Universität Münster Fachbereich Mathematik und Informatik

Masterarbeit Mathematik

Machine learning based surrogate modeling to accelerate parabolic PDE constrained optimization

von: Andy Kevin Wert Matrikelnummer: 461478

Erstgutachter: Prof. Dr. Mario Ohlberger Zweitgutachter: Dr. Stephan Rave

Contents

1	Intr	oduction	3
2	Parabolic optimal control problems		
	2.1	Introduction to the problem	4
	2.2	Finite element discretization	5
		2.2.1 Discretization in space	5
		2.2.2 Discretization in time	7
		2.2.3 Crank-Nicolson scheme	7
		2.2.4 Calculation of the objective function value	8
	2.3	Optimization of the control variable	9
3	Ense	emble-based optimization algorithm	10
4	Adaptive-ML-EnOpt algorithm		16
	4.1	Deep neural networks	16
	4.2	Modifying the EnOpt algorithm by using a neural network-based surrogate	20
5	Nun	nerical experiments	24
Bi	Bibliography		

1 Introduction

[1]

2 Parabolic optimal control problems

2.1 Introduction to the problem

Our optimization problem is based on the problem that is presented in [2]. We consider a state variable u and a control variable q, defined on $(0,T) \times \Omega$ with $T \in \mathbb{R}$ and $\Omega \subset \mathbb{R}^n$.

The goal of this thesis is to minimize the function

$$J(q, u) = \frac{1}{2} \int_0^T \int_{\Omega} (u(t, x) - \hat{u}(t, x))^2 dx dt + \frac{\alpha}{2} \int_0^T \int_{\Omega} q(t, x)^2 dx dt,$$
 (2.1a)

subject to the constraints

$$\partial_t u - \Delta u = f + q \quad \text{in } (0, T) \times \Omega,$$

 $u(0) = u_0 \quad \text{in } \Omega,$ (2.1b)

with homogeneous Dirichlet boundary conditions on $(0,T) \times \partial \Omega$.

Let $V = H_0^1(\Omega)$, $H = L^2(\Omega)$ and I = (0,T). We define our state space as

$$X := \{ v \mid v \in L^2(I, V) \text{ and } \partial_t v \in L^2(I, V^*) \}$$

and the control space as

$$Q := L^2(I, L^2(\Omega)).$$

The notion of the inner products and norms on $L^2(\Omega)$ and $L^2(I, L^2(\Omega))$ is introduced as

$$(v, w) := (v, w)_{L^{2}(\Omega)},$$

$$(v, w)_{I} := (v, w)_{L^{2}(I, L^{2}(\Omega))},$$

$$||v||_{I} := ||v||_{L^{2}(I, L^{2}(\Omega))}.$$

By using the inner product, the weak form of the state equations (2.1b) for $q, f \in Q$ and $u_0 \in V$ is given as

$$(\partial_t u, \phi) + (\nabla u, \nabla \phi) = (f + q, \phi) \quad \forall \phi \in X,$$

$$u(0) = u_0 \qquad \text{in } \Omega.$$
(2.2)

With the weak state equations (2.2), we define the weak formulation of the optimal control problem (2.1) as

Minimize
$$J(q, u) := \frac{1}{2} \|u - \hat{u}\|_{I}^{2} + \frac{\alpha}{2} \|q\|_{I}^{2}$$
 subject to (2.2) and $(q, u) \in Q \times X$. (2.3)

Now, we cite two results of the problems (2.2) and (2.3).

Proposition 2.1 ([2]). For fixed $q, f \in Q$, and $u_0 \in V$ there exists a unique solution $u \in X$ of problem (2.2). Moreover, the solution exhibits the improved regularity

$$u \in L^2(I, H^2(\Omega) \cap V) \cap H^1(I, L^2(\Omega)) \hookrightarrow C(\bar{I}, V).$$

It holds the stability estimate

$$\|\partial_t u\|_I + \|\nabla^2 u\|_I \le C\{\|f + q\|_I + \|\nabla u_0\|\}.$$

Proposition 2.2 ([2]). For given $f, \hat{u} \in L^2(I, H)$, $u_0 \in V$, and $\alpha > 0$, the optimal control Problem (2.3) admits a unique solution $(\bar{q}, \bar{u}) \in Q \times X$. The optimal control \bar{q} possesses the regularity

$$\bar{q} \in L^2(I, H^2(\Omega)) \cap H^1(I, L^2(\Omega)).$$

Due to the existence and uniqueness results from Proposition 2.1, we define u(q) as the unique solution of (2.2) with respect to some $q \in Q$. This enables us to define a reduced cost functional $j: Q \to \mathbb{R}$ that is only dependent on the control q as

$$j(q) := J(q, u(q)).$$

From now on, the optimal control problem that we examine is:

minimize
$$j(q)$$
 subject to $q \in Q$. (2.4)

2.2 Finite element discretization

In order to solve the optimization problem (2.4) numerically, the discretization of our model is now discussed. We begin with the presentation of the discretization in space with a n-D continuous Galerkin method. Then, we look at the discretization in time, which is done with a 1D continuous Galerkin method. From now on, we will also discuss some implementation details, so, in this chapter, how we handle the calculation of the objective function j. To solve the partial equations of (2.2), we use the Python package pyMOR.

2.2.1 Discretization in space

The discretization in space is shown on a 2-dimensional rectangular space $\Omega \subset \mathbb{R}^2$ with linear finite elements. We assume to have a vertex set $\mathcal{V} = (x_1, \dots, x_N) \in (\mathbb{R}^2)^N$ with a convex hull that is equal to $\bar{\Omega}$ and $x_i \neq x_j$ for all $i \neq j$ in $\{1, \dots, N\}$. Let $\hat{T} = \{(x, y) \in [0, 1]^2 \mid y \leq 1 - x\}$ be the reference triangle. Then,

$$\theta_l(\xi) = x_{l_1} + D\theta_l \begin{pmatrix} \xi_1 \\ \xi_2 \end{pmatrix} \text{ with } D\theta_l = (x_{l_2} - x_{l_1} \ x_{l_3} - x_{l_1})$$

is a transformation from the reference triangle \hat{T} to some other triangle T_l with the corners $x_{l_1}, x_{l_2}, x_{l_3} \in \mathcal{V}$.

We define now a mesh $\mathcal{T} = \{T_l\}$ which consists of triangles $T_l = \theta_l(\hat{T})$, where $T_l \cap T_m$ for $T_l, T_m \in \mathcal{T}$ is either a common side, a common corner, or empty, and where $\bar{\Omega} = \bigcup_{T_l \in \mathcal{T}} T_l$. We also assume that every vertex in \mathcal{V} is a corner of at least one triangle of \mathcal{T} .

In our implementation, we discretize a rectangular domain by specifying the number of grid intervals first. Then, we divide the domain into smaller rectangles of the same size, so that the number of rectangles along the x- and the y-axis is equal to the predefined number of grid intervals. Each smaller rectangular unit is then divided into four equally sized triangles by adding a vertex into the center of the rectangle which is connected with the corners of the unit. The vertex set of the whole domain is now given by the union of the corners of all triangles. As an example, if we have given a domain $\Omega = [a, a]$ with a > 0 and we define the number of grid intervals as 2, then our mesh would look like that:

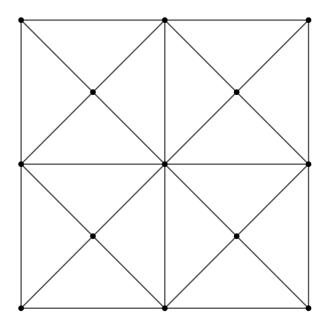


Figure 2.1: Example of a mesh with 2 grid intervals in a square shaped domain.

