, t_k, \mathbf{W}_k)$.

end

 Compute representations of the partial gradients $\nabla_{\mathbf{u}} \hat{F}_{\mu_k}^k = \nabla_{\mathbf{u}} \hat{F}_{\mu_k}(\mathbf{u}^k, t_k, \mathbf{W}_k)$, $\nabla_t \hat{F}_{\mu_k}^k = \nabla_t \hat{F}_{\mu_k}(\mathbf{u}^k, t_k, \mathbf{W}_k)$, and $\nabla_{\mathbf{W}} \hat{F}_{\mu_k}^k = \nabla_{\mathbf{W}} \hat{F}_{\mu_k}(\mathbf{u}^k, t_k, \mathbf{W}_k)$, see (9.23).

if $\mu_k = \mu$ and $\|\nabla_{\mathbf{u}} \hat{F}_{\mu_k}^k\|_U^2 + \|\nabla_t \hat{F}_{\mu_k}^k\|^2 + \|\nabla_{\mathbf{W}} \hat{F}_{\mu_k}^k\|_{L_{\mathbb{P}}^2(\Xi)}^2 < \varepsilon_{\text{grad}}^2$, **then**

 | STOP and return $(\bar{\mathbf{u}}, \bar{t}, \bar{\mathbf{W}}) = (\mathbf{u}^k, t_k, \mathbf{W}_k)$.

end

 Compute the tensors corresponding to $Y_{k,i}^2$ and to $Z_{k,i}^2$.

 To compute the update direction $(s^k, \tau_k, \mathbf{S}_k) \in U \times \mathbb{R} \times \mathcal{P}_d(\Xi)$, solve the barrier-Newton system (9.26) with $(u, t, \mathbf{W}) = (\mathbf{u}^k, t_k, \mathbf{W}_k)$ and $\mu = \mu_k$ approximately by applying a forward-backward block Gauss-Seidel iteration.

 Set $\mu_{k+1} := \max\{\mu, \mu_{\text{fac}} \cdot \mu_k\}$.

 Compute stepsizes $\sigma_u^k, \sigma_t^k, \sigma_W^k \in (0, 1]$ such that $(\mathbf{u}^k + \sigma_u^k s^k, t_k + \sigma_t^k \tau_k, \mathbf{W}_k + \sigma_W^k \mathbf{S}_k)$ is approximately feasible for (9.18) w. r. t. c_μ and μ_{k+1} .

if such stepsizes cannot be computed or are too small, **then**

 | STOP and return $(\bar{\mathbf{u}}, \bar{t}, \bar{\mathbf{W}}) = (\mathbf{u}^k, t_k, \mathbf{W}_k)$.

end

 Set $\mathbf{u}^{k+1} := \mathbf{u}^k + \sigma_u^k s^k$, $t_{k+1} := t_k + \sigma_t^k \tau_k$, and $\mathbf{W}_{k+1} := \mathbf{W}_k + \sigma_W^k \mathbf{S}_k$.

end

The object $\mathbf{u}^k \in U$ is a finite element function represented by a vector $\mathbf{u}^k \in \mathbb{R}^{d_u}$. When computing the tensors corresponding to $\hat{J}(\mathbf{u}^k, \cdot)$ and $\nabla_{\mathbf{u}} \hat{J}(\mathbf{u}^k, \cdot)$, we solve the discretized state and adjoint equation by AMEn and store the low-rank tensors $\mathbf{y}_k, \mathbf{z}_k \in \mathbb{R}^{d_0 \times \dots \times d_m}$ representing the discrete state and adjoint state, respectively, to initialize AMEn in the next iteration. The $\mathbb{R}^{d_1 \times \dots \times d_m}$ -tensor representing $\hat{J}(\mathbf{u}^k, \cdot)$ is

$$(\mathbb{1}^\top \circ_1 [(M_H \circ_1 (\mathbf{Q}\mathbf{y}_k - \hat{\mathbf{q}})) \odot (\mathbb{1} \otimes \boldsymbol{\omega}) \odot (\mathbf{Q}\mathbf{y}_k - \hat{\mathbf{q}})])(1, \bullet, \dots, \bullet) + \left(\frac{\gamma}{2}(\mathbf{u}^k)^\top \tilde{\mathbf{M}}\mathbf{u}^k\right) \boldsymbol{\omega}^{-1}, \quad (9.33)$$

with $\mathbb{1} \in \mathbb{R}^{d_H}$, cf. (6.18), and the one representing $\nabla_{\mathbf{u}} \hat{J}(\mathbf{u}^k, \cdot)$ is given by

$$-(\tilde{\mathbf{M}}^{-1}\mathbf{B}^\top) \circ_1 \mathbf{z}_k + \gamma \mathbf{u}^k \otimes \boldsymbol{\omega}^{-1} \in \mathbb{R}^{d_u \times d_1 \times \dots \times d_m},$$

cf. (6.22). In (9.33) and in the following we use interpolation, i. e., we construct a tensor of function evaluations in the Gaussian quadrature nodes \mathbf{a}_l multiplied by the weights $\boldsymbol{\omega}(l)^{-1}$.

The tensors $\mathbf{Y}_{k,i} \in \mathbb{R}^{d_1 \times \dots \times d_m}$ and $\mathbf{Z}_{k,i} \in \mathbb{R}^{d_1 \times \dots \times d_m}$ shall represent the functions $Y_{k,i}$ and $Z_{k,i}$, cf. (6.10). To compute $\mathbf{Z}_{k,i} \approx \boldsymbol{\omega}^{-1} \odot (\boldsymbol{\omega} \odot \mathbf{Y}_{k,i})^{-1}$, we first apply the Newton-Schulz method, see Subsection 2.1.3. If the relative error $\frac{\|(\boldsymbol{\omega} \odot \mathbf{Z}_{k,i}) \odot (\boldsymbol{\omega} \odot \mathbf{Y}_{k,i}) - \mathbb{1}\|_{\mathbb{F}}}{\|\mathbb{1}\|_{\mathbb{F}}}$ is greater than 10^{-4} , we perform additional AMEn iterations to compute the elementwise reciprocal, cf. [46]. Algorithm 2 is stopped if the relative error in the elementwise reciprocal exceeds 0.1. Note that in theory we have $Y_{k,i}(\xi) \geq c_\mu \mu_k$ for a. e. $\xi \in \Xi$ because $(\mathbf{u}^k, t_k, \mathbf{W}_k)$ is approximately feasible for (9.18) w. r. t. c_μ and μ_k . Therefore, the functions $Z_{k,i} \in L_{\mathbb{P}}^\infty(\Xi)$ are well-defined and fulfill $Z_{k,i}(\cdot) \in (0, \frac{1}{c_\mu \mu_k}]$ almost surely.

The evaluation of the gradient (9.23) becomes

$$\begin{aligned}\nabla_u \hat{F}_{\mu_k}(\mathbf{u}^k, t_k, \mathbf{W}_k) &= \mu_k \sum_{i=1}^2 a_i \mathbb{E}[\nabla_u \hat{J}(\mathbf{u}^k, \cdot) Z_{k,i}] \\ \nabla_t \hat{F}_{\mu_k}(\mathbf{u}^k, t_k, \mathbf{W}_k) &= 1 - \mu_k \sum_{i=1}^2 a_i \mathbb{E}[Z_{k,i}] \\ \nabla_W \hat{F}_{\mu_k}(\mathbf{u}^k, t_k, \mathbf{W}_k) &= \mathbb{1} - \mu_k \sum_{i=1}^2 Z_{k,i}.\end{aligned}$$

The expectations can be evaluated using tensor contractions. To evaluate the gradient norm, we use the discrete $L_{\mathbb{P}}^2(\Xi)$ -norm because it is also given as a tensor contraction.

The tensors corresponding to $Y_{k,i}^2$ are $\mathbf{Y}_{k,i} \odot (\boldsymbol{\omega} \odot \mathbf{Y}_{k,i})$ and are computed by truncated componentwise multiplication and i -mode matrix products, cf. (6.25). Analogously, we compute the tensors corresponding to $Z_{k,i}^2$.

To solve the barrier-Newton system (9.26) approximately, we perform at most 5 iterations of the described forward-backward block Gauss-Seidel method. We stop this iteration earlier if the relative residual is smaller than 0.1. In (9.28), we use an approximation $\tilde{\nabla}_{uu}^2 \hat{F}_{\mu_k}(\mathbf{u}^k, t_k, \mathbf{W}_k)$ of the Hessian involving a constant (w. r. t. ξ) approximation $\tilde{\nabla}_{uu}^2 \hat{J}(\mathbf{u}^k)$ of the Hessian $\nabla_{uu}^2 \hat{J}(\mathbf{u}^k, \cdot)$, cf. (6.26). This leads to

$$\begin{aligned}& \mu_k \sum_{i=1}^2 \left(a_i \mathbb{E}[Z_{k,i}] \tilde{\nabla}_{uu}^2 \hat{J}(\mathbf{u}^k) \tilde{\mathfrak{s}}^{k,\ell+1} + a_i^2 \mathbb{E}[(\nabla_u \hat{J}(\mathbf{u}^k, \cdot), \tilde{\mathfrak{s}}^{k,\ell+1})_U \nabla_u \hat{J}(\mathbf{u}^k, \cdot) Z_{k,i}^2] \right) = \\ & - \mu_k \sum_{i=1}^2 a_i \mathbb{E}[\nabla_u \hat{J}(\mathbf{u}^k, \cdot) Z_{k,i}] + \mu_k \tau_{k,\ell} \sum_{i=1}^2 a_i^2 \mathbb{E}[\nabla_u \hat{J}(\mathbf{u}^k, \cdot) Z_{k,i}^2] \\ & + \mu \sum_{i=1}^2 a_i^2 \mathbb{E}[\nabla_u \hat{J}(\mathbf{u}^k, \cdot) \mathbf{S}_{k,\ell} Z_{k,i}^2].\end{aligned}$$

A PCG method is used to solve this equation because the discrete reference Hessian $\tilde{\nabla}_{uu}^2 \hat{J}(\mathbf{u}^k)$ can only be applied efficiently, but should not be formed as an explicit, typically dense matrix. The expectation $\mathbb{E}[Z_{k,i}]$ and the right-hand side are precomputed before solving w. r. t. $\tilde{\mathfrak{s}}^{k,\ell+1}$. Expected values are computed by tensor contractions and componentwise

multiplication. (9.32) is solved analogously. Equation (9.29) becomes

$$\begin{aligned}\tilde{\tau}_{k,\ell+1} &= \frac{-\nabla_t \hat{F}_{\mu_k}(u^k, t_k, W_k) - \nabla_{ut}^2 \hat{F}_{\mu_k}(u^k, t_k, W_k) \tilde{s}^{k,\ell+1} - \nabla_{Wt}^2 \hat{F}_{\mu_k}(u^k, t_k, W_k) S_{k,\ell}}{\nabla_{tt}^2 \hat{F}_{\mu_k}(u^k, t_k, W_k)} \\ &= \left(\mu_k \sum_{i=1}^2 a_i^2 \mathbb{E}[Z_{k,i}^2] \right)^{-1} \left(-1 + \mu_k \sum_{i=1}^2 a_i \mathbb{E}[Z_{k,i}] \right. \\ &\quad \left. + \mu_k \sum_{i=1}^2 a_i^2 \mathbb{E}[(\nabla_u \hat{J}(u^k, \cdot), \tilde{s}^{k,\ell+1})_U Z_{k,i}^2] - \mu_k \sum_{i=1}^2 a_i \mathbb{E}[Z_{k,i}^2 S_{k,\ell}] \right)\end{aligned}$$

and (9.31) is formulated analogously. We transform (9.30) to

$$\begin{aligned}\tilde{S}_{k,\ell+1} &:= \nabla_{WW}^2 \hat{F}_{\mu_k}(u^k, t_k, W_k)^{-1} \left(-\nabla_W \hat{F}_{\mu_k}(u^k, t_k, W_k) \right. \\ &\quad \left. - \nabla_{uW}^2 \hat{F}_{\mu_k}(u^k, t_k, W_k) \tilde{s}^{k,\ell+1} - \nabla_{tW}^2 \hat{F}_{\mu_k}(u^k, t_k, W_k) \tilde{\tau}_{k,\ell+1} \right) \\ &= \frac{1}{\mu_k} (Z_{k,1}^2 + Z_{k,2}^2)^{-1} \left(-\mathbb{1} + \mu_k (Z_{k,1} + Z_{k,2}) \right. \\ &\quad \left. + \mu_k (a_1 Z_{k,1}^2 + a_2 Z_{k,2}^2) (\nabla_u \hat{J}(u^k, \cdot), \tilde{s}^{k,\ell+1})_U - \mu_k \tilde{\tau}_{k,\ell+1} (a_1 Z_{k,1}^2 + a_2 Z_{k,2}^2) \right) \\ &= \frac{1}{\mu_k} (Y_{k,2}^2 + Y_{k,1}^2)^{-1} \left(-Y_{k,1}^2 Y_{k,2}^2 + \mu_k (Y_{k,1} Y_{k,2}^2 + Y_{k,1}^2 Y_{k,2}) \right. \\ &\quad \left. + \mu_k (a_1 Y_{k,2}^2 + a_2 Y_{k,1}^2) (\nabla_u \hat{J}(u^k, \cdot), \tilde{s}^{k,\ell+1})_U - \mu \tilde{\tau}_{k,\ell+1} (a_1 Y_{k,2}^2 + a_2 Y_{k,1}^2) \right).\end{aligned}$$

The advantage of the latter formulation is that it does not depend on the possibly inexact computed $Z_{k,i}$. The tensor representing $\tilde{S}_{k,\ell+1}$ is computed by AMEn with componentwise multiplication operators.

We propose a stepsize selection strategy, which uses the fact that (9.25) is a box constraint w. r. t. W and t and that it can be rewritten as

$$\operatorname{ess\,inf}_{\xi \in \Xi} W(\xi) - a_i (\hat{J}(u, \xi) - t) \geq c_\mu \tilde{\mu} \quad \text{for } i \in \{1, 2\}.$$

It tries to avoid too many evaluations of the random variable objective function $\hat{J}(u, \cdot)$ for different u and works as follows:

- S1. Choose factors $\sigma_{\text{fac},u} \in (0, 1)$, $\hat{\sigma}_{\text{fac}} \in (0, 1]$, and initial stepsizes $\sigma_u, \hat{\sigma} \in (0, 1]$.
- S2. Compute the tensor corresponding to $\hat{J}(u^k + \sigma_u s^k, \cdot) \in L_{\mathbb{P}}^\infty(\Xi)$ and store this quantity and the corresponding state for the next iteration.
- S3. Compute $\delta_{k,i} := \min_{\xi \in \Xi} W_k(\xi) + \hat{\sigma} S_k(\xi) - a_i \hat{J}(u^k + \sigma_u s^k, \xi)$ for $i \in \{1, 2\}$. This is done by the multilevel coordinate search (MCS) method [65], a non-rigorous global optimization routine, which works with evaluations of the respective polynomial only, see Subsection 2.1.3.
- S4. If $\delta_{k,i} \geq c_\mu \mu_{k+1} - \hat{\sigma} a_i \tau_k - a_i t_k$ for $i \in \{1, 2\}$, choose $\sigma_u^k = \sigma_u$ and $\sigma_t^k = \sigma_W^k = \hat{\sigma}$ and STOP. The new iterate $(u^k + \sigma_u^k s^k, t_k + \sigma_t^k \tau_k, W_k + \sigma_W^k S_k)$ is then approximately feasible for (9.18) w. r. t. c_μ and μ_{k+1} .
Otherwise, compute $\nu_{k,i} := \min_{\xi \in \Xi} W_k(\xi) - a_i \hat{J}(u^k + \sigma_u s^k, \xi)$ for $i \in \{1, 2\}$ by MCS.

S5. If $\nu_{k,i} < c_\mu \mu_{k+1} - a_i t_k$ for some $i \in \{1, 2\}$, decrease σ_u by the factor $\sigma_{\text{fac},u}$ and go to S2 because $(\mathbf{u}^k + \sigma_u \mathbf{s}^k, t_k, \mathbf{W}_k)$ is not approximately feasible.

Otherwise, i. e., $\nu_{k,i} \geq c_\mu \mu_{k+1} - a_i t_k$ for all $i \in \{1, 2\}$ and $c_\mu \mu_{k+1} - a_i t_k > \delta_{k,i} + \hat{\sigma} a_i \tau_k$ for at least one $i \in \{1, 2\}$ giving $\nu_{k,i} > \delta_{k,i} + \hat{\sigma} a_i \tau_k$ for at least one $i \in \{1, 2\}$, compute

$$\sigma^k := \min_{i: \nu_{k,i} > \delta_{k,i} + \hat{\sigma} a_i \tau_k} \frac{c_\mu \mu_{k+1} - \nu_{k,i} - a_i t_k}{-\frac{\nu_{k,i}}{\hat{\sigma}} + \frac{\delta_{k,i}}{\hat{\sigma}} + a_i \tau_k}.$$

This stepsize is non-negative because the denominator is negative and the numerator is non-positive.

S6. If $\sigma^k < \hat{\sigma}_{\text{fac}} \sigma_u$, decrease σ_u by the factor $\sigma_{\text{fac},u}$ and go to S2 to have stepsizes of approximately the same magnitude.

