
ENHANCED UNIVERSAL KRIGING FOR TRANSFORMED INPUT PARAMETER SPACES

• **Matthias Fischer**

Institute of Engineering Mechanics
Karlsruhe Institute of Technology
76131 Karlsruhe, Germany
matthias.fischer@kit.edu

• **Carsten Proppe**

Institute of Engineering Mechanics
Karlsruhe Institute of Technology
76131 Karlsruhe, Germany

ABSTRACT

With computational models becoming more expensive and complex, surrogate models have gained increasing attention in many scientific disciplines and are often necessary to conduct sensitivity studies, parameter optimization etc. In the scientific discipline of uncertainty quantification (UQ), model input quantities are often described by probability distributions. For the construction of surrogate models, space-filling designs are generated in the input space to define training points, and evaluations of the computational model at these points are then conducted. The physical parameter space is often transformed into an i. i. d. uniform input space in order to apply space-filling training procedures in a sensible way. Due to this transformation surrogate modeling techniques tend to suffer with regard to their prediction accuracy. Therefore, a new method is proposed in this paper where input parameter transformations are applied to basis functions for universal kriging. To speed up hyperparameter optimization for universal kriging, suitable expressions for efficient gradient-based optimization are developed. Several benchmark functions are investigated and the proposed method is compared with conventional methods.

Manuscript accepted in: Probabilistic Engineering Mechanics on June 27, 2023
<https://doi.org/10.1016/j.pro bengmech.2023.103486>

©2023. This manuscript version is made available under the CC-BY-NC-ND 4.0 license
<https://creativecommons.org/licenses/by-nc-nd/4.0/>

Keywords universal kriging · transformed input parameter space · basis functions · hyperparameter optimization

1 Introduction

Surrogate modeling methods are widely used in various scientific disciplines, such as engineering, meteorology or physics. Simulation runs for complex computational models are often computationally expensive. Surrogate modeling techniques provide the opportunity to construct a surrogate model with relatively few evaluations of the computational model.

For instance, in meteorology, the German Meteorological Service (DWD) currently uses a global grid with a mesh size of 13 km to achieve sufficiently accurate weather predictions. Even at these high resolutions, many atmospheric processes occur on a sub-grid scale and still need to be parameterized, since they cannot be explicitly resolved. Consequently, these models are computationally expensive to run which makes further studies, such as parameter identification or sensitivity studies, infeasible.

By using surrogate models, a relationship between model inputs and outputs can be determined with relatively low computational cost. Model parameters of the computational model are often considered as input parameters of the surrogate model in the context of uncertainty quantification (UQ). Epistemic uncertainties can arise from a lack of knowledge of these parameters. These uncertainties are often described by probability distributions that are determined

based on measurements, expert knowledge or parameter identification studies. Before applying surrogate modeling techniques, input parameters are often transformed to independent and identically distributed (i. i. d.) uniform variables. In this unit hypercube, space-filling designs, such as Latin hypercube sampling, can be applied in a meaningful way. However, prediction accuracy of surrogate models may suffer from such transformations as will be shown in this paper. The aim of this work is to develop a strategy to overcome such prediction accuracy losses as a result of input space transformations.

Although the exact relationship between input and output quantities is not known in advance, it may be possible to assume certain underlying trends in order to improve prediction accuracy. Such trends can be incorporated as basis functions by various surrogate modeling techniques, e. g. universal kriging. Usually, simple basis functions (often low order polynomials) are used for the sake of simplicity and to encounter the risk of overfitting. In a high dimensional input space, higher order multivariate polynomials up to a certain polynomial degree would lead to an excessive increase in the number of polynomials, commonly known as the "curse of dimensionality". To overcome this problem, methods for polynomial basis selection can be applied, such as least-angle regression. Kersaudy et al. [1] proposed the method *LARS-Kriging-PC modeling* where explicit multivariate basis functions for universal kriging are obtained by least-angle regression based on polynomial chaos expansions. The purpose of the method is to select polynomials that bring the most relevant information to the kriging model. The results show that the benefits of both Gaussian process regression and polynomial chaos expansion are combined. The concept is similar to *blind kriging* where the underlying trend is identified from data using a Bayesian variable selection technique (Joseph et al. [2]). As with the case of *LARS-Kriging-PC modeling*, a set of candidate basis functions has to be defined beforehand. In both methods, the choice of polynomials is not based on prior knowledge of the problem, but on the regression technique itself. However, Oakley [3] emphasizes that the choice of basis functions should be chosen to incorporate any beliefs regarding the problem, e. g. the physical evolution of the output variable with respect to the input parameters (see e. g. Kersaudy et al. [1]).

None of the above mentioned methods consider the input parameter transformations to the i. i. d. uniform input space which is performed in this study. However, the transformation may highly impact the shape of the model function. Although it is straightforward to define simple polynomial basis functions in the i. i. d. uniform input space, i. e. the space in which the surrogate model is built, this is not the most meaningful choice. More precisely, defining basis functions in the original, physical parameter space is considered more sensible. In this paper, a method for defining basis functions for universal kriging is proposed that takes input parameter transformations into account. In particular, isoprobabilistic transformation is applied to the construction of basis functions.

In this study it will be assumed that stationarity holds in the i. i. d. uniform space in reference to the uniform density of training points in this space. This can generally be considered to be a reasonable first approach if nothing is known about the physical problem. Accordingly, this corresponds to non-stationary basis functions in the physical input space. The idea of finding transformations between spaces so that stationarity and isotropy holds in the deformed space was already used by Sampson and Guttorp [4] where thin-plate splines were applied to achieve such a mapping. Schmidt and O'Hagan [5] enhanced this method by using a Bayesian approach where the mapping is described by a function with a Gaussian process prior. Both methods have been developed for a setting with only two input dimensions in the context of geostatistics.

A related approach of incorporating non-stationary kernel functions has been proposed by Gibbs [6] and a simplified version has been demonstrated by Xiong et al. [7], where a density function is constructed that aims at describing the smoothness of function value changes with respect to input parameters. The density function is then used to describe a non-linear mapping to an input space where uniform smoothness and therefore a stationary kernel function can be assumed. Using this mapping, non-stationary kernel functions can be defined in the original input space. The difference compared to this paper is that in the approach by Gibbs [6] the transformation is based on the estimated smoothness of function values, whereas here it is based on the joint PDF of input parameters and thus the density of training points.

The core of Gaussian process regression methods lies in the estimation of hyperparameters. These are generally determined by likelihood maximization. With a high number of training points and input dimensions, the estimation of hyperparameters may be computationally expensive. Therefore, gradient-based hyperparameter estimation for Gaussian process regression is frequently applied to speed up hyperparameter optimization. The derivative of the log marginal likelihood with respect to the hyperparameters is determined for universal kriging in this work.

The remainder of this paper is organized as follows. In Section 2 applied methods are depicted. First, the applied training procedure and fundamentals of Gaussian process regression are presented. Then, a technique for

the construction of basis functions for transformed input parameter spaces is proposed. Suitable expressions for gradient-based hyperparameter optimization for universal kriging are developed. Model validation measures are specified. In Section 3 the model setup and the set of benchmark problems for investigation of proposed methods are presented. Validation results are shown in Section 4 and discussed in Section 5. Finally, Section 6 concludes this work.

2 Universal kriging with basis functions for transformed input parameter spaces

In this section, the applied training procedure, simple kriging and universal kriging are briefly presented. A new way of defining basis functions for transformed input parameter spaces is then proposed. Equations for gradient-based hyperparameter optimization for universal kriging are developed. Finally, the applied model validation strategy is described.

2.1 Training procedure

In this study, problems with non-uniform input parameter distributions are considered. In terms of training strategies, it is reasonable to generate a space-filling design that takes the PDFs of input parameters into account, i.e. to generate more training points in regions associated with higher probability and vice versa. This generally yields higher overall prediction accuracy, because the resulting surrogate model is then more accurate in regions associated with higher probabilities. In particular, lower validation errors are obtained, if the model is validated by Monte Carlo sampling with respect to parameter PDFs. Lu et al. [8] described a procedure to find such space-filling designs. However, surrogate modeling methods, such as Gaussian process regression, may struggle with inhomogeneous density of training points in the input space, because optimal covariance parameters are usually strongly dependent on the density of training points. However, if specific problems are considered where a higher accuracy is desired in certain regions of the input space, as it is the case in reliability analysis (see e. g. Bourinet [9]), adjusted methods are preferred that assign a higher training point density to such regions.

Therefore, the focus of this study is the case where input parameter spaces are transformed to i. i. d. uniform parameter spaces. Isoprobabilistic transformation, in particular Rosenblatt transformation [10], is applied to transform the physical input vector $\mathbf{x} = (x_1, x_2 \dots x_p)^\top$ to the i. i. d. uniform input vector $\mathbf{u} = (u_1, u_2 \dots u_p)^\top$ with input space dimension p . This makes the application of space-filling designs more sensible, because a homogeneous space-filling design can then be obtained in the uniform input space. Thus, the training point density induces a model accuracy which corresponds to the probability of parameter values. Furthermore, in the case of highly different scales between the input parameters, surrogate modeling methods may suffer without such a transformation. Let $\hat{\mathbf{x}} = (\hat{x}_1, \hat{x}_2 \dots \hat{x}_p)^\top$ be a random vector in the physical input space with joint cumulative distribution function $F_{\hat{\mathbf{x}}}(\mathbf{x})$. The Rosenblatt transformation $\mathbf{u} = \mathcal{T}_{\text{ros}}(\mathbf{x})$ is then given by

$$\begin{aligned} u_1 &= P\{\hat{x}_1 \leq x_1\} = F_{\hat{x}_1}(x_1) \\ u_2 &= P\{\hat{x}_2 \leq x_2 \mid \hat{x}_1 = x_1\} = F_{\hat{x}_2|\hat{x}_1}(x_2 \mid x_1) \\ &\vdots \\ u_p &= P\{\hat{x}_p \leq x_p \mid \hat{x}_{p-1} = x_{p-1} \dots \hat{x}_1 = x_1\} = F_{\hat{x}_p|\hat{x}_{p-1} \dots \hat{x}_1}(x_p \mid x_{p-1} \dots x_1), \end{aligned}$$

where $F_{\square}(\cdot)$ denote respective conditional cumulative distribution functions. It can be shown that transformed random vector $\hat{\mathbf{u}} = (\hat{u}_1, \hat{u}_2 \dots \hat{u}_p)^\top = \mathcal{T}_{\text{ros}}(\hat{\mathbf{x}})$ is then i. i. d. uniformly distributed on the p -dimensional unit hypercube. $F_{\square}(\cdot)$ can be determined from the joint PDF $f_{\hat{\mathbf{x}}}(\mathbf{x})$ which is assumed to be known in this study (see e. g. Melchers and Beck [11]).