Now, let $\mathcal{P}_1(\hat{T}, \mathbb{R})$ be the space of polynomials up to order 1 in \hat{T} . Then, $\{\psi_1, \psi_2, \psi_3\}$ with $\psi_1(\xi) = 1 - \xi_1 - \xi_2, \psi_2(\xi) = \xi_1, \psi_3(\xi) = \xi_2$ defines a basis of $\mathcal{P}_1(\hat{T}, \mathbb{R})$. Using this basis, we set

$$V_h = \operatorname{span}\{\phi_i, i = 0, \dots, N\} \cap V$$

as the finite element space of our state variables with

$$\phi_i|_{T_l} = \begin{cases} 0 & \text{if } x_i \notin T_l \\ \psi_1 \circ \theta_l^{-1} & \text{if } \theta_l \begin{pmatrix} 0 \\ 0 \end{pmatrix} = x_i \\ \psi_2 \circ \theta_l^{-1} & \text{if } \theta_l \begin{pmatrix} 1 \\ 0 \end{pmatrix} = x_i \\ \psi_3 \circ \theta_l^{-1} & \text{if } \theta_l \begin{pmatrix} 0 \\ 1 \end{pmatrix} = x_i \end{cases}$$

for all $T_l \in \mathcal{T}$ and i = 1, ..., N.

By construction, every $u \in V_h$ is uniquely defined by

$$u = \sum_{i=1}^{N} U_i \phi_i$$

with $U_i = u(x_i)$.

Now, we want to calculate $\int_{\Omega} u \cdot v \, dx$ and $\int_{\Omega} \nabla u \cdot \nabla v \, dx$ for all $u, v \in V_h$. In order to do that, we set the mass matrix $M_n = \left(\int_{\Omega} \phi_i \cdot \phi_j \, dx\right)_{i,j=1,\dots,N}$ and the stiffness matrix

$$L_n = \left(\int_{\Omega} \nabla \phi_i \cdot \nabla \phi_j \, \mathrm{d}x \right)_{i,j=1,\dots,N}$$
. Let

$$U = \begin{pmatrix} U_1 \\ \vdots \\ U_n \end{pmatrix} \text{ and } V = \begin{pmatrix} V_1 \\ \vdots \\ V_n \end{pmatrix}.$$

Then we have

$$\int_{\Omega} u \cdot v \, \mathrm{d}x = U^T M_n V \text{ and } \int_{\Omega} \nabla v \cdot \nabla u \, \mathrm{d}x = U^T L_n V.$$

2.2.2 Discretization in time

At first, we partition the time interval $\bar{I} = [0, T]$ as

$$\bar{I} = \{0\} \cup I_1 \cup I_2 \cup \cdots \cup I_M$$

with subintervals $I_m = (t_{m-1}, t_m]$, where $t_m = m \frac{T}{M}$ for m = 0, ..., M and $M \in \mathbb{N}$. We want that the discretizations of our functions are continuous in \bar{I} and piecewise polynomial of order 1 in all subintervals I_m , so the discretization space of our state variables is

$$X_{k,h} := \{ v \in C(\bar{I}, V_h) \mid v |_{I_m} \in \mathcal{P}_1(I_m, V_h), m = 1, 2, \dots, M \},$$

where $\mathcal{P}_1(I_m, V_h)$ denotes the space of polynomials up to order 1, defined on I_m with values in V_h . Similarly, we define the time-discretized space of our control variables as

$$Q_d := \{ v \in C(\bar{I}, H) \mid v|_{I_m} \in \mathcal{P}_1(I_m, H), m = 1, 2, \dots, M \} \supset X_{k,h}.$$

By using the Lagrange basis of $\mathcal{P}_1(I_m,\mathbb{R})$, we can write every function $v \in Q_d$ as

$$v(t,\cdot) = \left(m - t\frac{M}{T}\right)v_{m-1}(\cdot) + \left(t\frac{M}{T} - m + 1\right)v_m(\cdot) \text{ for } t \in I_m,$$

where $v_m(\cdot) = v(t_m, \cdot)$.

2.2.3 Crank-Nicolson scheme

Now, we solve the weak state equations (2.2) for the state $u \in X_{k,h}$, the control $q \in Q_d$, and $f \in Q$ numerically. For m = 0, we set

$$U_0 = \begin{pmatrix} U_{0,1} \\ \vdots \\ U_{0,n} \end{pmatrix},$$

where $U_{0,i} = u_0(x_i)$ for i = 1, ..., N.

For m = 1, ..., M, we get with the Crank-Nicolson scheme that for all $v \in V_h$:

$$(u_m, v) + \frac{T}{2M}(\nabla u_m, \nabla v) = (u_{m-1}, v) - \frac{T}{2M}(\nabla u_{m-1}, \nabla v) + \frac{T}{2M}(f_{m-1} + q_{m-1}, v) + \frac{T}{2M}(f_m + q_m, v),$$

where u_m is a time discretization of u at the time step t_m , while $f_m = f(t_m, \cdot)$ and $q_m = q(t_m, \cdot)$. To solve the above equation, we define the matrix $\tilde{M}_n \in \mathbb{R}^{N \times N}$ as

$$\left(\tilde{M}_n\right)_{i,j} = \begin{cases}
0 & \text{if } x_i \text{ or } x_j \text{ in } \partial\Omega \text{ and } i \neq j \\
1 & \text{if } x_i \text{ or } x_j \text{ in } \partial\Omega \text{ and } i = j \\
(M_n)_{i,j} & \text{else}
\end{cases}$$

and the matrix $\tilde{L}_n \in \mathbb{R}^{N \times N}$ as

$$\left(\tilde{L}_n\right)_{i,j} = \begin{cases} 0 & \text{if } x_j \text{ in } \partial\Omega\\ (L_n)_{i,j} & \text{else,} \end{cases}$$

so that $(u_m, v) = U_m^T \tilde{M}_n V$ and $(\nabla u_m, \nabla v) = U_m^T \tilde{L}_n V$ for all $m = 0, \dots, M$, which is giving us

$$V^{T}\tilde{M}_{n}^{T}U_{m} + \frac{T}{2M}V^{T}\tilde{L}_{n}^{T}U_{m} = V^{T}\tilde{M}_{n}^{T}U_{m-1} - \frac{T}{2M}V^{T}\tilde{L}_{n}^{T}U_{m-1} + \frac{T}{2M}(f_{m-1} + q_{m-1}, v) + \frac{T}{2M}(f_{m} + q_{m}, v).$$

In the pyMOR implementation, vectors F_m for m = 0, ..., M are defined such that $V^T F_m \approx (f_m + q_m, v)$ for all $v \in V_h$ and $(F_m)_i = 0$ if the *i*-th entry in the vertex set \mathcal{V} lies on the boundary of Ω . By using these vectors, we get the equation

$$\left(\tilde{M}_{n}^{T} + \frac{T}{2M}\tilde{L}_{n}^{T}\right)U_{m} = \tilde{M}_{n}^{T}U_{m-1} - \frac{T}{2M}\tilde{L}_{n}^{T}U_{m-1} + \frac{T}{2M}F_{m-1} + \frac{T}{2M}F_{m}, \quad (2.5)$$

which is solved after U_m with functions from the Python package SciPy.

2.2.4 Calculation of the objective function value

For fixed $\hat{u}, f \in Q$, we define u = u(q) for all $q \in Q_d$, so that it satisfies (2.5). We calculate j(q) now in the following way:

$$j(q) \approx \frac{1}{2} \sum_{m=1}^{M} \int_{t_{m-1}}^{t_m} \left(\left(m - t \frac{M}{T} \right) (u_{m-1} - \hat{u}_{m-1}) + \left(t \frac{M}{T} - m + 1 \right) (u_m - \hat{u}_m), \right.$$

$$\left. \left(m - t \frac{M}{T} \right) (u_{m-1} - \hat{u}_{m-1}) + \left(t \frac{M}{T} - m + 1 \right) (u_m - \hat{u}_m) \right) dt$$

$$+ \frac{\alpha}{2} \sum_{m=1}^{M} \int_{t_{m-1}}^{t_m} \left(\left(m - t \frac{M}{T} \right) q_{m-1} + \left(t \frac{M}{T} - m + 1 \right) q_m, \right.$$

$$\left. \left(m - t \frac{M}{T} \right) q_{m-1} + \left(t \frac{M}{T} - m + 1 \right) q_m \right) dt.$$

Integration by substitution yields

$$j(q) \approx \frac{T}{6M} \sum_{m=1}^{M} (u_{m-1} - \hat{u}_{m-1}, u_{m-1} - \hat{u}_{m-1}) + (u_{m-1} - \hat{u}_{m-1}, u_m - \hat{u}_m)$$

$$+ (u_m - \hat{u}_m, u_m - \hat{u}_m)$$

$$+ \frac{\alpha T}{6M} \sum_{m=1}^{M} (q_{m-1}, q_{m-1}) + (q_{m-1}, q_m) + (q_m, q_m)$$

$$\approx \frac{T}{6M} \sum_{m=1}^{M} (U_{m-1} - \hat{U}_{m-1}) M_n (U_{m-1} - \hat{U}_{m-1})$$

$$+ (U_{m-1} - \hat{U}_{m-1}) M_n (U_m - \hat{U}_m)$$

$$+ (U_m - \hat{U}_m) M_n (U_m - \hat{U}_m)$$

$$+ \frac{\alpha T}{6M} \sum_{m=1}^{M} Q_{m-1} M_n Q_{m-1} + Q_{m-1} M_n Q_m + Q_m M_n Q_m,$$

where $\hat{U}_m = (\hat{u}(t_m, x_i))_{i=1,...,N}$ and $Q_m = (q(t_m, x_i))_{i=1,...,N}$ for m = 0,...,M.

