# Polynomial Chaos Expanded Gaussian Process

**Dominik Polke\*, Tim Kösters, Elmar Ahle**
Electrical Engineering and Computer Science
University of Applied Sciences Niederrhein
Krefeld, Germany
{dominik.polke, elmar.ahle}@hs-niederrhein.de
tim.koesters@stud.hn.de

**Dirk Söffker**
Dynamics and Control
University of Duisburg-Essen
Duisburg, Germany
soeffker@uni-due.de

## Abstract

In complex and unknown processes, global models are initially generated over the entire experimental space, but they often fail to provide accurate predictions in local areas. Recognizing this limitation, this study addresses the need for models that effectively represent both global and local experimental spaces. It introduces a novel machine learning (ML) approach: Polynomial Chaos Expanded Gaussian Process (PCEGP), leveraging polynomial chaos expansion (PCE) to calculate input-dependent hyperparameters of the Gaussian process (GP). This approach provides a mathematically interpretable method that incorporates non-stationary covariance functions and heteroscedastic noise estimation to generate locally adapted models. The model performance is compared to different algorithms in benchmark tests for regression tasks. The results demonstrate low prediction errors of the PCEGP in these benchmark applications, highlighting model performance that is often competitive with or superior to previous methods. A key advantage of the presented model is the transparency and traceability in the calculation of hyperparameters and model predictions.

*Keywords* Machine learning, Gaussian process, polynomial chaos expansion, Bayesian hyperparameter optimization, non-stationary covariance functions

## 1 Introduction

System modeling of real processes occurs within a specific global experimental space. A widely used choice for modeling unknown system behavior is the Gaussian process (GP), presented in [1]. As modeling often aims to follow some objective, e.g., the optimization of a target variable, the model prediction is improved in areas related to that desired objective. Therefore, more experiments are placed in the area of interest, resulting in different densities of data in the whole experimental space [2]. It is also possible that even with uniformly distributed density of training data, the smoothness of the target variable changes in the experimental space and requires locally adapted models to best represent the system behavior [3].

Therefore, models which can handle different densities of training data and the variety of smoothness of the target variable over the whole experimental space are needed. In addition, these kind of models should also be able to model the heteroscedastic noise of data, as presented, e.g., in [4].

One approach to handle different densities of data across the entire design space with GPs is to use non-stationary covariance functions. Non-stationarity allows the modeling of the input-output relationship based on the input value range [5]. In [3], a non-stationary version of the Matérn covariance function is presented, which allows to model non-stationary behavior in data. A second level GP is used in [5] to calculate the lengthscale parameters of the first level GP in order to model non-stationary behavior. For the anisotropic case, this leads to a second level GP for each input, which can be computationally expensive due to the high dimensions.

In [6], the use of deep Gaussian processes (DGPs) with Bayesian optimization (BO) is investigated. The non-stationarity of data is modeled by considering a functional composition of stationary GPs, providing a multiple layer structure. However, the nesting of different models leads to a more difficult interpretability and is challenging for model training. Furthermore, there are approaches, [7], that use local experts to model non-stationary behavior. Such approaches require training and testing of several different models.

In [2], the non-stationarity of GPs is achieved by using a deep neural network (DNN) for hyperparameter estimation of the GP depending on the input data. Since DNNs are a type of black box learning, the training process of a DNN lacks transparency. This makes it difficult to understand how exactly the network is trained and how it makes its decisions [8].

These approaches result in either a structural limitation of the model, a computationally expensive model, the need for different local models, or a model that is difficult to interpret. In this work, a fully transparent and mathematically interpretable method for generating non-stationary covariance functions and heteroscedastic noise estimation is presented. Here, the novel combination of polynomial chaos expansion (PCE) and Gaussian processes (GPs), called Polynomial Chaos Expanded Gaussian Process (PCEGP) is presented. The PCE is utilized to calculate input dependent hyperparameters of the GP. This approach provides the ability to model the global and local behavior of a complex process with only one model.

Besides providing mathematical transparency and traceability, PCE offers several advantages. For example, the PCE is not numerically expensive to evaluate and works well for machine learning (ML) regression on small data sets [9]. In addition, the PCE representation allows the computation of statistical moments and sensitivities of the output with respect to the inputs as described in [10] and [11].

This paper is structured as follows. First, a short introduction and further references to GPs, PCE, and Bayesian hyperparameter optimization are given in Section 2. In Section 3, the ML approach with the novel combination of PCE and GP is introduced. Then, the method to generate non-stationary covariance functions and the heteroscedastic noise estimation are introduced. In Section 4, the Bayesian hyperparameter optimization for the PCEGP is presented. The evaluation of the new ML approach is carried out in Section 5, using various benchmark datasets. The results illustrate the PCEGP's ability to achieve low prediction errors across diverse applications, underscoring a model performance that competes with or surpasses previous methods. Finally, a summary and outlook for future work are given in Section 6.

## 2 Background and Related Work

In this section, the necessary components for the novel ML approach are presented. This includes the GP, the PCE, and the Bayesian optimization for hyperparameter optimization using the tree-structured Parzen estimator.

### 2.1 Gaussian Process

The GP is a probabilistic ML approach using kernel machines [1]. The principle advantage over other kernel-based approaches is given by an uncertainty measurement of the predictions.

A GP is a collection of random variables. Each finite number of these random variables have a joint Gaussian distribution. In the following the GP is considered for regression, where the training data are denoted as $\mathcal{D} = \{\boldsymbol{x}_i, y_i\}_{i=1}^N$ with $\boldsymbol{x}_i \in \mathbb{R}^{n_x}$ and $y_i \in \mathbb{R}$. The number of training points is denoted with $N$ and the number of inputs is denoted with $n_x$.

The mean function $m(\boldsymbol{x})$ and the covariance function $k(\boldsymbol{x}, \boldsymbol{x}')$, also called kernel function, completely describe the GP

$$f(\boldsymbol{x}) \sim \mathcal{GP}\left(m(\boldsymbol{x}), k(\boldsymbol{x}, \boldsymbol{x}')\right). \tag{1}$$

Most ML applications have a large number of input para-meters, e.g., in chemical engineering [12]. This results in a high-dimensional feature space for the learning algorithm. For distance-based ML algorithms, this leads to the so-called curse of dimensionality [13]. To avoid this, unimportant features can be removed with automatic relevance determination (ARD) [14] or by techniques, presented in [15]. For GPs, the ARD is realized by introducing a lengthscale parameter for each feature [16].

A common choice for the stationary covariance function of the GP is the squared exponential kernel with ARD, which is introduced in [1]. It is defined by

$$k\left(\boldsymbol{x}, \boldsymbol{x}'\right) = \sigma_{\mathrm{f}}^2 \exp\left(-\frac{1}{2}\left(\boldsymbol{x} - \boldsymbol{x}'\right)^{\mathsf{T}} \boldsymbol{L}\left(\boldsymbol{x} - \boldsymbol{x}'\right)\right), \tag{2}$$

with the lengthscale matrix $\boldsymbol{L} = \mathrm{diag}\left([l_1 \ldots l_{n_x}]\right)^{-2}$, which includes lengthscale parameter $l_i \in \mathbb{R}^+$ for each input, and the signal variance $\sigma_{\mathrm{f}}^2 \in \mathbb{R}^+$. The noise variance $\sigma_{\mathrm{n}}^2 \in \mathbb{R}^+$ is used to model the noise in the data and also serves as hyperparameter of the GP. These hyperparameters are used to fit the model to the given data. In the following all hyperparameters are summarized in $\boldsymbol{\theta} \in \mathbb{R}^{n_x+2}$, which leads to the covariance matrix

$$\boldsymbol{K}\left(\boldsymbol{X}, \boldsymbol{\theta}\right) = \boldsymbol{K}\left(\boldsymbol{X}, \boldsymbol{L}, \sigma_{\mathrm{f}}\right) + \sigma_{\mathrm{n}}^2 \boldsymbol{I}. \tag{3}$$

In the following, $\boldsymbol{K}\left(\boldsymbol{X}\right) \in \mathbb{R}^{N \times N}$ denotes the covariance matrix depending on the training data and $\boldsymbol{k}\left(\boldsymbol{X}, \boldsymbol{x}^*\right) \in \mathbb{R}^{N \times 1}$ the covariance vector depending on the training data and the next point $\boldsymbol{x}^*$ to predict. The input training points of the whole dataset $\mathcal{D}$ are denoted as $\boldsymbol{X} \in \mathbb{R}^{N \times n_x}$. The covariance matrix $\boldsymbol{K}\left(\boldsymbol{X}\right)$ is calculated by the covariance function $k\left(\boldsymbol{x}, \boldsymbol{x}'\right)$ and thus dependent on the training data $\mathcal{D}$ and the hyperparameters $\boldsymbol{L}$, $\sigma_{\mathrm{f}}$, and $\sigma_{\mathrm{n}}$.

Usually, these hyperparameters are optimized by maximizing the marginal log likelihood (MLL), which is given by

$$\log p\left(\boldsymbol{y} \mid \boldsymbol{X}\right) = -\frac{1}{2}\,\boldsymbol{y}^{\intercal}\boldsymbol{K}\left(\boldsymbol{X}\right)^{-1}\boldsymbol{y} - \frac{1}{2}\log|\boldsymbol{K}\left(\boldsymbol{X}\right)| - \frac{N}{2}\log\left(2\pi\right)$$

under the assumption that $\boldsymbol{y} \sim \mathcal{N}\left(0, \boldsymbol{K}\left(\boldsymbol{X}\right) + \sigma_{\mathrm{n}}^2\,\boldsymbol{I}\right)$ for the GP output. For the optimization of the hyperparameters with a gradient descent algorithm, the gradient of the marginal likelihood is used and given by

$$\frac{\partial}{\partial \boldsymbol{\theta}_i}\log p\left(\boldsymbol{y} \mid \boldsymbol{X}, \boldsymbol{\theta}\right) = \frac{1}{2}\,\boldsymbol{y}^{\intercal}\,\boldsymbol{K}\left(\boldsymbol{X}\right)^{-1}\frac{\partial}{\partial \boldsymbol{\theta}_i}\boldsymbol{K}\left(\boldsymbol{X}\right)^{-1}\boldsymbol{y} - \frac{1}{2}\mathrm{tr}\left(\boldsymbol{K}\left(\boldsymbol{X}\right)^{-1}\frac{\partial \boldsymbol{K}\left(\boldsymbol{X}\right)}{\partial \boldsymbol{\theta}_i}\right). \tag{4}$$

To predict an output based on a new input $\boldsymbol{x}^*$ and the given dataset $\mathcal{D}$, the posterior distribution is defined by

$$\mu\left(\boldsymbol{x}^*\right) = \boldsymbol{k}\left(\boldsymbol{X}, \boldsymbol{x}^*\right)^{\intercal}\boldsymbol{K}\left(\boldsymbol{X}\right)^{-1}\boldsymbol{y}, \tag{5}$$

$$\sigma^2\left(\boldsymbol{x}^*\right) = k\left(\boldsymbol{x}^*, \boldsymbol{x}^*\right) - \boldsymbol{k}\left(\boldsymbol{X}, \boldsymbol{x}^*\right)^{\intercal}\boldsymbol{K}\left(\boldsymbol{X}\right)^{-1}\boldsymbol{k}\left(\boldsymbol{X}, \boldsymbol{x}^*\right), \tag{6}$$

where $\mu\left(\boldsymbol{x}^*\right) \in \mathbb{R}$ denotes the mean prediction and $\sigma^2\left(\boldsymbol{x}^*\right) \in \mathbb{R}^+$ the uncertainty prediction of the point to evaluate $\boldsymbol{x}^* \in \mathbb{R}^{n_x}$.