Maximin Latin hypercube sampling (Morris and Mitchell [12]) for generating a set of n training points in the i. i. d. uniform parameter space $\mathbf{U} = \{\mathbf{u}_i, i = 1 \dots n\}$ is applied, where \mathbf{U} corresponds to the set of training points in the physical parameter space $\mathbf{X} = \{\mathbf{x}_i, i = 1 \dots n\}$ such that $\{\mathbf{u}_i = \mathcal{T}_{\text{ros}}(\mathbf{x}_i), i = 1 \dots n\}$. Evaluations of the computational model $f(\mathbf{x}_i)$ are conducted for all training points resulting in the output vector $\mathbf{y} = \{y_i = f(\mathbf{x}_i), i = 1 \dots n\}$.

The aim of kriging is to build a surrogate model \mathcal{M} for a scalar model output y , i.e. a quantity of interest (QoI), based on the experimental design and model evaluations $\{\mathbf{U}, \mathbf{y}\}$. Note that all equations in the context of kriging are represented with respect to i. i. d. uniform variables \mathbf{u} instead of physical variables \mathbf{x} , because surrogate models are built in the i. i. d. uniform parameter space in this study. Prediction mean and prediction variance at a set of

input points $\mathbf{U}_* = \{\mathbf{u}_{*i}, i = 1 \dots l\}$ are to be determined.

2.2 Simple kriging

For simple kriging, a zero-mean Gaussian process

$$g_{\text{SK}}(\mathbf{u}) \sim \mathcal{GP}(0, k(\mathbf{u}, \mathbf{u}')) \quad (1)$$

with covariance function $k(\mathbf{u}, \mathbf{u}')$, also known as kernel function, is assumed as prior. The kernel function describes the dependence structure between values of the stochastic process at different points, usually depending on their distance.

Here, the anisotropic form of the radial-basis function

$$k(\mathbf{u}, \mathbf{u}') = \theta_0 \exp\left(-\sum_{i=1}^p \left(\frac{|u_i - u'_i|}{\theta_i}\right)^2\right) \quad (2)$$

with respect to hyperparameters $\boldsymbol{\theta} = \{\theta_i, i = 0 \dots p\}$ is used to allow for different smoothness between input dimensions. Furthermore, i. i. d. Gaussian noise with variance σ_n^2 is added to allow for aleatoric uncertainties in the simulations of the computational model.

Let

$$\begin{aligned} \mathbf{K} &= \{K_{ij} = k(\mathbf{u}_i, \mathbf{u}_j), \quad i = 1 \dots n, \quad j = 1 \dots n\}, \\ \mathbf{k} &= \{k_{ij} = k(\mathbf{u}_i, \mathbf{u}_{*j}), \quad i = 1 \dots n, \quad j = 1 \dots l\} \quad \text{and} \\ \boldsymbol{\sigma}_0^2 &= \{\sigma_{0j}^2 = k(\mathbf{u}_{*j}, \mathbf{u}_{*j}), \quad j = 1 \dots l\} \end{aligned}$$

be the vectors and matrices of kernel function evaluations at training points \mathbf{U} and prediction points \mathbf{U}_* , respectively.

The best linear unbiased predictor (BLUP) and its prediction variance for the set of prediction points \mathbf{U}_* under the assumptions of simple kriging, as shown by Rasmussen and Williams [13], are then

$$\begin{aligned} \mathcal{M}(\mathbf{U}_*) &= \mathbf{k}^\top \mathbf{K}_y^{-1} \mathbf{y} \\ \sigma^2(\mathbf{U}_*) &= \sigma_0^2 - \mathbf{k}^\top \mathbf{K}_y^{-1} \mathbf{k} \end{aligned}$$

with $\mathbf{K}_y = \mathbf{K} + \sigma_n^2 \mathbf{1}$. The hyperparameters $\boldsymbol{\theta}$ and noise parameter σ_n are determined by maximum likelihood estimation.

2.3 Universal kriging

The theory of universal kriging was introduced by Matheron [14] in the field of geostatistics. A prior

$$g_{\text{UK}}(\mathbf{u}) = g_{\text{SK}}(\mathbf{u}) + \mathbf{h}(\mathbf{u})^\top \boldsymbol{\beta} \quad (3)$$

is used, with zero-mean Gaussian process $g_{\text{SK}}(\mathbf{u})$ (Eq. 1), vectors of known basis functions $\mathbf{h}(\mathbf{u}) = \{h_j(\mathbf{u}), j = 1 \dots q\}$ and coefficients $\boldsymbol{\beta} = \{\beta_j, j = 1 \dots q\}$. Coefficients $\boldsymbol{\beta}$ are unknown, but not required to be specified for the computation (see Rasmussen and Williams [13]).

Let

$$\mathbf{H} = \{H_{ij} = h_i(\mathbf{u}_j), \quad i = 1 \dots q, \quad j = 1 \dots n\} \quad \text{and} \quad (4)$$

$$\mathbf{H}_* = \{H_{*ij} = h_i(\mathbf{u}_{*j}), \quad i = 1 \dots q, \quad j = 1 \dots l\} \quad (5)$$

be the matrices of basis function evaluations at training points \mathbf{U} and prediction points \mathbf{U}_* , respectively.

The best linear unbiased predictor (BLUP) and its prediction variance for the set of prediction points \mathbf{U}_* under the assumptions of universal kriging, as shown by Rasmussen and Williams [13], are then

$$\mathcal{M}(\mathbf{U}_*) = \mathbf{H}_*^\top \boldsymbol{\mu} + \mathbf{k}^\top \mathbf{K}_y^{-1} (\mathbf{y} - \mathbf{H}^\top \boldsymbol{\mu}) \quad (6)$$

$$\sigma^2(\mathbf{U}_*) = \sigma_0^2 - \mathbf{k}^\top \mathbf{K}_y^{-1} \mathbf{k} + \mathbf{R}^\top (\mathbf{H} \mathbf{K}_y^{-1} \mathbf{H}^\top)^{-1} \mathbf{R} \quad (7)$$

with $\boldsymbol{\mu} = (\mathbf{H}\mathbf{K}_y^{-1}\mathbf{H}^\top)^{-1}\mathbf{H}\mathbf{K}_y^{-1}\mathbf{y}$ and $\mathbf{R} = \mathbf{H}_* - \mathbf{H}\mathbf{K}_y^{-1}\mathbf{k}$.

In the case of a constant scalar basis function $h(\mathbf{u}) = 1$, Eq. 6 and Eq. 7 become

$$\begin{aligned}\mathcal{M}(\mathbf{U}_*) &= \boldsymbol{\mu}\mathbf{I} + \mathbf{k}^\top\mathbf{K}_y^{-1}(\mathbf{y} - \boldsymbol{\mu}\mathbf{I}) \\ \sigma^2(\mathbf{U}_*) &= \sigma_0^2 - \mathbf{k}^\top\mathbf{K}_y^{-1}\mathbf{k} + \mathbf{R}^\top(\mathbf{H}\mathbf{K}_y^{-1}\mathbf{H}^\top)^{-1}\mathbf{R}\end{aligned}$$

where $\boldsymbol{\mu} = \frac{\mathbf{I}^\top\mathbf{K}_y^{-1}\mathbf{y}}{\mathbf{I}^\top\mathbf{K}_y^{-1}\mathbf{I}}$. This case is known as ordinary kriging, i. e. kriging with unknown constant mean.

2.4 Basis functions for transformed input parameter spaces

In many cases, a zero or constant mean function for the Gaussian process is sufficient which corresponds to scalar basis functions $\mathbf{h}(u) = 0$ or $\mathbf{h}(u) = 1$ in Eq. 3, respectively. However, incorporation of more sophisticated basis functions may significantly improve prediction accuracy. Kersaudy et al. [1] proposed the method *LARS-Kriging-PC modeling* where explicit basis functions for universal kriging are obtained by least-angle regression based on polynomial chaos expansions. A sparse polynomial basis is generated to tackle the curse of dimensionality. Similarly, for *blind kriging*, the underlying trend is identified from data using a Bayesian variable selection technique (Joseph et al. [2]). In both cases, basis functions are selected from a set of candidate functions based on the surrogate modeling technique itself rather than physical knowledge about the problem.