2.3 Optimization of the control variable

To optimize the control variable, we write every $q \in Q_d$, by using a fixed basis

$$\Phi = \{\phi_1, \dots, \phi_{N_b}\}\tag{2.6}$$

with $\phi_1, \ldots, \phi_{N_b} \in H$ and scalars $q_1^0, q_1^1, \ldots, q_1^M, \ldots, q_{N_b}^0, q_{N_b}^1, \ldots, q_{N_b}^M \in \mathbb{R}$, as

$$q(t,x) = \sum_{i=1}^{N_b} \alpha_i(t)\phi_i(x)$$
(2.7)

with

$$\alpha_i(t) = \begin{cases} q_i^{m-1} \left(m - t \frac{M}{T} \right) + q_i^m \left(t \frac{M}{T} - m + 1 \right) & \text{if } t \in I_m \text{ with } m = 1, \dots, M \\ q_i^0 & \text{if } t = 0 \end{cases}$$

Each control variable that is written in this form can be represented as a vector

$$\mathbf{q} = [q_1^0, q_1^1, \dots, q_1^M, \dots, q_{N_b}^0, q_{N_b}^1, \dots, q_{N_b}^M]^T \in \mathcal{D} := \mathbb{R}^{N_q},$$

with $N_q = (M+1) \cdot N_b$. Therefore, we write

$$j(\mathbf{q}) := j(q)$$

for each q that is defined like in (2.7).

In the next chapters, we present algorithms that minimize $j(\mathbf{q})$ with respect to its control vector \mathbf{q} .

3 Ensemble-based optimization algorithm

The adaptive ensemble-based algorithm (EnOpt) is usually used to maximize the net present value of oil recovery methods with respect to a control vector. Examples are presented in [1], [3], [4]. In this chapter, we want to utilize the EnOpt algorithm to optimize the objective function j. Our implementation is similar to that in [1].

We begin by describing this algorithm for a general function $F: \mathbb{R}^{N_{\mathbf{q}}} \to \mathbb{R}$ to iteratively solve the optimization problem

$$\underset{\mathbf{q} \in \mathcal{D}}{\text{maximize}} F(\mathbf{q}).$$

We start at an initialization \mathbf{q}_0 , which is updated iteratively with a preconditioned gradient ascent method that is given by

$$\mathbf{q}_{k+1} = \mathbf{q}_k + \beta_k \mathbf{d}_k,$$

$$\mathbf{d}_k pprox rac{\mathbf{C}_{\mathbf{q}_k}^k \mathbf{G}_k}{\|\mathbf{C}_{\mathbf{q}_k}^k \mathbf{G}_k\|}_{\infty},$$

where k = 0, 1, 2, ... denotes the optimization iteration. β_k with $\beta_k > 0$ is computed by using a line search. Furthermore, $\mathbf{C}_{\mathbf{q}_k}^k$ denotes the user-defined covariance matrix of the control variables at the k-th iteration and \mathbf{G}_k is the approximate gradient of F with respect to the control variables.

We define the initial covariance matrix $\mathbf{C}_{\mathbf{q}_0}^0$ so that the covariance between controls of different basis functions ϕ_i, ϕ_j is zero and

$$Cov(q_j^i, q_j^{i+h}) = \sigma_j^2 \rho^h \left(\frac{1}{1 - \rho^2} \right), \text{ for all } h \in \{0, \dots, M - i\},$$

where $\sigma_j^2 > 0$ is the variance for the basis function ϕ_j and $\rho \in (-1,1)$ the correlation coefficient.

That means that for $\mathbf{C}_j := \left(\operatorname{Cov}(q_j^i, q_j^k) \right)_{i,k}$ with $j = 1, \dots, N_b$, we set

$$\mathbf{C}_{\mathbf{q}_0}^0 = \begin{pmatrix} \mathbf{C}_1 & 0 & \cdots & 0 \\ 0 & \mathbf{C}_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \mathbf{C}_{N_b} \end{pmatrix}. \tag{3.1}$$

To compute the step direction \mathbf{d}_k at iteration step k, we sample $\mathbf{q}_{k,m} \in \mathcal{D}$ for $m = 1, \ldots, N$, with $N \in \mathbb{N}$, from a multivariate Gaussian distribution with mean \mathbf{q}_k and covariance $\mathbf{C}_{\mathbf{q}_k}^k$, and then we define

$$\mathbf{C}_{\mathbf{q}_{k},F}^{k} := \frac{1}{N-1} \sum_{m=1}^{N} (\mathbf{q}_{k,m} - \mathbf{q}_{k}) (F(\mathbf{q}_{k,m}) - F(\mathbf{q}_{k})). \tag{3.2}$$

Now, we set $\mathbf{d}_k = \frac{\mathbf{C}_{\mathbf{q}_k,F}^k}{\|\mathbf{C}_{\mathbf{q}_k,F}^k\|_{\infty}}$. This is valid since $\mathbf{C}_{\mathbf{q}_k,F}^k$ is an estimation of $\mathbf{C}_{\mathbf{q}_k}^k\mathbf{G}_k$, which can by shown like in [3]. Here, we begin with the Taylor expansion around \mathbf{q}_k and get

$$F(\mathbf{q}) = F(\mathbf{q}_k) + (\mathbf{q} - \mathbf{q}_k)^T \nabla F(\mathbf{q}_k) + O(\|\mathbf{q} - \mathbf{q}_k\|^2)$$

$$\Rightarrow F(\mathbf{q}) - F(\mathbf{q}_k) = (\mathbf{q} - \mathbf{q}_k)^T \mathbf{G}_k + O(\|\mathbf{q} - \mathbf{q}_k\|^2).$$

Multiplying both sides by $(\mathbf{q} - \mathbf{q}_k)$ and setting $\mathbf{q} = \mathbf{q}_{k,m}$ yields

$$(\mathbf{q}_{k,m} - \mathbf{q}_k)(F(\mathbf{q}_{k,m}) - F(\mathbf{q}_k))$$

$$= (\mathbf{q}_{k,m} - \mathbf{q}_k)(\mathbf{q}_{k,m} - \mathbf{q}_k)^T \mathbf{G}_k + O(\|\mathbf{q}_{k,m} - \mathbf{q}_k\|^3),$$

where $O(\|\mathbf{q}_{k,m} - \mathbf{q}_k\|^3)$ are the remaining terms containing order ≥ 3 of $(\mathbf{q}_{k,m} - \mathbf{q}_k)$. Neglecting $O(\|\mathbf{q}_{k,m} - \mathbf{q}_k\|^3)$ gives by summation over all samples and multiplication of both sides with $\frac{1}{N-1}$:

$$\frac{1}{N-1} \sum_{m=1}^{N} (\mathbf{q}_{k,m} - \mathbf{q}_{k}) (F(\mathbf{q}_{k,m}) - F(\mathbf{q}_{k}))$$

$$\approx \left(\frac{1}{N-1} \sum_{m=1}^{N} (\mathbf{q}_{k,m} - \mathbf{q}_{k}) (\mathbf{q}_{k,m} - \mathbf{q}_{k})^{T} \right) \mathbf{G}_{k}$$

$$\implies \mathbf{C}_{\mathbf{q}_{k},F}^{k} \approx \mathbf{C}_{\mathbf{q}_{k}}^{k} \mathbf{G}_{k},$$

since $\frac{1}{N-1} \sum_{m=1}^{N} (\mathbf{q}_{k,m} - \mathbf{q}_k) (\mathbf{q}_{k,m} - \mathbf{q}_k)^T$ is itself an approximation of $\mathbf{C}_{\mathbf{q}_k}^k$.