For GPs, methods exists to generate non-stationary covariance functions, e.g., [3] and [4]. In [2] the hyperparameters of the GP are calculated with a DNN, dependent on the training data. Thus, any stationary covariance function can be transformed to a non-stationary covariance function.

## 2.2 Polynomial Chaos Expansion

The polynomial chaos expansion, presented in [17], is a popular method in uncertainty quantification (UQ) of models which are parametrized by independent random variables with extensions to dependent variables [18]. By viewing the model as an input-output map, the effect of input to output uncertainties can be assessed in UQ [19]. The PCE allows the modeling of different data distributions with different kinds of polynomial bases.

In the case of UQ for high-fidelity models, a large number of simulations is required. To avoid the computational load, surrogate models can be used. One of the most widely used methods to build a surrogate model in UQ is the generalized polynomial chaos. The PCE is well suited for the approximation of functions with random variables, because of the orthogonality of the polynomial basis to the probability measure of the variables [20].

In general, the polynomial chaos expansion is defined by

$$f\left(x\right) = \sum_{i=0}^{\infty} \alpha_i\,\phi_i\left(x\right), \tag{7}$$

where $\alpha_i \in \mathbb{R}$ denotes the polynomial coefficients and $\phi_i\left(x\right)$ represents the polynomial basis. This infinite summation has to be truncated at the finite term $q \in \mathbb{N}_0$, which results in

$$f\left(x\right) = \sum_{i=0}^{q} \alpha_i\,\phi_i\left(x\right). \tag{8}$$

In this work, the coefficients of the PCE are optimized in a two-step optimization strategy, which consists of the Tree-structured Parzen Estimator (TPE) and the stochastic gradient descent algorithm Adam [21], presented in Section 3.

There exists a number of orthogonal polynomial bases and their corresponding probability density functions. For example, the basis of the Legendre-polynomials, which are defined by the polynomial function

$$Le_n(x) = \frac{1}{2^n}\binom{2n}{n}(1 - x^2)^{n/2} \tag{9}$$

and the probability density function

$$p\left(x\right) = \frac{1}{2}. \tag{10}$$

As demonstrated by [9] and [22], there are already studies showing the application of PCE beyond the typical use case for UQ. In [9], it is shown how the PCE can be used for ML regression. In [22], a combination of PCE and GP for ML is presented, in which the PCE is used to describe the global behavior of the computational model while the GP describes the local variability of the model. This combination of PCE and GP is completely different to the approach developed in this work.

### 2.3 Bayesian Hyperparameter Optimization Using the Tree-Structured Parzen Estimator

Bayesian optimization is a powerful method for efficient global optimization of expensive and unknown black-box functions [23]. This optimization approach consists of two parts. The first part is a probabilistic model used to model the black-box function. The second part of Bayesian optimization is the acquisition function. The acquisition function expresses the greatest possible improvement for the selected optimization goal, based on the model prediction.

In Bayesian optimization, the model quality improves over iterations, when successive measurements are incorporated in model training [23]. For TPE, the choice of the acquisition function is the expected improvement (EI) acquisition function. This function balances exploration with exploitation during the search [24]. Exploration involves investigating unexplored regions to discover new solutions, while exploitation focuses on optimizing known promising solutions. With EI, a balance between discovering potential improvements and optimizing known areas can be achieved and is described by

$$\mathrm{EI}\left(x,\,\xi\right) = \left(\mu - f\left(x^*\right) - \xi\right)\Phi\left(\frac{\mu - f\left(x^*\right) - \xi}{\sigma}\right) + \sigma\,\varphi\left(\frac{\mu - f\left(x^*\right) - \xi}{\sigma}\right), \tag{11}$$

with the mean $\mu$, the standard deviation $\sigma$, the so far observed minima $f\left(x^*\right)$, and a weight factor $\xi \in \mathbb{R}$ for fine tuning between exploitation and exploration. In addition, $\Phi$ denotes the probability density function and $\varphi$ denotes the cumulative distribution function.

The Bayesian optimization for hyperparameter optimization proposes the next set of hyperparameters and evaluates the objective function with these hyperparameters. These steps are carried out until a termination criterion is fulfilled. In this work, the Bayesian optimization is used for hyperparameter optimization of the ML model, which is generally introduced in [25] and [26].

How Bayesian optimization can be accelerated to hyperparameter optimization is shown in [27]. The tree-structured Parzen estimator [28] is used to model the objective function, which has to be optimized. The inputs of this function are the hyperparameters of the model and the output is the error of the ML model as an objective to be minimized.

## 3 Polynomial Chaos Expanded Gaussian Processes

In this section the novel ML approach as a combination of PCE and GP is presented, where the PCE is used to calculate the hyperparameters of the GP. First, the general architecture of this approach is presented. Then, the calculation of the GP's lengthscale parameters through the PCE is introduced. After that, the heteroscedastic noise estimation of this approach is presented. In Appendix D, the calculation of the hyperparameters for scalar inputs and multiple inputs is explained in detail.

### 3.1 Architecture of the Polynomial Chaos Expanded Gaussian Process

The general architecture of the PCEGP is shown in Figure 1. It consists of data scaling of the input $x^*$, data point dependent hyperparameter calculation, and the prediction of the output $\hat{y}_{\mathrm{PCEGP}}$ in $\mathbb{R}^{2\times 1}$, which consists of the mean and uncertainty prediction. First, the given input data point $x^*$ is scaled to $x_{\mathrm{s}}^*$ with a scaler, e.g., the min-max-scaler.

As illustrated in Figure 1, the scaled data point $x_{\mathrm{s}}^*$ is the input for the PCEs and also for the covariance functions of the GP. The PCEs predict the hyperparameters $\hat{l}_1\left(x_{\mathrm{s}}^*\right)\dots\hat{l}_{n_k}\left(x_{\mathrm{s}}^*\right)$ for each covariance funcion $k_1\dots k_{n_k}$, where $n_k \in \mathbb{N}$ denotes the number of used covariance functions. In Appendix E, the algorithm for model prediction is shown.

The choice of the PCE basis function depends on the input data. In Appendix C, different polynomial bases are shown. The approach of this study allows the combination of different polynomial bases for the calculation of the GP hyperparameters depending on the input data. This eliminates the choice of the polynomial basis and the weighting of these becomes part of the hyperparameter optimization. As shown in Figure 1, the combination of different PCE basis functions is achieved by a linear combination of these bases.

To be able to perform versatile modeling tasks with the novel approach, different covariance functions are combined and the influence of the respective covariance function is controlled via hyperparameter optimization. The covariance functions used in this work are shown in Appendix B and can be exchanged or extended by further ones as needed. The combination of the covariance functions is realized by summation of the covariance functions with

$$k_{\Sigma}\left(\boldsymbol{x},\boldsymbol{x}'\right) = k_1\left(\boldsymbol{x},\boldsymbol{x}'\right) + k_2\left(\boldsymbol{x},\boldsymbol{x}'\right) + \cdots + k_{n_k}\left(\boldsymbol{x},\boldsymbol{x}'\right). \tag{12}$$

In this work, the covariance functions and polynomial bases have to be chosen manually and are not part of the hyperparameter optimization. In Section 5, the setup of the model for the experiments carried out in this work is presented in detail.
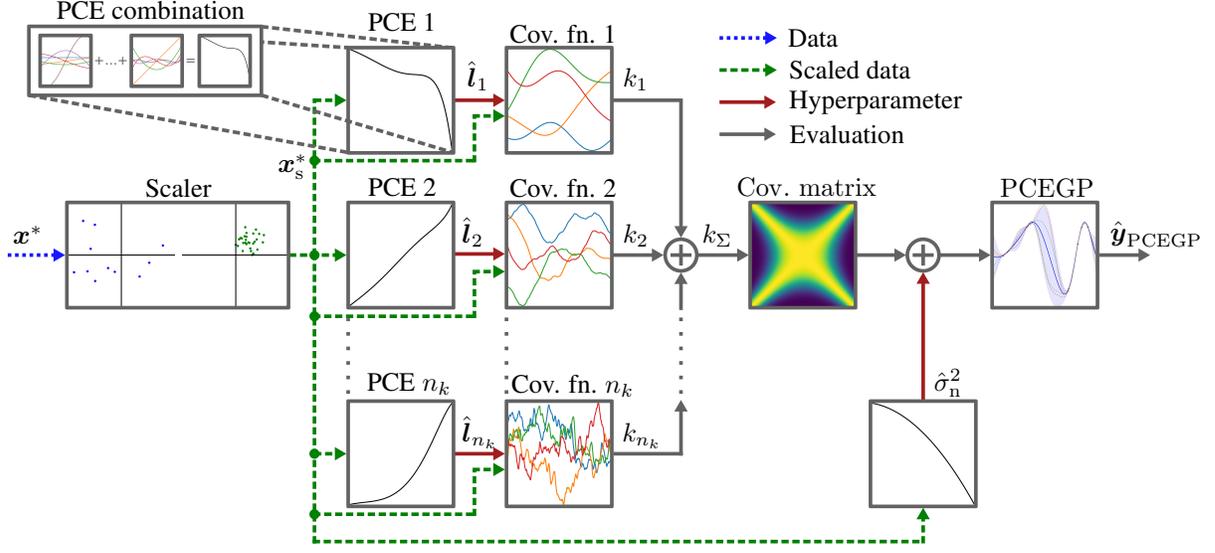
4

Figure 1: Evaluation of the Polynomial Chaos Expanded Gaussian Process with data scaling, point-dependent hyperparameter calculation, and output prediction

## 3.2 PCE for Non-Stationary Covariance Function

Non-stationary covariance functions enable the model to dynamically adjust to functions characterized by varying degrees of smoothness and training point densities across the input space. For GPs, this can be achieved by the estimation of an input dependent lengthscale parameter vector $l(x)$ of a chosen covariance function. In this work, PCE is used to calculate the lengthscale parameters as a function of the input data.