In particular, the selection of basis functions becomes more important when using transformed input parameter spaces. In this study, transformation of the physical input space to an i. i. d. uniform input space is carried out as described in Section 2.1. Kriging techniques are then applied in the i. i. d. uniform input space. Although it is then straightforward to define simple polynomial basis functions in the i. i. d. uniform input space, this is not the most meaningful choice. Defining basis functions in the original, physical parameter space is considered to be more sensible. This is because physical relationships that are attempted to be captured by trend functions are assumed to be related more directly to physical input parameters than to transformed i. i. d. uniform input parameters. Therefore, it is proposed to define polynomial basis functions $f(\mathbf{x})$ in the original, physical input space \mathbf{x} and express them as functions $h(\mathbf{u}) = f(\mathcal{T}_{\text{ros}}^{-1}(\mathbf{u}))$ in the i. i. d. uniform input space \mathbf{u} by means of the inverse Rosenblatt transformation $\mathcal{T}_{\text{ros}}^{-1}$ (see Section 2.1). In the one-dimensional case, the transformation becomes the quantile function or percent-point-function $x = \text{PPF}(u)$ (inverse cumulative distribution function) of the input parameter. The transformed basis function is then $h(u) = f(\text{PPF}(u))$. In case of the particular linear basis function $f(x) = x$, the transformed basis function results in the quantile function $h(u) = \text{PPF}(u)$. If there is no correlation between the input parameters, i. e. $\{\rho_{\hat{x}_i, \hat{x}_j} = 0 \quad \forall i, j \in \{1, \dots, p\}, i \neq j\}$, the transformation reduces to independent functions $\{h(u_i) = f(\text{PPF}(u_i)), i = 1 \dots p\}$ in all input parameters.

Fig. 1 illustrates the transformation of basis functions. Transformed basis functions $\mathbf{h}(\mathbf{u}) = \{h_j(\mathbf{u}), j = 1 \dots q\}$ can then be used as basis functions in the universal kriging method (Eq. 3), i. e. for computing \mathbf{H} and \mathbf{H}_* (Eq. 4 and Eq. 5).

2.5 Link to non-stationary kernel functions

In this paper, the procedure of building surrogate models is carried out in the i. i. d. uniform parameter space. However, by using the Rosenblatt transformation, uniform input variables can be transformed into the physical parameter space as described in Section 2.1. All equations from Section 2.2 and Section 2.3 can thus be expressed with respect to physical parameters \mathbf{x} . In the case of universal kriging, the original non-transformed definition of trend functions in the physical parameter space are used instead of the transformed basis functions in the i. i. d. uniform parameter space according to Section 2.4.

In particular, the description of equations results in non-stationary kernel functions, because the kernel functions are defined as stationary with respect to i. i. d. uniform variables. The anisotropic formulation of the radial-basis function kernel (Eq. 2) thus results in

$$k(\mathbf{x}, \mathbf{x}') = \theta_0 \exp \left(- \sum_{i=1}^p \left(\frac{|\mathcal{T}_{\text{ros}}(\mathbf{x})_i - \mathcal{T}_{\text{ros}}(\mathbf{x}')_i|}{\theta_i} \right)^2 \right).$$

It is emphasized that the definition of stationary kernel functions in the i. i. d. uniform parameter space is meaningful instead of using stationary kernel functions in the physical parameter space to account for the non-homogeneous density

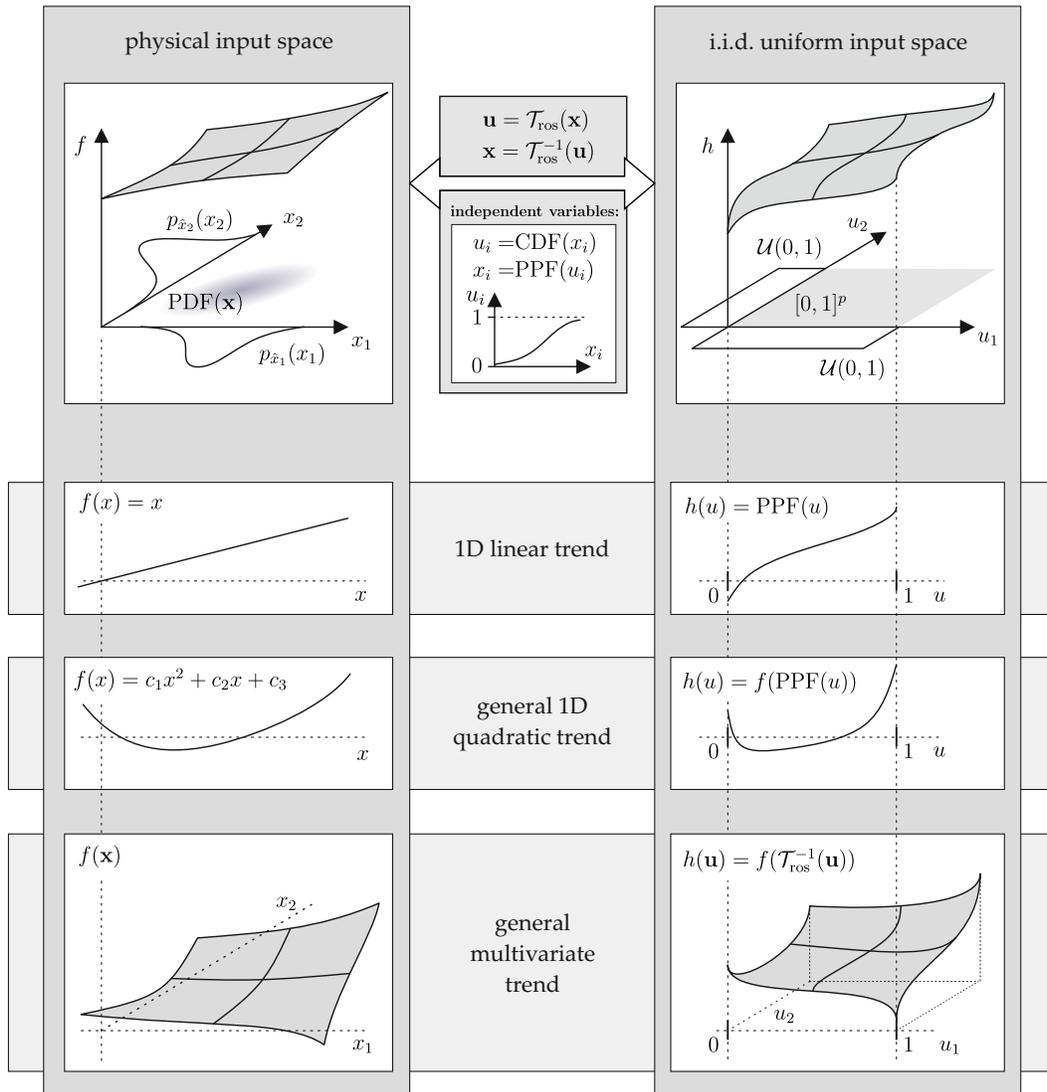


Figure 1: Visualization of model functions (top) and selected trend functions (three rows below) in the physical (left) and i. i. d. uniform (right) parameter space. The Rosenblatt transformation \mathcal{T}_{ros} (quantile functions PPF in case of independent variables) defines the relationship between both parameter spaces.

of training points. According to Section 2.1, the density of training points in the physical parameter space corresponds to the joint PDF of input parameters. Training points are sparse in the tails of the PDF and a larger length scale of the kernel function in these regions is desirable, whereas in regions with a higher density of training points, i. e. closer to the mean value, a smaller length scale is sensible.

2.6 Gradient-based hyperparameter estimation for universal kriging

Hyperparameter optimization for universal kriging can be computationally intensive, particularly with many input dimensions and a large number of training points. Optimization of the hyperparameters is most commonly conducted by maximizing the marginal likelihood which is a measure for the model evidence. However, other methods exist, such as maximization of the *pseudo-likelihood* which is obtained by cross-validation (see Rasmussen and Williams [13] for further discussion), which are not considered in this work. Preferably, gradient-based hyperparameter optimization with incorporation of gradient information is conducted to reduce computation time. Therefore, the gradient of the log marginal likelihood with respect to the hyperparameters θ has to be determined.

The log marginal likelihood for simple kriging is given as

$$\log p(\mathbf{y}|\mathbf{U}, \theta) = -\frac{1}{2}\mathbf{y}^\top \mathbf{K}_y^{-1} \mathbf{y} - \frac{1}{2} \log |\mathbf{K}_y| - \frac{p}{2} \log 2\pi. \quad (8)$$

Since only \mathbf{K}_y depends on θ , an expression for the gradient of the log marginal likelihood is aimed that only depends on $\frac{\partial \mathbf{K}_y}{\partial \theta}$. If the gradient of \mathbf{K}_y is available, as it is the case for widely used kernel functions, the gradient of the log marginal likelihood can be determined.

For simple kriging, the gradient of the log marginal likelihood by the use of matrix identities in A becomes

$$\begin{aligned} \frac{\partial}{\partial \theta_l} \log p(\mathbf{y}|\mathbf{U}, \theta) &= \frac{1}{2}\mathbf{y}^\top \mathbf{K}_y^{-1} \frac{\partial \mathbf{K}_y}{\partial \theta_l} \mathbf{K}_y^{-1} \mathbf{y} - \frac{1}{2} \text{tr} \left(\mathbf{K}_y^{-1} \frac{\partial \mathbf{K}_y}{\partial \theta_l} \right) \\ &= \frac{1}{2} \text{tr} \left((\boldsymbol{\alpha} \boldsymbol{\alpha}^\top - \mathbf{K}_y^{-1}) \frac{\partial \mathbf{K}_y}{\partial \theta_l} \right) \end{aligned}$$

with $\boldsymbol{\alpha} = \mathbf{K}_y^{-1} \mathbf{y} = \mathbf{L}_K^{-\top} (\mathbf{L}_K^{-1} \mathbf{y})$ and $\mathbf{L}_K = \text{cholesky}(\mathbf{K}_y)$.

Using the Einstein summation convention, the derivative of the log marginal likelihood given $\boldsymbol{\alpha}$ can be expressed as

$$\frac{\partial}{\partial \theta_l} \log p(\mathbf{y}|\mathbf{U}, \theta) = \frac{1}{2} \left((\alpha_i \alpha_j - \delta_{im} [\mathbf{K}_y^{-1}]_{mj}) \frac{\partial K_{y\,ji}}{\partial \theta_l} \right) \quad (9)$$

with Kronecker delta δ_{im} .