By using the samples $\{\mathbf{q}_{k-1,m}\}_{m=1}^{N}$ and the covariance matrix $\mathbf{C}_{\mathbf{q}_{k-1}}^{k-1}$ from the last iteration, we update $\mathbf{C}_{\mathbf{q}_{k-1}}^{k-1}$, like in [5], by setting

$$\mathbf{C}_{\mathbf{q}_{k}}^{k} = \mathbf{C}_{\mathbf{q}_{k-1}}^{k-1} + \tilde{\beta}_{k}\tilde{\mathbf{d}}_{k} \text{ with}$$
(3.3)

$$\tilde{\mathbf{d}}_{k} = N^{-1} \sum_{m=1}^{N} (F(\mathbf{q}_{k-1,m}) - F(\mathbf{q}_{k}))((\mathbf{q}_{k-1,m} - \mathbf{q}_{k})(\mathbf{q}_{k-1,m} - \mathbf{q}_{k})^{T} - \mathbf{C}_{\mathbf{q}_{k-1}}^{k-1}),$$
(3.4)

where $\tilde{\beta}_k$ is a step size that is chosen so that no entries of the diagonal of $\mathbf{C}_{\mathbf{q}_k}^k$ are negative. How we set $\tilde{\beta}_k$ is shown below at the implementation of the entire EnOpt algorithm.

Now that we have described the optimization steps of this algorithm, we iterate until $F(\mathbf{q}_k) < F(\mathbf{q}_{k-1}) + \varepsilon$, where $\varepsilon > 0$.

In our implementation, the EnOpt algorithm takes the objective function $F: \mathbb{R}^{N_{\mathbf{q}}} \to \mathbb{R}$, our initial iterate $\mathbf{q}_0 \in \mathbb{R}^{N_{\mathbf{q}}}$, the sample size $N \in \mathbb{N}$, the tolerance $\varepsilon > 0$, the maximum number of iterations $k^* \in \mathbb{N}$, the initial step size $\beta_1 > 0$ for the computation of the next iterate, the initial step size $\beta_2 > 0$ for the iteration of the covariance matrix, the step size contraction $r \in (0,1)$, the maximum number of step size trials $\nu^* \in \mathbb{N}$, the variance $\sigma^2 \in \mathbb{R}^{N_b}$ with positive elements, the correlation coefficient $\rho \in (-1,1)$, the number of time steps $N_t \in \mathbb{N}$, the number of basis functions $N_b \in \mathbb{N}$, a projection PR and an initial

covariance $C_{init} \in \mathbb{R}^{N_q}$. PR is set to the identity function lambda mu : mu and C_{init} is set to None by default.

PR is used to project the inputs onto a given set and \mathbf{C}_{init} is an alternative definition for the initialization of the covariance. In this chapter we do not need to specify these inputs, however PR could be used if we had an optimization problem where the iterates were restricted to a certain spatial domain, which is not the case.

The implementation in pseudo code is shown here:

Algorithm 1 EnOpt algorithm

```
1: function EnOpt(F, \mathbf{q}_0, N, \varepsilon, k^*, \beta_1, \beta_2, r, \nu^*, \sigma^2, \rho, N_t, N_b, PR \leftarrow
                                                                                                                                                lambda mu
     \text{mu}, \mathbf{C}_{\text{init}} \leftarrow \text{None})
            F_k^{\text{prev}} \leftarrow \text{F}(\mathbf{q}_0)
2:
           \mathbf{q}_k, T_k, \mathbf{C}_k, F_k \leftarrow \text{OptStep}(\mathbf{F}, \mathbf{q}_0, N, 0, [], \mathbf{C}_{\text{init}}, F_k^{\text{prev}}, \beta_1, \beta_2, r, \varepsilon, \nu^*, \sigma^2, \rho, N_t, N_b, PR)
3:
4:
           5:
6:
                   \mathbf{q}_k, T_k, \mathbf{C}_k, F_k \leftarrow \text{OptStep}(\mathbf{F}, \mathbf{q}_k, N, k, T_k, \mathbf{C}_k, F_k, \beta_1, \beta_2, r, \varepsilon, \nu^*, \sigma^2, \rho, N_t, N_b, PR)
7 \cdot
                   k \leftarrow k+1
8:
           return q_k, k
9:
```

We begin by initializing F_k^{prev} as $F(\mathbf{q}_0)$. The next iterate \mathbf{q}_k , along with the sample T_k , the covariance \mathbf{C}_k and the function value of \mathbf{q}_k , denoted by F_k , are computed by calling the initial optimization step OPTSTEP, which is shown in algorithm 2.

After we initialized k as 1, we loop until the stop criterion is satisfied. In each optimization loop, F_k^{prev} is set to the function value of the last iterate and $\mathbf{q}_k, T_k, \mathbf{C}_k$ and F_k are updated by calling OPTSTEP again.

The next algorithms use some functions from Python such as LEN for the length of a list or an array, as well as functions from the Python package NumPy, which are identified by starting with 'np.'.

Algorithm 2 OptStep algorithm

```
OPTSTEP(F, \mathbf{q}_k, N, k, T_k, \mathbf{C}_k, F_k, \beta_1, \beta_2, r, \varepsilon, \nu^*, \sigma^2, \rho, N_t, N_b, PR
  1: function
          lambda mu : mu)
                     \begin{aligned} N_{\mathbf{q}} \leftarrow & \text{len}(\mathbf{q}_k) \\ \mathbf{C}_k^{\text{next}} \leftarrow & \text{np.zeros}((N_{\mathbf{q}}, N_{\mathbf{q}})) \end{aligned}
  2:
  3:
                     if k = 0 then
  4:
                               if C_k is None then
  5:
                                         \mathbf{C}_k^{\text{next}} \leftarrow \text{INITCov}(\mathbf{q}_k, \sigma^2, \rho, N_t, N_b)
  6:
  7:
                                         \mathbf{C}_k^{	ext{next}} \leftarrow \mathbf{C}_k
  8:
  9:
                               \mathbf{C}_k^{\text{next}} \leftarrow \text{UPDATECov}(\mathbf{q}_k, T_k, \mathbf{C}_k, F_k, \beta_2)
10:
                    \begin{aligned} & \text{sample} \leftarrow \text{np.random.MULTIVARIATE\_NORMAL}(\mathbf{q}_k, \mathbf{C}_k^{\text{next}}, \text{size} \leftarrow N) \\ & T_k^{\text{next}} \leftarrow [\text{PR}(\text{sample}[j]), \text{F}(\text{PR}(\text{sample}[j]))]_{j=0}^{N-1} \end{aligned}
11:
12:
13:
                     \mathbf{C}_F^k \leftarrow \text{np.zeros}(N_{\mathbf{q}})
                     for m = 0, ..., N - 1 do
14:
                               \mathbf{C}_F^k \leftarrow \mathbf{C}_F^k + (T_k^{\text{next}}[m][0] - \mathbf{q}_k) \cdot (T_k^{\text{next}}[m][1] - F_k)
15:
                    \mathbf{C}_F^k \leftarrow 1/(N-1) \cdot \mathbf{C}_F^k
\mathbf{d}_k \leftarrow \mathbf{C}_F^k / \|\mathbf{C}_F^k\|_{\infty}
\mathbf{q}_k^{\text{next}}, F_k^{\text{next}} \leftarrow \text{LINESEARCH}(\mathbf{F}, \mathbf{q}_k, \mathbf{d}_k, \beta_1, r, \varepsilon, \nu^*, \text{PR})
\mathbf{return} \ \mathbf{q}_k^{\text{next}}, T_k^{\text{next}}, \mathbf{C}_k^{\text{next}}, F_k^{\text{next}}
16:
17:
18:
19:
```

We start the OPTSTEP algorithm by updating the covariance matrix \mathbf{C}_k or, if k is zero, initializing it. In the case that k is zero, we first check if there is a predefined covariance matrix \mathbf{C}_k . When there is one, $\mathbf{C}_k^{\text{next}}$ is set to that matrix. Otherwise, we define $\mathbf{C}_k^{\text{next}}$ by calling the function INITCOV, which sets the matrix like it is described in (3.1).

