

# Introduction to Uncertainty Quantification

## Module 3.4: Surrogate Models for Uncertainty Propagation

### 1 Introduction

Let us consider again the function  $Y = h(\mathbf{X})$  where  $\mathbf{X}$  is a random vector. In many applications, the function  $h(\mathbf{X})$  represents a complex numerical model of a physical system that comes at high computational cost (e.g. a finite element model). The Monte Carlo methods we discussed in the previous lesson are robust, but even with variance reduction they can be prohibitively expensive. We therefore seek alternatives that allow us to propagate uncertainty.

The most common approach in recent years has been to develop a so-called *surrogate model* (also sometimes referred to as a *metamodel* or *emulator*) that approximates the system  $h(\mathbf{X})$  with a simpler mathematical function,  $\hat{h}(\mathbf{X})$ , that can be learned from data generated by running a small number of full model evaluations and is computationally cheap to evaluate. The fact that the surrogate model is computationally inexpensive allows Monte Carlo simulation and related methods to be performed on the approximate function, and therefore facilitates uncertainty propagation.

Many different types of surrogate models are used for uncertainty propagation. The two most widely used are Gaussian process regression (GPR) with the popular Kriging method being a special case and polynomial chaos expansions (PCE). We will cover these two methods in some detail. There are several other methods that are also used including radial basis functions (RBFs), deep neural networks (DNNs), support vector machines, and other forms of regression such as polynomial response surfaces. We will also briefly touch on some of these methods, but will not cover them in detail here.

### 2 Gaussian Process Regression

Gaussian Process (GP) regression is a machine learning (ML) method that seeks to identify the best fit Gaussian stochastic process to a set of data points. GP regression is a widely-used ML method to construct surrogate models for complex computational models. To introduce GP regression, consider a computational model  $y(\mathbf{x})$  having input vector  $\mathbf{x} = [x_1, x_2, \dots, x_n]^T \in \mathbb{R}^n$ . The model  $y(\mathbf{x})$  is then approximated by

$$\mathcal{Y}(\mathbf{x}, \omega) = \mathcal{F}(\mathbf{x}) + Z(\mathbf{x}, \omega) \quad (1)$$

where  $\mathcal{F}(\cdot)$  is a regression model and  $Z(\cdot)$  is a zero-mean Gaussian random process having sample space indexed by  $\omega \in \Omega$ . We consider that the regression model  $\mathcal{F}(\cdot)$  is defined through a linear combination of basis functions  $\mathbf{f}(\mathbf{x}) = [f_1(\mathbf{x}), f_2(\mathbf{x}), \dots, f_p(\mathbf{x})]^T$  having coefficients  $\boldsymbol{\beta} = [\beta_1, \beta_2, \dots, \beta_p]^T$  as

$$\mathcal{F}(\mathbf{x}) = \boldsymbol{\beta}^T \mathbf{f}(\mathbf{x}) \quad (2)$$

This regression model represents the mean value or the *trend* in the model and is often modeled with polynomial functions – see Section 2.3. The Gaussian random process  $Z(\cdot)$  is considered to have zero mean and covariance:

$$\mathbb{E}[Z(\mathbf{x}_1)Z(\mathbf{x}_2)] = \sigma_z^2 \mathcal{R}(\mathbf{x}_1, \mathbf{x}_2; \boldsymbol{\theta}) \quad (3)$$

where  $\mathcal{R}(\cdot)$  is the autocorrelation function having hyperparameters  $\boldsymbol{\theta}$  and  $\sigma_z$  is the variance of the process. The regression function is modeled using a positive definite kernel function  $k(\mathbf{x}_1, \mathbf{x}_2)$ , several examples of which are provided in Section 2.3.

The goal of GP regression is then to learn the hyperparameters  $\boldsymbol{\theta}$  and regression coefficients  $\boldsymbol{\beta}$  such that the model provides the best fit to a set of training data corresponding to input realizations  $\mathbf{X} = \{\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(n)}\}$  and corresponding model outputs  $\mathbf{Y} = \{y^{(1)}, y^{(2)}, \dots, y^{(n)}\}$  evaluated by  $y^{(i)} = y(\mathbf{x}^{(i)})$ . We will consider two distinct cases. In the first case, the data points are assumed to be exact and the GP model is designed to interpolate between the training points. This is referred to as Kriging [cite]. In the second case, the data points are assumed to contain some noise and the model is fit as a regressor. Additional details on both cases and more information on the method can be found in the textbook by Rasmussen and Williams [1].

## 2.1 Kriging for Interpolation

The first case we'll consider is the case where the GP model serves as an interpolant between the existing data points, which we refer to as Kriging. Kriging was originally developed to interpolate geospatial data by Krige [2] and the methods was later formalized in this setting by Matheron [3]. Kriging was first used for surrogate modeling in the context of design of computer experiments by Sacks et al. [4] in 1989.

The Kriging method takes the input-output data,  $\mathbf{X}, \mathbf{Y}$ , and builds a probability model to interpolate the data by assuming that the interpolant is a Gaussian process. That is, the Gaussian assumption states that the joint probability distribution between the interpolant (predictor)  $\mathcal{Y}(\mathbf{x})$  and the data  $\mathbf{Y}$  follows the joint normal distribution with

$$\begin{Bmatrix} \mathcal{Y}(\mathbf{x}) \\ \mathbf{Y} \end{Bmatrix} \sim N \left( \begin{Bmatrix} \mathbf{f}(\mathbf{x})^T \boldsymbol{\beta} \\ \mathbf{F} \boldsymbol{\beta} \end{Bmatrix}, \sigma_z^2 \begin{Bmatrix} 1 & \mathbf{r}(\mathbf{x})^T \\ \mathbf{r}(\mathbf{x}) & \mathbf{R} \end{Bmatrix} \right) \quad (4)$$

where

- $\mathbf{F}$  is the matrix of basis function evaluations at the training points given by  $F_{ij} = f_j(\mathbf{x}^{(i)})$ ,  $i = 1, \dots, n$ ,  $j = 1, \dots, p$
- $\mathbf{r}(\mathbf{x})$  is the vector of correlations between the prediction point  $\mathbf{x}$  and the training points  $\mathbf{x}^{(i)}$  given by  $r_i = \mathcal{R}(\mathbf{x}, \mathbf{x}^{(i)}; \boldsymbol{\theta})$ ,  $i = 1, \dots, n$
- $\mathbf{R}$  is the correlation matrix of points in the training set given by  $R_{ij} = \mathcal{R}(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}; \boldsymbol{\theta})$ ,  $i, j = 1, \dots, n$

Given the joint normal probability density in Eq. (4), the conditional distribution of the prediction at point  $\mathbf{x}$  is also normal

$$\mathcal{Y}(\mathbf{x}) \mid \mathbf{X}, \mathbf{Y} \sim \mathcal{N}(\hat{y}(\mathbf{x}), \sigma_{\hat{y}}^2(\mathbf{x})), \quad (5)$$

having mean

$$\hat{y}(\mathbf{x}) = \mathbf{f}(\mathbf{x})^T \boldsymbol{\beta} + \mathbf{r}(\mathbf{x})^T \mathbf{R}^{-1}(\mathbf{Y} - \mathbf{F} \boldsymbol{\beta}) \quad (6)$$

and variance

$$\sigma_{\hat{y}}^2(\mathbf{x}) = \sigma_z^2 (1 - \mathbf{r}(\mathbf{x})^T \mathbf{R}^{-1} \mathbf{r}(\mathbf{x}) + \mathbf{t}(\mathbf{x})^T (\mathbf{F}^T \mathbf{R}^{-1} \mathbf{F})^{-1} \mathbf{t}(\mathbf{x})) \quad (7)$$

where

$$\mathbf{t}(\mathbf{x}) = \mathbf{F}^T \mathbf{R}^{-1} \mathbf{r}(\mathbf{x}) - \mathbf{f}(\mathbf{x}) \quad (8)$$

Eq. (5) is referred to as the *Kriging predictor*. The prediction takes the form of a normal random variable having mean in Eq. (6) and variance given in Eq. (7). We specifically note that the variance of the predictor in Eq. (7) is zero at the training points and therefore the Kriging predictor serves to *interpolate* between the training data.

To fit the GP model and obtain the Kriging predictor, it is necessary to determine the best fit regression parameters  $\beta$ , variance  $\sigma_z$ , and GP hyperparameters  $\theta$ . This can be done in several ways. Here, we will describe the two most commonly applied methods – maximum likelihood estimation (MLE) and cross-validation (CV).

### Maximum Likelihood Estimation

Using MLE, the predictor can be determined by identifying the hyperparameters  $\theta$ , variance  $\sigma_z$ , and the regression coefficients  $\beta$  such that the likelihood of the observations  $\mathbf{Y}$  is maximized. Since  $\mathbf{Y}$  is assumed to be drawn from a multivariate normal distribution, the likelihood function can be given by

$$\mathcal{L}(\theta, \sigma_z, \beta | \mathbf{Y}) = \frac{(\det \mathbf{R})^{-1/2}}{(2\pi\sigma_z^2)^{n/2}} \exp \left[ -\frac{1}{2\sigma_z^2} (\mathbf{Y} - \mathbf{F}\beta)^T \mathbf{R}^{-1} (\mathbf{Y} - \mathbf{F}\beta) \right] \quad (9)$$

Maximizing this likelihood function, we can obtain the regression coefficients  $\beta$  and variance  $\sigma_z^2$  analytically as

$$\hat{\beta} = \beta(\theta) = (\mathbf{F}^T \mathbf{R}^{-1} \mathbf{F})^{-1} \mathbf{F}^T \mathbf{R}^{-1} \mathbf{Y} \quad (10)$$

and

$$\hat{\sigma}_z^2 = \sigma_z^2(\theta) = \frac{1}{n} (\mathbf{Y} - \mathbf{F}\beta)^T \mathbf{R}^{-1} (\mathbf{Y} - \mathbf{F}\beta) \quad (11)$$

The proof can be found in [5].