In the following, the derivation of the gradient of the log marginal likelihood for universal kriging is carried out, where suitable abbreviations for efficient computation are introduced.

The log marginal likelihood for universal kriging, as shown by Rasmussen and Williams [13], is

$$\log p(\mathbf{y}|\mathbf{U}, \theta) = -\frac{1}{2}\mathbf{y}^\top \mathbf{K}_y^{-1} \mathbf{y} + \frac{1}{2}\mathbf{y}^\top \mathbf{C} \mathbf{y} - \frac{1}{2} \log |\mathbf{K}_y| - \frac{1}{2} \log |\mathbf{A}| - \frac{p-m}{2} \log 2\pi, \quad (10)$$

where $\mathbf{A} = \mathbf{H} \mathbf{K}_y^{-1} \mathbf{H}^\top$, $\mathbf{C} = \mathbf{K}_y^{-1} \mathbf{H}^\top \mathbf{A}^{-1} \mathbf{H} \mathbf{K}_y^{-1}$ and $m = \text{rank}(\mathbf{H}^\top)$.

Inserting \mathbf{A} and \mathbf{C} in Eq. 10 results in

$$\begin{aligned} \log p(\mathbf{y}|\mathbf{U}, \theta) &= -\frac{1}{2}\mathbf{y}^\top \mathbf{K}_y^{-1} \mathbf{y} + \frac{1}{2}\mathbf{y}^\top (\mathbf{K}_y^{-1} \mathbf{H}^\top (\mathbf{H} \mathbf{K}_y^{-1} \mathbf{H}^\top)^{-1} \mathbf{H} \mathbf{K}_y^{-1}) \mathbf{y} \\ &\quad - \frac{1}{2} \log |\mathbf{K}_y| - \frac{1}{2} \log |\mathbf{H} \mathbf{K}_y^{-1} \mathbf{H}^\top| - \frac{n-m}{2} \log 2\pi \\ &= -\frac{1}{2}\mathbf{y}^\top \boldsymbol{\alpha} + \frac{1}{2}\mathbf{y}^\top \boldsymbol{\gamma} \boldsymbol{\eta} \boldsymbol{\alpha} \\ &\quad - \frac{1}{2} \log |\mathbf{K}_y| - \frac{1}{2} \log |\mathbf{H} \boldsymbol{\gamma}| - \frac{n-m}{2} \log 2\pi, \end{aligned}$$

where α , γ and η are defined as

$$\alpha = \mathbf{K}_y^{-1} \mathbf{y} = \mathbf{L}_K^{-\top} (\mathbf{L}_K^{-1} \mathbf{y}), \quad (11)$$

$$\gamma = \mathbf{K}_y^{-1} \mathbf{H}^\top = \mathbf{L}_K^{-\top} (\mathbf{L}_K^{-1} \mathbf{H}^\top) \quad \text{and} \quad (12)$$

$$\eta = (\mathbf{H} \mathbf{K}_y^{-1} \mathbf{H}^\top)^{-1} \mathbf{H} = \mathbf{L}_\eta^{-\top} (\mathbf{L}_\eta^{-1} \mathbf{H}) \quad (13)$$

with $\mathbf{L}_K = \text{cholesky}(\mathbf{K}_y)$ and $\mathbf{L}_\eta = \text{cholesky}(\mathbf{H} \mathbf{K}_y^{-1} \mathbf{H}^\top)$. Cholesky decomposition is applied in order to efficiently determine the inverse.

Taking the derivative with respect to θ by the use of matrix identities in A results in

$$\begin{aligned} & \frac{\partial}{\partial \theta_l} \log p(\mathbf{y} | \mathbf{U}, \theta) \\ &= \frac{1}{2} \mathbf{y}^\top \mathbf{K}_y^{-1} \frac{\partial \mathbf{K}_y}{\partial \theta_l} \mathbf{K}_y^{-1} \mathbf{y} \\ & \quad - \frac{1}{2} \mathbf{y}^\top \left(\mathbf{K}_y^{-1} \frac{\partial \mathbf{K}_y}{\partial \theta_l} \mathbf{K}_y^{-1} \mathbf{H}^\top (\mathbf{H} \mathbf{K}_y^{-1} \mathbf{H}^\top)^{-1} \mathbf{H} \mathbf{K}_y^{-1} \right) \mathbf{y} \\ & \quad + \frac{1}{2} \mathbf{y}^\top \left(\mathbf{K}_y^{-1} \mathbf{H}^\top (\mathbf{H} \mathbf{K}_y^{-1} \mathbf{H}^\top)^{-1} \mathbf{H} \mathbf{K}_y^{-1} \frac{\partial \mathbf{K}_y}{\partial \theta_l} \mathbf{K}_y^{-1} \mathbf{H}^\top (\mathbf{H} \mathbf{K}_y^{-1} \mathbf{H}^\top)^{-1} \mathbf{H} \mathbf{K}_y^{-1} \right) \mathbf{y} \\ & \quad - \frac{1}{2} \mathbf{y}^\top \left(\mathbf{K}_y^{-1} \mathbf{H}^\top (\mathbf{H} \mathbf{K}_y^{-1} \mathbf{H}^\top)^{-1} \mathbf{H} \mathbf{K}_y^{-1} \frac{\partial \mathbf{K}_y}{\partial \theta_l} \mathbf{K}_y^{-1} \right) \mathbf{y} \\ & \quad - \frac{1}{2} \text{tr} \left(\mathbf{K}_y^{-1} \frac{\partial \mathbf{K}_y}{\partial \theta_l} \right) - \frac{1}{2} \text{tr} \left(-(\mathbf{H} \mathbf{K}_y^{-1} \mathbf{H}^\top)^{-1} \mathbf{H} \mathbf{K}_y^{-1} \frac{\partial \mathbf{K}_y}{\partial \theta_l} \mathbf{K}_y^{-1} \mathbf{H}^\top \right). \end{aligned}$$

In the next step, symmetry and positive definiteness of \mathbf{K}_y (by definition) and therefore symmetry of \mathbf{K}_y^{-1} are taken into account. Furthermore, the matrix identity in Eq. 18 is used to simplify the result. With α , η and γ , it follows

$$\begin{aligned} & \frac{\partial}{\partial \theta_l} \log p(\mathbf{y} | \mathbf{U}, \theta) \\ &= \frac{1}{2} \alpha^\top \frac{\partial \mathbf{K}_y}{\partial \theta_l} \alpha - \frac{1}{2} \alpha^\top \frac{\partial \mathbf{K}_y}{\partial \theta_l} \gamma \eta \alpha + \frac{1}{2} \alpha^\top \eta^\top \gamma^\top \frac{\partial \mathbf{K}_y}{\partial \theta_l} \gamma \eta \alpha \\ & \quad - \frac{1}{2} \alpha^\top \eta^\top \gamma^\top \frac{\partial \mathbf{K}_y}{\partial \theta_l} \alpha - \frac{1}{2} \text{tr} \left(\mathbf{K}_y^{-1} \frac{\partial \mathbf{K}_y}{\partial \theta_l} \right) + \frac{1}{2} \text{tr} \left(\gamma \eta \mathbf{K}_y^{-1} \frac{\partial \mathbf{K}_y}{\partial \theta_l} \right) \\ &= \frac{1}{2} \text{tr} \left(\alpha^\top \frac{\partial \mathbf{K}_y}{\partial \theta_l} \alpha - \alpha^\top \frac{\partial \mathbf{K}_y}{\partial \theta_l} \gamma \eta \alpha \right. \\ & \quad \left. + \alpha^\top \eta^\top \gamma^\top \frac{\partial \mathbf{K}_y}{\partial \theta_l} \gamma \eta \alpha - \alpha^\top \eta^\top \gamma^\top \frac{\partial \mathbf{K}_y}{\partial \theta_l} \alpha \right) \\ & \quad - \frac{1}{2} \text{tr} \left(\mathbf{K}_y^{-1} \frac{\partial \mathbf{K}_y}{\partial \theta_l} \right) + \frac{1}{2} \text{tr} \left(\gamma \eta \mathbf{K}_y^{-1} \frac{\partial \mathbf{K}_y}{\partial \theta_l} \right) \\ &= \frac{1}{2} \text{tr} \left(\left(\alpha \alpha^\top - \gamma \eta \alpha \alpha^\top + \gamma \eta \alpha \alpha^\top \eta^\top \gamma^\top \right. \right. \\ & \quad \left. \left. - \alpha \alpha^\top \eta^\top \gamma^\top - \mathbf{K}_y^{-1} + \gamma \eta \mathbf{K}_y^{-1} \right) \frac{\partial \mathbf{K}_y}{\partial \theta_l} \right). \end{aligned}$$

The abbreviations $\rho = \alpha \alpha^\top$, $\varepsilon = \gamma \eta$ and $\xi = \varepsilon \rho$ are introduced. It follows

$$\frac{\partial}{\partial \theta_l} \log p(\mathbf{y} | \mathbf{U}, \theta) = \frac{1}{2} \text{tr} \left(\left(\rho - \xi - \xi^\top + \xi \varepsilon^\top + (\varepsilon - \mathbb{1}) \mathbf{K}_y^{-1} \right) \frac{\partial \mathbf{K}_y}{\partial \theta_l} \right).$$

Using the Einstein summation convention, the following calculation steps are required to determine the derivative of the log marginal likelihood given α , γ and η from Eq. 11, Eq. 12 and Eq. 13:

$$\begin{aligned}
 \rho_{ij} &= \alpha_i \alpha_j, \\
 \varepsilon_{ij} &= \gamma_{ik} \eta_{kj}, \\
 \xi_{il} &= \varepsilon_{ij} \rho_{jl}, \\
 \frac{\partial}{\partial \theta_l} \log p(\mathbf{y}|\mathbf{U}, \boldsymbol{\theta}) &= \frac{1}{2} \left(\left(\rho_{ij} - \xi_{ij} - \xi_{ji} + \xi_{im} \varepsilon_{jm} + (\varepsilon_{im} - \delta_{im}) [\mathbf{K}_y^{-1}]_{mj} \right) \frac{\partial K_{y_{ji}}}{\partial \theta_l} \right)
 \end{aligned} \tag{14}$$

with Kronecker delta δ_{im} .