If k is not zero, we get the updated covariance matrix by calling UPDATECOV, which is described in (3.3) and algorithm 3.

This covariance matrix is now used to get a Gaussian distributed sample with N elements around \mathbf{q}_k . The projected samples PR(sample[j]) and their respective function values F(PR(sample[j])) for $j=0,\ldots,N-1$ are stored in the list T_k^{next} . The definition is here abbreviated as $[PR(\text{sample}[j]), F(PR(\text{sample}[j]))]_{j=0}^{N-1}$. In the code, this is done with an iteration through a for-loop.

After that, we set \mathbf{C}_F^k like it is defined in (3.2), and we set \mathbf{d}_k to \mathbf{C}_F^k divided by its sup norm. Here, $\|\mathbf{C}_F^k\|_{\infty}$ is an abbreviation of np.MAX(np.ABS(\mathbf{C}_F^k)).

The next iterate $\mathbf{q}_k^{\text{next}}$ and its function value F_k^{next} is now computed with the line search

The next iterate $\mathbf{q}_k^{\text{next}}$ and its function value F_k^{next} is now computed with the line search algorithm LineSearch, that is shown in algorithm 4.

We return $\mathbf{q}_k^{\text{next}}, T_k^{\text{next}}, \mathbf{C}_k^{\text{next}}$ and F_k^{next} .

Algorithm 3 Covariance matrix update

```
1: function UPDATECOV(\mathbf{q}_k, T_k, \mathbf{C}_k, F_k, \beta_2)
                 N_{\mathbf{q}} \leftarrow \text{LEN}(\mathbf{q}_k)
                 N \leftarrow \text{Len}(T_k)
  3:
                 \mathbf{d}_k^{\mathrm{cov}} \leftarrow \text{NP.ZEROS}((N_{\mathbf{q}}, N_{\mathbf{q}}))
  4:
                 for m = 0, ..., N - 1 do
  5:
                          \mathbf{d}_k^{\text{cov}} \leftarrow \mathbf{d}_k^{\text{cov}} + (T_k[m][1] - F_k) \cdot ((T_k[m][0] - \mathbf{q}_k).\text{RESHAPE}((N_{\mathbf{q}}, 1)) \cdot (T_k[m][0] - \mathbf{q}_k)
  6:
        \mathbf{q}_k).RESHAPE((1, N_{\mathbf{q}})) - \mathbf{C}_k)
                \mathbf{d}_k^{\text{cov}} \leftarrow \mathbf{d}_k^{\text{cov}}/N
  7:
                 \mathbf{C}_{\mathrm{diag}}, \mathbf{d}_{\mathrm{diag}} \leftarrow \mathrm{np.ZEROS}(N_{\mathbf{q}}), \mathrm{np.ZEROS}(N_{\mathbf{q}})
  8:
                 for i = 0, ..., N_{\mathbf{q}} - 1 do
  9:
                          \mathbf{C}_{\mathrm{diag}}[i] \leftarrow \mathbf{C}_k[i,i]
10:
                         \mathbf{d}_{\mathrm{diag}}[i] \leftarrow \mathbf{d}_k^{\mathrm{cov}}[i,i]
11:
                 \beta_2^{\text{iter}} \leftarrow \beta_2
12:
                 while np.MIN(\mathbf{C}_{\mathrm{diag}} + \beta_2^{\mathrm{iter}} * \mathbf{d}_{\mathrm{diag}}) \le 0 do
13:
                         \beta_2^{\text{iter}} \leftarrow \beta_2^{\text{iter}}/2
14:
                 return \mathbf{C}_k + \beta_2^{\text{iter}} * \mathbf{d}_k^{\text{cov}}
15:
```

For the update of the covariance matrix, we calculate the step direction $\mathbf{d}_k^{\text{cov}}$ like in 3.4 first. Then we want to make sure that the updated covariance matrix has only positive values on its diagonal.

For this purpose, $\mathbf{C}_{\mathrm{diag}}$ and $\mathbf{d}_{\mathrm{diag}}$ are defined as the vector of values on the diagonal of C_k and d_k^{cov} respectively. Then, the step size β_2^{iter} , initialized as β_2 , is halfed until np.MIN($\mathbf{C}_{\text{diag}} + \beta_2^{\text{iter}} * \mathbf{d}_{\text{diag}}$) is positive. Finally, the matrix $\mathbf{C}_k + \beta_2^{\text{iter}} * \mathbf{d}_k^{\text{cov}}$ is returned, like in (3.3).

Algorithm 4 Line search

```
1: function LineSearch(F, \mathbf{q}_k, F_k, \mathbf{d}_k, \beta_1, r, \varepsilon, \nu^*, PR)
                               \beta_1^{\text{iter}} \leftarrow \beta_1 \\ \mathbf{q}_k^{\text{next}} \leftarrow \text{PR}(\mathbf{q}_k + \beta_1^{\text{iter}} \mathbf{d}_k) \\ F_k^{\text{next}} \leftarrow F(\mathbf{q}_k^{\text{next}})
    2:
    3:
    4:
    5:
                                while F_k^{	ext{next}} - F_k \le \varepsilon and \nu < \nu^* do \beta_1^{	ext{iter}} \leftarrow r\beta_1^{	ext{iter}} \mathbf{q}_k^{	ext{next}} \leftarrow \text{PR}(\mathbf{q}_k + \beta_1^{	ext{iter}} \mathbf{d}_k)
    6:
    7:
    8:
                               \begin{array}{c} \nu \leftarrow \nu + 1 \\ F_k^{\text{next}} \leftarrow F(\mathbf{q}_k^{\text{next}}) \\ \mathbf{return} \ \mathbf{q}_k^{\text{next}}, F_k^{\text{next}} \end{array}
    9:
10:
```

At the start of the line search algorithm, we initialize the step size β_1^{iter} as β_1 . Then we repeatedly calculate $\mathbf{q}_k^{\mathrm{next}} = \mathrm{PR}(\mathbf{q}_k + \beta_1^{\mathrm{iter}} \mathbf{d}_k)$ and $F_k^{\mathrm{next}} = \mathrm{F}(\mathbf{q}_k^{\mathrm{next}})$ with simultaneous reduction of β_1^{iter} by multiplication with r until either $F_k^{\mathrm{next}} - F_k > \varepsilon$ or $\nu \geq \nu^*$. After the termination of the while-loop, $\mathbf{q}_k^{\mathrm{next}}$ and F_k^{next} are returned. The reason for stopping the while loop also shows when $\mathbf{q}_k^{\mathrm{next}}$ is the last iteration of the EnOpt algorithm,

as the EnOpt algorithm stops when the function value of the next iteration is not greater

than the function value of the last iteration plus ε , i.e. when $F_k^{\text{next}} > F_k + \varepsilon$. Now we use this algorithm to optimize our objective function j. Since this is a maximization procedure and j should be minimized, we apply -j to the EnOpt algorithm, which gives us:

Algorithm 5 FOM-EnOpt algorithm

- 1: function FOM-ENOPT($\mathbf{q}_0, N, \varepsilon, k^*, \beta_1, \beta_2, r, \nu^*, \sigma^2, \rho, N_t, \mathbf{q}$ -base)
- return EnOpt(-J, $\mathbf{q}_0, N, \varepsilon, k^*, \beta_1, \beta_2, r, \nu^*, \sigma^2, \rho, N_t, \text{LEN}(\mathbf{q}_{\text{base}}))$

 $\mathbf{q}_{\mathrm{base}}$ is here a list that consists of the basis functions, so Φ in (2.6). There are some more inputs that FOM-ENOPT requires, but we omit these as they are only needed for the calculation of J.

4 Adaptive-ML-EnOpt algorithm

In this chapter, we introduce the Adaptive-ML-EnOpt algorithm [1], which is a modified version of the EnOpt algorithm. This algorithm is supposed to reduce the number of FOM evaluations by using a machine learning-based surrogate function, which improves the computation speed with respect to the EnOpt algorithm. Therefore, we introduce deep neural networks (DNNs) next. After that, the Adaptive-ML-EnOpt-algorithm is presented.

4.1 Deep neural networks

This description of deep neural networks is based on the definitions in [1].