We can then obtain the hyperparameters  $\theta$  by solving the following optimization problem

$$\hat{\theta} = \arg \min_{\theta} (-\log(\mathcal{L}(\theta | \mathbf{Y}))) \quad (12)$$

Applying the relations in Eqs. (9), (10), and (11) this optimization problem can be expressed

$$\hat{\theta} = \arg \min_{\theta} \frac{1}{2} (\log(\det(\mathbf{R})) + n \log(2\pi\sigma_z^2) + n) \quad (13)$$

Solution of this optimization problem can be performed with your choice of optimizer. Performing this optimization is outside the scope of this document.

### Cross-Validation

An alternate approach to obtain the GP parameters is to perform cross-validation (CV). In the process of *K-fold cross-validation*, the training data  $\mathcal{D} = [\mathbf{X}, \mathbf{Y}]$  is split into  $K$  mutually exclusive and collectively complete subsets  $\mathcal{D}_k, k = 1, \dots, K$  such that

$$\begin{aligned} \mathcal{D}_i \cap \mathcal{D}_j &= \emptyset, \quad \forall i \neq j \quad \text{and} \\ \cup_{k=1}^K \mathcal{D}_k &= \mathcal{D} \end{aligned} \quad (14)$$

Then,  $K$  predictions are obtained by estimating the model parameters such that, in the  $k^{th}$  prediction all of the data are used for training except the  $k^{th}$  subset. The model is then used to predict the values from the  $k^{th}$  subset. The special case of *leave-one-out cross-validation* (LOO-CV) sets  $K = n$  so that each subset contains a single training point.

The CV error for the  $k^{th}$  subset is given by

$$\epsilon_{CV}^k = \sum_{(x^{(i)}, y^{(i)}) \in \mathcal{D}_k} \left( y^{(i)} - \hat{y}_{\mathcal{D} \setminus \mathcal{D}_k} \left( x^{(i)} \mid \beta, \sigma_z^2, \theta \right) \right)^2 \quad (15)$$

where  $\hat{y}_{\mathcal{D} \setminus \mathcal{D}_k}(\cdot)$  denotes the mean Kriging predictor using all training data except those in  $\mathcal{D}_k$ . The total CV error is then given by averaging the CV errors across all  $K$  subsets

$$\epsilon_{CV} = \frac{1}{K} \sum_{k=1}^K \epsilon_{CV}^k \quad (16)$$

To train the model using CV, we then aim to identify the parameters  $\beta$ ,  $\sigma_z^2$ , and  $\theta$  that minimize the total CV error where we can again take advantage of the analytical expression  $\beta$  in Eq. (10). We further note that  $\hat{y}(\mathbf{x})$  does not depend on  $\sigma_z^2$ . Therefore, the minimization reduces to

$$\hat{\theta} = \arg \min_{\theta} \epsilon_{CV}(\theta \mid \mathbf{Y}) \quad (17)$$

The variance is then computed as

$$\hat{\sigma}_z^2 = \hat{\sigma}_z^2(\hat{\theta}) = \frac{1}{K} \sum_{k=1}^K \sum_{(\mathbf{x}^{(i)}, y^{(i)}) \in \mathcal{D}_k} \frac{\left( y^{(i)} - \hat{y}_{\mathcal{D} \setminus \mathcal{D}_k}(\mathbf{x}^{(i)} \mid \hat{\theta}) \right)^2}{\sigma_{\hat{y}, \mathcal{D} \setminus \mathcal{D}_k}^2(\mathbf{x}^{(i)} \mid \hat{\theta}) / \sigma_z^2} \quad (18)$$

where  $\sigma_{\hat{y}, \mathcal{D} \setminus \mathcal{D}_k}^2(\cdot)$  is the variance of the GP predictor given from Eq. (7) estimated from all data except those in  $\mathcal{D}_k$ . A more detailed discussion of these estimates can be found in [6].

## 2.2 Regression of Noisy Functions

When using a GP for prediction of noisy function, we assume that the noisy function can be expressed as

$$y = h(\mathbf{X}) + \varepsilon \quad (19)$$

where  $\varepsilon$  is an additive noise this is assumed to follow a zero-mean Gaussian distribution  $\varepsilon \sim N(0, \Sigma_n)$  with covariance matrix  $\Sigma_n$ . Depending on the nature of the noise, three cases arise:

- *Homoscedastic noise*:  $\Sigma_n = \sigma_n^2 \mathbf{I}$  where  $\mathbf{I}$  is the identity matrix. Here, the variance is independent and constant across all observations.
- *Independent Heteroscedastic noise (Heterogeneous)*:  $\Sigma_n = \text{diag}(\sigma_n^2)$ . Here the noise variances differ at each point, but remain uncorrelated.
- *General Heteroscedastic noise*: Here  $\Sigma_n$  can be any valid covariance matrix allowing for differing noise at each point and correlations between the components at each point.

With the introduction of noise, the joint Gaussian distribution from Eq. (4) becomes

$$\begin{Bmatrix} \mathcal{Y}(\mathbf{x}) \\ \mathbf{Y} \end{Bmatrix} \sim N \left( \begin{Bmatrix} \mathbf{f}(\mathbf{x})^T \beta \\ \mathbf{F} \beta \end{Bmatrix}, \begin{Bmatrix} \sigma_z^2 & \sigma_z^2 \mathbf{r}(\mathbf{x})^T \\ \sigma_z^2 \mathbf{r}(\mathbf{x}) & \sigma_z^2 \mathbf{R} + \Sigma_n \end{Bmatrix} \right) \quad (20)$$

where the terms are defined exactly as above, except we now introduce the noise covariance  $\Sigma_n$ . This yields the conditional distribution of the predictor having the same Gaussian form

$$\mathcal{Y}(\mathbf{x}) \mid \mathbf{X}, \mathbf{Y} \sim \mathcal{N}(\hat{y}(\mathbf{x}), \sigma_{\hat{y}}^2(\mathbf{x})), \quad (21)$$

but with mean and variance modified to account for the noise as follows. Introducing the covariance matrix  $\mathbf{C} = \sigma_z^2 \mathbf{R} + \Sigma_n$  and the cross-covariance  $\mathbf{c}(\mathbf{x}) = \sigma_z^2 \mathbf{r}(\mathbf{x})$ , the mean of the predictor is given by

$$\hat{y}(\mathbf{x}) = \mathbf{f}(\mathbf{x})^T \beta + \mathbf{c}(\mathbf{x})^T \mathbf{C}^{-1} (\mathbf{Y} - \mathbf{F} \beta) \quad (22)$$

and variance is given by

$$\sigma_y^2(\mathbf{x}) = (\sigma_z^2 - \mathbf{c}(\mathbf{x})^T \mathbf{C}^{-1} \mathbf{c}(\mathbf{x}) + \mathbf{u}_c(\mathbf{x})^T (\mathbf{F}^T \mathbf{C}^{-1} \mathbf{F})^{-1} \mathbf{u}_c(\mathbf{x})) \quad (23)$$

where

$$\mathbf{u}_c(\mathbf{x}) = \mathbf{F}^T \mathbf{C}^{-1} \mathbf{c}(\mathbf{x}) - \mathbf{f}(\mathbf{x}) \quad (24)$$

One can clearly see that these equations are identical to the noise-free Eqs. (6)-(8) but accounting for the full covariance with noise.

In the following, we will deal only with the *homoscedastic* case for which  $\mathbf{C} = \sigma_z^2 \mathbf{R} + \sigma_n^2 \mathbf{I}$ . The total variance is therefore given by

$$\sigma_t^2 = \sigma_z^2 + \sigma_n^2 \quad (25)$$

Let us define the ratio of noise to total variance by

$$\tau = \frac{\sigma_n^2}{\sigma_t^2} \quad (26)$$

and perform the following factorizations

$$\tilde{\mathbf{r}}(\mathbf{x}) = (1 - \tau) \mathbf{r}(\mathbf{x}) \quad (27)$$

and

$$\tilde{\mathbf{R}} = (1 - \tau) \mathbf{R} + \tau \mathbf{I}. \quad (28)$$

We can then write the joint Gaussian distribution as

$$\left\{ \begin{array}{c} \mathcal{Y}(\mathbf{x}) \\ \mathbf{Y} \end{array} \right\} \sim N \left( \left\{ \begin{array}{c} \mathbf{f}(\mathbf{x})^T \boldsymbol{\beta} \\ \mathbf{F} \boldsymbol{\beta} \end{array} \right\}, \sigma_t^2 \left\{ \begin{array}{cc} 1 & \tilde{\mathbf{r}}^T(\mathbf{x}) \\ \tilde{\mathbf{r}}(\mathbf{x}) & \tilde{\mathbf{R}} \end{array} \right\} \right) \quad (29)$$

conditional distribution of the predictor as

$$\mathcal{Y}(\mathbf{x}) \mid \mathbf{X}, \mathbf{Y} \sim \mathcal{N}(\hat{y}(\mathbf{x}), \sigma_y^2(\mathbf{x})), \quad (30)$$

having mean

$$\hat{y}(\mathbf{x}) = \mathbf{f}(\mathbf{x})^T \boldsymbol{\beta} + \tilde{\mathbf{r}}(\mathbf{x})^T \tilde{\mathbf{R}}^{-1} (\mathbf{Y} - \mathbf{F} \boldsymbol{\beta}) \quad (31)$$

and variance

$$\sigma_y^2(\mathbf{x}) = \sigma_t^2 \left( 1 - \tilde{\mathbf{r}}(\mathbf{x})^T \tilde{\mathbf{R}}^{-1} \tilde{\mathbf{r}}(\mathbf{x}) + \mathbf{u}(\mathbf{x})^T (\mathbf{F}^T \tilde{\mathbf{R}}^{-1} \mathbf{F})^{-1} \mathbf{u}(\mathbf{x}) \right) \quad (32)$$

where

$$\mathbf{u}(\mathbf{x}) = \mathbf{F}^T \tilde{\mathbf{R}}^{-1} \tilde{\mathbf{r}}(\mathbf{x}) - \mathbf{f}(\mathbf{x}) \quad (33)$$

In the noisy case described here, the Gaussian predictor does not possess zero variance, and thus it does not provide an exact estimator at the existing data points. It therefore serves as a regressor, rather than an interpolant. Training the GP model follows the same process as in the noiseless case above where again we can determine the model variance  $\sigma_z^2$ , regression parameters  $\boldsymbol{\beta}$ , correlation function hyperparameters  $\boldsymbol{\theta}$  using either MLE or CV, but we now must also learn the noise variance  $\sigma_n^2$ . We review these two approaches next – again for the homoscedastic case.