The complexity of the computation of the gradient of the log marginal likelihood (Eq. 9 and Eq. 14) is dominated by the inverse of matrix \mathbf{K}_y which is of the computational complexity $\mathcal{O}(n^3)$. Once the inverse is determined, it can be used for the computation of all hyperparameters θ_l . In contrast, computation of the gradient of the log marginal likelihood based on Eq. 8 and Eq. 10, i. e. without using the gradient of matrix \mathbf{K}_y , is of the computational complexity $\mathcal{O}(p \cdot n^3)$, because the inverse of \mathbf{K}_y has to be determined towards all input dimensions p . Gradient-based hyperparameter optimization with incorporation of gradient information is therefore more beneficial.

2.7 Model validation

When surrogate models are constructed, it is essential to assess their quality based on model validation measures. Therefore, the generalization error

$$\text{SE} = \mathbb{E} [(f(\hat{\mathbf{x}}) - \mathcal{M}(\hat{\mathbf{x}}))^2]$$

is considered, where $\mathbb{E}[\cdot]$ is the mathematical expectation operator. It describes the squared difference between original physical model f and surrogate model prediction \mathcal{M} and is therefore denoted as SE (*squared error*). It can be expressed as

$$\begin{aligned}
 \text{SE} &= \int_{\mathcal{D}_{\hat{\mathbf{x}}}} (f(\mathbf{x}) - \mathcal{M}(\mathbf{x}))^2 f_{\hat{\mathbf{x}}}(\mathbf{x}) d\mathbf{x} \\
 &= \int_{\mathcal{D}_{\hat{\mathbf{u}}}} (f(\mathbf{u}) - \mathcal{M}(\mathbf{u}))^2 f_{\hat{\mathbf{u}}}(\mathbf{u}) d\mathbf{u} = \int_{\mathcal{D}_{\hat{\mathbf{u}}}} (f(\mathbf{u}) - \mathcal{M}(\mathbf{u}))^2 d\mathbf{u}
 \end{aligned} \tag{15}$$

where $f_{\hat{\mathbf{x}}}$ (resp. $\mathcal{D}_{\hat{\mathbf{x}}}$) is the PDF (resp. the support) of the random input vector $\hat{\mathbf{x}}$ and $f_{\hat{\mathbf{u}}}$ (resp. $\mathcal{D}_{\hat{\mathbf{u}}}$) is the PDF (resp. the support) of the random input vector $\hat{\mathbf{u}}$. The PDF of the i. i. d. uniform random input vector $\hat{\mathbf{u}}$ is $f_{\hat{\mathbf{u}}}(\mathbf{u}) = 1$.

The SE value (Eq. 15) is generally not known analytically, since the model function \mathcal{M} is assumed to be known only at certain points as is the case for complex computer models. Therefore, the generalization error is estimated by using a validation set $\{(\mathbf{u}_{\text{val},i}, y_{\text{val},i}), i = 1 \dots n_{\text{val}}\}$ with n_{val} validation points obtained from evaluations of the computer model. The validation points are obtained by Monte-Carlo sampling with respect to the parameter PDFs. The normalized mean squared error results in

$$\text{NMSE} = \frac{\hat{\text{SE}}}{\hat{\sigma}_y^2} = \frac{1}{\sigma_{y_{\text{val}}}^2} \frac{1}{n_{\text{val}}} \sum_{i=1}^{n_{\text{val}}} (y_{\text{val},i} - \mathcal{M}(\mathbf{u}_{\text{val},i}))^2.$$

Here,

$$\begin{aligned}
 \bar{y} &= \frac{1}{n_{\text{val}}} \sum_{i=1}^{n_{\text{val}}} y_{\text{val},i} \quad \text{and} \\
 \sigma_{y_{\text{val}}}^2 &= \frac{1}{n_{\text{val}} - 1} \sum_{i=1}^{n_{\text{val}}} (y_{\text{val},i} - \bar{y})^2
 \end{aligned}$$

are the mean and the variance of evaluations $y_{\text{val},i}$ as estimates for the output variable y , respectively.

Model accuracy is considered to be high if NMSE values are close to 0 and low if NMSE values are close to 1. By definition, values are non-negative and should not exceed 1 as the covariance between the surrogate model and data would then be higher than the variance of the data. The NMSE error is used since normalization allows comparison between different problems, i. e. different scales. Compared to the Q^2 error ($Q^2 = 1 - \text{NMSE}$, see e. g. Kersaudy et al. [1]), the NMSE is more sensible to visualize graphically on a logarithmic scale which is reasonable for small values.

In the case of expensive computational models, cross-validation techniques such as leave-one-out cross-validation can be applied instead of using a separate validation data set. If the cross-validation error is costly to compute, analytical expressions can be used to reduce computational cost arising by obtaining many separate leave-one-out surrogate models (see e. g. Dubrule [15]). Since computationally inexpensive test cases are used in this study to investigate the proposed method, separate validation data sets are used to ensure accurate validation results.

3 Test cases

The proposed method is applied to several benchmark problems that are shown in Table 1. Probability density functions are assigned to input parameters. All investigated functions contain a non-uniform distribution for at least one input dimension because otherwise the Rosenblatt transformation would not affect the trend function and the proposed method would not differ from the conventional methods.

Table 1: Benchmark functions including labels, mathematical expressions, number of input dimensions, probability density functions (PDF) and parameter correlations.

#	equation	dim.	input parameter PDFs* $f_{\hat{x}_i}$ and pairwise Pearson correlation coefficients $\rho_{\hat{x}_i, \hat{x}_j}$ ($\rho_{\hat{x}_i, \hat{x}_j} = 0$ if not stated)
1	Oakley & O'Hagan [16] $f(x) = 5 + x + \cos(x)$	1	$\hat{x} \sim \mathcal{N}(0, 4)$
2	Lognormal Ratio [17] $f(\mathbf{x}) = \frac{x_1}{x_2}$	2	$\hat{x}_{1,2} \sim \mathcal{LN}(1, 0.5)$ $\rho_{\hat{x}_1, \hat{x}_2} = 0.3$
3	Webster et al. [18] $f(\mathbf{x}) = x_1^2 + x_2^3$	2	$\hat{x}_1 \sim \mathcal{U}(1, 10)$ $\hat{x}_2 \sim \mathcal{N}(2, 1)$
4	Short Column [17] $f(\mathbf{x}) = 1 - \frac{4}{1125} \frac{x_2}{x_1} - \frac{1}{5625} \left(\frac{x_3}{x_1} \right)^2$	3	$\hat{x}_1 \sim \mathcal{LN}(5, 0.5)$ $\hat{x}_2 \sim \mathcal{N}(2000, 400)$ $\hat{x}_3 \sim \mathcal{N}(500, 100)$ $\rho_{\hat{x}_2, \hat{x}_3} = 0.5$
5	Cantilever Beam [17] $f(\mathbf{x}) = \frac{5 \cdot 10^5}{x_1} \sqrt{\left(\frac{x_2}{16} \right)^2 + \left(\frac{x_3}{4} \right)^2}$	3	$\hat{x}_1 \sim \mathcal{N}(2.9e7, 1.45e6)$ $\hat{x}_2 \sim \mathcal{N}(1000, 100)$ $\hat{x}_3 \sim \mathcal{N}(500, 100)$
6	Borehole [19, 20] $f(\mathbf{x}) = \frac{2\pi x_3(x_4 - x_6)}{\ln\left(\frac{x_2}{x_1}\right) \left(1 + \frac{2x_7 x_3}{\ln(x_2/x_1)x_7^2 x_8} + \frac{x_3}{x_5}\right)}$	8	$\hat{x}_1 \sim \mathcal{N}(0.1, 0.0162)$ $\hat{x}_2 \sim \mathcal{LN}(3700, 4890)$ $\hat{x}_3 \sim \mathcal{U}(63\ 070, 115\ 600)$ $\hat{x}_4 \sim \mathcal{U}(990, 1110)$ $\hat{x}_5 \sim \mathcal{U}(63.1, 116)$ $\hat{x}_6 \sim \mathcal{U}(700, 820)$ $\hat{x}_7 \sim \mathcal{U}(1120, 1680)$ $\hat{x}_8 \sim \mathcal{U}(9\ 855, 12\ 045)$

7	Steel Column [21] $f(\mathbf{x}) = x_1 - \frac{P}{2x_5x_6} - \frac{x_8PE_b}{x_5x_6x_7(E_b-P)},$ $P = x_2 + x_3 + x_4,$ $E_b = \frac{8\pi^2}{9 \cdot 10^8} x_5x_6x_7^2x_9$	9	$\hat{x}_1 \sim \mathcal{LN}(400, 35)$ $\hat{x}_2 \sim \mathcal{N}(5e5, 5e4)$ $\hat{x}_{3,4} \sim \mathcal{G}(6e5, 9e4)$ $\hat{x}_5 \sim \mathcal{LN}(300, 3),$ $\hat{x}_6 \sim \mathcal{LN}(20, 2)$ $\hat{x}_7 \sim \mathcal{LN}(300, 5)$ $\hat{x}_8 \sim \mathcal{N}(30, 10)$ $\hat{x}_9 \sim \mathcal{W}(2.1e5, 4200)$
8	Sulfur Model [22] $f(\mathbf{x}) = -5.488 \cdot 10^{-9}$ $\cdot x_1^2x_2x_3^2x_4x_5x_6x_7x_8x_9$	9	$\hat{x}_1 \sim \mathcal{LN}(0.76, 0.152)$ $\hat{x}_2 \sim \mathcal{LN}(0.39, 0.039)$ $\hat{x}_3 \sim \mathcal{LN}(0.85, 0.085)$ $\hat{x}_4 \sim \mathcal{LN}(0.3, 0.09)$ $\hat{x}_5 \sim \mathcal{LN}(5.0, 2.0),$ $\hat{x}_6 \sim \mathcal{LN}(1.7, 0.34)$ $\hat{x}_7 \sim \mathcal{LN}(71.0, 10.65)$ $\hat{x}_8 \sim \mathcal{LN}(0.5, 0.25)$ $\hat{x}_9 \sim \mathcal{LN}(5.5, 2.75)$
9	Oakley & O'Hagan [23] $f(\mathbf{x}) = \mathbf{a}_1^T \mathbf{x} + \mathbf{a}_2^T \sin(\mathbf{x})$ $+ \mathbf{a}_3^T \cos(\mathbf{x}) + \mathbf{x}^T \mathbf{M} \mathbf{x},$ $\mathbf{a}_i, \mathbf{M} \text{ according to [23]}$	15	$\hat{x}_i \sim \mathcal{N}(0, 1), \quad i = 1 \dots 15$

*PDF parameters correspond to mean μ and standard deviation σ for normal \mathcal{N} , log-normal \mathcal{LN} , Weibull \mathcal{W} and Gumbel \mathcal{G} distributions and to lower and upper limit for uniform \mathcal{U} distributions.