For this purpose, the used covariance functions in Appendix B have to be adapted to estimate a data point dependent lengthscale $l(x)$. For example, the squared exponential covariance function

$$k(x, x') = \sigma_{\mathrm{f}}^2 \exp\left(-\frac{(x - x')^2}{2\, l^2}\right), \tag{13}$$

is changed to

$$k(x, x') = \sigma_{\mathrm{f}}^2 \exp\left(-\frac{(l(x) \odot x - l(x') \odot x')^2}{2}\right), \tag{14}$$

which enables the possibility to calculate the lengthscale for each data point separately by PCE. Here, $\odot$ defines the Hadamard product. This approach is adapted from [2], where this is done via a DNN. This makes it possible to transform any stationary covariance function into a non-stationary one. Here, $l(x)$ is generally described as an input-dependent lengthscale. The lengthscale value calculated by PCE is denoted with $\hat{l}(x)$. In Figure 1, the case $\hat{l}(x_{\mathrm{s}}^*)$ is shown, where the data point $x_{\mathrm{s}}^*$ is already scaled.

The choice of the polynomial basis depends on the distribution of the input data. With

$$\hat{l}(x) = \sum_{i=0}^{q} \alpha_{l,i}\, \phi_i(x), \tag{15}$$

the lengthscale becomes an input dependent hyperparameter for the covariance functions, estimated by the PCE. From this point onward, the optimizer's task shifts from directly adjusting the lengthscale parameters to optimizing the $q + 1$ coefficients of the PCE. In the case that the lengthscale parameter is calculated by a linear combination of $b$ polynomial bases, the number of parameters to be optimized changes to $b(q + 1)$. The polynomial coefficients of the input dependent lengthscale $\hat{l}(x)$ are summarized in the vector $\alpha_l \in \mathbb{R}^{b(q+1) \times 1}$.

The different polynomial bases, normally used for PCE, are shown in Appendix C. The superposition of the polynomial bases is achieved by a linear combination of the polynomials, as demonstrated in Figure 1 with

$$\hat{l}(x) = \sum_{i=0}^{q_1} \alpha_{l,i,1}\, \phi_{i,1}(x) + \sum_{i=0}^{q_2} \alpha_{l,i,2}\, \phi_{i,2}(x) + \cdots + \sum_{i=0}^{q_b} \alpha_{l,i,b}\, \phi_{i,b}(x). \tag{16}$$

### 3.3 PCE for Heteroscedastic Noise Estimation

A separate PCE is also used to model the heteroscedastic noise parameter $\hat{\sigma}_{\mathrm{n}}^2(\boldsymbol{x})$ of the GP as a function of the input data. Since the noise variance has to be one point per data point, all inputs are weighted and summarized, which leads to

$$\hat{\sigma}_{\mathrm{n}}^2(\boldsymbol{x}) = \frac{1}{n_x} \sum_{i=0}^{r} \sum_{j=1}^{n_x} \alpha_{\sigma_{\mathrm{n}},i} \phi_{\sigma_{\mathrm{n}},i}(x_j) \tag{17}$$

for the heteroscedastic noise variance. As the noise of the signal can be dependent on the measurement range for real sensor data, but often does not reflect a complex function, truncation of the polynomial is possible for a significantly smaller value $r$. This also limits the number of parameters to be optimized for modeling the noise variance and results in $\boldsymbol{\alpha}_{\sigma_{\mathrm{n}}} \in \mathbb{R}^{(r+1) \times 1}$.

## 4  Hyperparameter Optimization

In this work, the hyperparameter optimization is done via the TPE and the stochastic gradient descent algorithm Adam [21]. The optimization loop is described in Appendix F. The first step is to investigate the data for the choice of a suitable scaler, covariance functions, and polynomial bases. After that, the data are scaled with the scaler. A termination criterion has to be selected for the optimization loop, i.e., the number of trials for hyperparameter optimization. Then, the hyperparameter optimization starts with the TPE as a Bayesian approach. First, initial hyperparameter combinations are placed with random search in the entire hyperparameter space. Then, the hyperparameters are suggested via the TPE algorithm.

The TPE algorithm uses an objective function for hyperparameter optimization. The inputs of this function are the hyperparameters of the model. This includes the selection of the polynomial degree $q$ for the lengthscale $\hat{l}(\boldsymbol{x})$ and the polynomial degree $r$ for the heteroscedastic noise variance $\hat{\sigma}_{\mathrm{n}}^2(\boldsymbol{x})$ calculation. In addition, the polynomial coefficients $\boldsymbol{\alpha}_l$ and $\boldsymbol{\alpha}_{\sigma_{\mathrm{n}}}$, as well as the output scales $\boldsymbol{\sigma}_{\mathrm{f}}^2$ of the covariance functions are suggested via TPE.

The output of the objective function is the loss of the model. The optimization objective is to minimize the loss metric. For this, a shuffled $k$-fold cross-validation is performed, in which $k$ different models are trained using a gradient descent method and the hyperparameters of the PCEGP are fine-tuned. The gradient is built over $n_{\mathrm{I}}$ iterations for hyperparameter fine-tuning.

The use of the TPE ensures that the gradient descent procedure is performed with different hyperparameter combinations and initial values for the polynomial coefficients. Compared to the grid search or random search, the TPE has the advantage that it learns from poor hyperparameter combinations and avoids these regions for new hyperparameter suggestions.

For the gradient descent based hyperparameter optimization, the covariance matrix depends now on the lengthscale PCE coefficients $\boldsymbol{\alpha}_l$, the noise PCE coefficients $\boldsymbol{\alpha}_{\sigma_{\mathrm{n}}}$, and the output scale $\sigma_{\mathrm{f}}$, which is described by

$$\boldsymbol{K}(\boldsymbol{X}) = \boldsymbol{K}\left(\hat{\boldsymbol{L}}(\boldsymbol{X}), \sigma_{\mathrm{f}}\right) + \hat{\boldsymbol{\sigma}}_{\mathrm{n}}^2(\boldsymbol{X})\,\boldsymbol{I}. \tag{18}$$

Here, $\hat{\boldsymbol{L}}(\boldsymbol{X}) \in \mathbb{R}^{n_x \times N}$ denotes the prediction of the lengthscale parameters for each data point, depending on the PCE coefficients $\boldsymbol{\alpha}_l$ and the input $\boldsymbol{X}$. With $\hat{\boldsymbol{\sigma}}_{\mathrm{n}}^2(\boldsymbol{X}) \in \mathbb{R}^{N \times 1}$ the prediction of the heteroscedastic noise, depending on the PCE coefficients $\boldsymbol{\alpha}_{\sigma_{\mathrm{n}}}$ and the input $\boldsymbol{X}$ is denoted. The task of the optimizer shifts from optimizing the lengthscale and noise directly to optimizing the polynomial coefficients. The gradient from Equation (4) is now built on the polynomial coefficients for hyperparameter optimization.

## 5  Experiments

In this section, the presented approach is applied to three different regression benchmarks with multiple inputs and are used to evaluate the models performance. This includes the Boston housing benchmark, energy efficiency benchmark, and concrete compressive benchmark.

The experiments, presented in the following, are carried out using the scikit-learn library [29] for data scaling, the GPyTorch library [30] for kernel implementations, and the chaospy library [31] for the implementation of the PCE. The TPE is implemented, using the Optuna library [32] and the gradient descent optimizer Adam is implemented using the torch library [33].

All experiments are carried out with the same model structure, which consists of the squared exponential, absolute exponential, Matérn$_{3/2}$, and rational quadratic covariance function for the GP. The fixed structure is illustrated in Appendix G. The lengthscale parameters $\hat{l}(\boldsymbol{x})$ of each covariance function are calculated by the Legendre PCE basis with a uniform distribution between $0$ and $1$. The polynomial degree is suggested by the TPE between $5$ and $10$. The noise variance $\hat{\sigma}_n^2(\boldsymbol{x})$ is fixed to $10^{-4}$ as the aim is to achieve the most accurate model predictions. All experiments are carried out with scaled inputs between $0$ and $1$. The output is normalized, which proved to be the best configuration.

The benchmark regression datasets and results are presented in the following. All experimental results for regression are generated using a ten-fold cross-validation. The partitioning of training and test datasets follows a ten-fold cross-validation approach, consistent with methodologies employed in prior studies with which the approach of this work is compared. Some of the approaches considered for comparison here repeat the cross-validation and some do not. For the approaches with repetition, the mean of the error metric is used to evaluate performance.

All datasets are segmented into ten equally sized, non-overlapping subsamples through random selection. One of these subsamples is assigned as the test dataset, while the remaining nine subsamples are utilized as training data. This process is iterated ten times, with each subsample taking on the role of the test set once throughout the iterations.

## 5.1 Boston Housing Benchmark

The approach of this work is tested on the Boston housing benchmark dataset [34], which summarizes the median house prices in Boston. The dataset consists of 506 data points with 13 inputs and one output.

Table 1 shows the results of this work on the Boston benchmark compared to seven different approaches. The Deep Gaussian Covariance Network (DGCN) algorithm [2] is included in the comparison.

The results on the Boston housing benchmark are compared to the best approach in [36], called tEMLAP, which is also a GP based approach. This approach consists of a Student $t$ noise likelihood, where the MLL is maximized with expectation maximization. The authors in [36] used the Laplace approximation for the Student's $t$ distribution. In addition, a comparison is made to the GPBoost approach [38], where the authors combine tree-boosting with a GP to improve the model prediction. In [38], the authors also used a standard GP for comparison.

The results of three non-GP-based approaches with good performances on the Boston housing benchmark are also considered in this work. These approaches consist of a Random Forest model [39], a XGBoost [37] based model, and a neural network ensemble [35]. In [37] and [35], the authors used the MSE for model evaluation. In this work, the RMSE is calculated with the MSE for comparison.

The neural network ensemble achieves the best results in the Boston housing benchmark. The approach of this work outperforms the GPBoost and standard GP on this benchmark dataset. In addition, the non-GP-based methods random forest and XGBoost are outperformed by the PCEGP. The DGCN and tEMLAP approaches show slightly better prediction performance than the PCEGP, but the PCEGP is competitive.

Table 1: Results on the Boston housing benchmark dataset

| Model | RMSE |
|---|---|
| Neural network ensemble [35] | 0.05 |
| DGCN [2] | 2.40 |
| tEMLAP [36] | 2.49 |
| **PCEGP** | **2.53** |
| Random Forest [37] | 2.59 |
| GPBoost [38] | 2.81 |
| GP [38] | 2.96 |
| XGBoost [39] | 4.62 |

## 5.2 Energy Efficiency Benchmark

In this section, the performance of the ML approach of this work is tested on the energy efficiency benchmark dataset. The dataset, initially introduced in [40], is designed to evaluate the energy efficiency of buildings by examining the heating load (HL) and cooling load (CL) requirements as dependent on various building parameters. The dataset includes eight input parameters and two target variables for prediction.

In Table 2, the results of the energy efficiency benchmark are shown. The approach of this work is compared to the DGCN algorithm, the GPBoost algorithm, and the standard GP. The GPBoost algorithm and the standard GP have been trained and evaluated for the prediction of HL. A comparison is also made to a fuzzy based approach. In [41], the authors used fuzzy inductive reasoning (FIR) and an adaptive neuro fuzzy inference system (ANFIS) for predicting the energy efficiency. Only the FIR model is used for comparison, since it is the best approach in [41].