## Maximum Likelihood Estimation

Once again, since  $\mathbf{Y}$  is assumed to be drawn from a multivariate normal distribution, the likelihood function can be given by

$$\mathcal{L}(\boldsymbol{\theta}, \sigma_z, \boldsymbol{\beta} | \mathbf{Y}) = \frac{(\det \tilde{\mathbf{R}})^{-1/2}}{(2\pi\sigma_t^2)^{n/2}} \exp \left[ -\frac{1}{2\sigma_t^2} (\mathbf{Y} - \mathbf{F}\boldsymbol{\beta})^T \tilde{\mathbf{R}}^{-1} (\mathbf{Y} - \mathbf{F}\boldsymbol{\beta}) \right] \quad (34)$$

which is identical to Eq. (9) except with  $\sigma_z^2$  replaced with  $\sigma_t^2$  and  $\mathbf{R}$  replaced with  $\tilde{\mathbf{R}}$ . Maximizing this likelihood function, we can obtain the regression coefficients  $\boldsymbol{\beta}$  and total variance  $\sigma_t^2$  analytically as

$$\hat{\boldsymbol{\beta}} = \boldsymbol{\beta}(\boldsymbol{\theta}) = \left( \mathbf{F}^T \tilde{\mathbf{R}}^{-1} \mathbf{F} \right)^{-1} \mathbf{F}^T \tilde{\mathbf{R}}^{-1} \mathbf{Y} \quad (35)$$

and

$$\hat{\sigma}_t^2 = \sigma_t^2(\boldsymbol{\theta}) = \frac{1}{n} (\mathbf{Y} - \mathbf{F}\boldsymbol{\beta})^T \tilde{\mathbf{R}}^{-1} (\mathbf{Y} - \mathbf{F}\boldsymbol{\beta}) \quad (36)$$

We can then obtain the hyperparameters  $\boldsymbol{\theta}$  by solving the following optimization problem

$$\hat{\boldsymbol{\theta}}, \hat{\tau} = \arg \min_{\boldsymbol{\theta}, \tau} \frac{1}{2} \left( \log(\det(\tilde{\mathbf{R}})) + n \log(2\pi\sigma_t^2) + n \right) \quad (37)$$

where we now also learn the additional parameter  $\tau \in (0, 1)$  that expresses the fraction of variance attributed to noise. Solution of this optimization problem can be performed with your choice of optimizer. Performing this optimization is outside the scope of this document.

### Cross-Validation

Using CV, we recognize that the CV error now depends on the additional parameter  $\tau$  because we use the revised estimator  $\hat{y}_{\mathcal{D} \setminus \mathcal{D}_k}(\mathbf{x} | \boldsymbol{\theta}, \tau)$  from Eq. (31) again estimated from all data except those in  $\mathcal{D}_k$ . Therefore, the minimization can be expressed as

$$\hat{\boldsymbol{\theta}}, \tau = \arg \min_{\boldsymbol{\theta}, \tau} \epsilon_{CV}(\boldsymbol{\theta}, \tau | \mathbf{Y}) \quad (38)$$

The total variance is then computed as

$$\hat{\sigma}_t^2 = \hat{\sigma}_t^2(\hat{\boldsymbol{\theta}}, \hat{\tau}) = \frac{1}{K} \sum_{k=1}^K \sum_{(\mathbf{x}^{(i)}, y^{(i)}) \in \mathcal{D}_k} \frac{\left( y^{(i)} - \hat{y}_{\mathcal{D} \setminus \mathcal{D}_k}(\mathbf{x}^{(i)} | \hat{\boldsymbol{\theta}}, \hat{\tau}) \right)^2}{\sigma_{\hat{y}, \mathcal{D} \setminus \mathcal{D}_k}^2(\mathbf{x}^{(i)} | \hat{\boldsymbol{\theta}}, \hat{\tau}) / \sigma_t^2} \quad (39)$$

where  $\sigma_{\hat{y}, \mathcal{D} \setminus \mathcal{D}_k}^2(\cdot)$  is the variance of the GP predictor given from Eq. (32) estimated from all data except those in  $\mathcal{D}_k$ .

## 2.3 Common Kernels and Regression Functions

In this section, we briefly review the commonly used correlation functions (kernels) and regression models in Gaussian process regression.

### Correlation Functions (Kernels)

The correlation function  $\mathcal{R}(\mathbf{x}_1, \mathbf{x}_2; \boldsymbol{\theta})$  is typically expressed through a symmetric positive semi-definite kernel used to define the functional form of the linear dependence between points in the GP. The form of this function can have profound impact on the resulting GP and its fit to the data. For example, the form of the kernel encodes assumptions such as continuity, differentiability, etc. and its hyperparameters  $\boldsymbol{\theta}$  typically dictate a length-scale that dictates how strongly correlated points in the GP are as a function of

the distance between them. In multi-dimensional cases having  $M$  dimensions the hyperparameter vector is typically also  $M$ -dimensional (i.e.  $\boldsymbol{\theta} \in \mathbb{R}^M$ ) such that there is one parameter (length-scale) per dimension.

Usually the multi-dimensional correlation function will be expressed as a function of one-dimensional kernels  $K(\mathbf{x}_1, \mathbf{x}_2; \boldsymbol{\theta})$  using either the *ellipsoidal form* where

$$\mathcal{R}(\mathbf{x}_1, \mathbf{x}_2; \boldsymbol{\theta}) = K(h) \quad (40)$$

where

$$h = \left[ \sum_{i=1}^M \left( \frac{x_{1i} - x_{2i}}{\theta_i} \right)^2 \right]^{0.5} \quad (41)$$

or in *separable form* where

$$\mathcal{R}(\mathbf{x}_1, \mathbf{x}_2; \boldsymbol{\theta}) = \prod_{i=1}^M K(x_{1i}, x_{2i}; \theta_i) \quad (42)$$

A special case of these correlation functions is the *isotropic* correlation function in which  $\theta_i = \theta, \forall i$ . That is, the correlation function has the same hyperparameter in each dimension. and the above expressions simplify as

$$\mathcal{R}(\mathbf{x}_1, \mathbf{x}_2; \theta) = K(h) \quad (43)$$

where

$$h = \frac{1}{\theta} \left[ \sum_{i=1}^M (x_{1i} - x_{2i})^2 \right]^{0.5} \quad (44)$$

or in *separable form* where

$$\mathcal{R}(\mathbf{x}_1, \mathbf{x}_2; \theta) = \prod_{i=1}^M K(x_{1i}, x_{2i}; \theta) \quad (45)$$

and the above optimizations reduce to identifying a scalar hyperparameter  $\theta$  rather than the full  $M$ -dimensional vector  $\boldsymbol{\theta}$ .

In all of the above cases (ellipsoidal, separable, isotropic), the kernel is typically selected from one of the following popular functional forms (all expressed in scalar form for integration in the ellipsoidal or separable forms):

- *Linear*: The linear kernel takes the following form:

$$K(x_1, x_2; \theta) = \max \left( 0, 1 - \frac{|x_1 - x_2|}{\theta} \right) \quad (46)$$

Here, the length-scale parameter dictates the slope of the function and the function implies a  $\mathcal{C}^0$  continuous but non-differentiable Gaussian process. Examples of this kernel are shown in 1a.

- *Exponential*: The exponential kernel takes the following form:

$$K(x_1, x_2; \theta) = \exp \left[ -\frac{|x_1 - x_2|}{\theta} \right] \quad (47)$$

Here, the length-scale parameter dictates the rate of exponential decay of the function and the function again implies a  $\mathcal{C}^0$  continuous but non-differentiable Gaussian process. Examples of this kernel are shown in 1b.

- *Gaussian*: (aka squared exponential or radial basis function kernel) The Gaussian kernel takes the following form:

$$K(x_1, x_2; \theta) = \exp \left[ -\frac{1}{2} \left( \frac{|x_1 - x_2|}{\theta} \right)^2 \right] \quad (48)$$

Here, the length-scale parameter correspond to the standard deviation of the Gaussian, which also dictate the rate of decay of the squared exponential. Unlike the exponential and linear, this kernel implies a continuous and infinitely differentiable Gaussian process – which is often desirable for practical applications and for encouraging smoothness of the resulting surrogate. Examples of this kernel are shown in 1c.

- *Matérn*: The Matérn kernel takes the following form:

$$K(x_1, x_2; \theta) = \frac{1}{2^{\nu-1}\Gamma(\nu)} \left( 2\sqrt{\nu} \frac{|x_1 - x_2|}{\theta} \right)^{\nu} \mathcal{K}_{\nu} \left( 2\sqrt{\nu} \frac{|x_1 - x_2|}{\theta} \right) \quad (49)$$

where  $\nu \geq 1/2$  is a shape parameter,  $\Gamma(\cdot)$  is the Gamma function, and  $\mathcal{K}_{\nu}$  is the modified Bessel function of the second kind. The resulting GP is  $\lceil \nu - 1 \rceil$  times differentiable where  $\lceil \cdot \rceil$  denotes the ceiling function. For  $\nu = 1/2$ , the Matérn corresponds to the exponential kernel and as  $\nu \rightarrow \infty$  it approaches the Gaussian kernel. The most commonly applied cases are the  $\nu = 3/2$  and  $\nu = 5/2$  cases which can be written as

$$K(x_1, x_2; \theta, \nu = 3/2) = \left( 1 + \sqrt{3} \frac{|x_1 - x_2|}{\theta} \right) \exp \left[ -\sqrt{3} \frac{|x_1 - x_2|}{\theta} \right] \quad (50)$$

and

$$K(x_1, x_2; \theta, \nu = 5/2) = \left( 1 + \sqrt{5} \frac{|x_1 - x_2|}{\theta} + \frac{5}{3} \frac{|x_1 - x_2|^2}{\theta^2} \right) \exp \left[ -\sqrt{5} \frac{|x_1 - x_2|}{\theta} \right] \quad (51)$$

Examples of these kernels are shown in Fig 1d and 1e.