Otherwise, choose $\sigma_u^k = \sigma_u$ and $\sigma_t^k = \sigma_W^k = \sigma^k$ and STOP. The new iterate $(\mathbf{u}^k + \sigma_u^k \mathbf{s}^k, t_k + \sigma_t^k \tau_k, \mathbf{W}_k + \sigma_W^k \mathbf{S}_k)$ is approximately feasible:

$$\begin{aligned} & \mathbf{W}_k(\cdot) + \sigma^k \mathbf{S}_k(\cdot) - a_i \hat{J}(\mathbf{u}^k + \sigma_u^k \mathbf{s}^k, \cdot) + a_i t_k + a_i \sigma^k \tau_k \\ &= \left(1 - \frac{\sigma^k}{\hat{\sigma}}\right) (\mathbf{W}_k(\cdot) - a_i \hat{J}(\mathbf{u}^k + \sigma_u^k \mathbf{s}^k, \cdot)) \\ & \quad + \frac{\sigma^k}{\hat{\sigma}} (\mathbf{W}_k(\cdot) + \hat{\sigma} \mathbf{S}_k(\cdot) - a_i \hat{J}(\mathbf{u}^k + \sigma_u \mathbf{s}^k, \cdot)) + a_i t_k + \sigma^k a_i \tau_k \\ &\geq \left(1 - \frac{\sigma^k}{\hat{\sigma}}\right) \nu_{k,i} + \frac{\sigma^k}{\hat{\sigma}} \delta_{k,i} + a_i t_k + \sigma^k a_i \tau_k \\ &= \sigma^k \left(-\frac{\nu_{k,i}}{\hat{\sigma}} + \frac{\delta_{k,i}}{\hat{\sigma}} + a_i \tau_k\right) + \nu_{k,i} + a_i t_k. \end{aligned}$$

If $-\frac{\nu_{k,i}}{\hat{\sigma}} + \frac{\delta_{k,i}}{\hat{\sigma}} + a_i \tau_k < 0$, we have

$$\sigma^k \left(-\frac{\nu_{k,i}}{\hat{\sigma}} + \frac{\delta_{k,i}}{\hat{\sigma}} + a_i \tau_k\right) + \nu_{k,i} + a_i t_k \geq c_\mu \mu_{k+1} - \nu_{k,i} - a_i \tau_k + \nu_{k,i} + a_i t_k = c_\mu \mu_{k+1}$$

by the definition of σ^k . Otherwise, we estimate

$$\underbrace{\sigma^k \left(-\frac{\nu_{k,i}}{\hat{\sigma}} + \frac{\delta_{k,i}}{\hat{\sigma}} + a_i \tau_k\right)}_{\geq 0} + \nu_{k,i} + a_i t_k \geq \nu_{k,i} + a_i t_k \geq c_\mu \mu_{k+1}.$$

Algorithm 2 is stopped if the stepsize σ_u becomes smaller than 10^{-3} during the iteration. In the numerical tests, we choose $\mu_0 = 10$, $\mu_{\text{fac}} = 0.8$, $c_\mu = 0.1$, $\varepsilon_{\text{grad}} = 0.01$, $\sigma_u = \hat{\sigma} = 0.7$, $\sigma_{\text{fac},u} = \hat{\sigma}_{\text{fac}} = 0.5$, and bound all appearing tensor ranks by 200. This rank bound is chosen quite large to be able to obtain accurate results for difficult setups. In easier cases such as $\beta = 0.5$, $\lambda = 1.00$, $\mu = 0.10$, see below, the largest used rank for the elementwise reciprocals $Z_{k,i}$ is 57 and the largest rank for the steps \mathbf{S}_k is 83, i. e., the rank bound is not restrictive. The initial control \mathbf{u}^0 is chosen to be the deterministic control and we take $t_0 = 0$ and $\mathbf{W}_0(\cdot) := a_2 \hat{J}(\mathbf{u}^0, \cdot) + \mu_0$. Since $\hat{J}(\mathbf{u}^0, \cdot) \geq 0$ and $a_2 > a_1 \geq 0$, this yields that $(\mathbf{u}^0, t_0, \mathbf{W}_0)$ is approximately feasible.

Clearly, Algorithm 2 is not as rigorous as, e. g., Algorithm 1 and the implementation presented in Chapter 8. The stopping criteria based on the quality of the elementwise reciprocal and the computed stepsize, for instance, come from practical experience with low-rank tensor implementations. In theory, developing a convergence theory for such an algorithm based

on results on inexact Newton methods [41] or inexact interior point methods [19] would be possible. But it is questionable if the requirements of such convergence results can be met in practice. Even a very small relative error in the tensor $\mathbf{Y}_{k,i}$ measured in a norm induced by an inner product can cause the reciprocal tensor $\mathbf{Z}_{k,i}$ to be not well-defined anymore because, e. g., one entry of $\mathbf{Y}_{k,i}$ could become zero, cf. the discussion in Subsection 2.1.3. Furthermore, the discretization relies on interpolation in the multivariate Gaussian quadrature nodes. As noted at the end of Section 6.2, this can be interpreted as a discretization based on samples in the quadrature nodes. This is a suitable interpretation in this case because the polynomials of arbitrary degree do not form a dense subset of $L_{\mathbb{P}}^{\infty}(\Xi)$ in general so that a discretization by polynomials is not appropriate.

In the following, we present results for different combinations of the quantile parameter $\beta \in \{0.5, 0.9\}$ in CVaR_{β} , the convex combination parameter $\lambda \in \{1.00, 0.75, 0.50\}$, see (9.5), and the log-barrier parameter $\mu \in \{0.1, 0.05, 0.02\}$. The results are summed up in Table 9.1.

β	λ	μ	number of iterations	achieved gradient norm	computing time (hours)
0.9	1.00	0.05	24	0.3012	3.83
0.9	1.00	0.10	25	0.0062	4.38
0.9	0.75	0.05	26	0.2869	1.97
0.9	0.50	0.05	28	0.0067	2.69
0.5	1.00	0.02	34	0.0071	5.56
0.5	1.00	0.05	28	0.0010	1.02
0.5	1.00	0.10	25	0.0083	0.62
0.5	0.75	0.02	33	0.0072	3.77
0.5	0.50	0.02	34	0.0068	2.51

Table 9.1.: Results of the implemented log-barrier method.

If the achieved gradient norm is greater than 0.01, the algorithm has stopped because the error in the elementwise reciprocal is too large. Stopping due to too small stepsizes has never been occurred although the full stepsize is not always taken, especially in the last iterations.

In general, it can be observed that it is easier to solve problems which are less risk-averse in the sense that smaller values of β or λ , or larger values of μ are chosen. In two cases, namely $\beta = 0.9$, $\mu = 0.05$, and $\lambda \in \{0.75, 1.00\}$, it is not possible to obtain a gradient norm smaller than 0.01. Nevertheless, we will see that the obtained solutions yield a smaller CVaR than the ones with $\lambda = 0.50$ or $\mu = 0.10$, for which a better accuracy is achieved. In contrast to the case $\beta = 0.9$, the algorithm computes a solution with a small gradient for $\beta = 0.5$ even with $\mu = 0.02$. The variability in the computing times is due to various influences: Depending on the obtained residual, possibly more Gauss-Seidel iterations have to be performed, especially for ill-conditioned systems. Furthermore, we use costly direct solves with full linear algebra in AMEn subproblems if the computation of the elementwise reciprocal with an iterative subsolver does not yield satisfactory results. The ranks needed to

compute, e. g., the reciprocal tensor accurately enough can vary depending on the problem data. Subsequent computations with this tensor are then more costly.

An interesting result, where the controls differ relatively much from each other, is obtained for $\beta = 0.5$, $\lambda = 1.00$, and $\mu = 0.02$, i. e., for minimization of the smoothed $\text{CVaR}_{0.5}$ with a rather small value of the barrier parameter μ . Figure 9.3 shows the distribution function of the “random variable objective function” [72] $\hat{J}(u, \cdot)$ for the deterministic, the risk-neutral, and the respective risk-averse control. In all distribution function plots, the expected value $\mathbb{E}[\hat{J}(u, \cdot)]$ of the respective random variable is marked by “*”. $\text{CVaR}_{0.5}[\hat{J}(u, \cdot)]$ and $\text{CVaR}_{0.9}[\hat{J}(u, \cdot)]$ are marked by “+” and “×”, respectively. As depicted in Figure 9.3, the risk-neutral control reduces the expectation as well as the $\text{CVaR}_{0.5}$ and the $\text{CVaR}_{0.9}$ of the random variable objective function compared to the values achieved by the deterministic control. Using the risk-averse approach, it is possible to decrease the $\text{CVaR}_{0.5}$ further while the expectation increases only a bit. Even the $\text{CVaR}_{0.9}$ is decreased in this setting although it is not minimized explicitly. The obtained controls and their differences are plotted in Figure 9.4. Especially the difference plots show that the uncertainty in the system is increasing from left to right. Additionally, the strips Ω_i can be recognized in these plots.

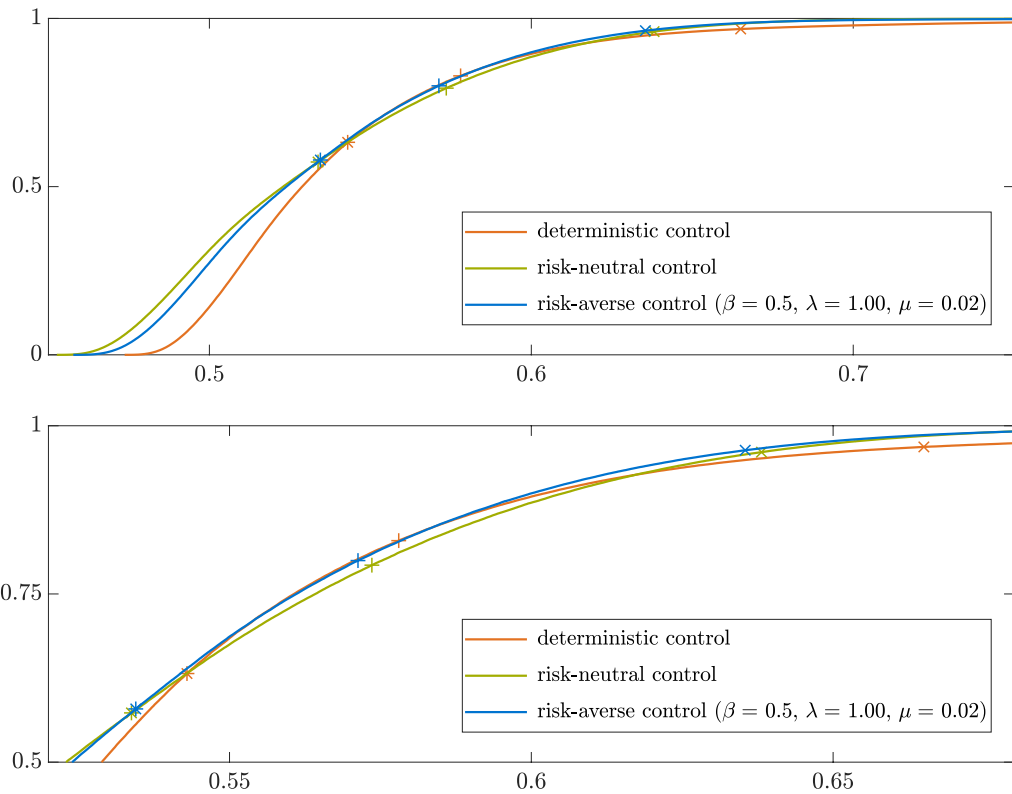


Figure 9.3.: Distribution function of the random variable objective function for the deterministic, the risk-neutral, and the risk-averse control with $\beta = 0.5$, $\lambda = 1.00$, and $\mu = 0.02$, details at the bottom

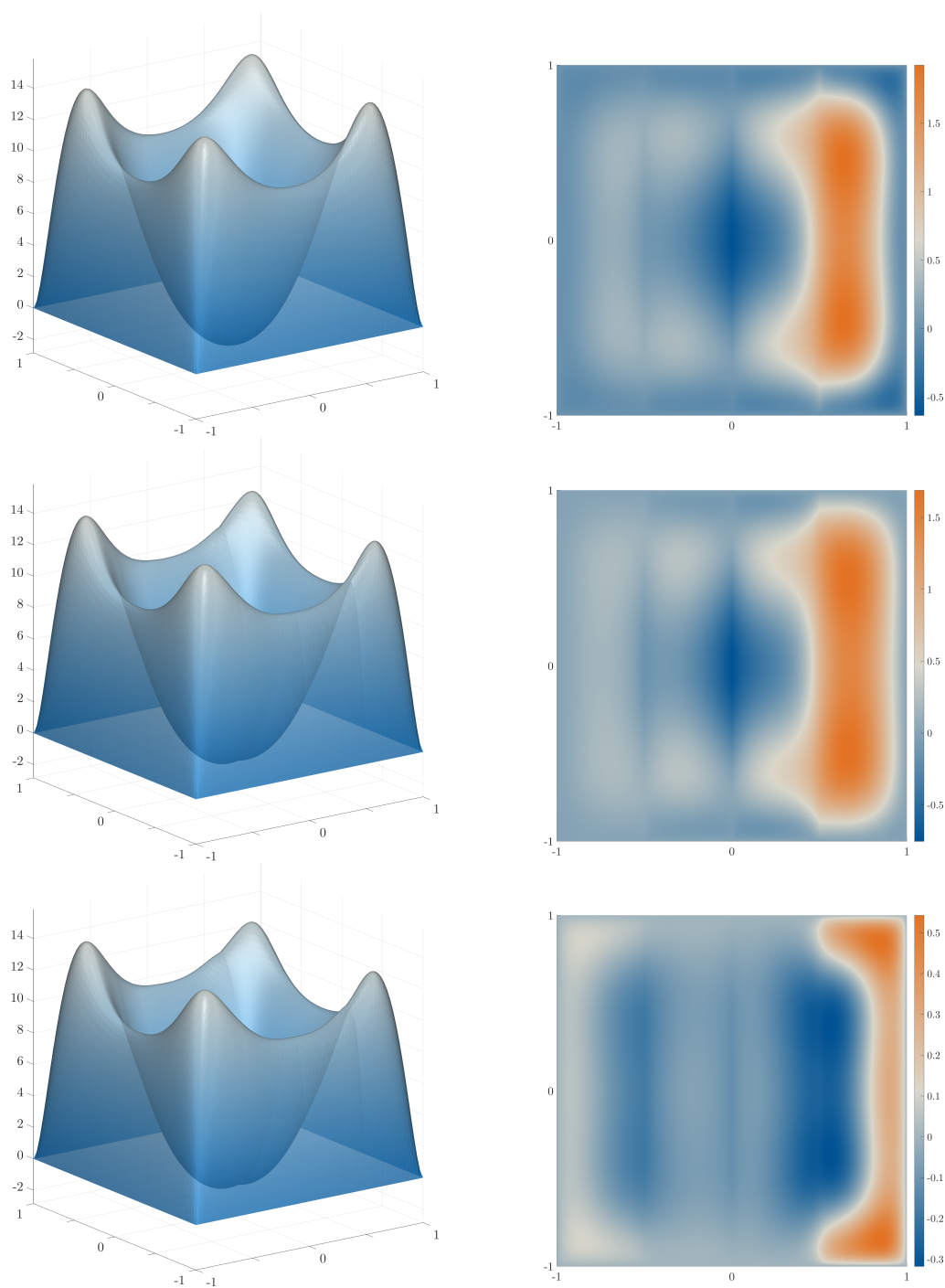


Figure 9.4.: Plots of the deterministic, the risk-neutral, and the risk-averse control (from top left to bottom left) for $\beta = 0.$, $\lambda = 1.00$, and $\mu = 0.02$. Differences between the deterministic and the risk-neutral, the deterministic and the risk-averse, and the risk-neutral and the risk-averse control (from top right to bottom right).

In the following, we investigate how the choice of one of the parameters β , λ , μ influences the result. Sometimes, the differences between the distribution functions are small so that we show regions of interest in the respective plots.

We start with the role of the log-barrier parameter μ . Figure 9.5 shows the distribution function of the random variable objective function for the risk-neutral and the risk-averse controls with $\beta = 0.5$ and $\lambda = 1.00$, i.e., we aim for minimizing smoothed versions of $\text{CVaR}_{0.5}$. The log-barrier parameter is decreased from 0.10 to 0.02. This setup is chosen because the algorithm is capable of computing a solution with small gradient, see Table 9.1, which is not the case for $\beta = 0.9$. We see that the $\text{CVaR}_{0.5}[\hat{J}(u, \cdot)]$ increases and approaches the value achieved by the risk-neutral control for larger values of μ while the expected value decreases slightly. This becomes clear by taking a look at Figure 9.1. For larger values of μ , the function \hat{v}_μ gets locally closer to the identity $v(t) = t$, which induces $\mathcal{R}_v = \mathbb{E}$.

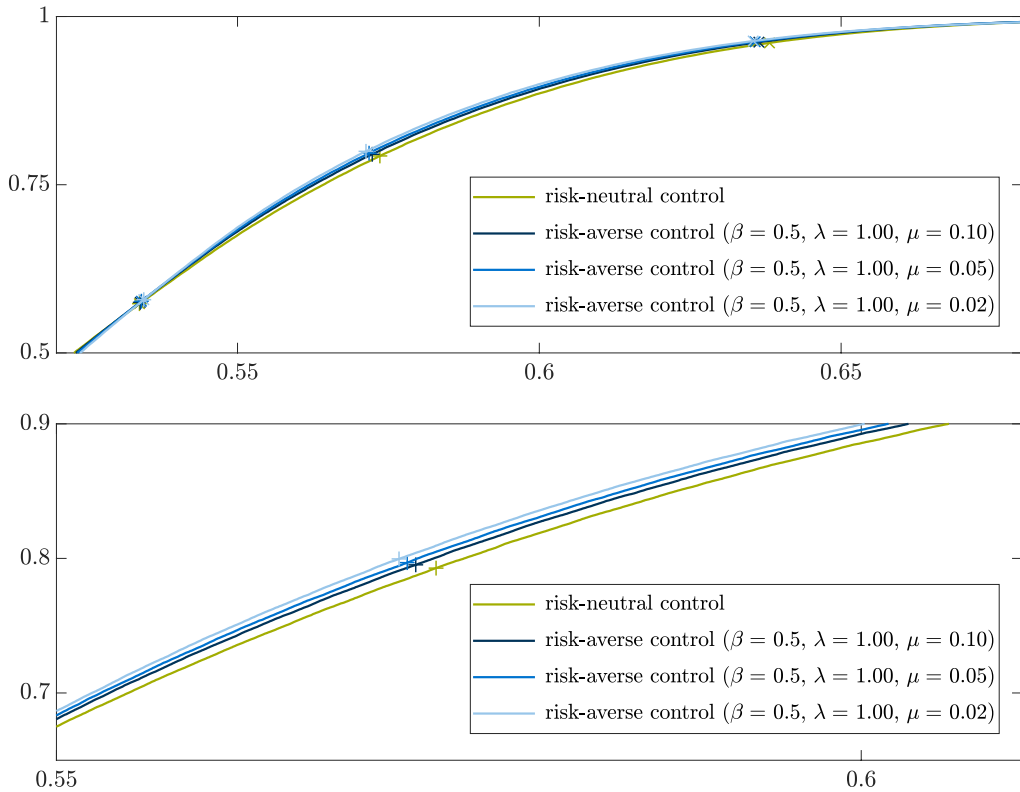


Figure 9.5.: Distribution function of the random variable objective function for the risk-neutral and the risk-averse controls with $\mu \in \{0.10, 0.05, 0.02\}$, $\beta = 0.5$, $\lambda = 1.00$

Figures 9.6 and 9.6 show the distribution functions as the combination parameter λ varies. Here we have two different setups, namely $(\beta, \mu) = (0.5, 0.02)$ and $(\beta, \mu) = (0.9, 0.05)$. As expected, the respective CVaR_β is smaller for larger values of λ and the expectation increases. It is remarkable that we make this observation also in the case $\beta = 0.9$ although the algorithm has to stop with a gradient norm of about 0.3 for $\lambda = 1.00$ and $\lambda = 0.75$ because the componentwise reciprocal cannot be computed accurately enough, see Table 9.1.