DNNs are used here to approximate a function $f: \mathbb{R}^{N_{\text{in}}} \to \mathbb{R}^{N_{\text{out}}}$ with $N_{\text{in}}, N_{\text{out}} \in \mathbb{N}$. We call $L \in \mathbb{N}$ the number of layers and $N_{\text{in}} = N_0, N_1, \dots, N_{L-1}, N_L = N_{\text{out}}$ the number of neurons in each layer. $W_i \in \mathbb{R}^{N_i \times N_{i-1}}$ denotes the weights in layer $i \in \{1, \dots, L\}$ and $b_i \in \mathbb{R}^{N_i}$ the biases of the layer $i \in \{1, \dots, L\}$. These are composed as $\mathbf{W} = ((W_1, b_1), \dots, (W_L, b_L))$, which is a tuple of pairs of corresponding weights and biases.

 $\rho: \mathbb{R} \to \mathbb{R}$ is the so-called activation function. A popular example is the rectified linear unit funtion $\rho(x) = \max(x, 0)$, however we will use the hyperbolic tangent funtion:

$$\rho(x) = \tanh(x) = \frac{\exp(2x) - 1}{\exp(2x) + 1}.$$

 $\rho_n^*: \mathbb{R}^n \to \mathbb{R}^n$ is now defined as the component-wise application of ρ onto a vector of dimension n, so $\rho_n^*(x) = \left[\rho(x_1), \dots, \rho(x_n)\right]^T$ for $x \in \mathbb{R}^n$. To calculate the output $\Phi_{\mathbf{W}}(x) \in \mathbb{R}^{N_{\mathrm{out}}}$ of a DNN for an input $x \in \mathbb{R}^{N_{\mathrm{in}}}$, we apply

To calculate the output $\Phi_{\mathbf{W}}(x) \in \mathbb{R}^{N_{\text{out}}}$ of a DNN for an input $x \in \mathbb{R}^{N_{\text{in}}}$, we apply the weights, biases, and activation function multiple times onto the input. It is calculated iteratively as shown here:

$$r_0(x) := x,$$
 $r_i(x) := \rho_{N_i}^*(W_i r_{i-1}(x) + b_i) \text{ for } i = 1, \dots, L-1,$
 $r_L(x) := W_L r_{L-1}(x) + b_L,$
 $\Phi_{\mathbf{W}}(x) := r_L(x).$

Now we try to optimize the parameters in **W** such that $\Phi_{\mathbf{W}} \approx f$. To achieve this, we sample a set that consists of inputs $x_i \in X \subset \mathbb{R}^{N_{\text{in}}}$ and corresponding outputs $f(x_i) \in \mathbb{R}^{N_{\text{out}}}$ and assemble them in the training set

$$T_{\text{train}} = \{(x_1, f(x_1)), \dots, (x_{N_{\text{train}}}, f(x_{N_{\text{train}}}))\} \subset X \times \mathbb{R}^{N_{\text{out}}}.$$

$$(4.1)$$

To evaluate the performance of our chosen \mathbf{W} , we use the mean squared error loss $\mathcal{L}(\Phi_{\mathbf{W}}, T_{\text{train}})$ to measure the distance between $\Phi_{\mathbf{W}}$ and f on a training set. The mean squared error loss is defined as

$$\mathscr{L}(\Phi_{\mathbf{W}}, T_{\text{train}}) := \frac{1}{|T_{\text{train}}|} \sum_{(x,y) \in T_{\text{train}}} \|\Phi_{\mathbf{W}}(x) - y\|_2^2.$$

Since we want $\Phi_{\mathbf{W}}$ to be close to f, we minimize the loss function with respect to \mathbf{W} . For that, we use some gradient-based optimization method. By the structure of the DNN, we can use the chain rule multiple times to divide the gradient of \mathcal{L} into much simpler gradient computations.

We want that $\Phi_{\mathbf{W}}$ is close to f on X but we train it only on a sample set of X, so we achieve that $\Phi_{\mathbf{W}}$ is only on T_{train} close to f. While we train, the mean squared error loss will eventually get better and better on the training set, but at some point the error on different samples will get worse [6]. We call that overfitting.

To prevent overfitting, we use early stopping. For early stopping, we evaluate the loss function on a validation set $T_{\text{val}} \subset X \times \mathbb{R}^{N_{\text{out}}}$, where usually $T_{\text{val}} \cap T_{\text{train}} = \emptyset$. Our algorithm for early stopping looks like this:

- let $\mathbf{W}^{(k)}$ be the weights in epoch k
- compute $\mathcal{L}(\Phi_{\mathbf{W}^{(k)}}, T_{\text{val}})$ in each epoch
- save $\mathbf{W}^{(k^*)}$ at iteration k^* if it is the minimizer over all previous weights
- if $\mathscr{L}(\Phi_{\mathbf{W}^{(k^*+i)}}, T_{\text{val}}) \ge \mathscr{L}(\Phi_{\mathbf{W}^{(k^*)}}, T_{\text{val}})$ for all i from 0 to a prescribed number: abort the training and use $\mathbf{W}^{(k^*)}$

So we abort the training if the minimum loss is not decreasing over a prescribed number of consecutive epochs. Our reasoning behind that is that the loss on the validation set is not srictly decreasing and can even increase over some epochs, but that is fine for us as long as we can decrease the loss over time.

To summarize, the training of one neural network is shown in algorithm 6. It takes the initialization of the neural network (DNN), the inputs and outputs of the training set $(x_{\text{train}}, y_{\text{train}})$, the inputs and outputs of the testing set $(x_{\text{test}}, y_{\text{test}})$, the loss function (loss_fn), the optimizer (optimizer), the number of training epochs (epochs) and the number (earlyStop) that describes after how many iterations without improvement of the test loss the training gets aborted. In our case, the loss function is chosen as the mean squared error loss and we use the L-BFGS optimizer with strong Wolfe line-search as our optimizer. The number of training epochs is only the maximum number of iterations since we apply early stopping to our training algorithm. Usually, the training terminates earlier because the loss over the testing set is not decreasing further.

The function 'testDNN' in algorithm 6 returns the loss of the function loss in between the output of the DNN with the current parameters and the output of the objective function over the testing set which are saved as y_{train} . So if loss in \mathcal{L} , we have testDNN(DNN, x_{test} , y_{test} , loss in \mathcal{L}) with

$$T_{\text{test}} = \{((x_{\text{test}})_1, (y_{\text{test}})_1), \dots, ((x_{\text{test}})_{N_{\text{test}}}, (y_{\text{test}})_{N_{\text{test}}})\}.$$

Algorithm 6 DNN training

```
1: function TRAINDNN(DNN, x_{\text{train}}, y_{\text{train}}, x_{\text{test}}, y_{\text{test}}, LOSS_FN, OPTIMIZER, epochs, earlyStop)
 2:
        minimalTestLoss \leftarrow TESTDNN(DNN, x_{\text{test}}, y_{\text{test}}, loss_fn)
 3:
        TORCH.SAVE(DNN.STATE_DICT(), 'checkpoint.pth')
 4:
        for epoch = 1, \dots, epochs do
 5:
            DNN.TRAIN()
 6:
            function CLOSURE()
 7:
                y_pred \leftarrow DNN(x_train).Reshape(len(y_train))
 8:
                loss \leftarrow Loss\_FN(y\_pred, y\_train)
 9:
                OPTIMIZER.ZERO_GRAD()
10:
                loss.backward()
11:
                return loss
12:
            OPTIMIZER.STEP(CLOSURE)
13:
            testLoss \leftarrow TESTDNN(DNN, x_{test}, y_{test}, Loss_{FN})
14:
            if testLoss < minimalTestLoss then
15:
                wait \leftarrow 0
16:
                minimalTestLoss \leftarrow testLoss
17:
                TORCH.SAVE(DNN.STATE_DICT(), 'checkpoint.pth')
18:
19:
            else
                wait \leftarrow wait +1
20:
            if wait \geq \text{earlyStop then}
21:
                DNN.LOAD_STATE_DICT(TORCH.LOAD('checkpoint.pth'))
22:
                return
23:
```

Since we search for local minima of the loss function, the initial value $\mathbf{W}^{(0)}$ of our iteration effects the local optimum that we get and therefore the performance. We use Kaiming initialization [7] to set our initial value $\mathbf{W}^{(0)}$. With Kaiming initialization, the starting values are initialized randomly since the elements of the weights W_i are sampled from a zero-mean Gaussian distribution whose standard deviation is $\sqrt{2/N_{i-1}}$ for $i \in \{1, \ldots, L\}$. The biases b_i are set to zero for $i \in \{1, \ldots, L\}$. The idea behind the random sampling is that the specified standard deviation prevents the exponential increase/ reduction of the input as shown in [7].