## Regression Functions

The regression function defines the underlying trend of the model. It is typically defined, as shown in Eq. (2) by a set of basis functions  $\mathbf{f}(\mathbf{x})$  and corresponding coefficients  $\boldsymbol{\beta}$ . These are typically defined through simple polynomials and the following summarizes the most commonly used regression model types.

- *Simple Kriging*: In simple Kriging, the basis functions are defined arbitrarily but the coefficients are pre-determined and are all set equal to 1, i.e.  $\beta_j = 1, \forall j$ . That is

$$\boldsymbol{\beta}^T \mathbf{f}(\mathbf{x}) = \sum_{j=1}^P f_j(\mathbf{x}) \quad (52)$$

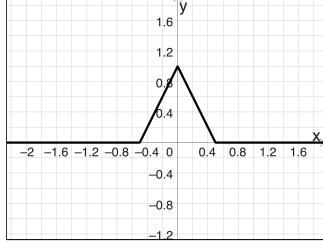
- *Ordinary Kriging*: In ordinary Kriging, the regression is simply set equal to a constant. That is  $f_0(\mathbf{x}) = 1$  and  $f_j(\mathbf{x}) = 0, \forall j > 0$  thus

$$\boldsymbol{\beta}^T \mathbf{f}(\mathbf{x}) = \beta_0 f_0(\mathbf{x}) = \beta_0 \quad (53)$$

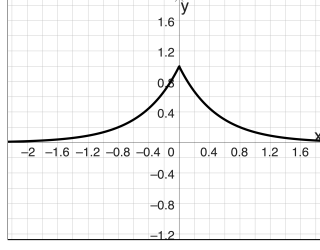
- *Universal Kriging*: Universal Kriging is the most general form of the Kriging regression model and generally prescribes the trend in terms of an arbitrary set of functions, usually specified as polynomials in the form

$$\boldsymbol{\beta}^T \mathbf{f}(\mathbf{x}) = \sum_{j=0}^P \beta_j f_j(\mathbf{x}) \quad (54)$$

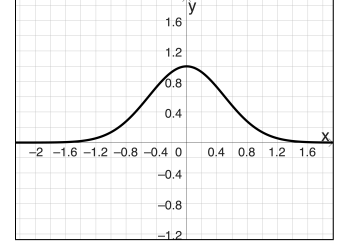




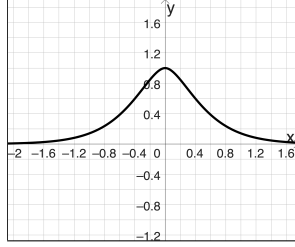
(a) A linear kernel



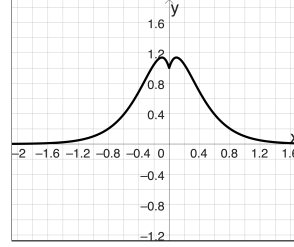
(b) An exponential kernel



(c) A Gaussian kernel



(d) Matérn kernel with  $\nu = 3/2$



(e) Matérn kernel with  $\nu = 5/2$

Figure 1: Several examples of common kernels.

Some of the common forms for universal Kriging include the linear and quadratic models given by

$$\boldsymbol{\beta}^T \mathbf{f}(\mathbf{x}) = \beta_0 + \sum_{i=1}^M \beta_i x_i \quad (55)$$

and

$$\boldsymbol{\beta}^T \mathbf{f}(\mathbf{x}) = \beta_0 + \sum_{i=1}^M \beta_i x_i + \sum_{i=1}^M \sum_{j=1}^M \beta_{ij} x_i x_j, \quad (56)$$

respectively.

Another novel form of Kriging trend uses a set of orthogonal polynomial basis functions, corresponding to a polynomial chaos expansion (PCE) – a commonly used form of surrogate model for uncertainty quantification discussed next. This PCE Kriging method is described in detail in [7].

### 3 Polynomial Chaos Expansions

The polynomial chaos expansion (PCE) aims to develop a surrogate model for the transformation  $y = h(\mathbf{x})$  by constructing an approximation through a series expansion in orthogonal polynomials. The following will build the basic concept and show how it is used for uncertainty quantification.

#### 3.1 General Setting

Consider an  $m$ -dimensional random vector  $\mathbf{X}$  having independent components defined through the joint probability density function  $f_{\mathbf{X}}(\mathbf{x}) = \prod_{i=1}^m f_{X_i}(x_i)$ . The PCE approximates the transformation  $Y = h(\mathbf{X})$  as

$$Y = h(\mathbf{X}) = \sum_{\boldsymbol{\alpha} \in \mathbb{N}^m} \beta_{\boldsymbol{\alpha}} \Psi_{\boldsymbol{\alpha}}(\boldsymbol{\xi}) \quad (57)$$

where  $\beta_{\alpha}$  are a set of deterministic coefficients,  $\alpha$  is a set of integers called the *multi-index* (more on this later), and  $\Psi_{\alpha}(\xi)$  are multivariate orthogonal polynomials expressed as a function of a set of independent random variables  $\xi$  referred to as the *germ*. These multivariate polynomials are orthogonal with respect to the probability density of the germ  $p_{\xi}(\xi)$  as

$$\langle \psi_j, \psi_k \rangle = \int \psi_j(\xi) \psi_k(\xi) p_{\xi}(\xi) d\xi = \delta_{jk} \quad (58)$$

where  $\delta_{jk}$  is the Kronecker delta with  $\delta_{jk} = 1$  if  $j = k$  and  $\delta_{jk} = 0$  if  $j \neq k$ . The multi-variate polynomial basis  $\Psi(\xi)$  are then constructed as a tensor product of the univariate orthogonal polynomials as

$$\Psi(\xi) = \prod_{i=1}^M \psi_{\alpha_i}(\xi_i) \quad (59)$$

To make the construction of the PCE practical, the summation in Eq. (57) needs to be truncated to include only a finite number of terms  $P$ . This is most commonly achieved by truncating such that only terms whose total degree is less than or equal to a value  $p$ , i.e.  $|\alpha| \leq p$ . This truncation defines a multi-index set

$$\mathcal{A}^{M,p} = \left\{ \alpha \in \mathbb{N}^M : |\alpha| = \sum_{i=1}^M \alpha_i \leq p \right\}. \quad (60)$$

whose cardinality is given by

$$\text{card } \mathcal{A}^{M,p} = \frac{(M+p)!}{M! p!} \equiv P. \quad (61)$$

This basis set is referred to as the *total degree basis*. Here, we see that the number of basis functions (and hence the number of coefficients to solve) grows extremely fast with the dimension  $M$  and the polynomial order  $p$ . To reduce the number of basis functions and improve scalability, Blatman and Sudret [8] proposed the so-called *hyperbolic truncation* scheme such that the basis set is defined by

$$\mathcal{A}^{M,p,q} = \left\{ \alpha \in \mathbb{N}^M : \|\alpha\|_q \equiv \left( \sum_{i=1}^M \alpha_i^q \right)^{1/q} \leq p \right\}. \quad (62)$$

In this scheme, only the basis vectors are retained having a specified  $L_q$  norm where  $q$  is a parameter of the scheme. In the case where  $q = 1$ , the truncation corresponds to the total degree basis. When  $q < 1$ , the scheme truncates terms that have higher-order interactions. This can dramatically reduce the cardinality of the basis set for high-dimensional problems with high polynomial order (i.e.  $M$  and  $p$  large). This scheme is illustrated in Figure [Create a figure like Figures 1 and 2 in \[8\] - Figures 2 and 3 are based on this reference.](#)

With this truncated basis set, we can now express the PCE in terms of a finite number of terms with some truncation error  $\epsilon$  as:

$$Y = h(\mathbf{X}) = \sum_{\alpha \in \mathcal{A}} \beta_{\alpha} \Psi_{\alpha}(\xi) + \epsilon. \quad (63)$$

The objective then in PCE is to determine the coefficients  $\beta$  that minimize the truncation error  $\epsilon$ . This can be viewed as a classical regression problem, and approaches to solve this regression problem will be discussed below. However, first we will discuss the selection of orthogonal polynomials for problems with different distribution types.

[Add single index notation.](#)

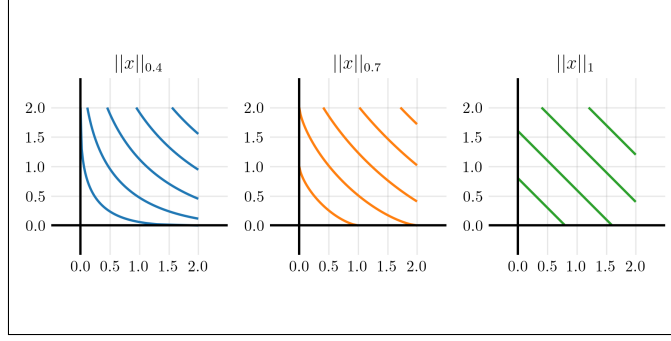


Figure 2: The level sets of the 0.4-norm, 1-norm, and 2-norm. Updated norms to 0.4, 0.7, and 1