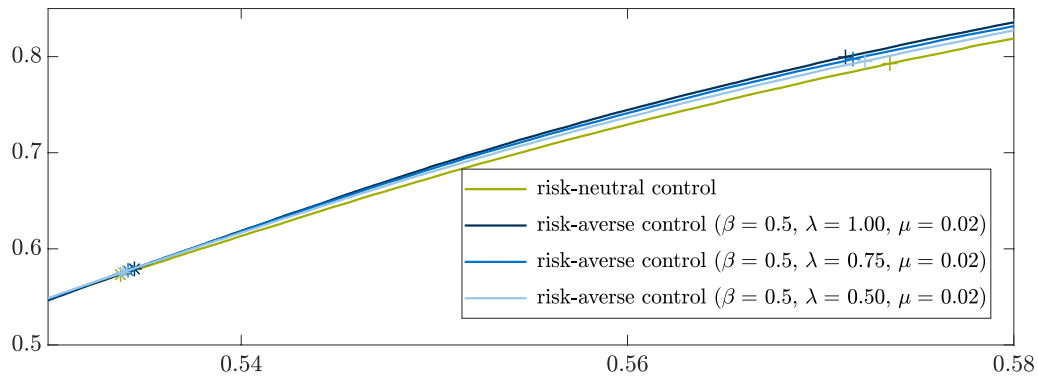


Figure 9.6.: Distribution function of the random variable objective function for the risk-neutral and the risk-averse controls with $\lambda \in \{1.00, 0.75, 0.50\}$ and $(\beta, \mu) = (0.5, 0.02)$

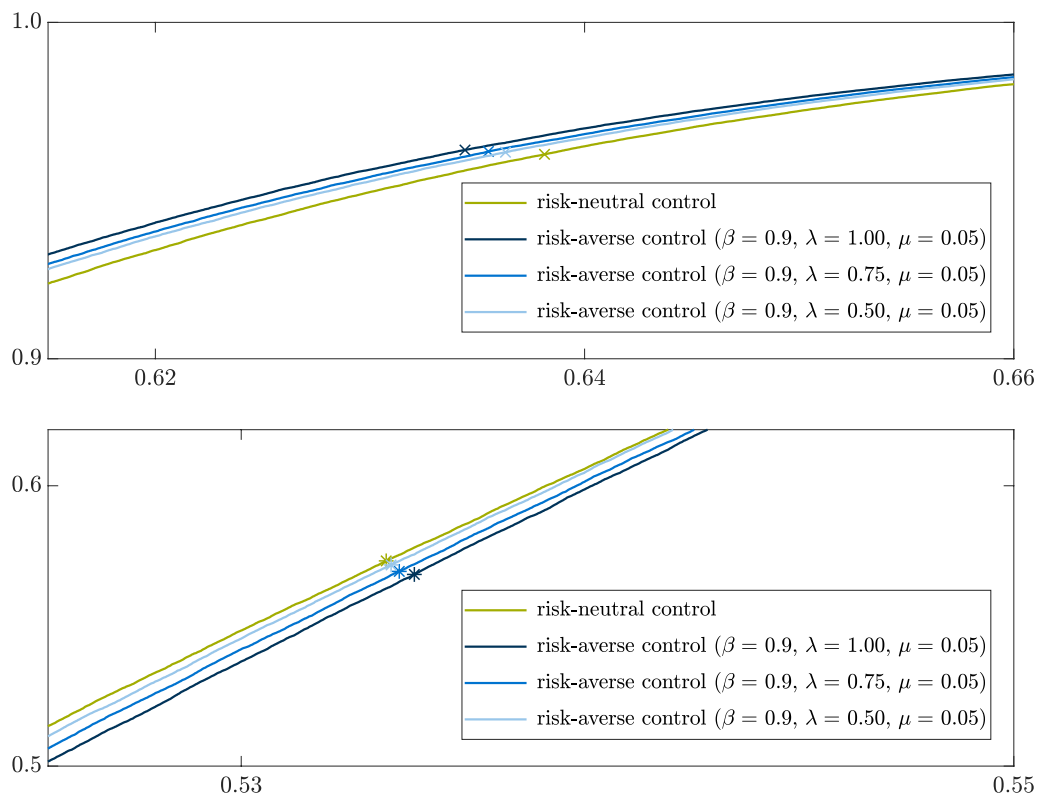


Figure 9.7.: Distribution function of the random variable objective function for the risk-neutral and the risk-averse controls with $\lambda \in \{1.00, 0.75, 0.50\}$ and $(\beta, \mu) = (0.9, 0.05)$

In Figure 9.8 we show the distribution function for different quantile parameters $\beta \in \{0.5, 0.9\}$ with $\lambda = 1.00$, i.e., we minimize a smoothed version of CVaR_β with $\mu = 0.10$ and $\mu = 0.05$, respectively. In both cases, the control resulting from a minimization of the smoothed $\text{CVaR}_{0.9}$ yields a smaller $\text{CVaR}_{0.9}$, but a larger $\text{CVaR}_{0.5}$ than the one which minimizes the the smoothed $\text{CVaR}_{0.5}$. Both are more risk-averse in terms of $\text{CVaR}_{0.5}$ and $\text{CVaR}_{0.9}$ than the risk-neutral control. In the case $\beta = 0.9$, $\mu = 0.05$ this behavior can be observed although the gradient norm is rather large in the computed solution. This solution achieves a smaller $\text{CVaR}_{0.9}$ than the one computed with $\beta = 0.9$ and $\mu = 0.10$.

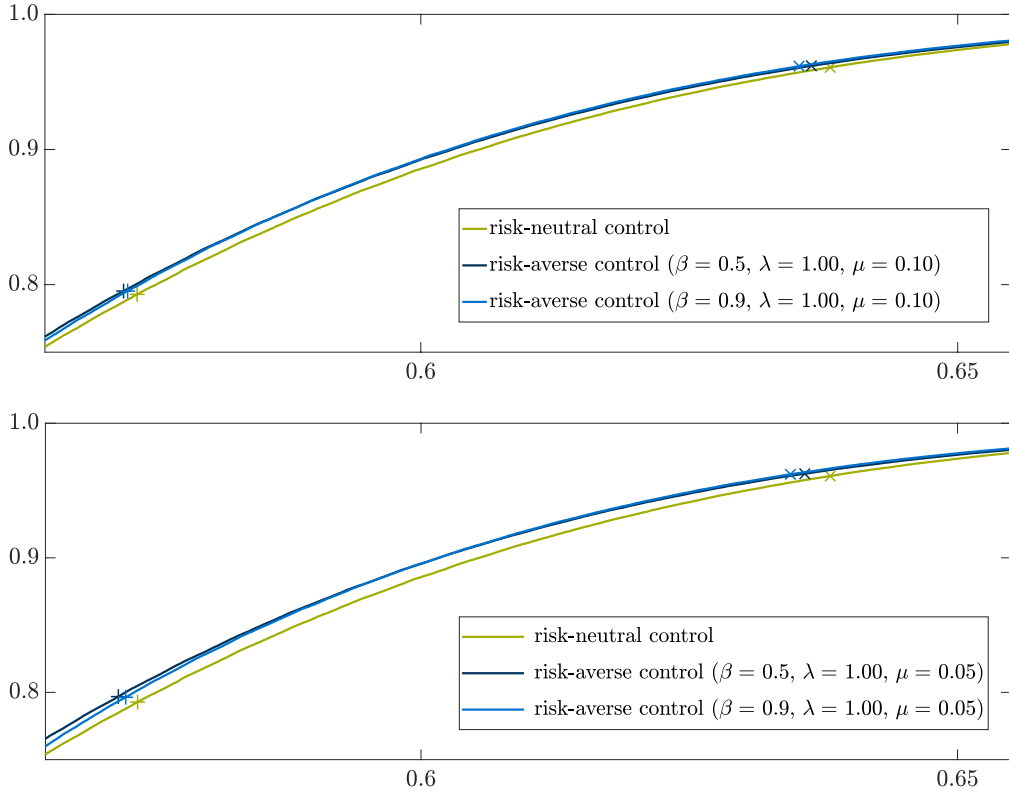


Figure 9.8.: Distribution function of the random variable objective function for the risk-neutral and the risk-averse controls with $\beta \in \{0.5, 0.9\}$, $\lambda = 1.00$, $\mu = 0.10$ (top), $\mu = 0.05$ (bottom)

In conclusion, we have proposed a log-barrier method implemented with low-rank tensors for solving risk-averse optimal control problems under uncertainty with a convex combination of the mean and the CVaR as risk measure. We have shown that keeping the barrier parameter $\mu > 0$ on a fixed level results in the minimization of a smoothed, monotonic, regular risk measure. We have formulated an existence result for convex problems as well as an optimality condition and the barrier-Newton equation in function space. The numerical results have shown that this is a hard problem class if implemented with low-rank tensors as already observed in [46]. Nevertheless, we are capable of computing risk-averse controls by this method, which yield a random variable objective with smaller CVaR than the risk-neutral control.

10. Conclusions and Perspectives

In this thesis, we have discussed optimal control problems of semilinear, elliptic PDEs with uncertain inputs and different risk measures. We have provided the necessary theory to formulate the PDEs in weak form in a tensor Banach space and have derived adjoint-based expressions for the derivatives of the reduced objective function in the risk-neutral and in the mean-variance case. A trust-region framework allowing for inexact objective function, gradient, and criticality measure evaluations has been established. It features global convergence while error bounds up to unknown multiplicative constants are sufficient. In the risk-neutral and in the mean variance case, we have investigated how the required error tolerances can be fulfilled based on the error in the state, in the adjoint state, and in the inexact projection onto the set of admissible controls. If the risk measure is the expected value and the problem data is essentially bounded w. r. t. the random parameters, error control of the state and adjoint state in the $L_{\mathbb{P}}^2(\Xi; H_0^1(\Omega))$ -norm is sufficient. We have derived an a posteriori estimator for this error based on a stochastic Galerkin discretization with polynomials in tensor product form. It has been discussed how the PDEs with uncertain inputs can be solved adaptively using low-rank tensor calculus and methods. This is necessary to make the computation efficient. Numerical results have shown the adaptivity of the proposed approach to different problem data in the sense that the constructed FE mesh resolves the PDE solution and the active set well and the polynomial grades are chosen dependent on the influence of the respective uncertain parameter on the controlled system.

The error estimation theory is rigorous such that global convergence of the algorithm can be established, but the practical implementation lacks this property at some points. We have used the AMEn method to solve tensor equations because of its efficiency in practice although its convergence theory is still limited to our knowledge. Additionally, it is not clear a priori if approximate solutions to the respective equations with moderate tensor ranks exist and can be computed by AMEn. In our experiments, the required small error tolerance has not been met by AMEn sometimes in the last iterations of the trust-region algorithm. Furthermore, we have bounded the number of used FE nodes due to efficiency reasons. Some terms in the error estimation—appearing only in the case that a nonlinear PDE is considered—have been neglected, namely the interpolation error and the $L_{\mathbb{P}}^{\infty}(\Xi; H_0^1(\Omega))$ -norm of the computed state.

In addition to the mentioned smooth risk measures, we also have discussed risk measures which are convex combinations of the mean and the conditional value-at-risk. They have many favorable properties, but are nonsmooth. We have proposed to apply a log-barrier method to a smooth reformulation of optimal control problems involving such risk measures and have investigated how this procedure affects the underlying risk measures if the barrier parameter is not driven to zero, but kept on a fixed level. In fact, useful properties, such as convexity and monotonicity, are preserved. Numerical results of a low-rank tensor implementation of this method working on a fixed discretization have been shown. The barrier-Newton

system has been solved approximately by a forward-backward block Gauss-Seidel method, which makes use of the structure of the system. The results reveal how the distribution function of the random variable objective function of the optimal control problem can be shaped by this procedure. Although the convergence of Newton's method in the proposed function space setting is well-known, it is hard to derive a practically relevant convergence result for a low-rank tensor implementation of it. Even very small perturbations due to truncation may cause the barrier terms to be not well-defined anymore so that iterative solvers for computing them may fail to converge.

Based on this work, several directions of future research can be pursued. The proposed approach for the analysis (Chapter 3) and the error estimation (Chapters 5 and 7) can be adapted to different settings of optimal control problems with semilinear, elliptic PDEs such as boundary control or different boundary conditions. Furthermore, low-rank methods are also suitable for time-dependent PDEs [105], where the time yields an additional tensor mode. The a posteriori error estimation procedure from Chapter 7 can be generalized to the 3D case and to more general coefficient functions. Additionally, the convergence of the discretization for the considered class of semilinear PDEs should be analyzed. Regarding the mean-variance risk measure, the question arises whether error estimation in $L_{\mathbb{P}}^q(\Xi; H_0^1(\Omega))$ with $q > 2$ is possible. For the ease of implementation, we have bounded all errors based on the $L_{\mathbb{P}}^2(\Xi; H_0^1(\Omega))$ -error in the state and adjoint state. By this procedure, we certainly overestimate the error in some cases. As discussed, weaker error control can be sufficient, for instance, an $L_{\mathbb{P}}^2(\Xi; L^2(\Omega))$ -error bound for the adjoint state to control the gradient error. In addition, alternative a posteriori error estimation techniques, such as the dual-weighted-residual method [18] for the objective function evaluation, can be employed, but have to be adapted to the stochastic setting.

In our numerical experiments, we have refrained from testing the scaling of the proposed algorithms w. r. t. to the number of the parameters because such tests can be found in [46] and our tests have been dedicated to adaptivity, which itself has required to run different setups. Regarding computing time and exactness of the obtained results, a comprehensive comparison to different treatments of the stochasticity, such as the use of adaptive sparse grids or multilevel Monte Carlo methods, would be interesting. For a fair comparison, one should find a setup for which the exact solution can be constructed, and the sample-based methods should also use adaptive grid refinement. Since they can benefit from parallelization, parallel codes for low-rank tensors, which are investigated currently [52], should be applied.

Based on the derived barrier-Newton system from Section 9.3 or a similar primal-dual system, one can derive an interior-point method for such problems. As shown, the necessary derivatives are Fréchet derivatives if $L_{\mathbb{P}}^{\infty}(\Xi)$ is considered as underlying function space. If one wants to design an adaptive algorithm for such problems, this space would have to be discretized suitably, which is a hard task because, e. g., the set of polynomials of arbitrary degree is not dense in $L_{\mathbb{P}}^{\infty}(\Xi)$. To overcome this issue, the space of continuous functions could be used alternatively.

To sum up, it shall be mentioned that many different methods and concepts from, e. g., multilinear algebra, functional analysis and PDE theory, numerical optimization and optimization in Banach spaces, numerical analysis and approximation theory as well as probability theory and uncertainty quantification, have been used in this dissertation. The interplay of them makes the topic very interesting and motivates future research in various directions.

Danksagung

Ich möchte an dieser Stelle mehreren Personen meinen Dank aussprechen, die zum Gelingen dieser Arbeit beigetragen haben.

Zuerst gilt mein Dank Prof. Dr. Michael Ulbrich, der mich im Rahmen des „TopMath“-Studiengangs bereits seit meinem dritten Studienjahr betreut hat und mich schon im Bachelor-Studium an Themen der Unsicherheitsquantifizierung und der Optimierung mit partiellen Differentialgleichungen heranführte. Sein Vorschlag, Niedrigrangtensormethoden im Bereich der Optimalsteuerung unter Unsicherheit und der Unsicherheitsquantifizierung zu verwenden, führte zu einem spannenden, vielseitigen Promotionsthema. Durch die Anstellung an seinem Lehrstuhl bot sich mir die Möglichkeit, neben der Forschungserfahrung auch umfangreiche Lehrerfahrung zu sammeln.

Prof. Dr. Thomas Surowiec danke ich für die spannende Zusammenarbeit im Bezug auf alternative Risikomaße, was die Basis für Kapitel 9 dieser Arbeit gebildet hat.

Außerdem freue ich mich, dass sich Prof. Dr. Christian Clason und Prof. Dr. Matthias Heinkenschloss bereiterklärt haben, diese Arbeit zu begutachten.

Für die Förderung ideeller und finanzieller Art möchte ich mich beim Elitestudienprogramm „TopMath – Mathematik mit Promotion“ und bei der International Research Training Group „IGDK Munich – Graz: Optimization and Numerical Analysis for Partial Differential Equations with Nonsmooth Structures“ und allen Verantwortlichen bedanken. Durch TopMath konnte ich schon früh sehr frei und forschungsorientiert studieren. Beide Programme ermöglichten mir die Teilnahme an interessanten Konferenzen, Workshops und Kursen.

Außerdem sollen sämtliche Kollegen und Kommilitonen aus dem Studium, von TopMath, dem IGDK und vom Lehrstuhl für Mathematische Optimierung nicht unerwähnt bleiben, ohne dass hier jeder persönlich genannt wird. Ich bedanke mich für hilfreiche fachliche Diskussionen und für die gute Zeit und Zusammenarbeit. Gesondert erwähnen will ich Johannes Milz und Johannes Haubner, die Teile dieser Arbeit genau gelesen und mit Anmerkungen versehen haben, und meinen Mentor Dr. Andre Milzarek, der mir gerade zu Beginn der Promotionszeit mit hilfreichen Ratschlägen zur Seite stand.

Zuletzt gebührt mein Dank meiner Familie, meinen Freunden und Gela. Danke, dass ihr für mich da seid, mich unterstützt und dass wir so viel gemeinsam erleben dürfen.

A. Appendix

A.1. Tensor Spaces

We give the proof of Proposition 2.6 here since it does not contribute to the main content of the thesis.

Proposition A.1 (Characterization of equivalence classes). *Let V and W be vector spaces over a field \mathbb{K} and let $v, \tilde{v} \in V$ and $w, \tilde{w} \in W$ be given. Consider the algebraic tensor space $V \otimes_a W$. The pair (\tilde{v}, \tilde{w}) belongs to the equivalence class $v \otimes w$ if and only if the following holds:*

$$(v = 0 \vee w = 0) \wedge (\tilde{v} = 0 \vee \tilde{w} = 0) \quad (\text{A.1})$$

or

$$(v \neq 0 \wedge w \neq 0) \vee (\tilde{v} \neq 0 \wedge \tilde{w} \neq 0) \quad \text{and} \quad \exists c \in \mathbb{K} \setminus \{0\} \text{ s. t. } \tilde{v} = cv \text{ and } \tilde{w} = c\tilde{w}. \quad (\text{A.2})$$

Case (A.1) holds exactly for pairs belonging to the equivalence class $0 \otimes 0$.

Proof. First we show that (A.1) holds if and only if both tensors $v \otimes w$ and $\tilde{v} \otimes \tilde{w}$ are in fact zero, i. e., they belong to the equivalence class $0 \otimes 0$: For some $v \in V$ we have $(v, 0) \in \mathcal{N}$ with \mathcal{N} defined as in Definition 2.5 taking $m = n = 1$, $\alpha_1 = 1$, $\beta_1 = 0$, $v_1 = v$, and $w_1 = 0$, and multiplying the resulting tensor by -1 . Analogously we get $(0, w) \in \mathcal{N}$ for any $w \in W$.