In [42], the authors developed a gated recurrent unit (GRU) to predict the energy efficiency. According to [42], this approach is the most common sequence learning model of Recurrent Neural Networks (RNN). In [43], the authors investigated the effectiveness of various learning methods in estimating the HL and CL of a building, based on the benchmark data. A comparison is also made to a weighted combination of random forests (RF) and backpropagation neural networks (BPNN), presented in [44]. The best model in [43] for both, HL and CL, is the random forest, which is also used for comparison.

The approaches from [42] and [43] have been developed with a focus on energy predictions. Especially the approach presented in [42] outperforms all other approaches. In both studies, additional data preprocessing methods are used to improve the model performance and to adapt the data to the model approaches in the best possible way. In [42] this is done by applying a polynomial equation to increase the number of features. Feature reduction is carried out in [43] to improve the prediction quality. This is not done in the other methods, but the best models should be mentioned.

The approach of this work shows a good model performance for HL and CL in this benchmark compared to other state of the art approaches. Especially for HL, the algorithm of this work is competitive to the DGCN and GPBoost approaches, which are the two best approaches without extra data preprocessing. In addition, the PCEGP outperforms the FIR approach as well as the standard GP. For CL, the RF-BNN and FIR approaches are outperformed by the PCEGP.

Table 2: Results on the energy efficiency benchmark dataset

| Model | RMSE HL | RMSE CL |
|---|---|---|
| GRU [42] | 0.01 | 0.01 |
| Random forest [43] | 0.24 | 0.62 |
| DGCN [2] | 0.37 | 0.54 |
| GPBoost [38] | 0.393 | - |
| **PCEGP** | **0.46** | **1.01** |
| RF-BPNN [44] | 0.46 | 1.16 |
| FIR [41] | 0.49 | 1.72 |
| GP [38] | 1.32 | - |

### 5.3 Concrete Compressive Benchmark

The last regression benchmark dataset considered in this work is the concrete compressive dataset [45]. This dataset is utilized to model the compressive strength of concrete using waste materials. The dataset includes eight inputs, and the target output is the compressive strength observed across 1030 data points.

Again, the DGCN algorithm, GPBoost algorithm, and the standard GP are considered to compare the performance of the PCEGP. In addition, an approach called least squares support vector machines (LSSVM) [46] is used for comparison. In [47], several ML algorithms are tested on the concrete compressive dataset. The bagging model and AdaBoost

Table 3: Results on the concrete compressive benchmark dataset

| Model | RMSE |
|---|---|
| DGCN [2] | 3.67 |
| GPBoost [38] | 3.96 |
| **PCEGP** | **4.17** |
| Random Forest Ensemble Learner [48] | 4.60 |
| GP [38] | 5.21 |
| LSSVM [46] | 6.69 |
| Bagging [47] | 6.91 |
| AdaBoost [47] | 7.53 |

model performed best and are considered in this work. The bagging and the AdaBoost approaches are based on ensemble techniques. In [48], the authors used an ensemble of random forests, which is also used for comparison of the model performance in this work.

In Table 3 the results again show a good model performance of the PCEGP. The RMSE of the PCEGP is better in five out of seven comparisons. The results show that the standard GP as well as various state of the art approaches are outperformed by the PCEGP. The GPBoost and DGCN have a smaller error here, but the PCEGP is again competitive. This once again underlines the broad applicability of the PCEGP.

## 6   Summary and Outlook

This paper presents a novel ML approach for modeling regression tasks. The Polynomial Chaos Expanded Gaussian Process is introduced to model the global and local behavior of data in one model through non-stationarity. This is achieved by calculating the lengthscale parameters of the GP's covariance functions as a function of input data using PCE. In addition the heteroscedastic noise estimation is introduced, using PCE.

The novel approach proved to be competitive on various benchmark datasets for regression and outperformed previous ML approaches. The hyperparameters obtained from PCE are expressed analytically as polynomials. This characteristic enhances model interpretability and traceability, especially when compared with other commonly used ML regression methods.

The algorithm currently also has limitations, as expert knowledge is required to search for the optimal parameters in order to restrict the search space of these parameters. It is also necessary to manually select a suitable data scaler, polynomial basis, and covariance function for the data.

Further research is needed on the automation of the selection of the data scaler, polynomial bases, covariance functions, and how to narrow the search range of the parameters to improve the convergence behavior. Since the PCEs are also suitable for small amounts of data, so an investigation and comparison of the model performance on small amounts of data is another point for future work. In addition, future work will investigate the performance of the PCEGP for time series forecasting and modeling of dynamic systems.

## References

[1] C. Rasmussen and C. Williams, *Gaussian process for machine learning*. London, England: The MIT Press, 2006.

[2] K. Cremanns and D. Roos, "Deep gaussian covariance network," 2017.

[3] C. J. Paciorek and M. J. Schervish, "Nonstationary covariance functions for gaussian process regression," in *Proceedings of the 16th International Conference on Neural Information Processing Systems*, NIPS'03, (Cambridge, MA, USA), pp. 273–280, MIT Press, 2003.

[4] M. Heinonen, H. Mannerström, J. Rousu, S. Kaski, and H. Lähdesmäki, "Non-stationary gaussian process regression with hamiltonian monte carlo," in *Proceedings of the 19th International Conference on Artificial Intelligence and Statistics* (A. Gretton and C. C. Robert, eds.), vol. 51 of *Proceedings of Machine Learning Research*, (Cadiz, Spain), pp. 732–740, PMLR, 2016.

[5] C. Plagemann, K. Kersting, and W. Burgard, "Nonstationary gaussian process regression using point estimates of local smoothness," in *Machine Learning and Knowledge Discovery in Databases* (W. Daelemans, B. Goethals, and K. Morik, eds.), (Berlin, Heidelberg), pp. 204–219, Springer Berlin Heidelberg, 2008.

[6] A. Hebbal, L. Brevault, M. Balesdent, E.-G. Talbi, and N. Melab, "Bayesian optimization using deep gaussian processes with applications to aerospace system design," *Optimization and Engineering*, vol. 22, no. 1, pp. 321–361, 2021.

[7] M. Trapp, R. Peharz, F. Pernkopf, and C. E. Rasmussen, "Deep structured mixtures of gaussian processes," in *Proceedings of the Twenty Third International Conference on Artificial Intelligence and Statistics* (S. Chiappa and R. Calandra, eds.), vol. 108 of *Proceedings of Machine Learning Research*, pp. 2251–2261, PMLR, 2020.

[8] Y. Liang, S. Li, C. Yan, M. Li, and C. Jiang, "Explaining the black-box model: A survey of local interpretation methods for deep neural networks," *Neurocomputing*, vol. 419, pp. 168–182, 2021.

[9] E. Torre, S. Marelli, P. Embrechts, and B. Sudret, "Data-driven polynomial chaos expansion for machine learning regression," *Journal of Computational Physics*, vol. 388, no. 4, pp. 601–623, 2019.

[10] H.-P. Wan, W.-X. Ren, and M. D. Todd, "Arbitrary polynomial chaos expansion method for uncertainty quantification and global sensitivity analysis in structural dynamics," *Mechanical Systems and Signal Processing*, vol. 142, p. 106732, 2020.

[11] T. A. Mara and W. E. Becker, "Polynomial chaos expansion for sensitivity analysis of model output with dependent inputs," *Reliability Engineering & System Safety*, vol. 214, p. 107795, 2021.

[12] J. Li, L. Pan, M. Suvarna, and X. Wang, "Machine learning aided supercritical water gasification for h2-rich syngas production with process optimization and catalyst screening," *Chemical Engineering Journal*, vol. 426, p. 131285, 2021.

[13] E. Keogh and A. Mueen, "Curse of dimensionality," *Encyclopedia of Machine Learning*, pp. 257–258, 2011.

[14] R. M. Neal, *Bayesian Learning for Neural Networks*, vol. 118 of *Lecture Notes in Statistics*. New York, NY: Springer New York, 1996.

[15] T. N. Varunram, M. B. Shivaprasad, K. H. Aishwarya, A. Balraj, S. V. Savish, and S. Ullas, "Analysis of different dimensionality reduction techniques and machine learning algorithms for an intrusion detection system," in *2021 IEEE 6th International Conference on Computing, Communication and Automation (ICCCA)*, pp. 237–242, 2021.

[16] T. Paananen, J. Piironen, M. R. Andersen, and A. Vehtari, "Variable selection for gaussian processes via sensitivity analysis of the posterior predictive distribution," in *Proceedings of the Twenty-Second International Conference on Artificial Intelligence and Statistics* (K. Chaudhuri and M. Sugiyama, eds.), vol. 89 of *Proceedings of Machine Learning Research*, pp. 1743–1752, PMLR, 2019.

[17] D. Shen, H. Wu, B. Xia, and D. Gan, "Polynomial chaos expansion for parametric problems in engineering systems: A review," *IEEE Systems Journal*, vol. 14, no. 3, pp. 4500–4514, 2020.

[18] S. Rahman, "A polynomial chaos expansion in dependent random variables," *Journal of Mathematical Analysis and Applications*, vol. 464, no. 1, pp. 749–775, 2018.

[19] B. Sudret, S. Marelli, and J. Wiart, "Surrogate models for uncertainty quantification: An overview," in *2017 11th European Conference on Antennas and Propagation (EUCAP)*, pp. 793–797, IEEE, 2017.

[20] J. D. Jakeman, F. Franzelin, A. Narayan, M. Eldred, and D. Plfüger, "Polynomial chaos expansions for dependent random variables," *Computer Methods in Applied Mechanics and Engineering*, vol. 351, pp. 643–666, 2019.

[21] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," *3rd International Conference for Learning Representations*, 2015.

[22] R. Schöbi, B. Sudret, and J. Wiart, "Polynomial-chaos-based kriging," *International Journal for Uncertainty Quantification*, vol. 5, no. 2, pp. 171–193, 2015.

[23] S. Greenhill, S. Rana, S. Gupta, P. Vellanki, and S. Venkatesh, "Bayesian optimization for adaptive experimental design: A review," *IEEE Access*, vol. 8, pp. 13937–13948, 2020.

[24] C. White, W. Neiswanger, and Y. Savani, "Bananas: Bayesian optimization with neural architectures for neural architecture search," *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 35, no. 12, pp. 10293–10301, 2021.

[25] J. Wu, X.-Y. Chen, H. Zhang, L.-D. Xiong, H. Lei, and S.-H. Deng, "Hyperparameter optimization for machine learning models based on bayesian optimization," *Journal of Electronic Science and Technology*, vol. 17, no. 1, pp. 26–40, 2019.

[26] V. Nguyen, "Bayesian optimization for accelerating hyper-parameter tuning," in *2019 IEEE Second International Conference on Artificial Intelligence and Knowledge Engineering (AIKE)*, pp. 302–305, IEEE, 2019.

[27] A. H. Victoria and G. Maragatham, "Automatic tuning of hyperparameters using bayesian optimization," *Evolving Systems*, vol. 12, no. 1, pp. 217–223, 2021.