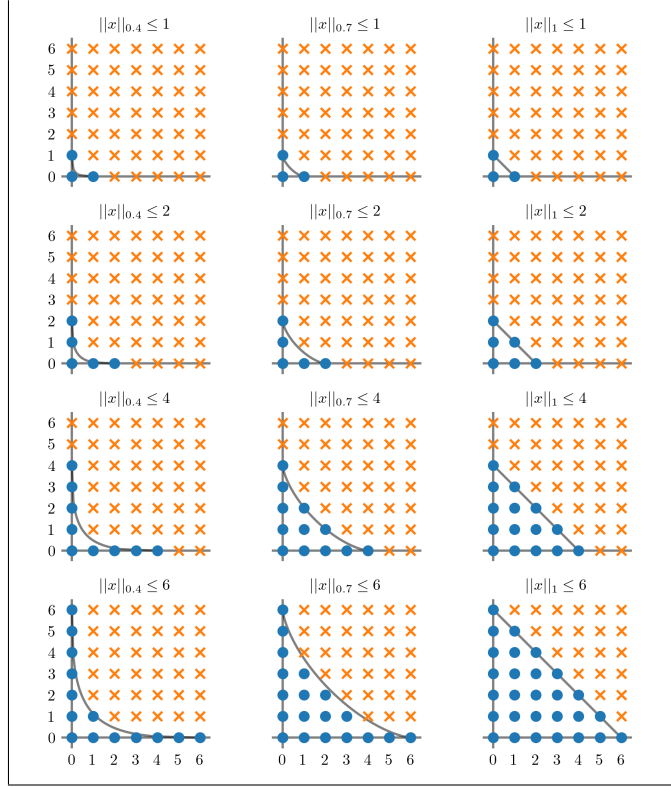


Figure 3: Truncation schemes for various norms and thresholds. Updated norms to 0.4, 0.7, and 1

### 3.1.1 Wiener-Askey Scheme

The specific orthogonal polynomials are selected according the Wiener-Askey scheme as detailed in [9]. The Wiener-Askey scheme specifically specifies the following polynomials corresponding to germs of different distributions:

- *Gaussian-Hermite*: When the germ  $\xi$  follows a standard normal distribution, the orthogonal polynomials are the (probabilists) Hermite polynomials  $H_n(x)$ . The Hermite polynomials are orthogonal

with respect to the following standard normal weight function

$$w(x) = e^{-\frac{x^2}{2}} \quad (64)$$

satisfying

$$\int_{-\infty}^{\infty} H_m(x) H_n(x) w(x) dx = 0 \quad \forall m \neq n. \quad (65)$$

We therefore have the following orthogonality condition

$$\int_{-\infty}^{\infty} H_m(x) H_n(x) e^{-\frac{x^2}{2}} dx = \sqrt{2\pi} n! \delta_{nm}, \quad (66)$$

Hence, the Hermite polynomials are orthogonal with respect to the standard normal probability density function. Additional details on the Hermite Polynomials can be found in Appendix A.1.

- *Beta-Jacobi* When the germ  $\xi$  follows the beta distribution, the orthogonal polynomials are the Jacobi polynomials,  $P_n^{(\alpha, \beta)}(x)$ . The Jacobi polynomials are orthogonal with respect to the following weight function

$$w(x) = (1-x)^\alpha (1+x)^\beta \quad (67)$$

on the interval  $[-1, 1]$  satisfying

$$\int_{-1}^1 P_m^{(\alpha, \beta)}(x) P_n^{(\alpha, \beta)}(x) w(x) dx = 0 \quad \forall m \neq n. \quad (68)$$

We have the following orthogonality condition

$$\int_{-1}^1 (1-x)^\alpha (1+x)^\beta P_m^{(\alpha, \beta)}(x) P_n^{(\alpha, \beta)}(x) dx = \frac{2^{\alpha+\beta+1}}{2n+\alpha+\beta+1} \frac{\Gamma(n+\alpha+1)\Gamma(n+\beta+1)}{\Gamma(n+\alpha+\beta+1)n!} \delta_{nm}, \quad \alpha, \beta > -1. \quad (69)$$

The Beta distribution has probability density function given by

$$f(x; \alpha, \beta) = \frac{1}{(b-a)^{\alpha+\beta+1} B(\alpha+1, \beta+1)} (x-a)^\beta (b-x)^\alpha \quad (70)$$

where  $B(p, q)$  is the *beta function* defined by

$$B(p, q) = \frac{\Gamma(p)\Gamma(q)}{\Gamma(p+q)} \quad (71)$$

where  $\Gamma(\cdot)$  is the Gamma function.

From these expressions, we can see that the weight function  $w(x)$  in Eq. (67) is a special scaled case of the Beta distribution defined on the interval  $[-1, 1]$ . Hence, the Jacobi polynomials are orthogonal with respect to the Beta distribution on  $[-1, 1]$ .

- *Uniform-Legendre* The special case of the Beta distribution with  $\alpha = \beta = 0$  corresponds to the uniform distribution and the Jacobi polynomials become the Legendre Polynomials,  $P_n(x)$ . Therefore, when the germ  $\xi$  follows the uniform distribution  $\sim U(-1, 1)$ , the orthogonal polynomials are the Legendre polynomials, which are orthogonal with respect to the following weight function

$$w(x) = 1 \quad (72)$$

satisfying

$$\int_{-1}^1 P_m(x) P_n(x) w(x) dx = 0 \quad \forall m \neq n. \quad (73)$$

We therefore have the following orthogonality condition

$$\int_{-1}^1 P_m(x) P_n(x) w(x) dx = \frac{2}{2n+1} \delta_{nm} \quad (74)$$

- *Gamma-Generalized Laguerre*

When the germ  $\xi$  follows the Gamma distribution, the orthogonal polynomials are the Generalized Laguerre polynomials,  $L_n^{(\alpha)}(x)$ . The Generalized Laguerre polynomials are orthogonal with respect to the following weight function

$$w(x) = x^\alpha e^{-x} \quad (75)$$

on the interval  $[0, \infty)$  satisfying

$$\int_0^\infty L_m^{(\alpha)}(x) L_n^{(\alpha)}(x) w(x) dx = 0 \quad \forall m \neq n. \quad (76)$$

We have the following orthogonality condition

$$\int_0^\infty L_m^{(\alpha)}(x) L_n^{(\alpha)}(x) x^\alpha e^{-x} dx = \frac{\Gamma(n + \alpha + 1)}{n!} \delta_{mn}, \quad \alpha > -1 \quad (77)$$

The Gamma distribution has probability density function given by

$$f(x; \alpha, \beta) = \frac{x^\alpha e^{-x/\beta}}{\beta^{\alpha+1} \Gamma(\alpha + 1)}, \quad \alpha > -1, \beta > 0 \quad (78)$$

where  $\Gamma(\cdot)$  is the Gamma function.

Regardless of the scale parameter  $\beta$  and constant  $\Gamma(\alpha + 1)$ , this density has the same form as the weight function and therefore the Generalized Laguerre polynomials are orthogonal with respect to the Gamma distribution.

- *Exponential-Laguerre*

The special case of the Gamma distribution with  $\alpha = 0$  corresponds to the exponential distribution. When the germ  $\xi$  follows the exponential distribution, the orthogonal polynomials are the Laguerre polynomials,  $L_n(x)$ . The Laguerre polynomials are orthogonal with respect to the following weight function

$$w(x) = e^{-x} \quad (79)$$

on the interval  $[0, \infty)$  satisfying

$$\int_0^\infty L_m(x) L_n(x) w(x) dx = 0 \quad \forall m \neq n. \quad (80)$$

We have the following orthogonality condition

$$\int_0^\infty L_m(x) L_n(x) e^{-x} dx = (n!)^2 \delta_{mn} \quad (81)$$

The Gamma distribution has probability density function given by

$$f(x; \alpha, \beta) = \frac{x^\alpha e^{-x/\beta}}{\beta^{\alpha+1} \Gamma(\alpha+1)}, \quad \alpha > -1, \beta > 0 \quad (82)$$

where  $\Gamma(\cdot)$  is the Gamma function.

Hence, the Laguerre polynomials are orthogonal with respect to the exponential probability density function. Additional details on the Laguerre Polynomials can be found in Appendix A.5.

## 3.2 Regression Methods

### Projection-based Approaches

Starting from the PCE in Eq. (57) and allowing that the germ corresponds to the vector of input random variables,  $\mathbf{X} = \boldsymbol{\xi}$ , we multiply by  $\Psi_{\alpha'}(\mathbf{X})$  and take the expectation, which yields

$$\beta_{\alpha} = E[\Psi_{\alpha}(\mathbf{X}) \cdot h(\mathbf{X})] \quad (83)$$

by orthonormality of the polynomials. This corresponds to the projection of the model  $h(\mathbf{X})$  onto the set of orthonormal polynomials  $\Psi_{\alpha}(\mathbf{X})$  and simplifies the estimation of the coefficients to the estimation of an expected value. Monte Carlo estimation would be inefficient. Instead, projection-based methods employ numerical integration schemes.

*Gaussian Quadrature:*

As the benchmark approach for numerical integration, Gaussian quadrature expresses the integral as a weighted sum of integrand evaluations at carefully selected quadrature points,  $\mathbf{x}^{(i)}$ . That is, the expectation is approximated by:

$$\begin{aligned} \beta_{\alpha} &= \mathbb{E}[\Psi_{\alpha}(\mathbf{X}) \cdot h(\mathbf{X})] \\ &= \int_{\Omega_{\mathbf{X}}} \Psi_{\alpha}(\mathbf{x}) h(\mathbf{x}) f_{\mathbf{X}}(\mathbf{x}) d\mathbf{x} \\ &\approx \sum_{i=1}^N w^{(i)} \Psi_{\alpha}(\mathbf{x}^{(i)}) h(\mathbf{x}^{(i)}) \end{aligned} \quad (84)$$

where the weights  $w^{(i)}$  and quadrature points  $\mathbf{x}^{(i)}$  are determined by the chosen polynomials  $\Psi_{\alpha}(x)$  and their order, and hence by the probability distributions of  $X_i$  following the Weiner-Askey scheme discussed above. In particular, the quadrature points  $x^{(i)}$  in each dimension correspond to the  $(p+1)$  roots of the  $p^{th}$ -order polynomial. Practically, the quadrature points and the weights are computed numerically using schemes such as the Golub-Welsch algorithm [10] or Newton's method for solving  $\Psi_p(x) = 0$ .

To perform this integration in multiple dimensions requires a tensor product grid of quadrature points. Since each dimension requires  $(p+1)$  quadrature points, an  $M$ -dimensional integral requires  $(p+1)^M$  quadrature points, which grows exponentially with dimension. Considering that the model evaluation  $h(\mathbf{x})$  is required at each quadrature point, this becomes computationally very expensive as the number of model evaluations grows exponentially. This is a consequence of the so-called *curse of dimensionality*.