Now we prove that the condition $v = 0 \vee w = 0$, cf. (A.1), is also necessary for a pair (v, w) to belong to $0 \otimes 0 = \mathcal{N}$. First observe that it is enough to allow only vectors from given bases $B \subset V$ and $C \subset W$ in the definition of the set \mathcal{N} : We have $\mathcal{N} = \tilde{\mathcal{N}}$ with

$$\tilde{\mathcal{N}} := \text{span} \left\{ \sum_{i=1}^m \sum_{j=1}^n \alpha_i \beta_j (v_i, w_j) - \left(\sum_{i=1}^m \alpha_i v_i, \sum_{j=1}^n \beta_j w_j \right) : \right. \\ \left. m, n \in \mathbb{N}, \alpha_i, \beta_j \in \mathbb{K}, v_i \in B, w_j \in C \right\}$$

We only show that linear combinations of pairs of the form

$$\sum_{i=1}^m \sum_{j=1}^n \alpha_i \beta_j (v_i, w_j) - \left(\sum_{i=1}^m \alpha_i v_i, \sum_{j=1}^n \beta_j w_j \right)$$

with general $v_i \in V$ and $w_i \in W$ belong to the set $\tilde{\mathcal{N}}$ because the inclusion $\tilde{\mathcal{N}} \subset \mathcal{N}$ is obvious. We insert the representations $v_i = \sum_{v \in B} \lambda_v^i v$ and $w_j = \sum_{w \in C} \mu_w^j w$, which are in fact finite

sums, into the formula and compute

$$\begin{aligned}
& \sum_{i=1}^m \sum_{j=1}^n \alpha_i \beta_j (v_i, w_j) - \left(\sum_{i=1}^m \alpha_i v_i, \sum_{j=1}^n \beta_j w_j \right) \\
&= \sum_{i=1}^m \sum_{j=1}^n \alpha_i \beta_j \left(\sum_{v \in B} \lambda_v^i v, \sum_{w \in C} \mu_w^j w \right) - \left(\sum_{i=1}^m \alpha_i \sum_{v \in B} \lambda_v^i v, \sum_{j=1}^n \beta_j \sum_{w \in C} \mu_w^j w \right) \\
&= \sum_{i=1}^m \sum_{j=1}^n \alpha_i \beta_j \left(\sum_{v \in B} \lambda_v^i v, \sum_{w \in C} \mu_w^j w \right) - \sum_{i=1}^m \sum_{j=1}^n \sum_{v \in B} \sum_{w \in C} \alpha_i \beta_j \lambda_v^i \mu_w^j (v, w) \\
&\quad + \sum_{i=1}^m \sum_{j=1}^n \sum_{v \in B} \sum_{w \in C} \alpha_i \beta_j \lambda_v^i \mu_w^j (v, w) - \left(\sum_{i=1}^m \sum_{v \in B} \alpha_i \lambda_v^i v, \sum_{j=1}^n \sum_{w \in C} \beta_j \mu_w^j w \right).
\end{aligned}$$

This is a linear combination belonging to the set $\tilde{\mathcal{N}}$. Now let $(v, w) \in \mathcal{N}$, i. e.,

$$\begin{aligned}
(v, w) &= \sum_{k=1}^K \gamma_k \left(\sum_{i=1}^m \sum_{j=1}^n \alpha_i^k \beta_j^k (v_i^k, w_j^k) - \left(\sum_{i=1}^m \alpha_i^k v_i^k, \sum_{j=1}^n \beta_j^k w_j^k \right) \right) \\
&= \sum_{k=1}^K \sum_{i=1}^m \sum_{j=1}^n \gamma_k \alpha_i^k \beta_j^k (v_i^k, w_j^k) - \sum_{k=1}^K \gamma_k \left(\sum_{i=1}^m \alpha_i^k v_i^k, \sum_{j=1}^n \beta_j^k w_j^k \right) \tag{A.3}
\end{aligned}$$

for some $m, n, K \in \mathbb{N}$, $\alpha_i^k, \beta_j^k, \gamma_k \in \mathbb{K}$, $v_i^k \in B$, $w_j^k \in C$. We assume that $v \neq 0$ and $w \neq 0$ holds. Thus, we can take the bases B and C such that $v \in B$ and $w \in C$. Furthermore, we take the shortest possible representation, i. e., K to be minimal. That gives us that $\gamma_k \neq 0$ and that α^k and β^k cannot both be unit vectors for all k because both cases would lead to a zero summand and the sum could be shortened. Therefore, (v, w) is not contained in the second sum over k in (A.3), which yields $(v, w) = \sum_{k=1}^K \sum_{i=1}^m \sum_{j=1}^n \gamma_k \alpha_i^k \beta_j^k (v_i^k, w_j^k)$ and $\sum_{k=1}^K \gamma_k \left(\sum_{i=1}^m \alpha_i^k v_i^k, \sum_{j=1}^n \beta_j^k w_j^k \right) = 0$. Now assume that a non-zero pair is contained in the second sum. Since $\gamma_k \neq 0$, it has to be contained again at least once, such that it cancels out. Because of the unique basis representation the respective terms would also cancel in the first sum over k , which is a contradiction to the minimality of K . Therefore we have $K = 1$ and w. l. o. g. $\gamma_1 = 1$. This gives $\left(\sum_{i=1}^m \alpha_i^1 v_i^1, \sum_{j=1}^n \beta_j^1 w_j^1 \right) = 0$ and thus $\sum_{i=1}^m \alpha_i^1 v_i^1 = 0$ and $\sum_{j=1}^n \beta_j^1 w_j^1 = 0$ and therefore $\alpha^1 = 0$ and $\beta^1 = 0$ due to the linear independence. We obtain $(v, w) = 0$, which is a contradiction to the assumption that $v \neq 0$ and $w \neq 0$.

In the second case (A.2) we have that $(v, w) - (\tilde{v}, \tilde{w}) = (v, c\tilde{w}) - (cv, \tilde{w}) \in \mathcal{N}$ by taking $m = n = 1$, $\alpha_1 = c$, $\beta_1 = c^{-1}$, $v_1 = v$ and $w_1 = c\tilde{w}$.

It remains to prove the “only if” part. Let $(\tilde{v}, \tilde{w}) \in v \otimes w$, i. e., $(v, w) - (\tilde{v}, \tilde{w}) \in \mathcal{N}$. Assume that (A.1) does not hold meaning that the first part of (A.2) holds and we have to show the second part. If $v = \tilde{v}$ and $w = \tilde{w}$, we can take $c = 1$ and are done; now we assume $c \notin \{0, 1\}$.

First observe that if $\tilde{v} = cv$ holds, we get $w = c\tilde{w}$: Represent $(v, w) - (cv, \tilde{w})$ as a shortest linear combination of the form (A.3) with $v \in B$, $w \in C$. Since $cv \notin B$ we can argue as above that (v, w) is contained in the first part and (cv, \tilde{w}) is contained in the second part. Assuming that other pairs than (cv, \tilde{w}) are contained in the second sum, we see again that the terms containing these would cancel in both sums and that we can take $K = 1$ and

$\gamma_1 = 1$. We have $(cv, \tilde{w}) = \left(\sum_{i=1}^m \alpha_i^1 v_i^1, \sum_{j=1}^n \beta_j^1 w_j^1 \right)$ and thus $\alpha^1 = ce^1$ if we take w.l.o.g. $v_1^1 = v$, where e^i is the i -th unit vector. Since (v, w) is contained in the first part, we get $(v, w) = \sum_{j=1}^n c\beta_j^1 (v, w_j^1)$ and therefore $\beta^1 = c^{-1}e^1$ if $w_1^1 = w$. This gives us finally that $(cv, \tilde{w}) = (cv, c^{-1}w)$ and hence $w = c\tilde{w}$. Analogously it follows from $(v, w) - (\tilde{v}, \tilde{w}) \in \mathcal{N}$ and $w = c\tilde{w}$ that $\tilde{v} = cv$.

Now we assume that the second part of (A.2) does not hold. That means that w.l.o.g. v and \tilde{v} are linearly independent. We get from the considerations above that also w and \tilde{w} are linearly independent. In the representation of $(v, w) - (\tilde{v}, \tilde{w})$ as in (A.3) we can take bases B and C containing both vectors, respectively: $v, \tilde{v} \in B$, $w, \tilde{w} \in C$. Again we argue that $(v, w) - (\tilde{v}, \tilde{w})$ is contained in the first sum and that we get $K = 1$ and $\gamma_1 = 1$ leading to $0 = \left(\sum_{i=1}^m \alpha_i^1 v_i^1, \sum_{j=1}^n \beta_j^1 w_j^1 \right)$. From that we obtain $\alpha^1 = 0$ and $\beta^1 = 0$, which gives $(v, w) - (\tilde{v}, \tilde{w}) = (0, 0)$, which contradicts the assumption that v and \tilde{v} are linearly independent. Since v and \tilde{v} are linearly dependent and both non-zero, there exists a constant $c \in \mathbb{K} \setminus \{0\}$ such that $\tilde{v} = cv$ holds. As shown above, we get that $w = c\tilde{w}$ is also true. \square

A.2. General L^p Spaces and Operator Theory

We discuss some general results about embeddings and interpolation of L^p spaces, which are used frequently for error and regularity estimates. Additionally, we state some results on operators between Banach spaces used for existence theory and error estimation of the considered PDEs.

Proposition A.2 (Estimating L^p -norms on finite measure spaces). *Let $(\Omega, \mathcal{A}, \mu)$ be a measure space with $0 < \mu(\Omega) < \infty$. Let $p \in [1, \infty]$ and $q \in [1, p]$ be given and let $v \in L_\mu^p(\Omega)$. Then, $v \in L_\mu^q(\Omega)$ and we have the estimate*

$$\|v\|_{L_\mu^q(\Omega)} \leq C \|v\|_{L_\mu^p(\Omega)}$$

with $C = \mu(\Omega)^{1/q-1/p}$. In case that μ is a probability measure, $\|v\|_{L_\mu^q(\Omega)} \leq \|v\|_{L_\mu^p(\Omega)}$ holds.

Proof. By Hölder's inequality we get

$$\|v\|_{L_\mu^q(\Omega)}^q = \|1 \cdot |v|^q\|_{L_\mu^1(\Omega)} \leq \|1\|_{L_\mu^{(1-q/p)^{-1}}(\Omega)} \cdot \| |v|^q \|_{L_\mu^{p/q}(\Omega)} = \mu(\Omega)^{1-q/p} \cdot \|v\|_{L_\mu^p(\Omega)}^q$$

and therefore the desired result. As usual, we define $0^{-1} = \infty$ in the case $q = p$ and get $\|1\|_{L_\mu^\infty(\Omega)} = 1 = \mu(\Omega)^0$. If μ is a probability measure, $\mu(\Omega) = 1$ and thus $C = 1$ holds. \square

Proposition A.3. *Let $A, A_0 : Y \rightarrow Y^*$ be bounded, self-adjoint and boundedly invertible linear operators between a Hilbert space Y and its dual space. We define the inner product $(y, v)_A := \langle Ay, v \rangle_{Y^*, Y}$ and norm $\|y\|_A := \sqrt{(y, y)_A}$ for $y, v \in Y$, and analogously the inner products induced by A_0 on Y and by A^{-1}, A_0^{-1} on Y^* . The norms induced by A and A_0 shall be equivalent with the constants $\sqrt{\lambda}$ and $\sqrt{\Lambda}$ for $0 < \lambda \leq \Lambda$, i. e.,*

$$\lambda \|y\|_A^2 \leq \|y\|_{A_0}^2 \leq \Lambda \|y\|_A^2 \tag{A.4}$$

holds for every $y \in Y$.

Then, the estimate

$$\frac{1}{\Lambda} \|b\|_{A^{-1}}^2 \leq \|b\|_{A_0^{-1}}^2 \leq \frac{1}{\lambda} \|b\|_{A^{-1}}^2 \quad (\text{A.5})$$

holds for every $b \in Y^*$.

Proof. From the first inequality in (A.4), it follows that

$$\lambda \|y\|_A^2 \leq \langle A_0 y, y \rangle_{Y^*, Y} = \langle A A^{-1} A_0 y, y \rangle_{Y^*, Y} \leq \|A^{-1} A_0 y\|_A \|y\|_A = \|A_0 y\|_{A^{-1}} \|y\|_A$$

and therefore $\|y\|_A \leq \frac{1}{\lambda} \|A_0 y\|_{A^{-1}}$ holds for all $y \in Y$. This induces with $b = A_0 y$, that $\|b\|_{A_0^{-1}}^2 = \|y\|_{A_0}^2 \leq \frac{1}{\lambda} \|A_0 y\|_{A^{-1}}^2 = \frac{1}{\lambda} \|b\|_{A^{-1}}^2$ holds for every $b \in Y^*$, which is the second inequality in (A.5). We have used the Cauchy-Schwarz inequality and that the operators A and A_0 are invertible. The rest of the prerequisites is only needed to make sure that the respective operators induce inner products and that (A.4) holds.

In the same fashion, the second inequality in (A.4), written as $\frac{1}{\Lambda} \|y\|_{A_0}^2 \leq \|y\|_A^2$, induces the estimate $\|b\|_{A^{-1}}^2 \leq \Lambda \|b\|_{A_0^{-1}}^2$, which is equivalent to the first inequality in (A.5). \square

Proposition A.4. *Let $N : Y \rightarrow Y^*$ be a monotone operator between a real, reflexive Banach space Y and its dual space. Let N be Gâteaux-differentiable at $y \in Y$. Then, the derivative $N'(y) : Y \rightarrow Y^*$ is also a monotone operator.*

Proof. Since $N'(y)$ is linear, it is sufficient to prove that $\langle N'(y)v, v \rangle_{Y^*, Y} \geq 0$ holds for all $v \in Y$. By the definition of the Gâteaux derivative and the continuity of $N'(y)$, we have

$$\begin{aligned} \langle N'(y)v, v \rangle_{Y^*, Y} &= \lim_{t \rightarrow 0^+} \frac{1}{t} \langle N(y + tv) - N(y), v \rangle_{Y^*, Y} \\ &= \lim_{t \rightarrow 0^+} \frac{1}{t^2} \underbrace{\langle N(y + tv) - N(y), y + tv - y \rangle_{Y^*, Y}}_{\geq 0} \geq 0, \end{aligned}$$

which is non-negative because of the monotonicity of N . \square

A.3. Superposition Operators between L^p Spaces

In this section, we discuss the properties of nonlinear superposition operators between general L^p spaces. We use this theory later to show required properties of the superposition operators between the state space $L_{\mathbb{P}}^p(\Xi; H_0^1(\Omega))$ and its dual.

Theorem A.5 (Well-definedness and continuity of superposition operators between L^p spaces). *Let $(\Omega, \mathcal{A}, \mu)$ be a measure space with a σ -finite measure μ and let $p, q \in [1, \infty)$. Let $\varphi : \Omega \times \mathbb{R} \rightarrow \mathbb{R}$ be a Carathéodory function, i. e., the function $\varphi(\cdot, t)$ is measurable for all $t \in \mathbb{R}$ and the function $\varphi(\omega, \cdot)$ is continuous for a. e. $\omega \in \Omega$. Furthermore, let φ fulfill the growth condition*

$$|\varphi(\omega, t)| \leq c_q(\omega) + c_\varphi |t|^{p/q} \quad \text{for all } t \in \mathbb{R} \text{ and a. e. } \omega \in \Omega$$

with some function $c_q \in L_\mu^q(\Omega)$ and a constant $c_\varphi \geq 0$.

Then, the superposition operator

$$N : L^p_\mu(\Omega) \rightarrow L^q_\mu(\Omega), \quad N(y)(\omega) := \varphi(\omega, y(\omega))$$

is well-defined and continuous.

Proof. The theorem is a consequence of [5, Thm. 3.1, Lem. 1.5, Thm. 1.1, Thm. 3.7]. [5, Thm. 3.1] provides a necessary and sufficient condition for the operator N to be well-defined. We use only the sufficient part and tighten the conditions a bit by requiring the growth conditions for (almost) every ω , t and using the Carathéodory property of φ and [5, Lem. 1.5, Thm. 1.1]. [5, Thm. 3.7] provides the continuity of the operator (and also some equivalence relation in a more general framework). \square

Remark A.6. Note that [5] is a suitable reference for this quite known statement although it covers the topic in a very general setting so that it is necessary to combine various theorems and lemmas from this book. Other references such as [48] restrict the discussion to domains equipped with the Lebesgue measure so that, e. g., the atomic part of the probability measure \mathbb{P} is not taken into account. This is especially important when providing necessary conditions for, e. g., boundedness, which can be seen in the following theorem, whose analog [48, Thm. 3] lacks this generality.

Theorem A.7 (Boundedness of superposition operators between L^p spaces). *Let $(\Omega, \mathcal{A}, \mu)$ be a measure space with a σ -finite measure μ and let $p, q \in [1, \infty)$. Let $\varphi : \Omega \times \mathbb{R} \rightarrow \mathbb{R}$ be a sup-measurable function and let $N : L^p_\mu(\Omega) \rightarrow L^q_\mu(\Omega)$ (defined as in Theorem A.5) be well-defined.*

Then, N is locally bounded if and only if the function $\varphi(\omega, \cdot)$ is bounded for each $\omega \in \Omega_d$, where $\Omega_d \subset \Omega$ is the purely atomic part of the measure μ .

Proof. See [5, Thm. 3.2]. \square

Remark A.8.

- If $\varphi : \Omega \times \mathbb{R} \rightarrow \mathbb{R}$ is a Carathéodory function, it is sup-measurable [5, Thm. 1.1].
- The fact that N is locally bounded means that $\limsup_{v \rightarrow y} \|N(v)\|_{L^q_\mu(\Omega)} < \infty$ holds for all $y \in L^p_\mu(\Omega)$.
- If μ is non-atomic, Theorem A.7 gives that the superposition operator N is automatically locally bounded if it is well-defined.

Theorem A.9 (Differentiability of superposition operators between L^p spaces). *Let $(\Omega, \mathcal{A}, \mu)$ be a measure space with a σ -finite measure μ and let $p, q \in [1, \infty)$ with $p > q$. Let $\varphi : \Omega \times \mathbb{R} \rightarrow \mathbb{R}$ be a Carathéodory function (see Theorem A.7) and let the restriction $\varphi(\omega, \cdot)$ be continuously differentiable for a. e. $\omega \in \Omega$. Let the function φ fulfill the growth condition*

$$|\varphi(\omega, t)| \leq c_q(\omega) + c_\varphi |t|^{p/q} \quad \text{for all } t \in \mathbb{R} \text{ and a. e. } \omega \in \Omega \quad (\text{A.6})$$

with some function $c_q \in L_\mu^q(\Omega)$ and a constant $c_\varphi \geq 0$. Furthermore, let the partial derivative $\varphi_t : \Omega \times \mathbb{R} \rightarrow \mathbb{R}$ fulfill the growth condition

$$|\varphi_t(\omega, t)| \leq c'_{pq/(p-q)}(\omega) + c'_{\varphi_t} |t|^{(p-q)/q} \quad \text{for all } t \in \mathbb{R} \text{ and a. e. } \omega \in \Omega \quad (\text{A.7})$$

with some function $c'_{pq/(p-q)} \in L_\mu^{pq/(p-q)}(\Omega)$ and a constant $c'_{\varphi_t} \geq 0$.