For the training of the DNN, we perform multiple restarts of the training algorithm with different initializations of $\mathbf{W}^{(0)}$ which minimizes the dependence of our neural network from the initial values. After we have trained enough DNNs, we select the neural network $\Phi_{\mathbf{W}^*}$ that has the smallest evaluation loss $\mathcal{L}(\Phi_{\mathbf{W}^*}, T_{\text{val}})$ over all restarts.

Before the whole algorithm for the construction of the DNN is presented, we look at the data that we use for the training. If we sample the inputs in a small area, it is likely that the corresponding outputs are also close to each other. Since we convert the values for the training of the neural network from double to float, it can even happen that the converted inputs or outputs are constant. In that case, the digits of these values that differ from each other get cut off at the conversion.

We want the values of the inputs and outputs to be distributed in such a way that significant differences are correctly represented. For that, the inputs $x \in \mathbb{R}^n$ and outputs $y \in \mathbb{R}$ are scaled to $\tilde{x} \in [0,1]^n$ and $\tilde{y} \in [0,1]$.

Let

$$T = \{(x_1, y_1), \dots, (x_N, y_N)\}\tag{4.2}$$

be a sample set of size N and

$$T_x = \{x_1, \dots, x_N\},$$
 $T_y = \{y_1, \dots, y_N\}$

the sets that contain the inputs/ output of that sample.

We define $x^{\text{low}}, x^{\text{upp}} \in \mathbb{R}^{\bar{n}}$ and $y^{\text{low}}, y^{\text{upp}} \in \mathbb{R}$ as

$$\begin{aligned} x_i^{\text{low}} &:= & \min\{x_i \mid x \in T_x\} \text{ for } i = 1, \dots, n, \\ x_i^{\text{upp}} &:= & \max\{x_i \mid x \in T_x\} \text{ for } i = 1, \dots, n, \\ y^{\text{low}} &:= & \min T_y, \\ y^{\text{upp}} &:= & \max T_y. \end{aligned}$$

Now, \tilde{x} and \tilde{y} are calculated as

$$\tilde{x}_i = \frac{x_i - x_i^{\text{low}}}{x_i^{\text{upp}} - x_i^{\text{low}}} \text{ for } i = 1, \dots, n,$$

$$\tilde{y} = \frac{y - y^{\text{low}}}{y^{\text{upp}} - y^{\text{low}}}.$$

After we have trained the neural network, we need to rescale the DNN outputs so that we have a proper approximation of the function f. The output $\Phi(\tilde{x})$ of the DNN Φ is rescaled with the calculation $\Phi(\tilde{x}) \cdot (y^{\text{upp}} - y^{\text{low}}) + y^{\text{low}}$.

To summarize, we present now the construction of a DNN as pseudo code. Training parameters like the neural network structure, the activation function, the loss function, the optimizer and the number of training epochs are stored in $V_{\rm DNN}$. We denote $x^{\rm low}$ as minIn, $x^{\rm upp}$ as maxIn, $y^{\rm low}$ as minOut and $y^{\rm upp}$ as maxOut. minIn and maxIn are calculated before the construction of the DNN and are taken as an input.

Algorithm 7 DNN construction

```
1: function CONSTRUCTDNN(sample, V_{\text{DNN}}, minIn, maxIn)
         normSample \leftarrow NP.ZEROS(LEN(sample), LEN(sample[0][0]))
         normVal \leftarrow NP.ZEROS(LEN(sample))
 3:
         for i = 0, \dots, \text{Len(sample)} - 1 \text{ do}
 4:
             normSample[i,:] \leftarrow sample[i][0]
 5:
             normVal[i] \leftarrow sample[i][1]
 6:
         minOut \leftarrow NP.MIN(normVal)
 7:
         maxOut \leftarrow NP.MAX(normVal)
 8:
         SCALEINPUT \leftarrow lambda mu : (mu - minIn)/(maxIn - minIn)
 9:
         SCALEOUTPUT \leftarrow lambda mu : (mu - minOut)/(maxOut - minOut)
10:
         RESCALEOUTPUT \leftarrow lambda mu : mu · (maxOut - minOut) + minOut
11:
12:
         normSample \leftarrow SCALEINPUT(normSample)
         normVal \leftarrow SCALEOutput(normVal)
13:
         x \leftarrow \text{TORCH.FROM\_NUMPY}(\text{normSample}).\text{TO}(\text{torch.float32})
14:
         y \leftarrow \text{TORCH.FROM\_NUMPY}(\text{normVal}).\text{TO}(\text{torch.float32})
15:
         train\_split \leftarrow INT(0.8 * LEN(x))
16:
         x_{\text{train}}, y_{\text{train}} \leftarrow x[: \text{train\_split}], y[: \text{train\_split}]
17:
         x_{\text{test}}, y_{\text{test}} \leftarrow x[\text{train\_split}:], y[\text{train\_split}:]
18:
         DNN \leftarrow FullyConnectedNN(V_{\text{DNN}}[0], activation_function = V_{\text{DNN}}[1])
19:
         LOSS\_FN \leftarrow NN.MSELOSS()
20:
         OPTIMIZER
                               \leftarrow
                                          TORCH.OPTIM.LBFGS(DNN.PARAMETERS(), lr
21:
    learning_rate, line_search_fn \leftarrow 'strong_wolfe')
22:
         TRAINDNN(DNN, x_{\text{train}}, y_{\text{train}}, x_{\text{test}}, y_{\text{test}}, LOSS_FN, OPTIMIZER, epochs)
         EVALDNN \leftarrow TESTDNN(DNN, x_{\text{test}}, y_{\text{test}}, \text{loss\_fn})
23:
         for i = 1, \ldots, \text{numberOfRestarts do}
24:
             DNN_i \leftarrow FullyConnectedNN(V_{DNN}[0], ACTIVATION\_FUNCTION = V_{DNN}[1])
25:
                                           TORCH.OPTIM.LBFGS(DNN_i.PARAMETERS( ), lr
             OPTIMIZER
                                  \leftarrow
26:
    learning_rate, line_search_fn \leftarrow 'strong_wolfe')
             TRAINDNN(DNN<sub>i</sub>, x_{\text{train}}, y_{\text{train}}, x_{\text{test}}, y_{\text{test}}, LOSS_FN, OPTIMIZER, epochs)
27:
             \text{EVALDNN}_i \leftarrow \text{TESTDNN}(\text{DNN}_i, x_{\text{test}}, y_{\text{test}}, \text{Loss\_fn})
28:
             if evalDNN_i < evalDNN then
29:
                  evalDNN \leftarrow evalDNN_i
30:
                  DNN \leftarrow DNN_i
31:
32:
         F_{\text{ML}} \leftarrow \text{lambda} \text{ mu} : \text{RESCALEOUTPUT}(\text{DNN}(\text{SCALEINPUT}(\text{mu})))
         return F_{\rm ML}
33:
```

4.2 Modifying the EnOpt algorithm by using a neural network-based surrogate

For the next step, we use neural networks like in [1] to get an improved version of the EnOpt algorithm. Here, we want to replace calls of the FOM function by a surrogate that is a neural network. The surrogate is supposed to be a local approximation of the FOM function around the current iterate. Doing that globally with a sufficient precision would be computationally too expensive.

Since the surrogate is only locally exact, we will introduce a trust region method. For that, we use an algorithm that projects inputs into the trust region.