*Sparse Quadrature:*

To mitigate the curse of dimensionality, quadrature can be conducted on so-called sparse quadrature grids whose weights are derived by hierarchically combining lower-order univariate quadrature terms. First derived by Smolyak [11], numerous such methods have been developed including those that use piecewise linear interpolation and those using polynomial interpolation including the popular Clenshaw-Curtis grids

that use the Chebyshev Gauss-Lobatto quadrature points and schemes using the Gauss-Patterson – both of which use nested grids of sparse quadrature points. A detailed discussion of sparse quadrature is beyond the scope of these notes.

## Least Squares

Introduced by Berveiller et al. [12], the most common approach to solving the PCE coefficients is by setting up the following least-squares problem. Recall the truncated PCE from Eq. (63) where the truncation error  $\epsilon$  can be expressed as

$$\epsilon = h(\mathbf{X}) - \boldsymbol{\beta}^T \Psi(\mathbf{X}) \quad (85)$$

where  $\boldsymbol{\beta} = [\beta_0, \beta_1, \dots, \beta_{P-1}]^T$ ,  $P = \text{card } \mathcal{A}^{M,p}$  and  $\Psi(\mathbf{X}) = \{[Psi_0(\mathbf{X}), \dots, \Psi_{P-1}(\mathbf{X})]^T$  are vectors containing the PCE coefficients and orthogonal polynomials, respectively. In least-squares minimization, we aim to find the coefficients  $\hat{\boldsymbol{\beta}}$  that minimize the expected mean square error of this truncation as

$$\hat{\boldsymbol{\beta}} = \arg \min_{\boldsymbol{\beta}} \mathbb{E} \left[ (h(\mathbf{X}) - \boldsymbol{\beta}^T \Psi(\mathbf{X}))^2 \right] \quad (86)$$

The standard approach to solving this minimization is referred to as *Ordinary Least Squares* (OLS). In OLS, we begin by evaluating the model  $Y = h(\mathbf{X})$  and the basis polynomials  $\Psi(\mathbf{X})$  at a set of  $N$  samples  $\mathcal{X} = [\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(N)}]^T$  (often referred to as the experimental design) to obtain the model responses  $\mathcal{Y} = [y^{(1)}, y^{(2)}, \dots, y^{(N)}]^T$  and the matrix of polynomial evaluations  $\Psi$  having terms  $\Psi_{ij} = \Psi_j(\mathbf{x}^{(i)})$ . The OLS problem can then be formulated practically as

$$\hat{\boldsymbol{\beta}} = \arg \min_{\boldsymbol{\beta}} \frac{1}{N} \sum_{i=1}^N \left[ \boldsymbol{\beta}^T \Psi(\mathbf{x}^{(i)}) - h(\mathbf{x}^{(i)}) \right]^2 \quad (87)$$

The closed-form solution to the minimization problem is then determined by

$$\hat{\boldsymbol{\beta}} = (\Psi^T \Psi)^{-1} \Psi^T \mathcal{Y} \quad (88)$$

The advantage of OLS is that it allows the coefficients to be determined from an arbitrary set of samples of the random vector  $\mathbf{X}$  as long as  $N \geq P$ .

## Sparse Regression Methods

As summarized by Luthen et al. [13], there are several ways to formulate the regression problem to promote a sparse solution; that is the solution that requires the fewest number of polynomials,  $P$ . Strictly speaking, to promote sparsity we aim to minimize the number of non-zero coefficients  $\boldsymbol{\beta}$ , which is commonly denoted by  $\|\boldsymbol{\beta}\|_0 = \sum_i I(\beta_i \neq 0)$  where  $I(\cdot)$  is the identity function that takes value 1 if the condition is satisfied and zero otherwise. We can integrate this condition as a penalty into the optimizer by reformulating the minimization as

$$\hat{\boldsymbol{\beta}} = \arg \min_{\boldsymbol{\beta}} \mathbb{E} \left[ (h(\mathbf{X}) - \boldsymbol{\beta}^T \Psi(\mathbf{X}))^2 \right] + \lambda \|\boldsymbol{\beta}\|_0, \quad (89)$$

which is referred to as  $l_0$ -minimization. However, this optimization is non-convex and requires a combinatorial solution that is computationally intractable for most practical problems.

Instead, it is common to formulate the problem as an  $l_1$ -minimization problem to achieve convexity, wherein  $\|\boldsymbol{\beta}\|_0$  is replaced with  $\|\boldsymbol{\beta}\|_1 = \sum_i |\beta_i|$ . There are several ways to formulate this  $l_1$ -minimization problem. Perhaps the most common is the *Lagrange formulation* given by

$$\hat{\boldsymbol{\beta}} = \arg \min_{\boldsymbol{\beta}} \mathbb{E} \left[ (h(\mathbf{X}) - \boldsymbol{\beta}^T \Psi(\mathbf{X}))^2 \right] + \lambda \|\boldsymbol{\beta}\|_1, \quad (90)$$

Alternatively, the  $l_1$ -minimization can be formulated as a constrained optimization problem using either the *basis pursuit* method where the  $l_1$  norm itself is minimized and constrained such that the least squares problem produces a sufficiently small error as

$$\begin{aligned}\hat{\boldsymbol{\beta}} &= \arg \min_{\boldsymbol{\beta}} \|\boldsymbol{\beta}\|_1 \\ \text{s.t. } &\mathbb{E} \left[ (h(\mathbf{X}) - \boldsymbol{\beta}^T \Psi(\mathbf{X}))^2 \right] \leq \sigma,\end{aligned}\tag{91}$$

where  $\sigma$  is an error threshold, or the *least absolute shrinkage and selection operator* (LASSO) where the  $l_2$  norm is again minimized subject to a constraint that the  $l_1$  norm must be smaller than a threshold  $\tau$  as

$$\begin{aligned}\hat{\boldsymbol{\beta}} &= \arg \min_{\boldsymbol{\beta}} \mathbb{E} \left[ (h(\mathbf{X}) - \boldsymbol{\beta}^T \Psi(\mathbf{X}))^2 \right] \\ \text{s.t. } &\|\boldsymbol{\beta}\|_1 \leq \tau,\end{aligned}\tag{92}$$

Numerous numerical methods have been devised to solve these sparse regression problems. Among the most widely-used are *Least Angle Regression* (LARS) [8] and its extension termed *Hybrid LARS with LOO-CV* [8], the *Orthogonal Matching Pursuit* (OMP) algorithm, and the *Subspace Pursuit* (SP) method. All of these are greedy algorithms that iteratively add and remove basis vectors according to some criterion. Details of these algorithms will not be provided here. For an in-depth review, see Luthen et al. [13].

### Bayesian Compressive Sensing

In Bayesian compressive sensing (BCS), or sparse Bayesian learning (SBL), the regression problem is reformulated as a Bayesian inverse problem in a special way to promote sparsity of the resulting expansion. Since we haven't covered Bayesian inversion in detail yet, we will provide only a brief overview of the method and will return to this at a later time.

The method starts by treating the training data  $\mathcal{Y}$  as independent noisy realizations of the PCE model following a normal distribution with mean  $\Psi\boldsymbol{\beta}$  and identical variances  $\sigma^2$ , i.e.  $p(\mathcal{Y}) = N(\Psi\boldsymbol{\beta}, \sigma^2\mathbf{I})$  where  $\mathbf{I}$  is the identity matrix. The variance  $\sigma$  is either a specified fixed value or a random variable whose distribution must be specified. The parameters of the PCE model,  $\boldsymbol{\beta}$  are also treated as normal random variables with zero mean and variances  $\gamma_i$  (each coefficient has its own variance), i.e.  $p(\beta_i) = N(0, \gamma_i)$ . These variances are the so-called *hyperparameters* of the model and are also treated as random variables having specified distribution  $p(\gamma_i)$ . The goal of Bayesian inference in this setting is to learn the distribution  $p(\boldsymbol{\beta}|\mathcal{Y})$  where the prior parameter distribution is given by  $p(\boldsymbol{\beta}) = \int p(\boldsymbol{\beta}|\boldsymbol{\gamma})p(\boldsymbol{\gamma})d\boldsymbol{\gamma}$ . In BCS, certain forms of the parameter distributions,  $p(\gamma_i)$ , are selected such that the effective prior distribution  $p(\boldsymbol{\beta})$  promotes sparsity. That is, it has a high weight near zero and heavy tails. Moreover, rather than seeking the full posterior distribution of the parameters  $p(\boldsymbol{\beta}|\mathcal{Y})$ , BCS seeks to identify only the *maximum a posterior* (MAP) estimate of the parameters  $\boldsymbol{\beta}^{\text{MAP}}$ , that is the estimate of  $\boldsymbol{\beta}$  that maximizes  $p(\boldsymbol{\beta}|\mathcal{Y})$ . Thanks to these simplifications and the assumption that  $p(\boldsymbol{\beta}|\boldsymbol{\gamma})$  and  $p(\mathcal{Y}|\boldsymbol{\beta}, \sigma^2)$  are normal, many of the requisite calculations can be performed analytically.