Then, the superposition operator $N : L_\mu^p(\Omega) \rightarrow L_\mu^q(\Omega)$ induced by φ (defined as in Theorem A.5) is continuously Fréchet-differentiable with derivative

$$[N'(y)v](\omega) = \varphi_t(\omega, y(\omega))v(\omega).$$

Proof. We adapt the proof from the very similar Proposition A.11 in [111].

Since φ is a Carathéodory function and condition (A.6) holds, the operator N is well-defined and continuous by Theorem A.5.

Moreover, φ_t is a Carathéodory function: Since $\varphi(\omega, \cdot)$ is continuously differentiable for a. e. $\omega \in \Omega$, the partial derivative $\varphi_t(\omega, \cdot)$ is continuous for a. e. $\omega \in \Omega$. By definition, $\varphi_t(\omega, t) = \lim_{h \rightarrow 0} \frac{1}{h}(\varphi(\omega, t+h) - \varphi(\omega, t))$ holds. Now, for a sequence $(h_n)_{n \in \mathbb{N}} \subset \mathbb{R} \setminus \{0\}$ converging to zero, we have that the difference quotient $\varphi_t^{(n)}(\omega, t) := \frac{1}{h_n}(\varphi(\omega, t+h_n) - \varphi(\omega, t))$ is measurable w. r. t. ω for all $t \in \mathbb{R}$ as sum of measurable functions. The pointwise limit $\varphi_t(\cdot, t)$ of these measurable functions is thus also measurable.

Due to the Carathéodory property and condition (A.7), the superposition operator $N_{\varphi_t} : L_\mu^p(\Omega) \rightarrow L_\mu^{pq/(p-q)}(\Omega)$ generated by φ_t is also well-defined and continuous by Theorem A.5. Note that the identification $\mathcal{L}(L_\mu^p(\Omega), L_\mu^q(\Omega)) \cong L_\mu^{pq/(p-q)}(\Omega)$ can be made by the Riesz representation theorem and Hölder's inequality. This gives $N'(y) \cong N_{\varphi_t}(y)$ and that the derivative $N' : L_\mu^p(\Omega) \rightarrow \mathcal{L}(L_\mu^p(\Omega), L_\mu^q(\Omega))$ as defined in the theorem is well-defined and continuous, cf. the proof of [111, Prop. A.11] for more details.

Moreover,

$$\begin{aligned} & \|N(y+v) - N(y) - N'(y)v\|_{L_\mu^q(\Omega)} \\ &= \|\varphi(\cdot, y(\cdot) + v(\cdot)) - \varphi(\cdot, y(\cdot)) - \varphi_t(\cdot, y(\cdot))v(\cdot)\|_{L_\mu^q(\Omega)} \\ &= \left\| \int_0^1 [\varphi_t(\cdot, y(\cdot) + \sigma v(\cdot)) - \varphi_t(\cdot, y(\cdot))]v(\cdot) \, d\sigma \right\|_{L_\mu^q(\Omega)} \\ &\leq \int_0^1 \|\varphi_t(\cdot, y(\cdot) + \sigma v(\cdot)) - \varphi_t(\cdot, y(\cdot))\|_{L_\mu^{pq/(p-q)}(\Omega)} \|v(\cdot)\|_{L_\mu^p(\Omega)} \, d\sigma \\ &\leq \left(\int_0^1 \|\varphi_t(\cdot, y(\cdot) + \sigma v(\cdot)) - \varphi_t(\cdot, y(\cdot))\|_{L_\mu^{pq/(p-q)}(\Omega)} \, d\sigma \right) \|v(\cdot)\|_{L_\mu^p(\Omega)} \\ &\leq \left(\sup_{\sigma \in [0,1]} \|\varphi_t(\cdot, y(\cdot) + \sigma v(\cdot)) - \varphi_t(\cdot, y(\cdot))\|_{L_\mu^{pq/(p-q)}(\Omega)} \right) \|v(\cdot)\|_{L_\mu^p(\Omega)} \\ &= o(\|v\|_{L_\mu^p(\Omega)}) \quad \text{as } \|v\|_{L_\mu^p(\Omega)} \rightarrow 0 \end{aligned}$$

follows from Hölder's inequality and the fact that N_{φ_t} is continuous. This shows that N is Fréchet-differentiable with the given derivative N' . \square

Alternatively, if μ is a finite measure, it is sufficient to require $\varphi(\cdot, 0) \in L_\mu^q(\Omega)$ and a suitable growth condition only on $\varphi_t(\omega, \cdot)$. Then, condition (A.6) follows directly:

Lemma A.10. *Let $(\Omega, \mathcal{A}, \mu)$ be a measure space with $\mu(\Omega) < \infty$ and let $p, q \in [1, \infty)$ with $p > q$. Let $\varphi : \Omega \times \mathbb{R} \rightarrow \mathbb{R}$ be a Carathéodory function, see Theorem A.7, and let the restriction $\varphi(\omega, \cdot)$ be continuously differentiable for a. e. $\omega \in \Omega$. Furthermore, let the partial derivative $\varphi_t : \Omega \times \mathbb{R} \rightarrow \mathbb{R}$ fulfill the ω -independent growth condition*

$$|\varphi_t(\omega, t)| \leq a'_{\varphi_t} + c'_{\varphi_t} |t|^{(p-q)/q} \quad \text{for all } t \in \mathbb{R} \text{ and a. e. } \omega \in \Omega,$$

with two constants $a'_{\varphi_t}, c'_{\varphi_t} \geq 0$ and let $\varphi(\cdot, 0) \in L^q_{\mu}(\Omega)$. Then, the function $\varphi(\omega, \cdot)$ fulfills

$$|\varphi(\omega, t)| \leq c_q(\omega) + c_{\varphi} |t|^{p/q}$$

for a. e. $\omega \in \Omega$ with $c_q \in L^q(\Omega)$ and a constant $c_{\varphi} \in \mathbb{R} \geq 0$.

Proof. We use a part of the proof of [111, Prop. A.11]. For a. e. $\omega \in \Omega$, we have

$$\begin{aligned} |\varphi(\omega, t)| &\leq |\varphi(\omega, 0)| + \int_0^1 |\varphi_t(\omega, \sigma t) t| \, d\sigma \\ &\leq |\varphi(\omega, 0)| + |t| \int_0^1 a'_{\varphi_t} + c'_{\varphi_t} |\sigma t|^{(p-q)/q} \, d\sigma \\ &\leq |\varphi(\omega, 0)| + a'_{\varphi_t} |t| + \frac{qc'_{\varphi_t}}{p} |t|^{p/q} \leq c_q(\omega) + c_{\varphi} |t|^{p/q} \end{aligned}$$

with some $c_{\varphi} \geq 0$ and $c_q \in L^q_{\mu}(\Omega)$ because $\varphi(\cdot, 0) \in L^q(\Omega)$ and constant functions belong to every $L^q_{\mu}(\Omega)$ since μ is finite. \square

A.4. Superposition Operators from $L^p_{\mathbb{P}}(\Xi; H_0^1(\Omega))$ to Its Dual

For an open, bounded Lipschitz domain $\Omega \subset \mathbb{R}^n$ ($n \in \{2, 3\}$) equipped with the Lebesgue measure λ , we define the deterministic state space $Y := H_0^1(\Omega)$ and its dual $Y^* = H^{-1}(\Omega)$. Furthermore, let $\Xi \subset \mathbb{R}^m$ be measurable and equipped with the probability measure \mathbb{P} . For $p \in [2, \infty)$ (if $n = 2$) and $p \in [2, 6]$ (if $n = 3$) we define the Bochner space $\mathbf{Y} := L^p_{\mathbb{P}}(\Xi; Y)$. We have the Sobolev embedding $H_0^1(\Omega) \hookrightarrow L^q(\Omega)$, i. e., $\|y\|_{L^q(\Omega)} \leq C_q \|y\|_{H_0^1(\Omega)}$ with some constant $C_q > 0$ for $q \in [1, \infty)$ if $n = 2$ and $q \in [1, 6]$ if $n = 3$. Now, for $p = q$ we can embed the space $\mathbf{Y} \hookrightarrow L^p_{\mathbb{P}}(\Xi; L^p(\Omega)) \cong L^p_{\lambda \otimes \mathbb{P}}(\Omega \times \Xi)$ and then use existing L^p -theory. The respective dual spaces are identified with $L^{p^*}_{\lambda \otimes \mathbb{P}}(\Omega \times \Xi) \cong L^{p^*}_{\mathbb{P}}(\Xi; L^{p^*}(\Omega)) \hookrightarrow L^{p^*}_{\mathbb{P}}(\Xi; H^{-1}(\Omega)) = \mathbf{Y}^*$ with the conjugate exponent $p^* = \frac{p}{p-1} \in (1, 2]$ if $n = 2$ and $p^* \in [\frac{1}{6}, 1]$ if $n = 3$.

Now we consider a function $\varphi : \Omega \times \Xi \times \mathbb{R} \rightarrow \mathbb{R}$ with some of the following properties:

Assumption A.11.

1. φ satisfies the Carathéodory property, i. e., $\varphi(\cdot, \cdot, t)$ is measurable for all $t \in \mathbb{R}$ and $\varphi(x, \xi, \cdot)$ is continuous for a. e. $x \in \Omega$, $\xi \in \Xi$.
2. φ satisfies the growth condition $|\varphi(x, \xi, t)| \leq c_{p^*}(x, \xi) + c_{\varphi} |t|^{p-1}$ for all $t \in \mathbb{R}$ and a. e. $x \in \Omega$, $\xi \in \Xi$ with some constant $c_{p^*} \geq 0$ and a function $c_{p^*} \in L^{p^*}_{\lambda \otimes \mathbb{P}}(\Omega \times \Xi)$.

3. The function $\varphi(x, \xi, \cdot)$ is monotonically increasing for a. e. $x \in \Omega$, $\xi \in \Xi$.
4. For a. e. $x \in \Omega$, $\xi \in \Xi$, the function $\varphi(x, \xi, \cdot)$ is continuously differentiable with partial derivative $\varphi_t : \Omega \times \Xi \times \mathbb{R} \rightarrow \mathbb{R}$.
5. It holds that $p > 2$ and the first partial derivative satisfies the growth condition

$$|\varphi_t(x, \xi, t)| \leq c'_{p/(p-2)}(x, \xi) + c'_{\varphi_t} |t|^{p-2} \text{ for all } t \in \mathbb{R} \text{ and a. e. } x \in \Omega, \xi \in \Xi$$

with some constant $c'_{\varphi_t} \geq 0$ and a function $c'_{p/(p-2)} \in L^{p/(p-2)}_{\lambda \otimes \mathbb{P}}(\Omega \times \Xi)$.

6. For a. e. $x \in \Omega$, $\xi \in \Xi$, the function $\varphi(x, \xi, \cdot)$ is twice continuously differentiable with first partial derivative $\varphi_t : \Omega \times \Xi \times \mathbb{R} \rightarrow \mathbb{R}$ and second partial derivative $\varphi_{tt} : \Omega \times \Xi \times \mathbb{R} \rightarrow \mathbb{R}$.

7. It holds that $p > 3$ and the second partial derivative satisfies the growth condition

$$|\varphi_{tt}(x, \xi, t)| \leq c''_{p/(p-3)}(x, \xi) + c''_{\varphi_{tt}} |t|^{p-3} \text{ for all } t \in \mathbb{R} \text{ and a. e. } x \in \Omega, \xi \in \Xi$$

with some constant $c''_{\varphi_{tt}} \geq 0$ and a function $c''_{p/(p-3)} \in L^{p/(p-3)}_{\lambda \otimes \mathbb{P}}(\Omega \times \Xi)$.

Using φ , we define the nonlinear operator

$$\mathbf{N} : \mathbf{Y} \rightarrow \mathbf{Y}^*, \quad \langle \mathbf{N}(\mathbf{y}), \mathbf{v} \rangle_{\mathbf{Y}^*, \mathbf{Y}} := \int_{\Xi} \int_{\Omega} \varphi(x, \xi, \mathbf{y}(x, \xi)) \mathbf{v}(x, \xi) \, dx \, d\mathbb{P} \quad \forall \mathbf{y}, \mathbf{v} \in \mathbf{Y}.$$

We investigate its properties:

Proposition A.12. *Under Assumption A.11:1-2, the operator $\mathbf{N} : \mathbf{Y} \rightarrow \mathbf{Y}^*$ is well-defined and continuous.*

Proof. We have that \mathbf{N} is well-defined and continuous as an operator from $L^p_{\lambda \otimes \mathbb{P}}(\Omega \times \Xi)$ to $L^{p^*}_{\lambda \otimes \mathbb{P}}(\Omega \times \Xi)$ by Theorem A.5. Using the continuous embeddings $\mathbf{Y} \hookrightarrow L^p_{\lambda \otimes \mathbb{P}}(\Omega \times \Xi)$ and $\mathbf{Y}^* \hookrightarrow L^{p^*}_{\lambda \otimes \mathbb{P}}(\Omega \times \Xi)$ as described above, we get that it is also well-defined and continuous from \mathbf{Y} to \mathbf{Y}^* . \square

Proposition A.13. *Given Assumption A.11:1-3, the operator $\mathbf{N} : \mathbf{Y} \rightarrow \mathbf{Y}^*$ is monotone.*

Proof. By Proposition A.12, the operator \mathbf{N} is well-defined. Since $\varphi(x, \xi, \cdot)$ is increasing for a. e. x, ξ by Assumption A.11:3, $\varphi(x, \xi, t) \leq \varphi(x, \xi, \tilde{t})$ holds if $t \leq \tilde{t}$. This gives $\varphi(x, \xi, \tilde{t}) - \varphi(x, \xi, t) \geq 0$ if $\tilde{t} - t \geq 0$ and $\varphi(x, \xi, t) - \varphi(x, \xi, \tilde{t}) \leq 0$ if $t - \tilde{t} \geq 0$, giving $(\varphi(x, \xi, t) - \varphi(x, \xi, \tilde{t}))(t - \tilde{t}) \geq 0$ for all $t, \tilde{t} \in \mathbb{R}$ and a. e. x, ξ . Hence, we get

$$\begin{aligned} \langle \mathbf{N}(\mathbf{y}) - \mathbf{N}(\tilde{\mathbf{y}}), \mathbf{y} - \tilde{\mathbf{y}} \rangle_{\mathbf{Y}^*, \mathbf{Y}} &= \\ \int_{\Xi} \int_{\Omega} (\varphi(x, \xi, \mathbf{y}(x, \xi)) - \varphi(x, \xi, \tilde{\mathbf{y}}(x, \xi))) (\mathbf{y}(x, \xi) - \tilde{\mathbf{y}}(x, \xi)) \, dx \, d\mathbb{P} &\geq 0 \end{aligned}$$

for all $\mathbf{y}, \tilde{\mathbf{y}} \in \mathbf{Y}$, showing the monotonicity of \mathbf{N} . \square

Proposition A.14. *Under Assumption A.11:1-2,4-5, the operator $\mathbf{N} : \mathbf{Y} \rightarrow \mathbf{Y}^*$ is continuously Fréchet-differentiable with derivative*

$$[\mathbf{N}'(\mathbf{y})\mathbf{v}](x, \xi) = \varphi_t(x, \xi, \mathbf{y}(x, \xi))\mathbf{v}(x, \xi).$$

Proof. We have that \mathbf{N} is continuously Fréchet-differentiable with the given derivative as an operator from $L_{\lambda \otimes \mathbb{P}}^p(\Omega \times \Xi)$ to $L_{\lambda \otimes \mathbb{P}}^{p^*}(\Omega \times \Xi)$ by Theorem A.9, having $q = p^* = \frac{p}{p-1}$, and thus $\frac{pq}{p-q} = \frac{p}{p-2}$ and $\frac{p-q}{q} = p-2$. Note that we require $p > p^*$ which is fulfilled if and only if $p > 2$. Again, the continuous embeddings $\mathbf{Y} \hookrightarrow L_{\lambda \otimes \mathbb{P}}^p(\Omega \times \Xi)$ and $L_{\lambda \otimes \mathbb{P}}^{p^*}(\Omega \times \Xi) \hookrightarrow \mathbf{Y}^*$ and the chain rule yield continuous F-differentiability from \mathbf{Y} to \mathbf{Y}^* . \square

Remark A.15. For $p = 2$ we can consider the subspace of $L_{\mathbb{P}}^2(\Xi; Y)$ of functions of the form $\mathbf{y}(x, \xi) = y(x)v(\xi)$ with some fixed $y \in Y$ with $\|y\|_Y = 1$ and arbitrary $v \in L_{\mathbb{P}}^2(\Xi)$. This subspace is isomorphic to $L_{\mathbb{P}}^2(\Xi)$. [5, Thm. 3.13] states that if a Carathéodory function $\varphi : \Xi \times \mathbb{R} \rightarrow \mathbb{R}$ generates a well-defined, F-differentiable superposition operator acting from $L_{\mathbb{P}}^2(\Xi)$ into itself, the function φ has the form $\varphi(\xi, t) = a(\xi) + b(\xi)t$ for a. e. $\xi \in \Xi_c$ with $a \in L_{\mathbb{P}}^2(\Xi)$ and $b \in L_{\mathbb{P}}^\infty(\Xi)$. $\Xi_c \subset \Xi$ is the non-atomic part of the probability measure \mathbb{P} . Hence, we need to require $p > 2$ if we want to work with “true” nonlinearities.

Proposition A.16. *Under Assumption A.11:1-2,5-7, the operator $\mathbf{N} : \mathbf{Y} \rightarrow \mathbf{Y}^*$ is twice continuously Fréchet-differentiable with second derivative*

$$[[\mathbf{N}''(\mathbf{y})\mathbf{v}]\tilde{\mathbf{v}}](x, \xi) = \varphi_{tt}(x, \xi, \mathbf{y}(x, \xi))\mathbf{v}(x, \xi)\tilde{\mathbf{v}}(x, \xi).$$

Proof. It follows from Proposition A.14 that \mathbf{N} is continuously F-differentiable. We show that \mathbf{N}' is again continuously F-differentiable with the given derivative. For this purpose we identify $\mathcal{L}(L^p, L^{p^*}) \cong L^{p/(p-2)}$ and consider the continuously F-differentiable superposition operator $\mathbf{N}_{\varphi_t} : L_{\lambda \otimes \mathbb{P}}^p(\Omega \times \Xi) \rightarrow L_{\lambda \otimes \mathbb{P}}^{p/(p-2)}(\Omega \times \Xi)$ induced by the Carathéodory function $\varphi_t : \Omega \times \Xi \times \mathbb{R} \rightarrow \mathbb{R}$, cf. the proof of Theorem A.9. It is well-defined and continuous, see Proposition A.12. Assumption A.11:6-7 and Theorem A.9 yield that it is even continuously F-differentiable. There we use $q = \frac{p}{p-2}$ and thus $p > q$ if and only if $p > 3$, $\frac{p}{q} = p-2$, $\frac{pq}{p-q} = \frac{p}{p-3}$, and $\frac{p-q}{q} = p-3$. Its derivative is $[\mathbf{N}'_{\varphi_t}(\mathbf{y})\mathbf{v}](x, \xi) = \varphi_{tt}(x, \xi, \mathbf{y}(x, \xi))\mathbf{v}(x, \xi)$. It is identified with the derivative given above. \square

A.5. Alternative Approach for the Discussion of the Semilinear, Elliptic PDE

We present an alternative approach for the existence of a unique solution of the PDE discussed in Section 3.2 for the more concrete case $\varphi(t) = \hat{\varphi}(x, \xi, t) = t^3$ and the choice $p = 4$. Recall that the state space is $\mathbf{Y} = L_{\mathbb{P}}^4(\Xi; H_0^1(\Omega))$. Since $\hat{\varphi}$ fulfills Assumption A.11 with $c_{p^*} \equiv 0$, $c_\varphi = 1$, $c'_{p/(p-2)} \equiv 0$, $c'_{\varphi_t} = 3$, $c''_{p/(p-3)} \equiv 0$, $c''_{\varphi_{tt}} = 6$, the superposition operator $\mathbf{N} : \mathbf{Y} \rightarrow \mathbf{Y}^*$ is twice continuously differentiable and monotone by the considerations in Section A.4. It follows that the operator $\tilde{\mathbf{N}} \equiv \mathbf{A} + \mathbf{N} : \mathbf{Y} \rightarrow \mathbf{Y}^*$ is a strictly monotone and twice continuously differentiable.