[28] Y. Ozaki, Y. Tanigaki, S. Watanabe, and M. Onishi, "Multiobjective tree-structured parzen estimator for computationally expensive optimization problems," in *Proceedings of the 2020 Genetic and Evolutionary Computation Conference* (C. A. Coello Coello, ed.), (New York, NY, USA), pp. 533–541, ACM, 2020.

[29] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, A. Müller, J. Nothman, G. Louppe, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and É. Duchesnay, "Scikit-learn: Machine learning in python," *Journal of Machine Learning Research*, pp. 2825–2830, 2011.

[30] J. Gardner, G. Pleiss, K. Q. Weinberger, D. Bindel, and A. G. Wilson, "Gpytorch: Blackbox matrix-matrix gaussian process inference with gpu acceleration," in *Advances in Neural Information Processing Systems* (S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, eds.), vol. 31, 2018.

[31] J. Feinberg and H. P. Langtangen, "Chaospy: An open source tool for designing methods of uncertainty quantification," *Journal of Computational Science*, vol. 11, pp. 46–57, 2015.

[32] T. Akiba, S. Sano, T. Yanase, T. Ohta, and M. Koyama, "Optuna: A next-generation hyperparameter optimization framework," in *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining* (A. Teredesai, V. Kumar, Y. Li, R. Rosales, E. Terzi, and G. Karypis, eds.), (New York, NY, USA), pp. 2623–2631, ACM, 2019.

[33] M. Novik, "torch-optimizer – collection of optimization algorithms for pytorch," 2020.

[34] D. Harrison and D. L. Rubinfeld, "Hedonic housing prices and the demand for clean air," *Journal of Environmental Economics and Management*, vol. 5, no. 1, pp. 81–102, 1978.

[35] S. Soares, C. H. Antunes, and R. Araújo, "Comparison of a genetic algorithm and simulated annealing for automatic neural network ensemble development," *Neurocomputing*, vol. 121, pp. 498–511, 2013.

[36] R. Ranjan, B. Huang, and A. Fatehi, "Robust gaussian process modeling using em algorithm," *Journal of Process Control*, vol. 42, pp. 125–136, 2016.

[37] A. B. Adetunji, O. N. Akande, F. A. Ajala, O. Oyewo, Y. F. Akande, and G. Oluwadara, "House price prediction using random forest machine learning technique," *Procedia Computer Science*, vol. 199, pp. 806–813, 2022.

[38] F. Sigrist, "Gaussian process boosting," *J. Mach. Learn. Res.*, vol. 23, no. 1, 2022.

[39] M. Shahhosseini, G. Hu, and H. Pham, "Optimizing ensemble weights for machine learning models: A case study for housing price prediction," in *Smart Service Systems, Operations Management, and Analytics* (H. Yang, R. Qiu, and W. Chen, eds.), Springer Proceedings in Business and Economics, pp. 87–97, Cham: Springer International Publishing, 2020.

[40] A. Tsanas and A. Xifara, "Accurate quantitative estimation of energy performance of residential buildings using statistical machine learning tools," *Energy and Buildings*, vol. 49, pp. 560–567, 2012.

[41] À. Nebot and F. Mugica, "Energy performance forecasting of residential buildings using fuzzy approaches," *Applied Sciences*, vol. 10, no. 2, p. 720, 2020.

[42] M. Sajjad, S. U. Khan, N. Khan, I. U. Haq, A. Ullah, M. Y. Lee, and S. W. Baik, "Towards efficient building designing: Heating and cooling load prediction via multi-output model," *Sensors (Basel, Switzerland)*, vol. 20, no. 22, 2020.

[43] W. Gao, J. Alsarraf, H. Moayedi, A. Shahsavar, and H. Nguyen, "Comprehensive preference learning and feature validity for designing energy-efficient residential buildings using machine learning paradigms," *Applied Soft Computing*, vol. 84, p. 105748, 2019.

[44] X. Zhang, J. Zhang, J. Zhang, and Y. Zhang, "Research on the combined prediction model of residential building energy consumption based on random forest and bp neural network," *Geofluids*, vol. 2021, pp. 1–12, 2021.

[45] I.-C. Yeh, "Modeling of strength of high-performance concrete using artificial neural networks," *Cement and Concrete Research*, vol. 28, no. 12, pp. 1797–1808, 1998.

[46] C. Liu, L. Tang, and J. Liu, "Least squares support vector machine with self-organizing multiple kernel learning and sparsity," *Neurocomputing*, vol. 331, pp. 493–504, 2019.

[47] W. Ahmad, A. Ahmad, K. A. Ostrowski, F. Aslam, P. Joyklad, and P. Zajdel, "Application of advanced machine learning approaches to predict the compressive strength of concrete containing supplementary cementitious materials," *Materials (Basel, Switzerland)*, vol. 14, no. 19, 2021.

[48] F. Farooq, W. Ahmed, A. Akbar, F. Aslam, and R. Alyousef, "Predictive modeling for sustainable high-performance concrete from industrial wastes: A comparison and optimization of models using ensemble learners," *Journal of Cleaner Production*, vol. 292, 2021.

# A  Notations

Table 4: Notations

| Formula symbol | Description |
| --- | --- |
| $b$ | Number of polynomial bases |
| $\mathcal{D} = \{\boldsymbol{x}_i, y_i\}_{i=1}^{N}$ | Training data set with $N$ data points |
| EI | Expected improvement acquisition function |
| $f(x)$ | Function value of the polynomial chaos expansion |
| $f(x^*)$ | Observed minima in $\mathbb{R}$ |
| $\boldsymbol{I}$ | Identity matrix |
| $k(\boldsymbol{x}, \boldsymbol{x}')$ | Covariance function of a Gaussian process between $\boldsymbol{x}$ and $\boldsymbol{x}'$ |
| $k_{\Sigma}(\boldsymbol{x}, \boldsymbol{x}')$ | Summation of $n_k$ covariance functions |
| $\boldsymbol{K}(\boldsymbol{X})$ | Covariance matrix depending on the training data in $\mathbb{R}^{N \times N}$ |
| $\boldsymbol{k}(\boldsymbol{X}, \boldsymbol{x}^*)$ | Covariance vector depending on the training data and the next point to predict in $\mathbb{R}^{N \times 1}$ |
| $l$ | Lengthscale parameter of a Gaussian process in $\mathbb{R}^+$ |
| $l(x)$ | Data point dependent lengthscale parameter of a Gaussian process in $\mathbb{R}$ |
| $\hat{l}(x)$ | Prediction of the lengthscale parameter by the polynomial chaos expansion in $\mathbb{R}$ |
| $l_i$ | Lengthscale parameter of a Gaussian process at input $i$ in $\mathbb{R}^+$ |
| $\hat{\boldsymbol{l}}(\boldsymbol{x})$ | Prediction of the lengthscale parameter vector by the polynomial chaos expansion in $\mathbb{R}^{n_x \times 1}$ |
| $\boldsymbol{L}$ | Lengthscale parameter diagonal matrix of a Gaussian process which includes all $l_i$ lengthscale parameters in $\mathbb{R}^{n_x \times n_x}$ |
| $\hat{\boldsymbol{L}}(\boldsymbol{X})$ | Prediction of the lengthscale parameter matrix by the polynomial chaos expansion in $\mathbb{R}^{n_x \times N}$ |
| $m(\boldsymbol{x})$ | Mean function of a Gaussian process for input vector $\boldsymbol{x}$ |
| $n_{\mathrm{F}}$ | Number of folds |
| $n_{\mathrm{I}}$ | Number of iterations |
| $n_{\mathrm{IT}}$ | Number of initial trials |
| $n_k$ | Number of covariance functions |
| $n_{\mathrm{T}}$ | Number of trials |
| $n_x$ | Number of input parameters |
| $N$ | Number of training points |
| $p(x)$ | Probability density function |
| $q$ | Polynomial degree of the lengthscale parameter estimation |
| $r$ | Polynomial degree of the noise variance estimation |
| $x$ | Scalar training input in $\mathbb{R}$ |
| $x^*$ | Scalar test input in $\mathbb{R}$ |
| $x_{\mathrm{s}}$ | Scaled scalar training input in $\mathbb{R}$ |
| $x_{\mathrm{s}}^*$ | Scaled scalar test input in $\mathbb{R}$ |
| $\boldsymbol{x}$ | Training input vector in $\mathbb{R}^{n_x}$ |
| $\boldsymbol{x}^*$ | Test input vector in $\mathbb{R}^{n_x}$ |
| $\boldsymbol{x}_{\mathrm{s}}$ | Scaled training input vector in $\mathbb{R}^{n_x}$ |
| $\boldsymbol{x}_{\mathrm{s}}$ | Scaled test input vector in $\mathbb{R}^{n_x}$ |
| $\boldsymbol{x}_i$ | $i$-th Training input vector in $\mathbb{R}^{n_x}$ |
| $\boldsymbol{X}$ | Training input matrix in $\mathbb{R}^{N \times n_x}$ |
| $y$ | Training output value in $\mathbb{R}$ |
| $y_{\mathrm{s}}$ | Scaled training output value in $\mathbb{R}$ |
| $\boldsymbol{y}$ | Training output vector in $\mathbb{R}^{N}$ |
| $\boldsymbol{y}_{\mathrm{s}}$ | Scaled training output vector in $\mathbb{R}^{N}$ |
| $y_i$ | $i$-th training output value in $\mathbb{R}$ |
| $\hat{\boldsymbol{y}}_{\mathrm{PCEGP}}$ | PCEGP model mean and uncertainty prediction in $\mathbb{R}^{2 \times 1}$ |
| $\hat{y}_{\mathrm{PCEGPs}}$ | Scaled PCEGP model prediction in $\mathbb{R}$ |
| $\hat{y}_{\mathrm{PCEGPs,m}}$ | Scaled PCEGP model mean prediction in $\mathbb{R}$ |
| $\hat{y}_{\mathrm{PCEGPs,v}}$ | Scaled PCEGP model variance prediction in $\mathbb{R}$ |
| $\alpha$ | Polynomial coefficient for the polynomial chaos expansion in $\mathbb{R}$ |
| $\alpha_i$ | Polynomial coefficient for the $i$-th polynomial in $\mathbb{R}$ |
| $\alpha_l$ | Polynomial coefficient for the lengthscale parameter in $\mathbb{R}$ |
| $\alpha_{l,i}$ | Polynomial coefficient for the lengthscale parameter and $i$-th polynomial in $\mathbb{R}$ |