For each problem with number of input dimensions p , a number of $n = 10p$ training points \mathbf{x}_i is generated by maximin Latin hypercube sampling. This number is chosen based on the recommendation by Loepky et al. [24] for conducting initial experiments. Evaluations of model functions $y_i = f(\mathbf{x}_i)$ are conducted. Surrogate models are built based on described methods in Section 2. The following Gaussian process regression methods are compared: simple kriging, ordinary kriging, universal kriging with linear trend, universal kriging with quadratic trend, universal kriging with transformed linear trend and universal kriging with transformed quadratic trend. A linear trend indicates that linear terms w. r. t. all input parameters are included. A quadratic trend indicates that polynomial terms up to the order of 2 w. r. t. all input parameters are included. A transformed trend indicates that inverse Rosenblatt transformation is applied to the corresponding trend function as demonstrated in Section 2.4. Transformations are only applied to linear and quadratic trends because transformation of a zero (simple kriging) or constant (ordinary kriging) trend would not result in any change.

For all combinations between each surrogate method and each benchmark problem, the demonstrated methods for constructing a surrogate model are applied 10 times for better statistical validity. As there may be multiple optima or unsuited initial values in the hyperparameter optimization, each optimization run is repeated 20 times with randomized initial hyperparameters and the hyperparameter set which yields the lowest validation error is selected. Computation time for one experiment therefore includes all 20 hyperparameter optimization runs. Mean value and standard deviation are determined for the validation errors in all cases.

The surrogate models are validated according to the validation measures in Section 2.7. In this study, a set of $n_{\text{val}} = 1000$ validation points is used for each validation.

Furthermore, all experiments are conducted for the case where the gradient of the log marginal likelihood is not included and where it is included according to Eq. 9 and Eq. 14 for comparison of computation time. For optimization, the L-BFGS-B method is used. In the case without gradient information, the gradient is estimated by 2-point finite difference estimation.

4 Results

The validation errors for all combinations between benchmark functions and surrogate methods, each consisting of the 10 experiments, are illustrated in Fig. 2, including their mean values and standard deviations. The mean values are shown in B. The simple kriging model shows the highest validation errors on average. The ordinary kriging models yield significantly lower validation errors compared to the simple kriging models. The universal kriging models yield even higher prediction accuracy than simple kriging and ordinary kriging. Depending on the problem, linear or quadratic basis functions are superior with regard to the validation error. However, for most cases a quadratic trend leads to smaller errors. The incorporation of transformed basis functions significantly reduces the validation errors compared to the non-transformed case for both linear and quadratic trends in most cases.

In Fig. 3 the effect of incorporating transformed basis functions on the surrogate model is illustrated for the *short column function* (benchmark problem #4). For the purpose of illustration, only cuts through the hypersurface are shown where only one input parameter is changed at a time and the output quantity with respect to each input parameter is shown. The values of the function that is to be predicted as well as the surrogate model are shown in the physical input space and the i. i. d. uniform input space, respectively.

In Fig. 4 the computation time for all combinations between benchmark functions and surrogate methods are shown including all 10 experiments with mean values and standard deviations. The computation time is shown for the cases without and with incorporation of the gradient of the log marginal likelihood for hyperparameter estimation. In all cases, the computation time with incorporation of gradient information can be reduced significantly.

5 Discussion

On average, prediction accuracy increases from simple kriging, to ordinary kriging, to universal kriging, because the surrogate models become more flexible and can better adapt to the problem. The effectiveness of applying certain types of basis functions highly depends on the problem itself. Depending on the relationship between input and output quantities in a model, linear or quadratic basis functions in universal kriging, or both, may lead to an improvement of prediction accuracy. In general, quadratic basis functions offer a more flexible way for the surrogate model to adapt to the data which yields an improvement compared to linear functions. However, polynomials of higher degrees combined with a high number of input dimensions result in a very high number of basis functions, known as the curse of dimensionality, which may lead to overfitting and high generalization errors.

The linear and quadratic basis functions that have been transformed by the Rosenblatt transformation lead to a significant improvement compared to linear and quadratic basis functions without transformation in most cases. This is due to the polynomial basis functions being defined in the physical input parameter space rather than in the transformed input parameter space.

The validation errors are very different between the considered benchmark functions (from RMSE= 0.01 to 0.4 for simple kriging) due to different numbers of dimensions and training points and due to incomparable nonlinearities of the problems. For example, model function values of benchmark function #2 become extremely high for small input parameter x_2 which is not unlikely according to the PDF. Surrogate models may strongly struggle to capture such relationships as can be seen from large validation errors with wide spread, as shown in Fig. 2.

From Fig. 3 it becomes clear that linear or quadratic basis functions without proposed transformation yield a linear or quadratic trend in the i. i. d. uniform input space, respectively. On the other hand, transformed basis functions yield a linear or quadratic trend in the physical input space. As it is generally more natural to assume linear or quadratic trends in the physical parameter space, rather than in the transformed i. i. d. uniform parameter space, the surrogate models with transformed basis functions are generally more accurate approximations of the model function. This can be further emphasized by the fact that the input parameter transformation generally increases the nonlinearities of the problem. Surrogate models with non-transformed basis functions generally tend to systematically deviate from the model function. In particular, the shape of such surrogate models tends to flatten in regions with lower PDF values (close to the margins) and to steepen in region with high PDF values (close to the mean). This artefact does not occur if transformed basis functions are used that account for the input space transformation.

The high validation errors for benchmark problem #9 in the case of universal kriging with a quadratic trend

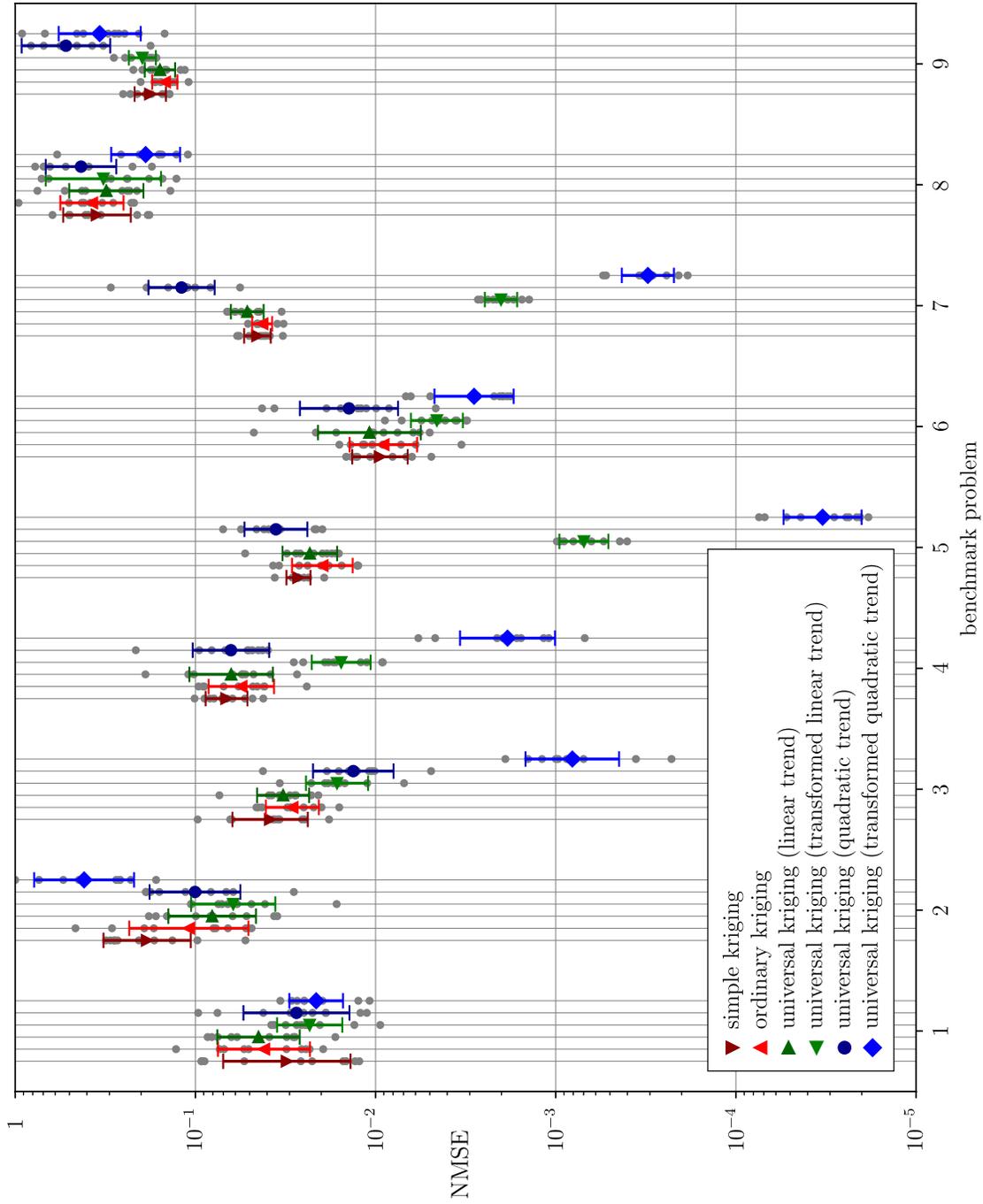


Figure 2: Model validation results: NMSE for all benchmark problems and investigated kriging methods, respectively. Mean values and standard deviations of the 10 experiments are shown for all cases.