Algorithm 8 Projection

```
1: function TR-PROJECTION(x, \mathbf{q}_k, \mathbf{d}_k)
```

- 2: upp $\leftarrow \mathbf{q}_k + \mathbf{d}_k$
- 3: $low \leftarrow \mathbf{q}_k \mathbf{d}_k$
- 4: **return** NP.MAXIMUM(NP.MINIMUM(x, upp), low)

We describe now the Adaptive-ML-EnOpt algorithm. This algorithm takes the initial guess $\mathbf{q}_0 \in \mathbb{R}^{N_u}$, the sample size $N \in \mathbb{N}$, the tolerances $\varepsilon_o, \varepsilon_i > 0$ for the outer and inner iterations, the maximum number of outer and inner iterations $k_o^*, k_i^* \in \mathbb{N}$, the DNN-specific variables V_{DNN} , the initial step size $\beta > 0$, the step size contraction $r \in (0,1)$ and the maximum number of step size trials $\nu^* \in \mathbb{N}$.

Algorithm 9 Adaptive-ML-EnOpt algorithm

```
1: function AML-ENOPT(F, \mathbf{q}_0, N, \varepsilon_o, \varepsilon_i, k_o^*, k_i^*, V_{\text{DNN}}, \beta_1, \beta_2, r, \nu^*, \sigma^2, \rho)
                F_k \leftarrow F(\mathbf{q}_0)
  2:
                F_k^{\text{next}} \leftarrow F_k
  3:
                \tilde{\mathbf{q}}_k, T_k, \mathbf{C}_k, \tilde{F}_k \leftarrow \text{OptStep}(F, \mathbf{q}_0, N, 0, [], \text{None}, F_k, \beta_1, \beta_2, r, \varepsilon_o, \nu^*, \sigma^2, \rho)
  4:
  5:
                \mathbf{q}_k \leftarrow \mathbf{q}_0
  6:
               \mathbf{q}_k^{	ext{next}} \leftarrow \mathbf{q}_k.	ext{copy}()
  7:
                \delta \leftarrow 40
  8:
                while \tilde{F}_k > F_k + \varepsilon_o and k < k_o^* do
  9:
                        T_k^x \leftarrow \text{NP.ZEROS}((N, nt + 1))
10:
                        for i = 0, ..., N - 1 do
11:
                               T_k^x[i,:] \leftarrow T_k[i][0]
12:
                        \min \text{In} \leftarrow \text{NP.ZEROS}(nt+1)
13:
                        \max In \leftarrow \text{NP.ZEROS}(nt+1)
14:
                        for i = 0, ..., nt - 1 do
15:
                                \min \operatorname{In}[i] \leftarrow \operatorname{NP.MIN}(T_k^x[:,i])
16:
                                \max \text{In}[i] \leftarrow \text{NP.MAX}(T_k^x[:,i])
17:
                        \mathbf{d}_k \leftarrow |\mathbf{q}_k - \tilde{\mathbf{q}}_k|
18:
                        while F_k^{\text{next}} \leq F_k + \varepsilon_o do
19:
                                F_{\text{ML}}^k \leftarrow \text{CONSTRUCTDNN}(T_k, V_{\text{DNN}}, \text{minIn}, \text{maxIn})
20:
                               F_k^{\text{ML}} \leftarrow F_{\text{ML}}^k(\mathbf{q}_k)
21:
                                flag_{TR} \leftarrow True
22:
                                while flag_{TR} do
23:
                                       \mathbf{d}_k^{\text{iter}} \leftarrow \delta \cdot \mathbf{d}_k
24:
                                       \mathbf{q}_k^{	ext{next}} \leftarrow 	ext{EnOpt}(F_{	ext{ML}}^k, \mathbf{q}_k, N, \varepsilon_i, k_i^*, \beta_1, \beta_2, r, \nu^*, \sigma^2, \rho, \mathbf{lambda} 	ext{ mu} :
25:
        TR-PROJECTION(mu, \mathbf{q}_k, \mathbf{d}_k^{\text{iter}}), \mathbf{C}_{\text{init}} \leftarrow \mathbf{C}_k)[0]
                                       F_k^{\text{next}} \leftarrow F(\mathbf{q}_k^{\text{next}})
\rho_k \leftarrow \frac{F_k^{\text{next}} - F_k}{F_{\text{ML}}^{\text{ML}}(\mathbf{q}_k^{\text{next}}) - F_k^{\text{approx}}}
26:
27:
                                       if \rho_k < 0.25 then
28:
                                               \delta \leftarrow 0.25 \cdot \delta
29:
                                       if \rho_k > 0.75 and NP.ANY(NP.ABS(\mathbf{q}_k - \mathbf{q}_k^{\text{next}}) -\mathbf{d}_k^{\text{iter}} == 0) then
30:
                                               \delta \leftarrow 2 \cdot \delta
31:
                                       if \rho_k > 0 then
32:
                                               \mathrm{flag}_{\mathrm{TR}} \leftarrow \mathbf{False}
33:
                        \tilde{\mathbf{q}}_k, T_k, \mathbf{C}_k, \tilde{F}_k \leftarrow \mathrm{OptStep}(F, \mathbf{q}_k^{\mathrm{next}}, N, k, T_k, \mathbf{C}_k, F_k^{\mathrm{next}}, \beta_1, \beta_2, r, \varepsilon, \nu^*, \sigma^2, \rho)
34:
                        F_k \leftarrow F_k^{\text{next}}
35:
                        \mathbf{q}_k \leftarrow \mathbf{q}_k^{\text{next}}.\text{COPY}()
36:
                        k \leftarrow k+1
37:
               return \mathbf{q}^* \leftarrow \mathbf{q}_k
38:
```

Algorithm 10 ROM-EnOpt algorithm

- 1: function ROM-ENOPT($\mathbf{q}_0, N, \varepsilon_o, \varepsilon_i, k_o^*, k_i^*, V_{\mathrm{DNN}}, \beta_1, \beta_2, r, \nu^*, \sigma^2, \rho$)
 2: return AML-ENOPT($-j, \mathbf{q}_0, N, \varepsilon_o, \varepsilon_i, k_o^*, k_i^*, V_{\mathrm{DNN}}, \beta_1, \beta_2, r, \nu^*, \sigma^2, \rho$)

5 Numerical experiments

Bibliography

- [1] T. Keil, H. Kleikamp, R. J. Lorentzen, M. B. Oguntola, and M. Ohlberger, "Adaptive machine learning-based surrogate modeling to accelerate PDE-constrained optimization in enhanced oil recovery," *Advances in Computational Mathematics*, vol. 48, no. 6, p. 73, Nov. 2022.
- [2] D. Meidner and B. Vexler, "A priori error estimates for space-time finite element discretization of parabolic optimal control problems part i: Problems without control constraints," SIAM Journal on Control and Optimization, vol. 47, no. 3, pp. 1150–1177, 2008. DOI: 10.1137/070694016. eprint: https://doi.org/10.1137/070694016. [Online]. Available: https://doi.org/10.1137/070694016.
- [3] M. B. Oguntola and R. J. Lorentzen, "Ensemble-based constrained optimization using an exterior penalty method," *Journal of Petroleum Science and Engineering*, vol. 207, p. 109165, 2021, ISSN: 0920-4105. DOI: https://doi.org/10.1016/j.petrol.2021. 109165. [Online]. Available: https://www.sciencedirect.com/science/article/pii/S0920410521008184.
- [4] Y. Zhang, A. S. Stordal, and R. J. Lorentzen, "A natural hessian approximation for ensemble based optimization," en, *Comput. Geosci.*, vol. 27, no. 2, pp. 355–364, Apr. 2023.
- [5] A. S. Stordal, S. P. Szklarz, and O. Leeuwenburgh, "A theoretical look at Ensemble-Based optimization in reservoir management," *Mathematical Geosciences*, vol. 48, no. 4, pp. 399–417, May 2016.
- [6] L. Prechelt, "Early stopping but when?" In Neural Networks: Tricks of the Trade: Second Edition, G. Montavon, G. B. Orr, and K.-R. Müller, Eds. Berlin, Heidelberg: Springer Berlin Heidelberg, 2012, pp. 53–67, ISBN: 978-3-642-35289-8. DOI: 10.1007/978-3-642-35289-8_5. [Online]. Available: https://doi.org/10.1007/978-3-642-35289-8_5.
- [7] K. He, X. Zhang, S. Ren, and J. Sun, "Delving deep into rectifiers: Surpassing human-level performance on imagenet classification," in 2015 IEEE International Conference on Computer Vision (ICCV), 2015, pp. 1026–1034. DOI: 10.1109/ICCV.2015.123.