Notations

| Formula symbol | Description |
|---|---|
| $\boldsymbol{\alpha}_l$ | Polynomial coefficient vector for the lengthscale parameter in $\mathbb{R}^{b(q+1)\times 1}$ |
| $\alpha_{\sigma_\mathrm{n}}$ | Polynomial coefficient for the noise variance $\sigma_\mathrm{n}$ in $\mathbb{R}$ |
| $\alpha_{\sigma_\mathrm{n},i}$ | Polynomial coefficient for the noise variance $\sigma_\mathrm{n}$ and $i$-th polynomial in $\mathbb{R}$ |
| $\boldsymbol{\alpha}_{\sigma_\mathrm{n}}$ | Polynomial coefficient vector for the noise variance $\sigma_\mathrm{n}$ in $\mathbb{R}^{(r+1)\times 1}$ |
| $\mu$ | Mean prediction of a Gaussian process |
| $\varphi$ | Cumulative distribution function |
| $\phi$ | Polynomial basis of the polynomial chaos expansion |
| $\phi_i$ | $i$-th polynomial of a polynomial basis |
| $\sigma$ | Standard deviation in $\mathbb{R}$ |
| $\sigma^2$ | Variance prediction of a Gaussian process |
| $\sigma_\mathrm{f}^2$ | Output scale variance of a Gaussian process in $\mathbb{R}^+$ |
| $\boldsymbol{\sigma}_\mathrm{f}^2$ | Output scale variance vector of multiple covariance functions of a Gaussian process in $\mathbb{R}^{n_k\times 1}$ |
| $\sigma_\mathrm{n}^2$ | Noise variance of a Gaussian process in $\mathbb{R}^+$ |
| $\hat{\sigma}_\mathrm{n}^2(x)$ | Prediction of the noise variance of the scalar input $x$ by the polynomial chaos expansion in $\mathbb{R}^+$ |
| $\hat{\sigma}_\mathrm{n}^2(\boldsymbol{x})$ | Prediction of the noise variance of the input vector $\boldsymbol{x}$ by the polynomial chaos expansion in $\mathbb{R}^+$ |
| $\hat{\boldsymbol{\sigma}}_\mathrm{n}^2(\boldsymbol{X})$ | Prediction of the noise variance vector by the polynomial chaos expansion in $\mathbb{R}^{N\times N}$ |
| $\boldsymbol{\theta}$ | Hyperparameter vector of a Gaussian process in $\mathbb{R}^{n_x+2}$ |
| $\boldsymbol{\theta}_i$ | $i$-th Hyperparameter of a Gaussian process in $\mathbb{R}$ |
| $\boldsymbol{\theta}_{\mathrm{PCEGP}}$ | Hyperparameter vector of the Polynomial Chaos Expanded Gaussian Process in $\mathbb{R}^{q+r+n_k+2\times 1}$ |
| $\xi$ | Weight factor for fine-tuning between exploitation and exploration |

# B Stationary and Non-stationary Covariance Functions

The following Table 5 shows four commonly used stationary covariance functions [1]. These four covariance functions are used in this work, which can represent a wide range of functional curves and thus ensure the variability of the model. These chosen covariance functions are the squared exponential, absolute exponential, Matérn$_{3/2}$, and rational quadratic. This covers function curves from very smooth to very wiggly. In Table 6, the non-stationary covariance versions of Table 5 are shown, as presented in [2]. Here, the lengthscale parameter $l\,(\boldsymbol{x})$ depends directly on the input $\boldsymbol{x}$. The product $l\,(\boldsymbol{x}) \odot \boldsymbol{x}$ denotes an elementwise multiplication of the elements in $l\,(\boldsymbol{x})$ and $\boldsymbol{x}$. In Figure 2, the function curve of the covariance functions are shown. These function curves are random drawn from the prior distribution of the GP.

Table 5: Stationary covariance functions

| Name | Formula |
|------|---------|
| Squared exponential | $k(\boldsymbol{x}, \boldsymbol{x}') = \sigma_{\mathrm{f}}^2 \exp\left(-\frac{\|\boldsymbol{x}-\boldsymbol{x}'\|^2}{2\,l^2}\right)$ |
| Absolute exponential | $k(\boldsymbol{x}, \boldsymbol{x}') = \sigma_{\mathrm{f}}^2 \exp\left(-\frac{\|\boldsymbol{x}-\boldsymbol{x}'\|}{l}\right)$ |
| Matérn$_{3/2}$ | $k(\boldsymbol{x}, \boldsymbol{x}') = \sigma_{\mathrm{f}}^2 \left(1 + \frac{\sqrt{3}\|\boldsymbol{x}-\boldsymbol{x}'\|}{l}\right) \exp\left(-\frac{\sqrt{3}\|\boldsymbol{x}-\boldsymbol{x}'\|}{l}\right)$ |
| Rational quadratic | $k(\boldsymbol{x}, \boldsymbol{x}') = \sigma_{\mathrm{f}}^2 \left(1 + \frac{\|\boldsymbol{x}-\boldsymbol{x}'\|^2}{2\,\alpha\,l^2}\right)^{-\alpha}$ |

Table 6: Non-stationary covariance functions

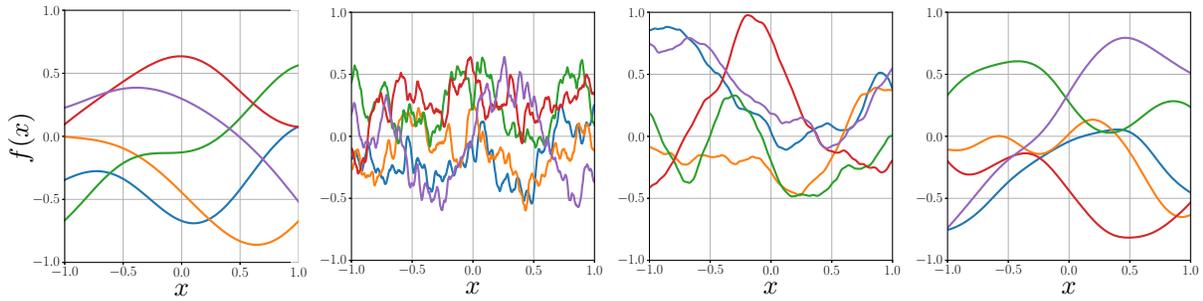| Name | Formula |
|------|---------|
| Squared exponential | $k(\boldsymbol{x}, \boldsymbol{x}') = \sigma_{\mathrm{f}}^2 \exp\left(-\frac{\|l(\boldsymbol{x})\odot\boldsymbol{x}-l(\boldsymbol{x}')\odot\boldsymbol{x}'\|^2}{2}\right)$ |
| Absolute exponential | $k(\boldsymbol{x}, \boldsymbol{x}') = \sigma_{\mathrm{f}}^2 \exp\left(-\|l\,(\boldsymbol{x}) \odot \boldsymbol{x} - l\,(\boldsymbol{x}') \odot \boldsymbol{x}'\|\right)$ |
| Matérn$_{3/2}$ | $k(\boldsymbol{x}, \boldsymbol{x}') = \sigma_{\mathrm{f}}^2 \left(1 + \sqrt{3}\,\|l\,(\boldsymbol{x}) \odot \boldsymbol{x} - l\,(\boldsymbol{x}') \odot \boldsymbol{x}'\|\right) \exp\left(-\sqrt{3}\,\|l\,(\boldsymbol{x}) \odot \boldsymbol{x} - l\,(\boldsymbol{x}') \odot \boldsymbol{x}'\|\right)$ |
| Rational quadratic | $k(\boldsymbol{x}, \boldsymbol{x}') = \sigma_{\mathrm{f}}^2 \left(1 + \frac{\|l(\boldsymbol{x})\odot\boldsymbol{x}-l(\boldsymbol{x}')\odot\boldsymbol{x}'\|^2}{2\alpha}\right)^{-\alpha}$ |



Figure 2: Covariance functions for GP from left to right: squared exponential, absolute exponential, Matérn$_{3/2}$, and rational quadratic

## C  Polynomial Chaos Expansion Bases

In Table 7, the formulas of four different types of commonly used PCE bases are shown. In addition, the distribution and probability density function are shown for each PCE basis. In Figure 3, the function curves of ten polynomials for each polynomial basis are shown. In uncertainty quantification, these polynomial bases are commonly used to model uncertain parameters, depending on the distribution of data. All polynomial bases of the PCE are orthogonal to their probability density function and build computational efficient surrogates [17]. In this work, the Legendre polynomials are used to calculate the hyperparameters of the GP. The concept of this work allows the combination of different polynomial bases to adapt the polynomials to better fit the hyperparameters to the data. With each additional polynomial bases the number of adjustable polynomial coefficients increases, which makes it more difficult for the optimizer to find the optimum. The number of polynomial coefficients for the lengthscale parameter can be calculated with $b\,(q+1)$, where $b$ denotes the number of different polynomial bases and $q$ denotes the polynomial degree.

Table 7: Polynomial bases for polynomial chaos expansion

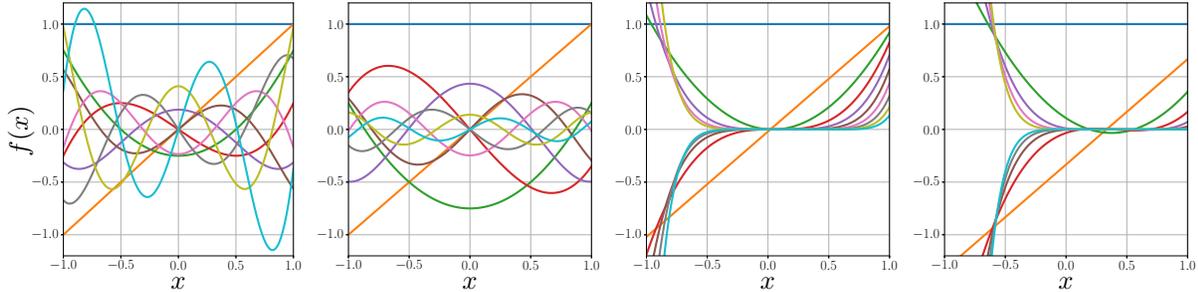| Name | Distribution | Density function | Formula |
|------|-------------|------------------|---------|
| Hermite | Normal | $\frac{1}{\sqrt{2\pi}}e^{-\frac{x^2}{2}}$ | $H_n(x) = (-1)^n e^{\frac{x^2}{2}} \frac{\mathrm{d}^n}{\mathrm{d}x^n}(e^{-\frac{x^2}{2}})$ |
| Legendre | Uniform | $\frac{1}{2}$ | $Le_n(x) = \frac{1}{2^n}\binom{2n}{n}(1-x^2)^{n/2}$ |
| Jacobi | Beta | $\frac{(1-x)^\alpha(1+x)^\beta}{2^{\alpha+\beta+1}B(\alpha+1,\beta+1)}$ | $J_n^{(\alpha,\beta)}(x) = \frac{1}{2^n n!}\sum_{k=0}^{n}(-1)^k\binom{n+\alpha}{k}\binom{n+\beta}{n-k}(1-x)^{n-k}(1+x)^k$ |
| Laguerre | Exponential | $e^{-x}$ | $La_n(x) = \frac{e^x}{n!}\frac{\mathrm{d}^n}{\mathrm{d}x^n}(x^n e^{-x})$ |



Figure 3: Polynomial bases for PCE from left to right: Hermite polynomials, Legendre polynomials, Laguerre polynomials, and Jacobi polynomials

## D  Non-stationary Lengthscale and Heterscedastic Noise Calculation via PCE

In the following, the calculation of the non-stationary lengthscale parameter and the heteroscedastic noise for single and for multiple inputs is presented. With

$$\hat{l}(x) = \sum_{i=0}^{q} \alpha_{l,i}\, \phi_i(x), \tag{19}$$

the calculation of the input dependent lengthscale $\hat{l}(x)$ parameter for a scalar $x$ via PCE is introduced in this work. Here, $\phi_i$ denotes the $i$-th polynomial of a polynomial basis which has to be chosen depending on the input data. These polynomials are weighted with the $i$-th polynomial coefficient $\alpha_{l,i}$. The polynomials are truncated at the polynomial degree $q$. In this work, the polynomial coefficients $\alpha_{l,i}$ are optimized to fit the non-stationary lengthscale $\hat{l}(x)$ to the data.