Unfortunately, the operator $\tilde{\mathbf{N}}$ is not coercive in \mathbf{Y} in the sense of Theorem 3.7: For a counterexample we choose $\Omega := (0, 1)^2 \subset \mathbb{R}^2$ and $\Xi := (0, 1) \subset \mathbb{R}$ with the uniform distribution. For $x \in \Omega$ we define

$$y_\varepsilon(x) := \begin{cases} \varepsilon^q x_1, & x_1 \leq x_2 \leq 2\varepsilon - x_1, \\ \varepsilon^q x_2, & x_2 < x_1 \text{ and } x_2 \leq 2\varepsilon - x_1, \\ -\varepsilon^q x_1 + 2\varepsilon^{q+1}, & 2\varepsilon - x_1 < x_2 \leq x_1 \text{ and } x_1 < 2\varepsilon, \\ -\varepsilon^q x_2 + 2\varepsilon^{q+1}, & x_1 < x_2 < 2\varepsilon \text{ and } 2\varepsilon - x_1 < x_2, \\ 0, & \text{otherwise} \end{cases}$$

for $\varepsilon \in (0, \frac{1}{2})$ and $q \in (-\frac{5}{4}, -\frac{3}{4})$, e. g., $q = -1$. Note that y_ε is continuous and belongs to $H_0^1(\Omega)$. In fact it is a multiple of a finite element ansatz function defined on triangles of equal area ε^2 . The weak derivative is, up to sets of measure zero,

$$\nabla y_\varepsilon(x) := \begin{cases} (\varepsilon^q, 0)^\top, & x_1 \leq x_2 \leq 2\varepsilon - x_1, \\ (0, \varepsilon^q)^\top, & x_2 < x_1 \text{ and } x_2 \leq 2\varepsilon - x_1, \\ (-\varepsilon^q, 0)^\top, & 2\varepsilon - x_1 < x_2 \leq x_1 \text{ and } x_1 < 2\varepsilon, \\ (0, -\varepsilon^q)^\top, & x_1 < x_2 < 2\varepsilon \text{ and } 2\varepsilon - x_1 < x_2, \\ (0, 0)^\top, & \text{otherwise.} \end{cases}$$

This yields

$$\|y_\varepsilon\|_{H_0^1(\Omega)}^2 = \int_\Omega \|\nabla y_\varepsilon(x)\|_{L^2(\Omega)^2}^2 dx = 4\varepsilon^{2q+2}$$

and

$$\begin{aligned} \|y_\varepsilon\|_{L^4(\Omega)}^4 &= \int_\Omega y_\varepsilon(x)^4 dx = 4 \int_0^\varepsilon \int_{x_1}^{2\varepsilon-x_1} \varepsilon^{4q} x_1^4 dx_2 dx_1 \\ &= 8\varepsilon^{4q} \int_0^\varepsilon x_1^4 (\varepsilon - x_1) dx_1 = \frac{4}{15} \varepsilon^{4q+6}. \end{aligned}$$

Furthermore, we define $v_\varepsilon(\xi) := \xi^{-1/2+\varepsilon}$ and get $\|v_\varepsilon\|_{L_{\mathbb{P}}^1(\Xi)} = \int_0^1 \xi^{-1/2+\varepsilon} d\xi = \frac{2}{1+2\varepsilon}$ as well as $\|v_\varepsilon\|_{L_{\mathbb{P}}^2(\Xi)}^2 = \int_0^1 \xi^{-1+2\varepsilon} d\xi = \frac{1}{2\varepsilon}$. For $\mathbf{y}_\varepsilon(x, \xi) := y_\varepsilon(x) v_\varepsilon(\xi)^{1/2}$ we obtain

$$\begin{aligned} \|\mathbf{y}_\varepsilon\|_{L_{\mathbb{P}}^2(\Xi; H_0^1(\Omega))}^2 &= \int_0^1 v_\varepsilon(\xi) \|y_\varepsilon\|_{H_0^1(\Omega)}^2 d\xi = \frac{8\varepsilon^{2q+2}}{1+2\varepsilon}, \\ \|\mathbf{y}_\varepsilon\|_{L_{\mathbb{P}}^4(\Xi; H_0^1(\Omega))}^4 &= \int_0^1 v_\varepsilon(\xi)^2 \|y_\varepsilon\|_{H_0^1(\Omega)}^4 d\xi = 8\varepsilon^{4q+3}, \\ \|\mathbf{y}_\varepsilon\|_{L_{\mathbb{P}}^4(\Xi; L^4(\Omega))}^4 &= \int_0^1 v_\varepsilon(\xi)^2 \|y_\varepsilon\|_{L^4(\Omega)}^4 d\xi = \frac{2}{15} \varepsilon^{4q+5}. \end{aligned}$$

This gives

$$\begin{aligned} \frac{\langle \mathbf{A}\mathbf{y}_\varepsilon + \mathbf{N}(\mathbf{y}_\varepsilon), \mathbf{y}_\varepsilon \rangle_{\mathbf{Y}^*, \mathbf{Y}}}{\|\mathbf{y}_\varepsilon\|_{\mathbf{Y}}} &\leq \frac{\bar{\kappa} \|\mathbf{y}_\varepsilon\|_{L_{\mathbb{P}}^2(\Xi; H_0^1(\Omega))}^2 + \|\mathbf{y}_\varepsilon\|_{L_{\mathbb{P}}^4(\Xi; L^4(\Omega))}^4}{\|\mathbf{y}_\varepsilon\|_{L_{\mathbb{P}}^4(\Xi; H_0^1(\Omega))}} \\ &= \frac{\bar{\kappa} \frac{8\varepsilon^{2q+2}}{1+2\varepsilon} + \frac{2}{15} \varepsilon^{4q+5}}{8^{1/4} \varepsilon^{q+3/4}} = \frac{8^{3/4} \bar{\kappa}}{1+2\varepsilon} \varepsilon^{q+5/4} + \frac{2^{1/4}}{15} \varepsilon^{3q+17/4}. \end{aligned}$$

By $q + \frac{3}{4} < 0$ we have $\|\mathbf{y}_\varepsilon\|_{L_{\mathbb{P}}^4(\Xi; H_0^1(\Omega))} = 8^{1/4} \varepsilon^{q+3/4} \xrightarrow{\varepsilon \rightarrow 0^+} \infty$, but

$$0 \leq \frac{\langle \mathbf{A}\mathbf{y}_\varepsilon + \mathbf{N}(\mathbf{y}_\varepsilon), \mathbf{y}_\varepsilon \rangle_{\mathbf{Y}^*, \mathbf{Y}}}{\|\mathbf{y}_\varepsilon\|_{\mathbf{Y}}} \leq \frac{8^{3/4} \varepsilon^{q+5/4}}{1+2\varepsilon} + \frac{2^{1/4}}{15} \varepsilon^{3(q+17/12)} \xrightarrow{\varepsilon \rightarrow 0^+} 0$$

holds by $q + \frac{5}{4} > 0$ and $q + \frac{17}{12} > q + \frac{15}{12} > 0$. Therefore, the operator $\tilde{\mathbf{N}}$ is *not* coercive on \mathbf{Y} .

To be able to apply Theorem 3.7 we show coercivity and the other prerequisites on a larger space $\tilde{\mathbf{Y}} \supset \mathbf{Y}$, namely $\tilde{\mathbf{Y}} := L_{\mathbb{P}}^2(\Xi; H_0^1(\Omega)) \cap L_{\mathbb{P}}^4(\Xi; L^4(\Omega))$, to infer that (3.12) has a unique solution in $\tilde{\mathbf{Y}}$. After that, we will see that this solution actually belongs to \mathbf{Y} . The space $\tilde{\mathbf{Y}}$ is equipped with the norm $\|\mathbf{y}\|_{\tilde{\mathbf{Y}}} := \|\mathbf{y}\|_{L_{\mathbb{P}}^2(\Xi; H_0^1(\Omega))} + \|\mathbf{y}\|_{L_{\mathbb{P}}^4(\Xi; L^4(\Omega))}$. By $L_{\mathbb{P}}^4(\Xi) \hookrightarrow L_{\mathbb{P}}^2(\Xi)$ (Proposition A.2) and $H_0^1(\Omega) \hookrightarrow L^4(\Omega)$ (Sobolev embedding) we have $\mathbf{Y} \subset \tilde{\mathbf{Y}}$.

The operator $\mathbf{A} : \hat{\mathbf{Y}} := L_{\mathbb{P}}^2(\Xi; H_0^1(\Omega)) \rightarrow L_{\mathbb{P}}^2(\Xi; H^{-1}(\Omega)) = \hat{\mathbf{Y}}^*$ is well-defined since $A \in L_{\mathbb{P}}^\infty(\Xi; \mathcal{L}(Y, Y^*))$, see Proposition 3.1, and thus also as an operator from $\tilde{\mathbf{Y}}$ to $\tilde{\mathbf{Y}}^*$: Since $\tilde{\mathbf{Y}} \subset \hat{\mathbf{Y}}$, meaning that $\|\mathbf{y}\|_{\tilde{\mathbf{Y}}} \leq c \|\mathbf{y}\|_{\hat{\mathbf{Y}}}$ for all $\mathbf{y} \in \tilde{\mathbf{Y}}$ (here with constant $c = 1$), we get

$$\|\mathbf{A}\|_{\mathcal{L}(\tilde{\mathbf{Y}}, \tilde{\mathbf{Y}}^*)} = \sup_{\substack{\mathbf{y}, \mathbf{v} \in \tilde{\mathbf{Y}}, \\ \|\mathbf{y}\|_{\tilde{\mathbf{Y}}} \leq 1, \|\mathbf{v}\|_{\tilde{\mathbf{Y}}} \leq 1}} \langle \mathbf{A}\mathbf{y}, \mathbf{v} \rangle_{\tilde{\mathbf{Y}}^*, \tilde{\mathbf{Y}}} \leq \sup_{\substack{\mathbf{y}, \mathbf{v} \in \hat{\mathbf{Y}}, \\ \|\mathbf{y}\|_{\hat{\mathbf{Y}}} \leq c, \|\mathbf{v}\|_{\hat{\mathbf{Y}}} \leq c}} \langle \mathbf{A}\mathbf{y}, \mathbf{v} \rangle_{\hat{\mathbf{Y}}^*, \hat{\mathbf{Y}}} = c^2 \|\mathbf{A}\|_{\mathcal{L}(\hat{\mathbf{Y}}, \hat{\mathbf{Y}}^*)}.$$

From the Theorem A.5 we have that $\mathbf{N} : L_{\mathbb{P}}^4(\Xi; L^4(\Omega)) \rightarrow L_{\mathbb{P}}^{4/3}(\Xi; L^{4/3}(\Omega))$ is well-defined and continuous and thus also well-defined and continuous as an operator from $\tilde{\mathbf{Y}}$ to $\tilde{\mathbf{Y}}^*$ by the same argument. This proves that $\tilde{\mathbf{N}} : \tilde{\mathbf{Y}} \rightarrow \tilde{\mathbf{Y}}^*$ is well-defined and continuous, in particular, hemicontinuous. Strict monotonicity of $\tilde{\mathbf{N}}$ is shown exactly as in Proposition A.13.

Now we prove coercivity. For $\|\mathbf{y}\|_{L_{\mathbb{P}}^2(\Xi; H_0^1(\Omega))} + \|\mathbf{y}\|_{L_{\mathbb{P}}^4(\Xi; L^4(\Omega))} \rightarrow \infty$ we have

$$\frac{\langle \mathbf{A}\mathbf{y} + \mathbf{N}(\mathbf{y}), \mathbf{y} \rangle_{\tilde{\mathbf{Y}}^*, \tilde{\mathbf{Y}}}}{\|\mathbf{y}\|_{\tilde{\mathbf{Y}}}} \geq \frac{\kappa \|\mathbf{y}\|_{L_{\mathbb{P}}^2(\Xi; H_0^1(\Omega))}^2 + \|\mathbf{y}\|_{L_{\mathbb{P}}^4(\Xi; L^4(\Omega))}^4}{\|\mathbf{y}\|_{L_{\mathbb{P}}^2(\Xi; H_0^1(\Omega))} + \|\mathbf{y}\|_{L_{\mathbb{P}}^4(\Xi; L^4(\Omega))}} \rightarrow \infty.$$

This can be shown as follows: Let $(a_n)_{n \in \mathbb{N}} \subset \mathbb{R}_{\geq 0}$ and $(b_n)_{n \in \mathbb{N}} \subset \mathbb{R}_{\geq 0}$ be two sequences such that $a_n + b_n > 0$ for all n , $a_n + b_n \xrightarrow{n \rightarrow \infty} \infty$ and let $c > 0$. Then, if $a_n \geq b_n$, giving $a_n > 0$, we have

$$d_n := \frac{ca_n^2 + b_n^4}{a_n + b_n} \geq \frac{ca_n^2}{2a_n} = \frac{ca_n}{2} \geq \frac{c}{4}(a_n + b_n),$$

and if $a_n \leq b_n$, giving $b_n > 0$, we have

$$d_n = \frac{ca_n^2 + b_n^4}{a_n + b_n} \geq \frac{b_n^4}{2b_n} = \frac{b_n^3}{2} \geq \frac{1}{16}(a_n + b_n)^3.$$

Overall, this gives

$$d_n \geq \min\left\{\frac{c}{4}(a_n + b_n), \frac{1}{16}(a_n + b_n)^3\right\} \geq \min\left\{\frac{c}{4}, \frac{1}{16}\right\}(a_n + b_n)$$

as long as $a_n + b_n \geq 1$. The assumption $a_n + b_n \xrightarrow{n \rightarrow \infty} \infty$ yields $d_n \xrightarrow{n \rightarrow \infty} \infty$.

Overall, we can apply Theorem 3.7 to deduce that (3.12) has a unique solution $\mathbf{y} \in \tilde{\mathbf{Y}} = L_{\mathbb{P}}^2(\Xi; H_0^1(\Omega)) \cap L_{\mathbb{P}}^4(\Xi; L^4(\Omega))$. By constructing and analyzing \mathbf{y} more explicitly as in Section 3.2, we see that this solution even belongs to $L_{\mathbb{P}}^4(\Xi; H_0^1(\Omega))$ if $r_f \geq 4$.

A.6. Concrete Computations for Section 9.3

In Section 9.3 we use some concretely computed results, which are presented in the following.

Computation of the Unique Stationary Point for Proposition 9.22

Let

$$\frac{\partial}{\partial \bar{w}} f_\mu(x, \bar{w}, t) = 1 - \frac{\mu}{\bar{w} - a_1(x-t)} - \frac{\mu}{\bar{w} - a_2(x-t)} = 0, \quad \bar{w} > a_i(x-t) \quad (i \in \{1, 2\}).$$

For $s_1 := a_1(x-t)$, $s_2 := a_2(x-t)$, i. e., $\bar{w} > s_1$ and $\bar{w} > s_2$, we perform the following equivalent transformations:

$$\begin{aligned} 1 - \frac{\mu}{\bar{w} - s_1} - \frac{\mu}{\bar{w} - s_2} &= 0, & \bar{w} > s_1, \bar{w} > s_2 \\ \Leftrightarrow (\bar{w} - s_1)(\bar{w} - s_2) - \mu(\bar{w} - s_2) - \mu(\bar{w} - s_1) &= 0, & \bar{w} > s_1, \bar{w} > s_2 \\ \Leftrightarrow \bar{w}^2 - (s_1 + s_2 + 2\mu)\bar{w} + s_1s_2 + \mu(s_1 + s_2) &= 0, & \bar{w} > s_1, \bar{w} > s_2 \\ \Leftrightarrow \bar{w} = \frac{s_1 + s_2 + 2\mu \pm \sqrt{(s_1 + s_2 + 2\mu)^2 - 4(s_1s_2 + \mu(s_1 + s_2))}}{2}, & \bar{w} > s_1, \bar{w} > s_2 \\ \Leftrightarrow \bar{w} = \mu + \frac{s_1 + s_2 + \sqrt{(s_1 - s_2)^2 + 4\mu^2}}{2}. \end{aligned}$$

In the last expression, taking “−” in front of the square root would result in $\bar{w} - s_2 = \frac{s_1 - s_2 + 2\mu - \sqrt{(s_1 - s_2)^2 + 4\mu^2}}{2} \leq \frac{s_1 - s_2}{2}$ and in the same manner $\bar{w} - s_1 \leq \frac{s_2 - s_1}{2}$. This gives that one of the two conditions $\bar{w} > s_1$ and $\bar{w} > s_2$ would be violated then which is why it is necessary to take “+”.