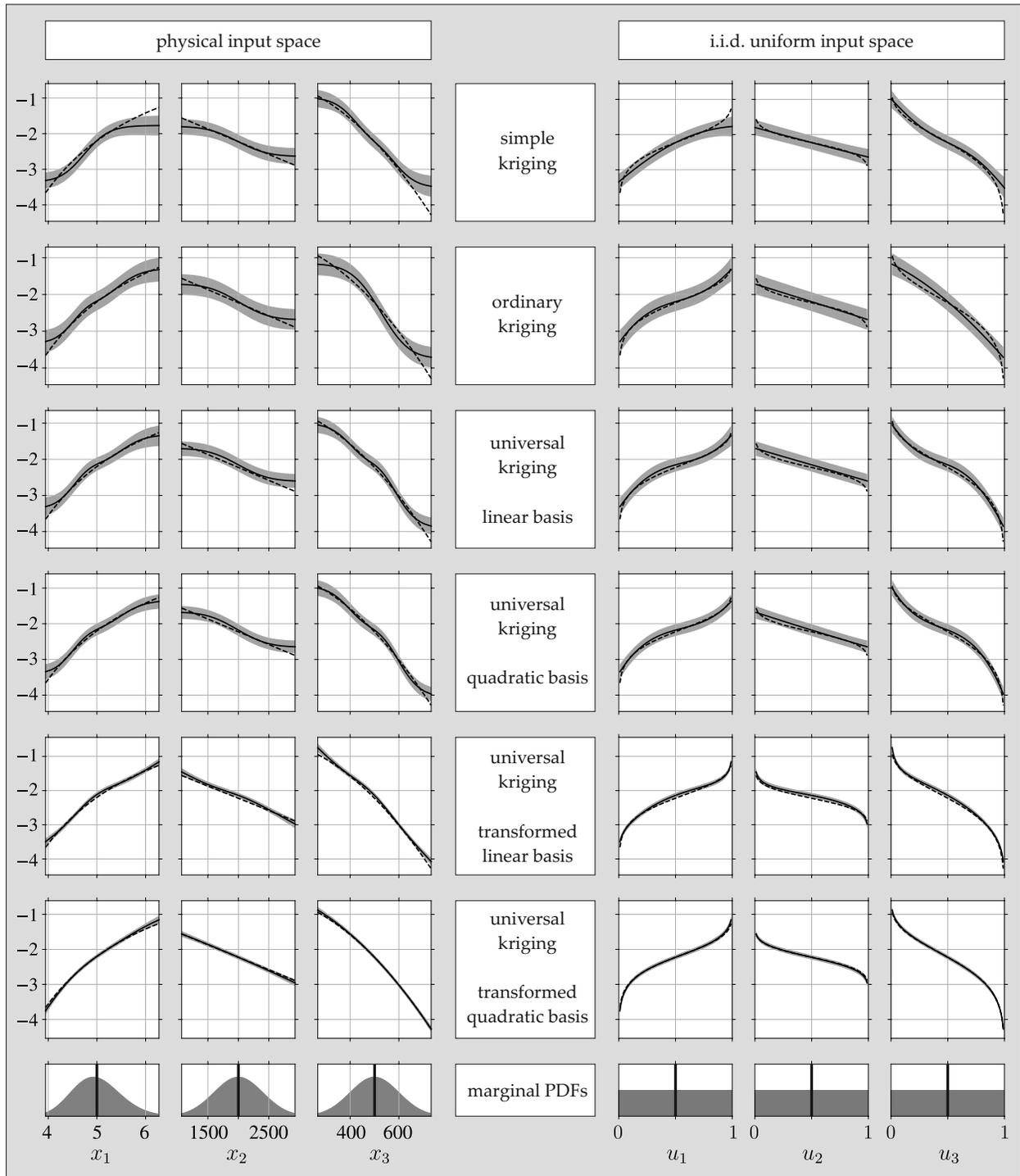


Figure 3: Visualization of benchmark problem #4 (short column function): True function (dashed line) and kriging prediction mean (solid line) and variance (grey shaded area) for investigated kriging methods (rows) and input parameters (columns), in the physical (left) and i. i. d. uniform (right) input space. Only cuts through the hypersurface are shown where one input parameter is changed whereas all other input parameters stay fixed at their mean value (shown in PDF on the bottom).

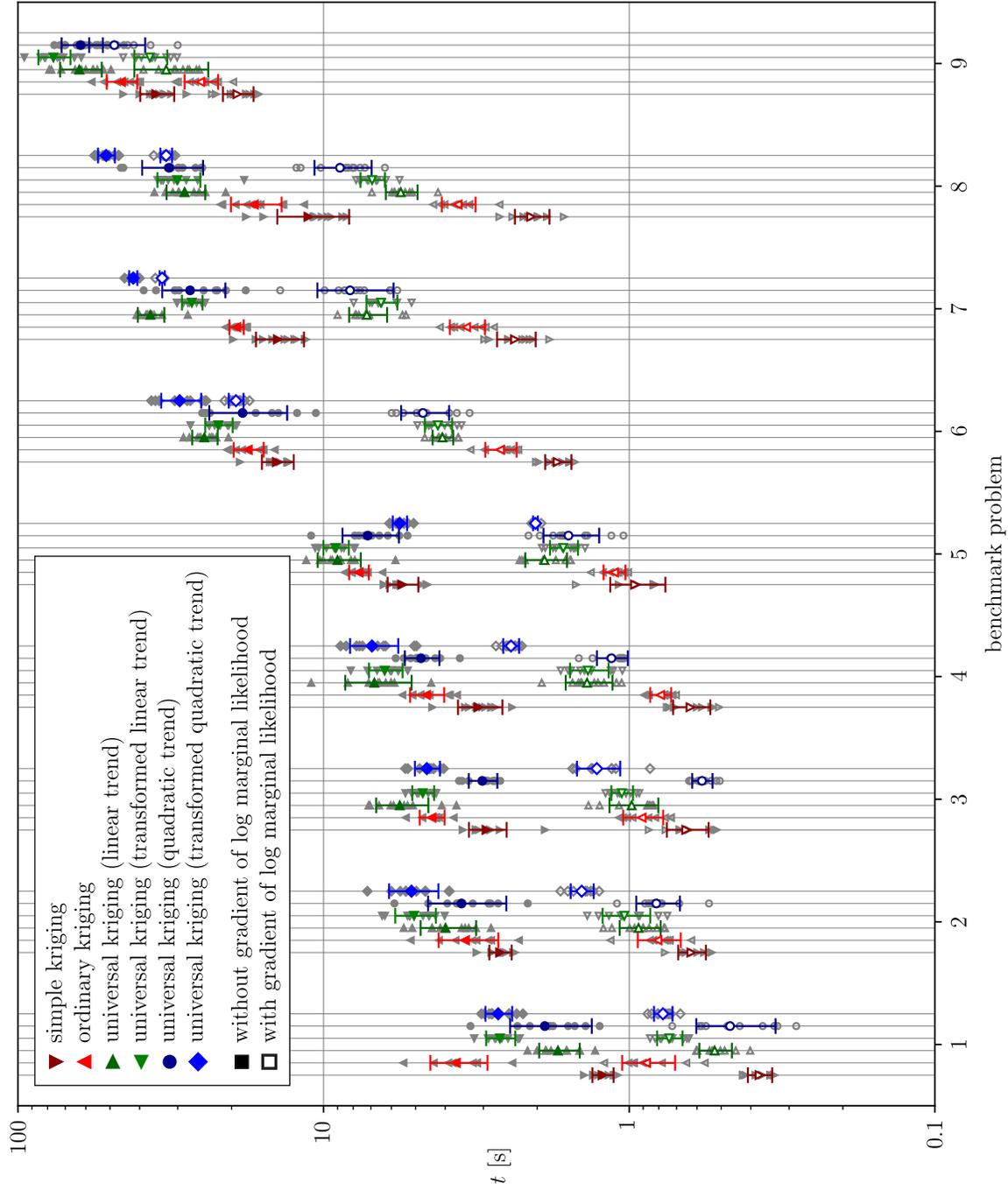


Figure 4: Runtime for obtaining a surrogate model for all benchmark problems and investigated kriging methods, respectively. Results are shown for the cases with (empty markers) and without (filled markers) incorporation of the gradient of the log marginal likelihood for hyperparameter optimization (see Section 2.6). Mean values and standard deviations of the 10 experiments are shown for all cases.

indicate the effect of overfitting. Using quadratic multivariate basis functions for a p -dimensional problem leads to $p(p+1)/2$ quadratic, p linear and 1 constant basis functions in a general case. Here, for $p = 15$ and given $n = 10p$, this means including 136 basis functions for a problem with 150 training points. In the context of polynomial chaos expansion, Sudret [25] argues that for polynomial fitting the number of training points should be at least 2-3 times as high as the number of basis functions in order to avoid overfitting. If this is not the case, a sparse set of basis functions could be determined, e. g. by means of least-angle regression (e. g. Kersaudy et al. [1]).