For multiple inputs the calculated lengthscale parameter $\hat{l}(x)$ changes to a lengthscale parameter vector $\hat{\boldsymbol{l}}(\boldsymbol{x})$. With an input data point of $n_x$ inputs, the calculation of the lengthscale parameter vector $\hat{\boldsymbol{l}}(\boldsymbol{x})$ is done with

$$\hat{\boldsymbol{l}}(\boldsymbol{x}) = \sum_{i=0}^{q} \alpha_{l,i} \begin{bmatrix} \phi_{l,i}(x_1) \\ \phi_{l,i}(x_2) \\ \vdots \\ \phi_{l,i}(x_{n_x}) \end{bmatrix} = \alpha_0 \begin{bmatrix} \phi_{l,0}(x_1) \\ \phi_{l,0}(x_2) \\ \vdots \\ \phi_{l,0}(x_{n_x}) \end{bmatrix} + \alpha_1 \begin{bmatrix} \phi_{l,1}(x_1) \\ \phi_{l,1}(x_2) \\ \vdots \\ \phi_{l,1}(x_{n_x}) \end{bmatrix} + \ldots \alpha_q \begin{bmatrix} \phi_{l,q}(x_1) \\ \phi_{l,q}(x_2) \\ \vdots \\ \phi_{l,q}(x_{n_x}) \end{bmatrix}. \tag{20}$$

The heteroscedastic noise estimation via PCE for a scalar input is defined by

$$\hat{\sigma}_{\mathrm{n}}^2(x) = \sum_{i=0}^{r} \alpha_{\sigma_{\mathrm{n}},i}\, \phi_i(x) \tag{21}$$

in this work. Here, $\hat{\sigma}_{\mathrm{n}}^2(\boldsymbol{x})$ denotes the input dependent noise variance of the model, which is calculated by the PCE with the polynomial degree $r$. In this context, $\phi_i$ again represents the $i$-th polynomial of a polynomial basis, the selection of which depends on the input data. The polynomials are multiplied by the corresponding $i$-th polynomial coefficient, denoted as $\alpha_{\sigma_{\mathrm{n}},i}$. Since the noise variance has to be one point per data point, all inputs are weighted and summarized, which leads to

$$\hat{\sigma}_{\mathrm{n}}^2(x) = \frac{1}{n_x} \sum_{i=0}^{r} \sum_{j=1}^{n_x} \alpha_{\sigma_{\mathrm{n}},i} \phi_{\sigma_{\mathrm{n}},i}(x_j) \tag{22}$$

to calculate the heteroscedastic noise.

# E  Algorithm for Model Prediction

The following Algorithm 1 shows, how the prediction of the PCEGP model works. The input is the test input data $\boldsymbol{x}^*$ with the point to predict. In the model, the covariance function $k_\Sigma(\boldsymbol{x}, \boldsymbol{x}')$, the polynomials for lengthscale $\hat{\boldsymbol{l}}(\boldsymbol{x})$, the polynomials for noise variance $\hat{\sigma}_n^2(\boldsymbol{x})$ are specified. In addition, the covariance matrix and noise depending on the training data $\boldsymbol{K}\left(\boldsymbol{X}_s, \hat{\boldsymbol{L}}(\boldsymbol{X}_s), \boldsymbol{\sigma}_f^2\right) + \hat{\boldsymbol{\sigma}}_n^2(\boldsymbol{X}_s)\,\boldsymbol{I}$, and the hyperparameters $\boldsymbol{\theta}_{\mathrm{PCEGP}} = \begin{bmatrix} q & r & \boldsymbol{\alpha}_l^\mathsf{T} & \boldsymbol{\alpha}_{\sigma_n}^\mathsf{T} & \boldsymbol{\sigma}_f^{2\mathsf{T}} \end{bmatrix}^\mathsf{T}$ are part of the trained model and needed to predict the output $\hat{\boldsymbol{y}}_{\mathrm{PCEGP}}$, consisting of the mean and uncertainty. First, the input data point $\boldsymbol{x}^*$ is scaled to $\boldsymbol{x}_s^*$. For each of the covariance functions, the lengthscale parameters are calculated using the PCE as a function of the input data. Then, the covariance $k_\Sigma(\boldsymbol{x}_i, \boldsymbol{x}^*)$ between the previous observed data points $\mathcal{D} = \{\boldsymbol{x}_{s,i}\}_{i=0}^N$ and the point to predict $\boldsymbol{x}^*$ is calculated. This results in the covariance vector $\boldsymbol{k}_\Sigma(\boldsymbol{X}_s, \boldsymbol{x}_s)$ which is then used to predict the scaled mean $\hat{y}_{\mathrm{PCEGPs,m}}$ and the scaled variance $\hat{y}_{\mathrm{PCEGPs,v}}$. The scaled prediction is then scaled back to the original value range $\hat{y}_{\mathrm{PCEGPm}}$ and $\hat{y}_{\mathrm{PCEGPv}}$.

---

**Algorithm 1** PCEGP model prediction

---

**Input:** Input data $\boldsymbol{x}^*$
**Output:** Model prediction $\hat{y}_{\mathrm{PCEGP}}$
**Specified in the model:** Covariance function $k_\Sigma(\boldsymbol{x}, \boldsymbol{x}')$, polynomials for lengthscale $\hat{\boldsymbol{l}}(\boldsymbol{x})$, polynomials for noise variance $\hat{\sigma}(\boldsymbol{x})$, covariance matrix and noise depending on the training data $\boldsymbol{K}\left(\boldsymbol{X}_s, \hat{\boldsymbol{L}}(\boldsymbol{X}_s), \boldsymbol{\sigma}_f^2\right) + \hat{\boldsymbol{\sigma}}_n^2(\boldsymbol{X}_s)\,\boldsymbol{I}$, hyperparameters $\boldsymbol{\theta}_{\mathrm{PCEGP}} = \begin{bmatrix} q & r & \boldsymbol{\alpha}_l^\mathsf{T} & \boldsymbol{\alpha}_{\sigma_n}^\mathsf{T} & \boldsymbol{\sigma}_f^{2\mathsf{T}} \end{bmatrix}^\mathsf{T}$

1: $\boldsymbol{x}_s^* \leftarrow$ Scale input data $\boldsymbol{x}^*$
2: **for** $k = 1 \ldots n_k$ **do**
3:    Calculate the point dependent lengthscale
   $\hat{l}_k(\boldsymbol{x}_s^*) \leftarrow \sum_{i=0}^{q_1} \alpha_{l,i,1}\,\phi_{i,1}(\boldsymbol{x}_s^*) + \sum_{i=0}^{q_2} \alpha_{l,i,2}\,\phi_{i,2}(\boldsymbol{x}_s^*) + \cdots + \sum_{i=0}^{q_b} \alpha_{l,i,b}\,\phi_{i,b}(\boldsymbol{x}_s^*)$
4: **end for**
5: **for** $i = 1 \ldots N$ **do**
6:    $k_\Sigma(\boldsymbol{x}_{s,i}, \boldsymbol{x}_s^*) \leftarrow$ Calculate point dependent covariance $\|\hat{\boldsymbol{l}}(\boldsymbol{x}_{s,i}) \odot \boldsymbol{x}_{s,i} - \hat{\boldsymbol{l}}(\boldsymbol{x}_s^*) \odot \boldsymbol{x}_s^*\|$
7: **end for**
8: $\boldsymbol{k}_\Sigma(\boldsymbol{X}_s, \boldsymbol{x}_s) \leftarrow$ Summarizing the computed covariances $k_\Sigma(\boldsymbol{x}_{s,i}, \boldsymbol{x}^*) \ldots k_\Sigma(\boldsymbol{x}_{s,N}, \boldsymbol{x}^*)$ to covariance vector
9: Compute the point dependent and noise variance $\hat{\sigma}_n^2(\boldsymbol{x}_s^*) \leftarrow \frac{1}{n_x} \sum_{i=0}^r \sum_{j=1}^{n_x} \alpha_{\sigma_n,i} \phi_{\sigma_n,i}(x_j^*)$
10: Compute the scaled output mean

   $\hat{y}_{\mathrm{PCEGPs,m}} \leftarrow \boldsymbol{k}_\Sigma(\boldsymbol{X}_s, \boldsymbol{x}_s^*)^\mathsf{T} \left(\boldsymbol{K}\left(\boldsymbol{X}_s, \hat{\boldsymbol{L}}(\boldsymbol{X}_s), \boldsymbol{\sigma}_f^2\right) + \hat{\boldsymbol{\sigma}}_n^2(\boldsymbol{X}_s)\,\boldsymbol{I}\right)^{-1} \boldsymbol{y}$

11: Compute the scaled output variance

   $\hat{y}_{\mathrm{PCEGPs,v}} \leftarrow k_\Sigma(\boldsymbol{x}_s^*, \boldsymbol{x}_s^*) + \hat{\sigma}_n^2(\boldsymbol{x}^*) - \boldsymbol{k}_\Sigma(\boldsymbol{X}_s, \boldsymbol{x}_s^*)^\mathsf{T} \left(\boldsymbol{K}\left(\boldsymbol{X}_s, \hat{\boldsymbol{L}}(\boldsymbol{X}_s), \boldsymbol{\sigma}_f\right) + \hat{\boldsymbol{\sigma}}_n^2(\boldsymbol{X}_s)\,\boldsymbol{I}\right)^{-1} \boldsymbol{k}_\Sigma(\boldsymbol{X}_s, \boldsymbol{x}_s^*)$

12: $\hat{\boldsymbol{y}}_{\mathrm{PCEGP}} \leftarrow$ Rescale the scaled output $\boldsymbol{y}_{\mathrm{PCEGP,s}}$ of the PCEGP model

---

## F    Algorithm for Hyperparameter Optimization

The following Algorithm 2 shows the hyperparameter optimization process of the PCEGP model. The input for the hyperparameter optimizaion are the data $\mathcal{D} = \{\boldsymbol{X}, \boldsymbol{y}\}$, the number of trials $n_\text{T}$, number of initial trials $n_\text{IT}$, number of iterations $n_\text{I}$, and the number of folds $n_\text{F}$. The number of initial trials $n_\text{IT}$ denotes the number of hyperparameter optimization runs with random search. The remaining $n_\text{T} - n_\text{IT}$ trials are carried out by the TPE algorithm. The number of iterations $n_\text{I}$ denotes the runs for parameter fine-tuning with the gradient descent optimizer Adam. With $n_\text{F}$, the number of folds for $k$-fold cross-validation is defined. The output of the hyperparameter optimization is the best set of hyperparameters $\boldsymbol{\theta}_\text{PCEGP} = \begin{bmatrix} q & r & \boldsymbol{\alpha}_l^\mathsf{T} & \boldsymbol{\alpha}_{\sigma_\text{n}}^\mathsf{T} & \boldsymbol{\sigma}_\text{f}^{2\mathsf{T}} \end{bmatrix}^\mathsf{T}$ for the PCEGP. In the model, the covariance function $k_\Sigma(\boldsymbol{x}, \boldsymbol{x}')$, the polynomials for lengthscale $\hat{\boldsymbol{l}}(\boldsymbol{x})$, and the polynomials for noise variance $\hat{\sigma}_\text{n}^2(\boldsymbol{x})$ are specified.