Computation of the Derivatives of v_μ

We compute the derivatives of the function v_μ , needed in the proof of Theorem 9.23. Recall that

$$v_\mu(s) = w_\mu(s) - \mu \ln(w_\mu(s) - a_1s) - \mu \ln(w_\mu(s) - a_2s) + \zeta(\mu)$$

holds with $w_\mu(s) = \mu + \frac{(a_1+a_2)s + \sqrt{(a_1-a_2)^2s^2 + 4\mu^2}}{2}$. Clearly, by the chain rule,

$$v'_\mu(s) = w'_\mu(s) - \mu \frac{w'_\mu(s) - a_1}{w_\mu(s) - a_1s} - \mu \frac{w'_\mu(s) - a_2}{w_\mu(s) - a_2s}$$

with $w'_\mu(s) = \frac{a_1+a_2}{2} + \frac{(a_1-a_2)^2s}{2\sqrt{(a_1-a_2)^2s^2 + 4\mu^2}}$. Indeed, the function w_μ is twice continuously differentiable on \mathbb{R} with the second derivative

$$\begin{aligned} w''_\mu(s) &= \frac{2\sqrt{(a_1-a_2)^2s^2 + 4\mu^2}(a_1-a_2)^2 - (a_1-a_2)^2s \frac{2(a_1-a_2)^2s}{\sqrt{(a_1-a_2)^2s^2 + 4\mu^2}}}{4((a_1-a_2)^2s^2 + 4\mu^2)} \\ &= \frac{4\mu^2(a_1-a_2)^2}{\sqrt{(a_1-a_2)^2s^2 + 4\mu^2}} = \frac{2\mu^2(a_1-a_2)^2}{((a_1-a_2)^2s^2 + 4\mu^2)^{3/2}} > 0 \quad \forall s \in \mathbb{R}. \end{aligned}$$

With this expression, the second derivative of v_μ can be computed defining $\beta := a_2 - a_1 > 0$ and $\gamma(s) := \sqrt{\beta^2 s^2 + 4\mu^2} \geq 2\mu$:

$$\begin{aligned} \tilde{v}_\mu''(s) &= w_\mu''(s) - \mu \frac{(w_\mu(s) - a_1 s) w_\mu''(s) - (w_\mu'(s) - a_1)^2}{(w_\mu(s) - a_1 s)^2} - \mu \frac{(w_\mu(s) - a_2 s) w_\mu''(s) - (w_\mu'(s) - a_2)^2}{(w_\mu(s) - a_2 s)^2} \\ &= \frac{2\mu^2 \beta^2}{\gamma^3(s)} - \mu \frac{(\mu + \frac{\beta s}{2} + \frac{\gamma(s)}{2}) \frac{2\mu^2 \beta^2}{\gamma^3(s)} - (\frac{\beta}{2} + \frac{\beta^2 s}{2\gamma(s)})^2}{(\mu + \frac{\beta s}{2} + \frac{\gamma(s)}{2})^2} - \mu \frac{(\mu - \frac{\beta s}{2} + \frac{\gamma(s)}{2}) \frac{2\mu^2 \beta^2}{\gamma^3(s)} - (-\frac{\beta}{2} + \frac{\beta^2 s}{2\gamma(s)})^2}{(\mu - \frac{\beta s}{2} + \frac{\gamma(s)}{2})^2} \\ &= \frac{\mu \beta^2}{\gamma^3(s) (\mu + \frac{\beta s}{2} + \frac{\gamma(s)}{2})^2 (\mu - \frac{\beta s}{2} + \frac{\gamma(s)}{2})^2} \left(2\mu (\mu + \frac{\beta s}{2} + \frac{\gamma(s)}{2})^2 (\mu - \frac{\beta s}{2} + \frac{\gamma(s)}{2})^2 \right. \\ &\quad \left. - 2\mu^2 (\mu + \frac{\beta s}{2} + \frac{\gamma(s)}{2}) (\mu - \frac{\beta s}{2} + \frac{\gamma(s)}{2})^2 + \gamma(s) (\frac{\gamma(s)}{2} + \frac{\beta s}{2})^2 (\mu - \frac{\beta s}{2} + \frac{\gamma(s)}{2})^2 \right. \\ &\quad \left. - 2\mu^2 (\mu - \frac{\beta s}{2} + \frac{\gamma(s)}{2}) (\mu + \frac{\beta s}{2} + \frac{\gamma(s)}{2})^2 + \gamma(s) (-\frac{\gamma(s)}{2} + \frac{\beta s}{2})^2 (\mu + \frac{\beta s}{2} + \frac{\gamma(s)}{2})^2 \right). \end{aligned}$$

The first factor is always positive and well-defined by $\frac{\pm\beta s + \gamma(s)}{2} > \frac{\pm\beta s + \beta|s|}{2} \geq 0$. Defining $\alpha_+(s) := \frac{\beta s + \gamma(s)}{2} > 0$ and $\alpha_-(s) := \frac{-\beta s + \gamma(s)}{2} > 0$, the second factor can be transformed to

$$\begin{aligned} &2\mu(\mu + \alpha_+(s))^2(\mu + \alpha_-(s))^2 - 2\mu^2[(\mu + \alpha_+(s))(\mu + \alpha_-(s))^2 + (\mu + \alpha_-(s))(\mu + \alpha_+(s))^2] \\ &\quad + \gamma(s)[\alpha_+(s)^2(\mu + \alpha_-(s))^2 + \alpha_-(s)^2(\mu + \alpha_+(s))^2] \\ &= 2\mu(\mu + \alpha_+(s))(\mu + \alpha_-(s)) [(\mu + \alpha_+(s))(\mu + \alpha_-(s)) - \mu(\mu + \alpha_-(s)) - \mu(\mu + \alpha_+(s))] \\ &\quad + \gamma(s)[\mu^2(\beta^2 s^2 + 2\mu^2) + 2\mu^3\gamma(s) + 2\mu^4] \\ &= 2\mu(\mu + \alpha_+(s))(\mu + \alpha_-(s)) \underbrace{[\alpha_+(s)(\mu + \alpha_-(s)) - \mu^2 - \mu\alpha_+(s)]}_{=0} \\ &\quad + \gamma(s)\mu^2[\beta^2 s^2 + 2\mu\gamma(s) + 4\mu^2] \\ &= \gamma(s)\mu^2[\gamma(s)^2 + 2\mu\gamma(s)] = \gamma(s)^2\mu^2(\gamma(s) + 2\mu) \geq 16\mu^5 > 0 \end{aligned}$$

using $\alpha_+(s) + \alpha_-(s) = \gamma(s)$, $\alpha_+(s)^2 + \alpha_-(s)^2 = \beta^2 s^2 + 2\mu^2$ and $\alpha_+(s)\alpha_-(s) = \mu^2$. We obtain a simplified version of the second derivative with $(\mu + \alpha_+(s))^2(\mu + \alpha_-(s))^2 = \mu^2(2\mu + \gamma(s))^2$:

$$v_\mu''(s) = \frac{\mu\beta^2}{\gamma(s)(2\mu + \gamma(s))} = \frac{\mu(a_2 - a_1)^2}{2\mu\sqrt{(a_2 - a_1)^2 s^2 + 4\mu^2} + (a_2 - a_1)^2 s^2 + 4\mu^2}.$$

Computation of Limits

For the proof of Theorem 9.23, we compute the limits $\lim_{s \rightarrow -\infty} v_\mu'(s)$ and $\lim_{s \rightarrow +\infty} v_\mu'(s)$. First, observe that

$$\lim_{s \rightarrow -\infty} w_\mu'(s) = \frac{a_1 + a_2}{2} - \frac{|a_2 - a_1|}{2} = a_1 \quad \text{and} \quad \lim_{s \rightarrow +\infty} w_\mu'(s) = \frac{a_1 + a_2}{2} + \frac{|a_2 - a_1|}{2} = a_2$$

holds ($a_2 > a_1$). Moreover,

$$\frac{w_\mu'(s) - a_1}{w_\mu(s) - a_1 s} = \frac{\frac{a_2 - a_1}{2} + \frac{(a_2 - a_1)^2 s}{2\sqrt{(a_2 - a_1)^2 s^2 + 4\mu^2}}}{\mu + \frac{(a_2 - a_1)s}{2} + \frac{\sqrt{(a_2 - a_1)^2 s^2 + 4\mu^2}}{2}} \rightarrow 0 \quad \text{as } s \rightarrow \pm\infty.$$

As $s \rightarrow +\infty$, the numerator tends to $a_2 - a_1$, but the denominator goes to $+\infty$. For $s \rightarrow -\infty$, the numerator becomes 0 and the denominator tends to μ . An analogous argumentation for the term $\frac{w'_\mu(s) - a_2}{w_\mu(s) - a_2 s}$ yields the limits

$$\lim_{s \rightarrow -\infty} v'_\mu(s) = a_1 \quad \text{and} \quad \lim_{s \rightarrow +\infty} v'_\mu(s) = a_2.$$

Computation of $\hat{v}_\mu(0)$ and $\hat{v}'_\mu(0)$ for Theorem 9.23

In the proof of Theorem 9.23, $d := \frac{2 - a_1 - a_2}{(1 - a_1)(a_2 - 1)}\mu = \frac{1}{a_2 - 1}\mu - \frac{1}{1 - a_1}\mu$ and $\hat{v}_\mu(s) = v_\mu(s + d) - d$ are defined. We compute

$$\begin{aligned} \hat{v}_\mu(0) &= v_\mu(d) - d = w_\mu(d) - \mu \ln(w_\mu(d) - a_1 d) - \mu \ln(w_\mu(d) - a_2 d) + \zeta(\mu) - d \\ &= \frac{a_2}{a_2 - 1}\mu - \frac{a_1}{1 - a_1}\mu - \mu \left(\ln\left(\frac{a_2 - a_1}{a_2 - 1}\mu\right) + \ln\left(\frac{a_2 - a_1}{1 - a_1}\mu\right) \right) \\ &\quad + \mu \left(\ln\left(\frac{a_2 - a_1}{a_2 - 1}\mu\right) + \ln\left(\frac{a_2 - a_1}{1 - a_1}\mu\right) - 2 \right) - \frac{1}{a_2 - 1}\mu + \frac{1}{1 - a_1}\mu \\ &= \frac{a_2}{a_2 - 1}\mu - \frac{a_1}{1 - a_1}\mu - 2\mu - \frac{1}{a_2 - 1}\mu + \frac{1}{1 - a_1}\mu = 0 \end{aligned}$$

using

$$\begin{aligned} \frac{2}{\mu} w_\mu(d) &= 2 + \frac{(a_1 + a_2)d + \sqrt{(a_1 - a_2)^2 d^2 + 4\mu^2}}{\mu} \\ &= 2 + \frac{(a_1 + a_2)(2 - a_1 - a_2)}{(1 - a_1)(a_2 - 1)} + \sqrt{\frac{(a_1 - a_2)^2 (2 - a_1 - a_2)^2}{(1 - a_1)^2 (a_2 - 1)^2} + 4} \\ &= \frac{2(1 - a_1)(a_2 - 1) + (a_1 + a_2)(2 - a_1 - a_2)}{(1 - a_1)(a_2 - 1)} + \frac{\sqrt{(a_2 - a_1)^2 (2 - a_1 - a_2)^2 + 4(1 - a_1)^2 (a_2 - 1)^2}}{(1 - a_1)(a_2 - 1)} \\ &= \frac{2(1 - a_1)(a_2 - 1 + a_1 + a_2) + (a_2 - 1)(1 - a_1 - a_1 - a_2)}{(1 - a_1)(a_2 - 1)} \\ &\quad + \frac{\sqrt{((1 - a_1) + (a_2 - 1))^2 ((1 - a_1) - (a_2 - 1))^2 + 4(1 - a_1)^2 (a_2 - 1)^2}}{(1 - a_1)(a_2 - 1)} \\ &= \frac{-(1 - a_1)^2 + 2a_2(1 - a_1) - (a_2 - 1)^2 - 2a_1(a_2 - 1) + \sqrt{((1 - a_1)^2 + (a_2 - 1)^2)^2}}{(1 - a_1)(a_2 - 1)} \\ &= \frac{2a_2(1 - a_1) - 2a_1(a_2 - 1)}{(1 - a_1)(a_2 - 1)} = \frac{2a_2}{a_2 - 1} - \frac{2a_1}{1 - a_1} \end{aligned}$$

and $w_\mu(d) - a_i d = \left(\frac{a_2 - a_i}{a_2 - 1} - \frac{a_1 - a_i}{1 - a_1}\right)\mu$. We see that the shift $\zeta(\mu)$ is defined such that the shifted function \hat{v}_μ has value 0 at $s = 0$.

Furthermore,

$$\begin{aligned} \hat{v}'_\mu(0) &= v'_\mu(d) = w'_\mu(d) - \mu \frac{w'_\mu(d) - a_1}{w_\mu(d) - a_1 d} - \mu \frac{w'_\mu(d) - a_2}{w_\mu(d) - a_2 d} \\ &= \frac{a_2}{(a_2 - 1)^2} + \frac{a_1}{(1 - a_1)^2} - \mu \frac{\frac{a_2}{(a_2 - 1)^2} + \frac{a_1}{(1 - a_1)^2} - a_1}{\frac{a_2 - a_1}{a_2 - 1}\mu} - \mu \frac{\frac{a_2}{(a_2 - 1)^2} + \frac{a_1}{(1 - a_1)^2} - a_2}{\frac{a_2 - a_1}{1 - a_1}\mu} \\ &= \left(\frac{a_2}{(a_2 - 1)^2} + \frac{a_1}{(1 - a_1)^2}\right) \left(\frac{a_2 - a_1 - (a_2 - 1) - (1 - a_1)}{a_2 - a_1}\right) + \frac{a_1(a_2 - 1)}{a_2 - a_1} + \frac{a_2(1 - a_1)}{a_2 - a_1} = 1 \end{aligned}$$

is computed using

$$\begin{aligned}
 w'_\mu(d) &= \frac{a_1+a_2}{2} + \frac{(a_1-a_2)^2 d}{2\sqrt{(a_1-a_2)^2 d^2 + 4\mu^2}} = \frac{a_1+a_2}{2} + \frac{\frac{(a_2-a_1)^2(2-a_1-a_2)}{(1-a_1)(a_2-1)}}{2\frac{(1-a_1)^2+(a_2-1)^2}{(1-a_1)(a_2-1)}} \\
 &= \frac{a_1+a_2}{2} + \frac{(a_2-a_1)^2((1-a_1)-(a_2-1))}{2((1-a_1)^2+(a_2-1)^2)} = \frac{a_1+a_2}{2} + \frac{(a_2-a_1)((1-a_1)^2-(a_2-1)^2)}{2((1-a_1)^2+(a_2-1)^2)} \\
 &= \frac{a_2(1-a_1)^2+a_1(a_2-1)^2}{(1-a_1)^2+(a_2-1)^2} = \frac{a_2}{(a_2-1)^2} + \frac{a_1}{(1-a_1)^2}.
 \end{aligned}$$

Now, we see what has been the idea behind the definition of the shift $\zeta(\mu)$ and the value of d : First, $d \in \mathbb{R}$ is chosen as the unique point where $v'_\mu(d) = 1$ holds. Then, the shift $\zeta(\mu)$ is defined such that $\mathcal{R}_\mu[0] = 0$ holds, which is provided by $\hat{v}_\mu(0) = 0$ and $\hat{v}'_\mu(0) = 1$.

Bibliography

- [1] P.-A. ABSIL, R. MAHONY, AND R. SEPULCHRE, *Optimization Algorithms on Matrix Manifolds*, Princeton University Press, Princeton, 2008.
- [2] M. AINSWORTH AND J. T. ODEN, *A Posteriori Error Estimation in Finite Element Analysis*, Pure and Applied Mathematics, John Wiley & Sons, New York, 2000.
- [3] A. ALEXANDERIAN, N. PETRA, G. STADLER, AND O. GHATTAS, *Mean-Variance Risk-Averse Optimal Control of Systems Governed by PDEs with Random Parameter Fields Using Quadratic Approximations*, SIAM/ASA J. Uncertain. Quantif., 5 (2017), pp. 1166–1192.
- [4] A. ALI, E. ULLMANN, AND M. HINZE, *Multilevel Monte Carlo Analysis for Optimal Control of Elliptic PDEs with Random Coefficients*, SIAM/ASA J. Uncertain. Quantif., 5 (2017), pp. 466–492.
- [5] J. APPELL AND P. ZABREĀKO, *Nonlinear Superposition Operators*, vol. 95 of Cambridge Tracts in Mathematics, Cambridge University Press, Cambridge, 1990.
- [6] P. ARTZNER, F. DELBAEN, J.-M. EBER, AND D. HEATH, *Coherent Measures of Risk*, Math. Finance, 9 (1999), pp. 203–228.
- [7] I. BABUŠKA, F. NOBILE, AND R. TEMPONE, *A stochastic collocation method for elliptic partial differential equations with random input data*, SIAM Rev., 52 (2010).
- [8] M. BACHMAYR, A. COHEN, AND W. DAHMEN, *Parametric PDEs: sparse or low-rank approximations?*, IMA J. Numer. Anal., (2017).
- [9] M. BACHMAYR AND W. DAHMEN, *Adaptive near-optimal rank tensor approximation for high-dimensional operator equations*, Found. Comput. Math., 15 (2015), pp. 839–898.
- [10] ———, *Adaptive low-rank methods for problems on Sobolev spaces with error control in L_2* , ESAIM Math. Model. Numer. Anal., 50 (2016), pp. 1107–1136.
- [11] M. BACHMAYR AND R. SCHNEIDER, *Iterative Methods Based on Soft Thresholding of Hierarchical Tensors*, Found. Comput. Math., 17 (2017), pp. 1037–1083.
- [12] M. BACHMAYR, R. SCHNEIDER, AND A. USCHMAJEV, *Tensor Networks and Hierarchical Tensors for the Solution of High-Dimensional Partial Differential Equations*, Found. Comput. Math., 16 (2016), pp. 1423–1472.
- [13] J. BALLANI AND L. GRASEDYCK, *A projection method to solve linear systems in tensor format*, Numer. Linear Algebra Appl., 20 (2013), pp. 27–43.

- [14] —, *Hierarchical tensor approximation of output quantities of parameter-dependent PDEs*, SIAM/ASA J. Uncertain. Quantif., 3 (2015), pp. 852–872.
- [15] J. BALLANI, L. GRASEDYCK, AND M. KLUGE, *Black box approximation of tensors in hierarchical tucker format*, Linear Algebra Appl., 438 (2013), pp. 639–657.
- [16] A. BARTH, C. SCHWAB, AND N. ZOLLINGER, *Multi-level Monte Carlo finite element method for elliptic PDEs with stochastic coefficients*, Numer. Math., 119 (2011), pp. 123–161.
- [17] H. BAUSCHKE AND P. COMBETTES, *Convex Analysis and Monotone Operator Theory in Hilbert Spaces*, Springer, New York, 2011.
- [18] R. BECKER AND R. RANNACHER, *An optimal control approach to a posteriori error estimation in finite element methods*, Acta Numer., 10 (2001), pp. 1–102.
- [19] S. BELLAVIA, *Inexact interior-point method*, J. Optim. Theory Appl., 96 (1998), pp. 109–121.
- [20] P. BENNER, S. DOLGOV, A. ONWUNTA, AND M. STOLL, *Low-rank solvers for unsteady Stokes–Brinkman optimal control problem with random data*, Comput. Methods Appl. Mech. Engrg., 304 (2016), pp. 26–54.
- [21] P. BENNER, A. ONWUNTA, AND M. STOLL, *Block-Diagonal Preconditioning for Optimal Control Problems Constrained by PDEs with Uncertain Inputs*, SIAM J. Matrix Anal. Appl., 37 (2013), pp. 491–518.
- [22] A. BESPALOV, C. E. POWELL, AND D. SILVESTER, *Energy norm a posteriori error estimation for parametric operator equations*, SIAM J. Sci. Comput., 36 (2014), pp. A339 – A363.
- [23] A. BESPALOV AND D. SILVESTER, *Efficient adaptive stochastic Galerkin methods for parametric operator equations*, SIAM J. Sci. Comput., 38 (2016), pp. A2118 – A2140.
- [24] A. BORZÌ, V. SCHULZ, C. SCHILLINGS, AND G. VON WINCKEL, *On the treatment of distributed uncertainties in PDE-constrained optimization*, GAMM-Mitt., 33 (2010), pp. 230–246.
- [25] A. BORZÌ AND G. VON WINCKEL, *A POD framework to determine robust controls in PDE optimization*, Comput. Vis. Sci., 14 (2011), pp. 91–103.
- [26] R. G. CARTER, *Numerical optimization in hilbert space using inexact function and gradient evaluations*, Technical report 89-45, ICASE, Langley, VA, 1989.
- [27] —, *On the Global Convergence of Trust Region Algorithms Using Inexact Gradient Information*, SIAM J. Numer. Anal., 28 (1991), pp. 251–265.
- [28] P. CHEN AND A. QUARTERONI, *Weighted Reduced Basis Method for Stochastic Optimal Control Problems with Elliptic PDE Constraint*, SIAM/ASA J. Uncertain. Quantif., 2 (2014), pp. 364–396.