### 3.3 Moment Estimation with PCE

The PCE has the added advantage that the moments of the function can be estimated directly from the coefficients. Consider the  $m^{\text{th}}$  moment about the origin given by

$$\mathbb{E}[Y^m] = \int [h(\mathbf{X})]^m p_{\mathbf{X}}(\mathbf{x})d\mathbf{x}\tag{93}$$



Applying the truncated PCE from Eq. (63) yields

$$\begin{aligned}
\mathbb{E}[Y^m] &\approx \int \left[ \sum_{\alpha \in \mathcal{A}} \beta_{\alpha} \Psi_{\alpha}(\mathbf{x}) \right]^m p_{\mathbf{X}}(\mathbf{x}) d\mathbf{x} \\
&\approx \int \sum_{\alpha_1 \in \mathcal{A}} \cdots \sum_{\alpha_m \in \mathcal{A}} \beta_{\alpha_1} \cdots \beta_{\alpha_m} \Psi_{\alpha_1}(\mathbf{x}) \cdots \Psi_{\alpha_m}(\mathbf{x}) p_{\mathbf{X}}(\mathbf{x}) d\mathbf{x} \\
&\approx \sum_{\alpha_1 \in \mathcal{A}} \cdots \sum_{\alpha_m \in \mathcal{A}} \beta_{\alpha_1} \cdots \beta_{\alpha_m} \int \Psi_{\alpha_1}(\mathbf{x}) \cdots \Psi_{\alpha_m}(\mathbf{x}) p_{\mathbf{X}}(\mathbf{x}) d\mathbf{x}
\end{aligned} \tag{94}$$

As we can see, moment estimation therefore requires only integration over the polynomial basis functions, which simplifies greatly thanks to their orthogonality with respect to the probability measure  $p_{\mathbf{X}}(\mathbf{x})$ . This can be seen in the estimation of the first two moments such that the mean value is computed by

$$\mathbb{E}[Y] \approx \sum_{\alpha \in \mathcal{A}} \beta_{\alpha} \int \Psi_{\alpha}(\mathbf{x}) p_{\mathbf{X}}(\mathbf{x}) d\mathbf{x} \tag{95}$$

Recall that  $\Psi_0(\mathbf{x}) = 1$ , this expression can be rewritten as

$$\mathbb{E}[Y] \approx \sum_{\alpha \in \mathcal{A}} \beta_{\alpha} \int \Psi_0(\mathbf{x}) \Psi_{\alpha}(\mathbf{x}) p_{\mathbf{X}}(\mathbf{x}) d\mathbf{x} \tag{96}$$

Thanks to the orthonormality of the polynomials, this simplifies to

$$\mathbb{E}[Y] \approx \beta_0 \tag{97}$$

Similarly, the second moment can be estimated by

$$\begin{aligned}
\mathbb{E}[Y^2] &\approx \int \left[ \sum_{\alpha \in \mathcal{A}} \beta_{\alpha} \Psi_{\alpha}(\mathbf{x}) \right]^2 p_{\mathbf{X}}(\mathbf{x}) d\mathbf{x} \\
&\approx \sum_{\alpha_1 \in \mathcal{A}} \sum_{\alpha_2 \in \mathcal{A}} \beta_{\alpha_1} \beta_{\alpha_2} \int \Psi_{\alpha_1}(\mathbf{x}) \Psi_{\alpha_2}(\mathbf{x}) p_{\mathbf{X}}(\mathbf{x}) d\mathbf{x} \\
&\approx \sum_{\alpha \in \mathcal{A}} \beta_{\alpha}^2 \int [\Psi_{\alpha}(\mathbf{x})]^2 p_{\mathbf{X}}(\mathbf{x}) d\mathbf{x} \\
&\approx \sum_{\alpha \in \mathcal{A}} \beta_{\alpha}^2
\end{aligned} \tag{98}$$

Recall that the variance is given by  $\sigma_Y^2 = \mathbb{E}[Y^2] - \mathbb{E}[Y]^2$ . Therefore, we have

$$\begin{aligned}
\sigma_Y^2 &\approx \sum_{\alpha \in \mathcal{A}} \beta_{\alpha}^2 - \beta_0^2 \\
&\approx \sum_{\substack{\alpha \in \mathcal{A} \\ \alpha \neq 0}} \beta_{\alpha}^2
\end{aligned} \tag{99}$$

That is, the variance can be obtained simply by the sum of the squares of all of the coefficients in the model excluding the intercept  $\beta_0$ .

Using Eq. (94), the higher-order moments of the model output can also be determined, but this requires the computation of higher-order products of the orthogonal polynomials, which are not trivial to evaluate and are therefore rarely used in practice. Nonetheless, expressions have been derived for third- and fourth-order relations for certain classes of polynomials (e.g. Hermite, Jacobi and generalized Laguerre) [14].

A closely related topic is the estimation of variance based sensitivity indices. This will be discussed in further detail in the following module on Global Sensitivity Analysis.

## 4 Other Surrogate Models

Numerous other surrogate models have been developed in practice, but are not as widely used as the GP and PCE surrogates. These include, but are not limited to, surrogate models based on radial basis functions, support vector regression, artificial neural networks, and polynomial regressions. These methods, though sometimes used in practice, generally do not have the same advantages as GPs and PCEs and are therefore less commonly used for UQ applications. We therefore will not explore them in further detail.

## A Definitions and Properties of Orthogonal Polynomials

Orthogonal polynomials in the Askey Scheme are defined through the *generalized hypergeometric series* denoted  ${}_pF_q(a_1, \dots, a_p; b_1, \dots, b_q; x)$ . First, denoting the rising factorial through the Pochhammer symbol

$$\begin{aligned} (a)_0 &= 1, \\ (a)_n &= a(a+1)(a+2) \cdots (a+n-1), \quad n \geq 1 \end{aligned} \tag{100}$$

the generalized hypergeometric series is expressed:

$${}_pF_q(a_1, \dots, a_p; b_1, \dots, b_q; x) = \sum_{k=0}^{\infty} \frac{(a_1)_k \cdots (a_p)_k}{(b_1)_k \cdots (b_q)_k} \frac{x^k}{k!}. \tag{101}$$

If any  $a_j$  is a non-positive integer, i.e.  $a_j = -n$ , then the series only has a finite number of terms and is, in fact, a polynomial of degree  $n$  given by

$${}_pF_q(a_1, \dots, a_p; b_1, \dots, b_q; x) = \sum_{k=0}^n \frac{(a_1)_k \cdots (a_p)_k}{(b_1)_k \cdots (b_q)_k} \frac{x^k}{k!}. \tag{102}$$

All orthogonal polynomials  $\{Q_n(x)\}$  also satisfy a three-term recurrence relation given by;

$$-xQ_n(x) = A_nQ_{n+1}(x) - (A_n + C_n)Q_n(x) + C_nQ_{n-1}(x), \quad n \geq 1 \tag{103}$$

where  $A_n, C_n \neq 0$  and  $C_n/A_{n-1} > 0$ . Given that  $Q_{-1}(x) = 0$  and  $Q_0(x) = 1$ , all  $Q_n(x)$  can be determined from the recurrence relation.

All orthogonal polynomials can also be determined by repeatedly applying the differential operator as:

$$Q_n(x) = \frac{1}{w(x)} \frac{d^n}{dx^n} [w(x)s^n(x)] \tag{104}$$

where  $w(x)$  is the orthogonal weighting function and  $s(x)$  is a polynomial of at most second order. This formula is referred to as the *generalized Rodrigues' formula*.

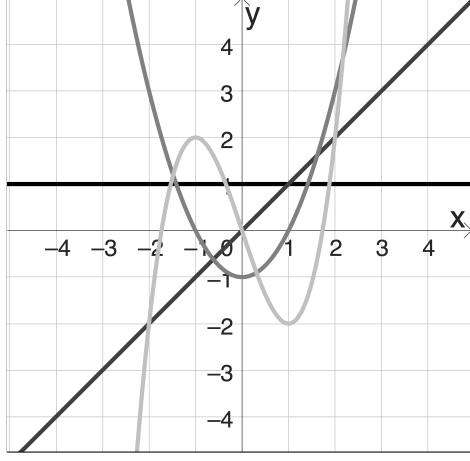
For each of the classes of orthogonal polynomials discussed above the in Askey scheme, we will provide the definition in terms of its generalized hypergeometric series, its recurrence relations, and its generalized Rodrigues' formula. We will also provide the first five orthogonal polynomials.

### A.1 Hermite Polynomials

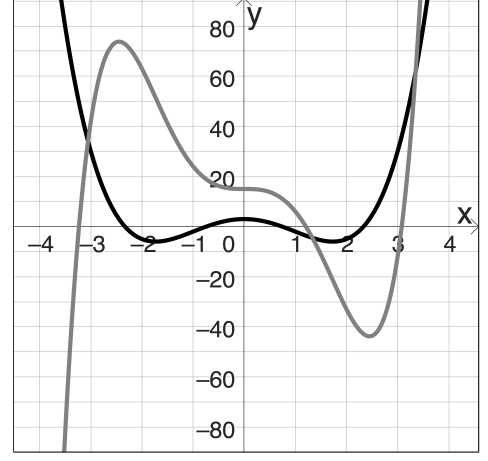
#### Hypergeometric Series Definition

The Hermite polynomials are defined in terms of the following hypergeometric series:

$$H_n(x) = (2x)^n {}_2F_0 \left( -\frac{n}{2}, -\frac{n-1}{2}; -\frac{1}{x^2} \right) \tag{105}$$



(a) The Hermite polynomials  $H_0, H_1, H_2, H_3$



(b) The Hermite polynomials  $H_4$  and  $H_5$

Figure 4: The first six Hermite polynomials Uploaded figure

### Rodrigues' Formula

$$H_n(x) = (-1)^n e^{\frac{x^2}{2}} \frac{d^n}{dx^n} e^{-\frac{x^2}{2}}, \quad (106)$$

where  $n = 0, 1, 2, \dots$  is the order of the polynomial.

### Recurrence Relations

$$H_{n+1}(x) = xH_n(x) - nH_{n-1}(x). \quad (107)$$

from which we can obtain the first six Hermite polynomials as

$$\begin{aligned} H_0(x) &= 1 \\ H_1(x) &= x \\ H_2(x) &= x^2 - 1 \\ H_3(x) &= x^3 - 3x \\ H_4(x) &= x^4 - 6x^2 + 3 \\ H_5(x) &= x^5 - 10x^3 + 15x \end{aligned} \quad (108)$$

which are plotted in Figure 4.