Incorporating the gradient of the log marginal likelihood for hyperparameter optimization leads to a significant computational speedup. The demonstrated equations in Section 2.6 have therefore been shown to be valid and useful in all cases.

6 Conclusion

An approach has been proposed where input parameter transformations are taken into account for the construction of basis functions for universal kriging. Since surrogate modeling methods are applied in the i. i. d. uniform input space to incorporate space-filling designs in a more meaningful way, basis functions in the universal kriging method also have to be defined in this space. It turned out to be more sensible to define basis functions in the physical input space. The inverse Rosenblatt transformation is applied to transform basis functions from the physical input space into the i. i. d. uniform input space which are then used for universal kriging. The transformed basis functions lead to a significant improvement in prediction accuracy compared to the case of non-transformed basis functions in most cases. Among the benchmark problems, the NMSE is in some cases lower by up to a factor of 10^3 with transformed basis functions. The authors of this paper highly recommend to consider proposed transformations to basis functions for universal kriging whenever non-uniform input parameter distributions are used and the surrogate model is constructed in the i. i. d. uniform parameter space. If the input parameters are uncorrelated, the inverse Rosenblatt transformation reduces to quantile functions of the input parameters which simplifies the transformation significantly. To speed up computation, demonstrated equations for gradient-based hyperparameter optimization can be applied.

The proposed method can be applied to other surrogate modeling methods that incorporate explicit basis functions where the input space is transformed to an i. i. d. uniform input space, e. g. to polynomial regression. In the latter case, the polynomial basis functions are transformed by the inverse Rosenblatt transformation before fitting the regression model to determine the regression coefficients. Eventually, the regression technique is no longer a polynomial regression, because transformations have been applied to the polynomials.

A Matrix identities

Derivatives of the elements of an inverse matrix:

$$\frac{\partial}{\partial x} \mathbf{M}(x)^{-1} = -\mathbf{M}(x)^{-1} \frac{\partial \mathbf{M}(x)}{\partial x} \mathbf{M}(x)^{-1} \quad (16)$$

Derivative of the log determinant of a positive definite symmetric matrix:

$$\frac{\partial}{\partial x} \log(|\mathbf{M}(x)|) = \text{tr} \left(\mathbf{M}(x)^{-1} \frac{\partial \mathbf{M}(x)}{\partial x} \right) \quad (17)$$

Cyclic permutation of matrices in the argument of a trace:

$$\text{tr}(\mathbf{M}_1 \mathbf{M}_2 \mathbf{M}_3) = \text{tr}(\mathbf{M}_2 \mathbf{M}_3 \mathbf{M}_1) = \text{tr}(\mathbf{M}_3 \mathbf{M}_1 \mathbf{M}_2) \quad (18)$$

B Validation errors

References

- [1] P. Kersaudy, B. Sudret, N. Varsier, O. Picon, and J. Wiart. A new surrogate modeling technique combining kriging and polynomial chaos expansions – application to uncertainty analysis in computational dosimetry. *Journal of Computational Physics*, 286:103–117, 2015. doi: 10.1016/j.jcp.2015.01.034.
- [2] V. R. Joseph, Y. Hung, and A. Sudjianto. Blind kriging: A new method for developing metamodels. *Journal of Mechanical Design*, 130, 2008. doi: 10.1115/1.2829873.

Table 2: Model validation results: NMSE for all benchmark problems and investigated kriging methods, respectively, as illustrated in Fig. 2.

#	simple kriging	ordinary kriging	universal kriging with			
			linear trend	quadratic trend	transf. linear trend	transf. quadratic trend
1	4.31E-02	4.99E-02	5.07E-02	2.50E-02	3.52E-02	2.25E-02
2	2.10E-01	1.49E-01	9.38E-02	6.93E-02	1.16E-01	5.10E-01
3	4.35E-02	3.06E-02	3.47E-02	1.76E-02	1.53E-02	9.41E-04
4	6.95E-02	6.04E-02	7.39E-02	1.66E-02	7.36E-02	2.26E-03
5	2.71E-02	2.13E-02	2.48E-02	7.30E-04	3.87E-02	3.76E-05
6	1.00E-02	9.81E-03	1.40E-02	4.85E-03	1.72E-02	3.27E-03
7	4.59E-02	4.29E-02	5.26E-02	2.05E-03	1.31E-01	3.26E-04
8	3.83E-01	4.11E-01	3.50E-01	4.41E-01	4.72E-01	2.13E-01
9	1.82E-01	1.50E-01	1.60E-01	2.00E-01	6.18E-01	3.93E-01

- [3] J. Oakley. Estimating percentiles of uncertain computer code outputs. *Journal of the Royal Statistical Society: Series C (Applied Statistics)*, 53(1):83–93, 2004. doi: 10.1046/j.0035-9254.2003.05044.x.
- [4] P. D. Sampson and P. Guttorp. Nonparametric estimation of nonstationary spatial covariance structure. *Journal of the American Statistical Association*, 87(417):108–119, 1992. doi: 10.2307/2290458.
- [5] A. M. Schmidt and A. O’Hagan. Bayesian inference for non-stationary spatial covariance structure via spatial deformations. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 65(3):743–758, 2003. doi: 10.1111/1467-9868.00413.
- [6] M. N. Gibbs. *Bayesian Gaussian processes for regression and classification*. PhD thesis, University of Cambridge, 1998.
- [7] Y. Xiong, W. Chen, D. Apley, and X. Ding. A non-stationary covariance-based kriging method for metamodelling in engineering design. *International Journal for Numerical Methods in Engineering*, 71(6):733–756, 2007. doi: <https://doi.org/10.1002/nme.1969>.
- [8] L. Lu, C. M. Anderson-Cook, and T. Ahmed. Non-uniform space filling (NUSF) designs. *Journal of Quality Technology*, 53(3):309–330, 2021. doi: 10.1080/00224065.2020.1727801.
- [9] J.-M. Bourinet. *Reliability analysis and optimal design under uncertainty - Focus on adaptive surrogate-based approaches*. Habilitation à diriger des recherches, Université Clermont Auvergne, 2018. tel-01737299v2.
- [10] M. Rosenblatt. Remarks on a multivariate transformation. *The Annals of Mathematical Statistics*, 23(3):470–472, 1952.
- [11] R. E. Melchers and A. T. Beck. *Rosenblatt and Other Transformations*, chapter B, pages 403–414. John Wiley & Sons, Ltd, 2017. ISBN 9781119266105.
- [12] M. D. Morris and T. J. Mitchell. Exploratory designs for computational experiments. *Journal of Statistical Planning and Inference*, 43(3):381–402, 1995. doi: 10.1016/0378-3758(94)00035-T.
- [13] C. E. Rasmussen and C. K. I. Williams. *Gaussian Processes for Machine Learning*. MIT Press, Cambridge, MA, USA, 2006. ISBN 026218253X.
- [14] G. Matheron. *Le krigeage universel*, volume 1. Cahiers du Centre de Morphologie Mathématique, École des Mines de Paris, Fontainebleau, 1969.
- [15] O. Dubrule. Cross validation of kriging in a unique neighborhood. *Journal of the International Association for Mathematical Geology*, 15:687–699, 1983. doi: 10.1007/bf01033232.
- [16] J. Oakley and A. O’Hagan. Bayesian inference for the uncertainty distribution of computer model outputs. *Biometrika*, 89(4):769–784, 2002. doi: 10.1093/biomet/89.4.769.
- [17] M. S. Eldred, H. Agarwal, V. M. Perez, S. F. Wojtkiewicz Jr., and J. E. Renaud. Investigation of reliability method formulations in DAKOTA/UQ. *Structure and Infrastructure Engineering*, 3(3):199–213, 2007. doi: 10.1080/15732470500254618.

- [18] M. D. Webster, M. A. Tatang, and G. J. McRae. Application of probabilistic collocation method for uncertainty analysis of a simple ocean model. *Joint Program Report Series, Report 4*, 1996. URL <http://globalchange.mit.edu/publication/15670>.
- [19] W. V. Harper and S. K. Gupta. Sensitivity/uncertainty analysis of a borehole scenario comparing latin hypercube sampling and deterministic sensitivity approaches. Technical report, Battelle Memorial Inst., Columbus, OH (USA). Office of Nuclear Waste Isolation, 1983.
- [20] M. D. Morris, T. J. Mitchell, and D. Ylvisaker. Bayesian design and analysis of computer experiments: Use of derivatives in surface prediction. *Technometrics*, 35(3):243–255, 1993. doi: 10.1080/00401706.1993.10485320.
- [21] M. Eldred, C. Webster, and P. Constantine. Evaluation of non-intrusive approaches for wiener-askey generalized polynomial chaos. *49th AIAA/ASME/ASCE/AHS/ASC Structures, Structural Dynamics, and Materials Conference*, 2008. doi: 10.2514/6.2008-1892.
- [22] M. Tatang, W. Pan, R. Prinn, and G. McRae. An efficient method for parametric uncertainty analysis of numerical geophysical model. *Journal of Geophysical Research*, 102:21925–21932, 1997. doi: 10.1029/97JD01654.
- [23] J. E. Oakley and A. O’Hagan. Probabilistic sensitivity analysis of complex models: a bayesian approach. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 66(3):751–769, 2004. doi: 10.1111/j.1467-9868.2004.05304.x.
- [24] J. L. Loepky, J. Sacks, and W. J. Welch. Choosing the sample size of a computer experiment: A practical guide. *Technometrics*, 51(4):366–376, 2009. doi: 10.1198/tech.2009.08040.
- [25] B. Sudret. Polynomial chaos expansions and stochastic finite element methods. In Kok-Kwang Phoon and Jianye Ching, editors, *Risk and Reliability in Geotechnical Engineering*, pages 265–300. CRC Press, 2014.