-
- [29] P. CHEN, A. QUARTERONI, AND G. ROZZA, *A Weighted Reduced Basis Method for Elliptic Partial Differential Equations with Random Input Data*, SIAM J. Numer. Anal., 51 (2013), pp. 3163–3185.
- [30] K. CLIFFE, M. GILES, R. SCHEICHL, AND A. TECKENTRUP, *Multilevel Monte Carlo methods and applications to elliptic PDEs with random coefficients*, Comput. Vis. Sci., 14 (2011), pp. 3–15.
- [31] A. CONN, N. GOULD, AND P. TOINT, *Trust-Region Methods*, SIAM and MOS, Philadelphia, 2000.
- [32] C. DA SILVA AND F. HERRMANN, *Optimization on the hierarchical Tucker manifold – Applications to tensor completion*, Linear Algebra Appl., 481 (2015), pp. 131–173.
- [33] V. DE SILVA AND L.-H. LIM, *Tensor rank and the ill-posedness of the best low-rank approximation problem*, SIAM J. Matrix Anal. Appl., 30 (2008), pp. 1084–1127.
- [34] A. DEFANT AND K. FLORET, *Tensor Norms and Operator Ideals*, Elsevier Science Publishers, Amsterdam, 1993.
- [35] S. V. DOLGOV, *Alternating minimal energy approach to ODEs and conservation laws in tensor product formats*. preprint, arXiv:1403.8085, 2014.
- [36] S. V. DOLGOV AND D. V. SAVOSTYANOV, *Alternating Minimal Energy Methods for Linear Systems in Higher Dimensions*, SIAM J. Sci. Comput., 36 (2014), pp. A2248–A2271.
- [37] W. DÖRFLER, *A Convergent Adaptive Algorithm for Poisson’s Equation*, SIAM J. Numer. Anal., 33 (1996), pp. 1106–1124.
- [38] M. EIGEL, C. J. GITTELSON, C. SCHWAB, AND E. ZANDER, *Adaptive stochastic Galerkin FEM*, Comput. Methods Appl. Mech. Engrg., 270 (2014), pp. 247–269.
- [39] ———, *A convergent adaptive stochastic Galerkin finite element method with quasi-optimal spatial meshes*, ESAIM Math. Model. Numer. Anal., 49 (2015), pp. 1367–1398.
- [40] M. EIGEL, M. PFEFFER, AND R. SCHNEIDER, *Adaptive stochastic Galerkin FEM with hierarchical tensor representations*, Numer. Math., 136 (2017), pp. 765–803.
- [41] S. C. EISENSTAT AND H. F. WALKER, *Globally convergent inexact Newton methods*, SIAM J. Optim., 4 (1994), pp. 393–422.
- [42] H. C. ELMAN AND Q. LIAO, *Reduced basis collocation methods for partial differential equations with random coefficients*, SIAM/ASA J. Uncertain. Quantif., 1 (2013), pp. 192–217.
- [43] M. ESPIG, W. HACKBUSCH, A. LITVINENKO, H. G. MATTHIES, AND E. ZANDER, *Efficient analysis of high dimensional data in tensor formats*, in Sparse Grids and Applications, J. Garcke and M. Griebel, eds., Springer, Berlin, 2013, pp. 31–56.

- [44] R. FORSTER AND R. KORNUBER, *A polynomial chaos approach to stochastic variational inequalities*, J. Numer. Math., 18 (2010), pp. 235–255.
- [45] J. GARCKE, *Sparse grids in a nutshell*, in Sparse Grids and Applications, J. Garcke and M. Griebel, eds., Springer, Berlin, 2013, pp. 57–80.
- [46] S. GARREIS AND M. ULBRICH, *Constrained Optimization with Low-Rank Tensors and Applications to Parametric Problems with PDEs*, SIAM J. Sci. Comput., 39 (2017), pp. A25–A54. The first version of this article has been acknowledged as Master’s thesis for Sebastian Garreis within the Master’s with Honours Examination for Mathematics at the Technical University of Munich.
- [47] T. GERSTNER AND M. GRIEBEL, *Numerical integration using sparse grids*, Numer. Algorithms, 18 (1998), pp. 209–232.
- [48] H. GOLDBERG, W. KAMPOWSKY, AND F. TRÖLTZSCH, *On NEMYTSKIJ Operators in l_p -Spaces of Abstract Functions*, Math. Nachr., 155 (1992), pp. 127–140.
- [49] G. H. GOLUB AND C. F. VAN LOAN, *Matrix Computations*, Johns Hopkins Univ. Press, Baltimore and London, 1996.
- [50] L. GRASEDYCK, *Hierarchical singular value decomposition of tensors*, SIAM J. Matrix Anal. Appl., 31 (2010), pp. 2029–2054.
- [51] L. GRASEDYCK, D. KRESSNER, AND C. TOBLER, *A literature survey of low-rank tensor approximation techniques*, GAMM-Mitt., 36 (2013), pp. 53–78.
- [52] L. GRASEDYCK AND C. LÖBBERT, *Distributed hierarchical SVD in the Hierarchical Tucker format*, Numer. Linear Algebra Appl., (2018), p. e2174.
- [53] M. GUNZBURGER, C. WEBSTER, AND G. ZHANG, *Stochastic finite element methods for partial differential equations with random input data*, Acta Numer., 23 (2014), pp. 521–650.
- [54] W. HACKBUSCH, *Tensor Spaces and Numerical Tensor Calculus*, Springer, Heidelberg, 2012.
- [55] W. HACKBUSCH, B. N. KHOROMSKIJ, AND E. E. TYRTYSHNIKOV, *Approximate iterations for structured matrices*, Numer. Math., 109 (2008), pp. 365–383.
- [56] W. HACKBUSCH AND S. KÜHN, *A new scheme for the tensor representation*, J. Fourier Anal. Appl., 15 (2009), pp. 706–722.
- [57] M. HEINKENSCHLOSS AND L. N. VICENTE, *Analysis of inexact trust-region SQP algorithms*, SIAM J. Optim., 12 (2001), pp. 283–302.
- [58] M. HINTERMÜLLER, K. ITO, AND K. KUNISCH, *The Primal-Dual Active Set Strategy as a Semismooth Newton Method*, SIAM J. Optim., 13 (2002), pp. 865–888.
- [59] M. HINZE, *A Variational Discretization Concept in Control Constrained Optimization: The Linear-Quadratic Case*, Comput. Optim. Appl., 30 (2005), pp. 45–61.

-
- [60] M. HINZE, R. PINNAU, M. ULBRICH, AND S. ULBRICH, *Optimization with PDE Constraints*, Springer, New York, 2009.
- [61] F. L. HITCHCOCK, *The expression of a tensor or a polyadic as a sum of products*, J. Math. Phys., 6 (1927), pp. 164–189.
- [62] S. HOLTZ, T. ROHWEDDER, AND R. SCHNEIDER, *The alternating linear scheme for tensor optimization in the tensor train format*, SIAM J. Sci. Comput., 34 (2012), pp. A683–A713.
- [63] L.S. HOU, J. LEE, AND H. MANOUZI, *Finite element approximations of stochastic optimal control problems constrained by stochastic elliptic PDEs*, J. Math. Anal. Appl., 384 (2011), pp. 87–103.
- [64] T. HUCKLE, K. WALDHERR, AND T. SCHULTE-HERBRÜGGEN, *Computations in Quantum Tensor Networks*, Linear Algebra Appl., 438 (2012), pp. 750–781.
- [65] W. HUYER AND A. NEUMAIER, *Global optimization by multilevel coordinate search*, J. Global Optim., 14 (1999), pp. 331–355.
- [66] T. HYTÖNEN, M. VERAAR, J. VAN NEERVEN, AND L. WEIS, *Analysis in Banach Spaces*, vol. I, Springer, Cham, 2016.
- [67] T. G. KOLDA AND B. W. BADER, *Tensor Decompositions and Applications*, SIAM Rev., 51 (2009), pp. 455–500.
- [68] D. P. KOURI, *A Multilevel Stochastic Collocation Algorithm for Optimization of PDEs with Uncertain Coefficients*, SIAM/ASA J. Uncertain. Quantif., 2 (2014), pp. 55–81.
- [69] D. P. KOURI, M. HEINKENSCHLOSS, D. RIDZAL, AND B. G. VAN BLOEMEN WAANDERS, *A Trust-Region Algorithm with Adaptive Stochastic Collocation for PDE Optimization under Uncertainty*, SIAM J. Sci. Comput., 35 (2013), pp. A1847–A1879.
- [70] ———, *Inexact Objective Function Evaluations in a Trust-Region Algorithm for PDE-Constrained Optimization under Uncertainty*, SIAM J. Sci. Comput., 36 (2014), pp. A3011–A3029.
- [71] D. P. KOURI AND T. M. SUROWIEC, *Risk-Averse PDE-Constrained Optimization Using the Conditional Value-At-Risk*, SIAM J. Optim., 26 (2016), pp. 365–396.
- [72] ———, *Existence and Optimality Conditions for Risk-Averse PDE-Constrained Optimization*, SIAM/ASA J. Uncertain. Quantif., 6 (2018), pp. 787–815.
- [73] D. KRESSNER, M. STEINLECHNER, AND B. VANDEREYCKEN, *Low-rank tensor completion by Riemannian optimization*, BIT, 54 (2014), pp. 447–468.
- [74] ———, *Preconditioned Low-rank Riemannian Optimization for Linear Systems with Tensor Product Structure*, SIAM J. Sci. Comput., 38 (2016), pp. A2018–A2044.
- [75] D. KRESSNER AND C. TOBLER, *Low-rank tensor Krylov subspace methods for parametrized linear systems*, SIAM J. Matrix Anal. Appl., 32 (2011), pp. 1288–1316.

- [76] ———, *Algorithm 941: htucker—a Matlab toolbox for tensors in hierarchical Tucker format*, ACM Trans. Math. Software, 40 (2014).
- [77] P. KROKHMAL, J. PALMQUIST, AND S. URYASEV, *Portfolio Optimization With Conditional Value-at-Risk Objective and Constraints*, Journal of Risk, 4 (2002), pp. 43–68.
- [78] A. KUNOTH AND C. SCHWAB, *Analytic Regularity and GPC Approximation for Control Problems Constrained by Linear Parametric Elliptic and Parabolic PDEs*, SIAM J. Control Optim., 51 (2013), pp. 2442–2471.
- [79] ———, *Sparse Adaptive Tensor Galerkin Approximations of Stochastic PDE-Constrained Control Problems*, SIAM/ASA J. Uncertain. Quantif., 4 (2016), pp. 1034–1059.
- [80] F. Y. KUO, C. SCHWAB, AND I. H. SLOAN, *Quasi-Monte Carlo Finite Element Methods for a Class of Elliptic Partial Differential Equations with Random Coefficients*, SIAM J. Numer. Anal., 50 (2012), pp. 3351–5574.
- [81] L. DE LATHAUWER, B. DE MOOR, AND J. VANDEWALLE, *A multilinear singular value decomposition*, SIAM J. Matrix Anal. Appl., 21 (2000), pp. 1253–1278.
- [82] H.-C. LEE AND J. LEE, *A Stochastic Galerkin Method for Stochastic Control Problems*, Commun. Comput. Phys., 14 (2013), pp. 77–106.
- [83] V. MURG, F. VERSTRAETE, Ö. LEGEZA, AND R. M. NOACK, *Simulating strongly correlated quantum systems with tree tensor network*, Phys. Rev. B, 82 (2010).
- [84] H. NIEDERREITER, *Random Number Generation and Quasi-Monte Carlo Methods*, vol. 63 of CBMS-NSF Regional Conf. Ser. in Appl. Math., SIAM, Philadelphia, 1992.
- [85] J. NOCEDAL AND S. J. WRIGHT, *Numerical Optimization*, Springer, New York, second ed., 2006.
- [86] E. NOVAK AND K. RITTER, *High dimensional integration of smooth functions over cubes*, Numer. Math., 75 (1996), pp. 79–97.
- [87] I. V. OSELEDETS, *DMRG approach to fast linear algebra in the TT-format*, Comput. Methods Appl. Math., 11 (2011), pp. 382–393.
- [88] ———, *Tensor-Train Decomposition*, SIAM J. Sci. Comput., 33 (2011), pp. 2295–2317.
- [89] ———, *TT-toolbox 2.2: Fast multidimensional array operations in TT format*, 2012. <https://github.com/oseledets/TT-Toolbox>.
- [90] I. V. OSELEDETS AND S. V. DOLGOV, *Solution of Linear Systems and Matrix Inversion in the TT-format*, SIAM J. Sci. Comput., 34 (2012), pp. A2718–A2739.
- [91] I. V. OSELEDETS AND E. E. TYRTYSHNIKOV, *TT-cross approximation for multidimensional arrays*, Linear Algebra Appl., 432 (2010), pp. 70–88.
- [92] L. QI AND J. SUN, *A nonsmooth version of Newton’s method*, Math. Program., 58 (1993), pp. 353–367.

-
- [93] M. SØRENSEN, L. DE LATHAUWER, P. COMON, S. ICART, AND L. DENEIRE, *Canonical polyadic decomposition with a columnwise orthonormal factor matrix*, SIAM J. Matrix Anal. Appl., 33 (2012), pp. 1190–1213.
- [94] R. T. ROCKAFELLAR AND S. URYASEV, *Optimization of conditional value-at-risk*, Journal of Risk, 2 (2000), pp. 21–41.
- [95] ———, *Conditional value-at-risk for general loss distributions*, J. Bank. Finance, 26 (2002), pp. 1443–1471.
- [96] ———, *The fundamental risk quadrangle in risk management, optimization and statistical estimation*, Surv. Oper. Res. Manag. Sci., 18 (2013), pp. 33–53.
- [97] R. T. ROCKAFELLAR AND R. J.-B. WETS, *Variational Analysis*, Springer, Berlin Heidelberg, 1998.
- [98] T. ROHWEDDER AND A. USCHMAJEV, *On local convergence of alternating schemes for optimization of convex problems in the tensor train format*, SIAM J. Numer. Anal., (2013).
- [99] A. RUSZCZYŃSKI AND A. SHAPIRO, *Optimization of Convex Risk Functions*, Math. Oper. Res., 31 (2006), pp. 433–452.
- [100] B. SAVAS AND L.-H. LIM, *Quasi-Newton Methods on Grassmannians and Multilinear Approximation of Tensors*, SIAM J. Sci. Comput., 32 (2010), pp. 3352–3393.
- [101] G. SCHULZ, *Iterative Berechnung der reziproken Matrix*, ZAMM, 13 (1933), pp. 57–59.
- [102] C. SCHWAB AND C. J. GITTELSON, *Sparse tensor discretizations of high-dimensional parametric and stochastic PDEs*, Acta Numer., 20 (2011), pp. 291–467.
- [103] A. SHAPIRO, D. DENTCHEVA, AND A. RUSZCZYŃSKI, *Lectures on Stochastic Programming: Modeling and Theory*, SIAM, Philadelphia, 2009.
- [104] M. STEINLECHNER, *Riemannian optimization for high-dimensional tensor completion*, SIAM J. Sci. Comput., 38 (2016), pp. S461–S484.
- [105] M. STOLL AND T. BREITEN, *A Low-Rank in Time Approach to PDE-Constrained Optimization*, SIAM J. Sci. Comput., 37 (2015), pp. B1–B29.
- [106] T. J. SULLIVAN, *Introduction to Uncertainty Quantification*, Springer, Cham, 2015.
- [107] A. TECKENTRUP, R. SCHEICHL, M. GILES, AND E. ULLMANN, *Further analysis of multilevel monte carlo methods for elliptic pdes with random coefficients*, Numer. Math., 125 (2013), pp. 569–600.
- [108] H. TIESLER, R. M. KIRBY, D. XIU, AND T. PREUSSER, *Stochastic Collocation for Optimal Control Problems with Stochastic PDE Constraints*, SIAM J. Control Optim., 50 (2012), pp. 2659–2682.

- [109] L. TUCKER, *Some mathematical notes on three-mode factor analysis*, Psychometrika, 31 (1966).
- [110] M. ULBRICH, *Semismooth Newton Methods for Operator Equations in Function Spaces*, SIAM J. Optim., 13 (2002), pp. 805–841.
- [111] M. ULBRICH, *Semismooth Newton Methods for Variational Inequalities and Constrained Optimization Problems in Function Spaces*, vol. 11 of MOS-SIAM Ser. Optim., SIAM, Philadelphia, 2011.
- [112] M. ULBRICH AND S. ULBRICH, *Primal-dual interior-point methods for PDE-constrained optimization*, Math. Program., Ser. B, 117 (2009), pp. 435–485.
- [113] A. USCHMAJEV AND B. VANDEREYCKEN, *The geometry of algorithms using hierarchical tensors*, Linear Algebra Appl., 439 (2013), pp. 133–166.
- [114] R. VERFÜRTH, *A Review of A Posteriori Error Estimation and Adaptive Mesh-Refinement Techniques*, Wiley-Teubner, Chichester and Stuttgart, 1996.
- [115] F. VERSTRAETE, V. MURG, AND J. I. CIRAC, *Matrix product states, projected entangled pair states, and variational renormalization group methods for quantum spin systems*, Adv. Phys., 57 (2008), pp. 143–224.
- [116] J. WEIDMANN, *Linear Operators in Hilbert Spaces*, Springer, New York, 1980.
- [117] E. ZEIDLER, *Nonlinear Functional Analysis and its Applications*, vol. II/B: Nonlinear Monotone Operators, Springer, New York, 1990.
- [118] J. C. ZIEMS, *Adaptive Multilevel Inexact SQP Methods for PDE-Constrained Optimization with Control Constraints*, SIAM J. Optim., 23 (2013), pp. 1257–1283.
- [119] J. C. ZIEMS AND S. ULBRICH, *Adaptive Multilevel Inexact SQP Methods for PDE-Constrained Optimization*, SIAM J. Optim., 21 (2011), pp. 1–40.