The algorithm works as follows. First the polynomial degree $q$ for the input-dependent lengthscale and polynomial degree $r$ for the heteroscedastic noise degree are suggested with random search. Then, the output scales for each covariance function are suggested. Depending on the polynomial degrees $q + 1$ and $r + 1$ the coefficients of the polynomials are suggested. In the case that $n_\text{IT} > trial$ is true, this is done with random search, otherwise with the TPE algorithm. Next, the setup of the model is done. The polynomial coefficients and output scales are fine tuned with the gradient descent algorithm Adam. If the loss improves, the best hyperparameter combination is saved. This is carried out until the number of trials $n_\text{T}$ is reached.

---

**Algorithm 2** Hyperparameter optimization

---

**Input:** Data $\mathcal{D} = \{\boldsymbol{X}, \boldsymbol{y}\}$, number of trials $n_{\mathrm{T}}$, number of initial trials $n_{\mathrm{IT}}$, number of iterations $n_{\mathrm{I}}$, number of folds $n_{\mathrm{F}}$

**Output:** Hyperparameters $\boldsymbol{\theta}_{\mathrm{PCEGP}} = \begin{bmatrix} q & r & \boldsymbol{\alpha}_l^{\mathsf{T}} & \boldsymbol{\alpha}_{\sigma_{\mathrm{n}}}^{\mathsf{T}} & \boldsymbol{\sigma}_{\mathrm{f}}^{2\mathsf{T}} \end{bmatrix}^{\mathsf{T}}$

**Specified in the model:** Covariance function $k_{\Sigma}(\boldsymbol{x}, \boldsymbol{x}')$, polynomials for lengthscale $\hat{\boldsymbol{l}}(\boldsymbol{x})$, polynomials for noise variance $\hat{\sigma}_{\mathrm{n}}^2(\boldsymbol{x})$

1:  $\boldsymbol{X}_{\mathrm{s}}, \boldsymbol{y}_{\mathrm{s}} \leftarrow$ Scale input data $\boldsymbol{X}$ and output data $\boldsymbol{y}$
2: **repeat**
3:    **if** $n_{\mathrm{IT}} > trial$ **then**
4:      $q, r \leftarrow$ Suggest lengthscale polynomial degree and noise variance polynomial degree with random search
5:      **for** $i = 1 \dots n_k$ **do**
6:        $\sigma_{\mathrm{f},i}^2 \leftarrow$ Suggest output scale with random search
7:        **for** $j = 0 \dots q$ **do**
8:          $\alpha_{l,i,j} \leftarrow$ Suggest lengthscale polynomial coefficient with random search
9:        **end for**
10:     **end for**
11:     **for** $i = 0 \dots r$ **do**
12:       $\alpha_{\sigma_{\mathrm{n}},i} \leftarrow$ Suggest noise variance polynomial coefficient with random search
13:     **end for**
14:    **else**
15:      $q, r \leftarrow$ Suggest lengthscale polynomial degree and noise variance polynomial degree with TPE
16:      **for** $i = 1 \dots n_k$ **do**
17:        $\sigma_{\mathrm{f},i}^2 \leftarrow$ Suggest output scale with TPE
18:        **for** $j = 0 \dots q$ **do**
19:          $\alpha_{l,i,j} \leftarrow$ Suggest lengthscale polynomial coefficient with TPE
20:        **end for**
21:     **end for**
22:     **for** $i = 0 \dots r$ **do**
23:       $\alpha_{\sigma_{\mathrm{n}},i} \leftarrow$ Suggest noise variance polynomial coefficient with TPE
24:     **end for**
25:    **end if**
26:    **for** $folds = 1 \dots n_{\mathrm{F}}$ **do**
27:      Setup PCEGP Model$\left(\boldsymbol{X}, q, r, \boldsymbol{\alpha}_l, \boldsymbol{\alpha}_{\sigma_{\mathrm{n}}}, \boldsymbol{\sigma}_{\mathrm{f}}^2\right)$ with covariance function $k_{\Sigma}(\boldsymbol{x}, \boldsymbol{x}')$, polynomials for lengthscale $\hat{\boldsymbol{l}}(\boldsymbol{x})$, and polynomials for noise variance $\hat{\sigma}_{\mathrm{n}}^2(\boldsymbol{x})$
28:      **for** $iterations = 1 \dots n_{\mathrm{I}}$ **do**
29:        $\boldsymbol{\alpha}_l, \boldsymbol{\alpha}_{\sigma_{\mathrm{n}}}, \boldsymbol{\sigma}_{\mathrm{f}}^2 \leftarrow$ Adjust the hyperparameters with gradient descent optimizer Adam
30:        $\hat{\boldsymbol{L}}(\boldsymbol{X}_{\mathrm{s}}), \hat{\boldsymbol{\sigma}}_{\mathrm{n}}^2(\boldsymbol{X}_{\mathrm{s}}) \leftarrow$ Compute the point dependent lengthscale and noise variance
31:        $\hat{\boldsymbol{y}}_{\mathrm{PCEGP,s}} \leftarrow$ Compute the scaled output of the PCEGP model for the given input $\boldsymbol{X}_{\mathrm{s}}$
32:        $\hat{\boldsymbol{y}}_{\mathrm{PCEGP}} \leftarrow$ Rescale the scaled output $\boldsymbol{y}_{\mathrm{PCEGP,s}}$ of the PCEGP model
33:        $loss \leftarrow$ Compute $loss$ with negative log likelihood of model prediction $\hat{\boldsymbol{y}}_{\mathrm{PCEGP}}$ and output data $\boldsymbol{y}$
34:      **end for**
35:      $mean\ loss \leftarrow$ Compute $mean\ loss$ of all folds $n_{\mathrm{F}}$
36:    **end for**
37:    **if** actual $mean\ loss < best\ mean\ loss$ **then**
38:      Safe the hyperparameter combination $\boldsymbol{\theta}_{\mathrm{PCEGP}} = \begin{bmatrix} q & r & \boldsymbol{\alpha}_l^{\mathsf{T}} & \boldsymbol{\alpha}_{\sigma_{\mathrm{n}}}^{\mathsf{T}} & \boldsymbol{\sigma}_{\mathrm{f}}^{2\mathsf{T}} \end{bmatrix}^{\mathsf{T}}$
39:    **end if**
40:    $trial \leftarrow trial + 1$
41: **until** $trial = n_{\mathrm{T}}$

# G  PCEGP Model Setup for Experiments

In this section, the PCEGP model setup for experiments is described. For all experiments, the same model setup is used, which is shown in Figure 4. Here, the inputs are scaled between $0$ and $1$ with the min-max scaler. The output is normalized, which proved to be the best configuration. To compute the input-dependent lengthscale parameters $\hat{l}(x^*)$, the Legendre PCE basis is used for each covariance function, with a polynomial degree $q$ to be optimized and a uniform distribution between $0$ and $1$. The polynomial degree is suggested by the TPE betweeen $5$ and $10$. With the squared exponential, absolute exponential, Matérn$_{3/2}$, and rational quadratic, four different covariance functions are used. These covariance functions vary from very sharp variations to smooth variations and are be combined to map a wide range of function curves. For the modeling tasks in this work, the noise variance is fixed to $10^{-4}$ as the aim is to achieve the most accurate model predictions in the benchmarks.
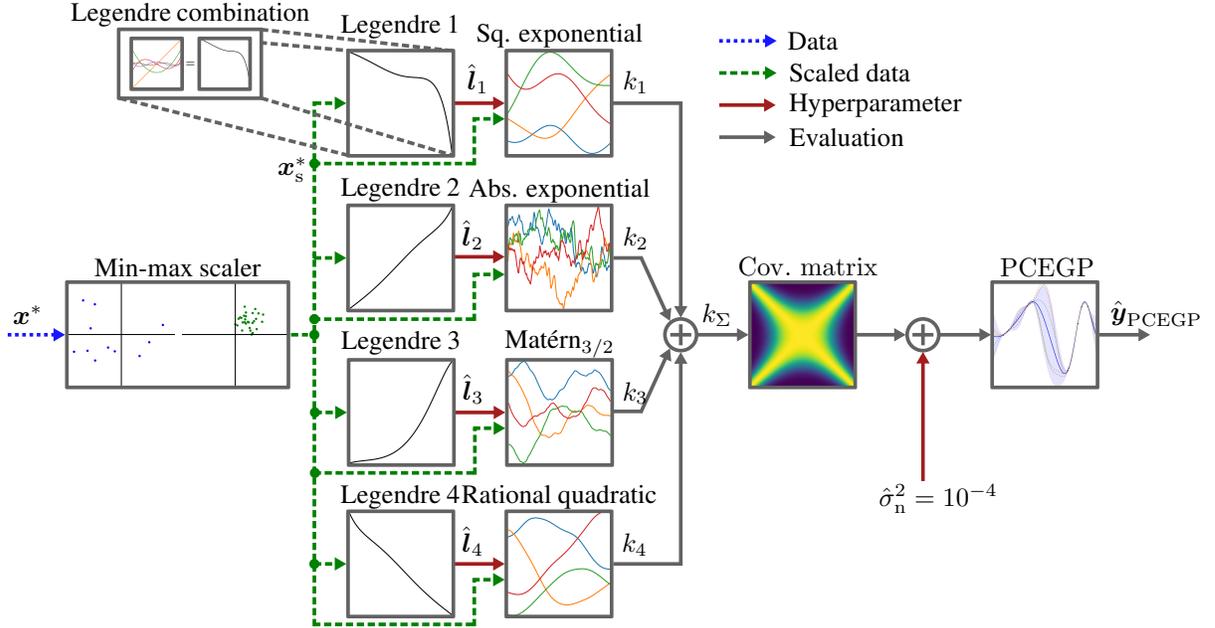


Figure 4: Polynomial Chaos Expanded Gaussian Process setup for experiments