## A.2 Jacobi Polynomials

### A.2.1 Hypergeometric Series Definition

The Jacobi polynomials are defined in terms of the following hypergeometric series:

$$P_n^{(\alpha, \beta)}(x) = \frac{(\alpha + 1)_n}{n!} {}_2F_1 \left( -n, 1 + \alpha + \beta + n; \alpha + 1; \frac{1}{2}(1 - x) \right), \quad (109)$$

### A.2.2 Rodrigues' Formula

$$P_n^{(\alpha, \beta)}(x) = \frac{(-1)^n}{2^n n!} (1 - x)^{-\alpha} (1 + x)^{-\beta} \frac{d^n}{dx^n} \left\{ (1 - x)^\alpha (1 + x)^\beta (1 - x^2)^n \right\}. \quad (110)$$

where  $\alpha$  and  $\beta$  are parameters of the model.

### A.2.3 Recurrence Relations

For a given  $\alpha$  and  $\beta$  the Jacobi polynomials can be determined from

$$2n(n + \alpha + \beta)(2n + \alpha + \beta - 2)P_n^{(\alpha, \beta)}(x) = (2n + \alpha + \beta - 1) \left\{ (2n + \alpha + \beta)(2n + \alpha + \beta - 2)x + \alpha^2 - \beta^2 \right\} P_{n-1}^{(\alpha, \beta)}(x) - 2(n + \alpha - 1)(n + \beta - 1)(2n + \alpha + \beta)P_{n-2}^{(\alpha, \beta)}(x), \quad (111)$$

from which we can obtain the first few Jacobi polynomials as

$$\begin{aligned} P_0^{(\alpha, \beta)}(x) &= 1 \\ P_1^{(\alpha, \beta)}(x) &= \frac{1}{2}[2(\alpha + 1) + (\alpha + \beta + 2)(x - 1)] \\ P_2^{(\alpha, \beta)}(x) &= \frac{1}{8}[4(\alpha + 1)(\alpha + 2) + 4(\alpha + \beta + 3)(\alpha + 2)(x - 1) + (\alpha + \beta + 3)(\alpha + \beta + 4)(x - 1)^2] \end{aligned} \quad (112)$$

which are plotted in Figure 5.

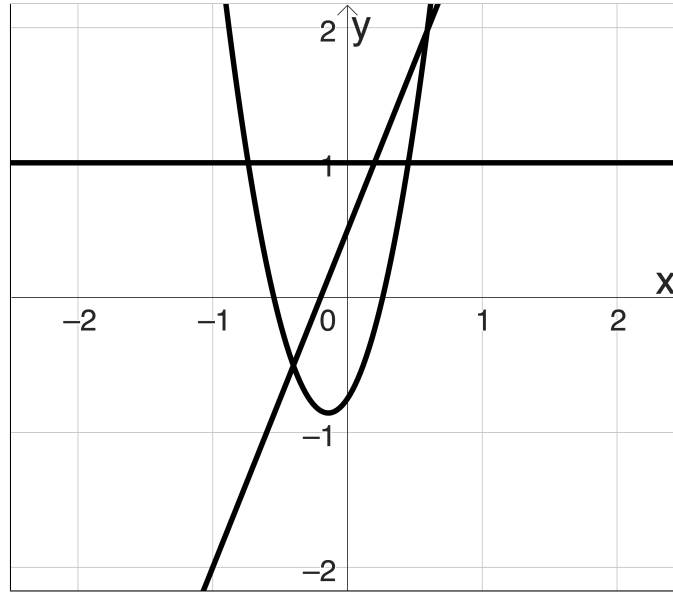


Figure 5: The first three Jacobi polynomials with  $\alpha = 2, \beta = 1$ . **Uploaded figure.**

## A.3 Legendre Polynomials

### A.3.1 Hypergeometric Series Definition

$$P_n(x) = P_n^{(0,0)}(x) = {}_2F_1\left(-n, n+1; 1; \frac{1}{2}(1-x)\right) \quad (113)$$

### A.3.2 Rodrigues' Formula

$$P_n(x) = \frac{1}{2^n n!} \frac{d^n}{dx^n} (x^2 - 1)^n. \quad (114)$$

### A.3.3 Recurrence Relations

$$(n+1)P_{n+1}(x) = (2n+1)xP_n(x) - nP_{n-1}(x) \quad (115)$$

from which we can obtain the first five Legendre polynomials as

$$\begin{aligned} P_0(x) &= 1 \\ P_1(x) &= x \\ P_2(x) &= \frac{1}{2}(3x^2 - 1) \\ P_3(x) &= \frac{1}{2}(5x^3 - 3x) \\ P_4(x) &= \frac{1}{8}(35x^4 - 30x^2 + 3) \\ P_5(x) &= \frac{1}{8}(63x^5 - 70x^3 + 15x) \end{aligned} \quad (116)$$

which are plotted in Figure 6.

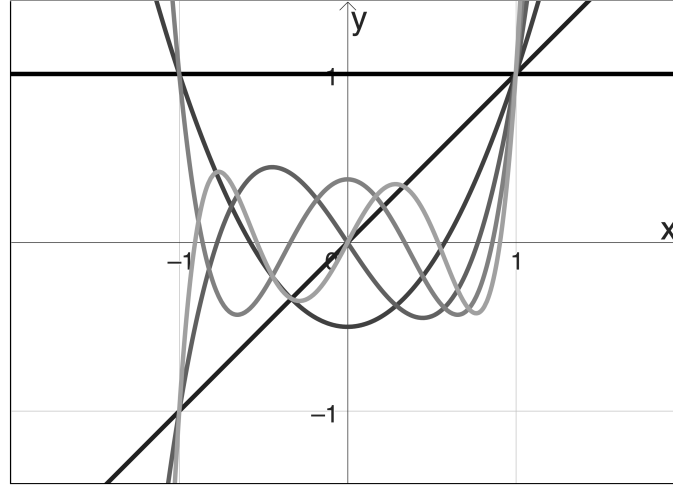


Figure 6: The first six Legendre Jacobi polynomials. Uploaded figure.

## A.4 Generalized Laguerre Polynomials

### A.4.1 Hypergeometric Series Definition

$$L_n^{(\alpha)}(x) = \binom{n+\alpha}{n} M(-n, \alpha+1, x) = \frac{(\alpha+1)_n}{n!} {}_1F_1(-n, \alpha+1, x) \quad (117)$$

### A.4.2 Rodrigues' Formula

$$L_n^{(\alpha)}(x) = \frac{x^{-\alpha} e^x}{n!} \frac{d^n}{dx^n} (e^{-x} x^{n+\alpha}). \quad (118)$$

### A.4.3 Recurrence Relations

$$(n+1)L_{n+1}^{(\alpha)}(x) = (2n+\alpha+1-x)L_n^{(\alpha)}(x) - (n+\alpha)L_{n-1}^{(\alpha)}(x) \quad (119)$$

from which we can obtain the first few Generalized Laguerre polynomials as

$$\begin{aligned}
L_0^{(\alpha)}(x) &= 1 \\
L_1^{(\alpha)}(x) &= -x + (\alpha + 1) \\
L_2^{(\alpha)}(x) &= \frac{x^2}{2} - (\alpha + 2)x + \frac{(\alpha + 1)(\alpha + 2)}{2} \\
L_3^{(\alpha)}(x) &= \frac{-x^3}{6} + \frac{(\alpha + 3)x^2}{2} - \frac{(\alpha + 2)(\alpha + 3)x}{2} + \frac{(\alpha + 1)(\alpha + 2)(\alpha + 3)}{6}
\end{aligned} \tag{120}$$

which are plotted in Figure 7.

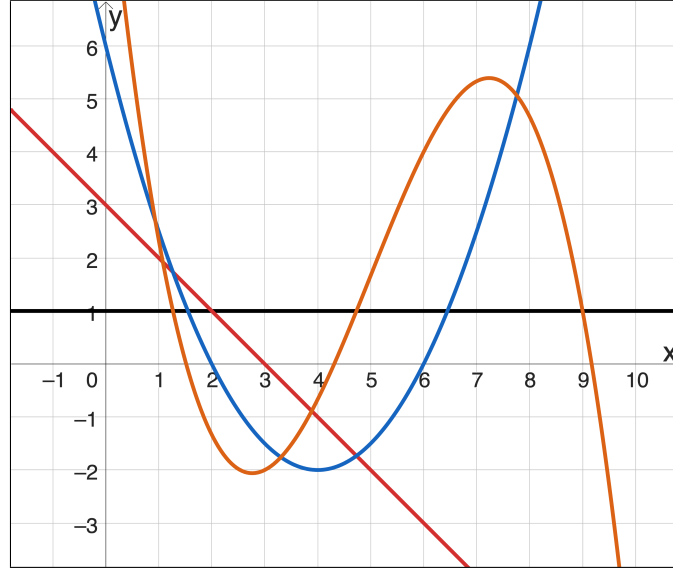


Figure 7: The first three Generalized Laguerre polynomials  $L_0^{(2)}(x)$ ,  $L_1^{(2)}(x)$ ,  $L_2^{(2)}(x)$  with  $\alpha = 2$ . **Uploaded figure.**

## A.5 Laguerre Polynomials

### A.5.1 Hypergeometric Series Definition

The Laguerre polynomials can be determined from the hypergeometric series in Eq. (117) by setting  $\alpha = 0$ .

### A.5.2 Rodrigues' Formula

$$L_n(x) = \frac{e^x}{n!} \frac{d^n}{dx^n} (e^{-x} x^n), \tag{121}$$

### A.5.3 Recurrence Relations

$$L_{k+1}(x) = \frac{(2k+1-x)L_k(x) - kL_{k-1}(x)}{k+1}, \tag{122}$$

from which we can obtain the first five Laguerre polynomials as

$$\begin{aligned}
L_0(x) &= 1 \\
L_1(x) &= 1 - x \\
L_2(x) &= \frac{1}{2}(x^2 - 4x + 2) \\
L_3(x) &= \frac{1}{6}(-x^3 + 9x^2 - 18x + 6) \\
L_4(x) &= \frac{1}{24}(x^4 - 16x^3 + 72x^2 - 96x + 24) \\
L_5(x) &= \frac{1}{120}(-x^5 + 25x^4 - 200x^3 + 600x^2 - 600x + 120)
\end{aligned} \tag{123}$$

which are plotted in Figure 8. Should all six of these be on one figure?

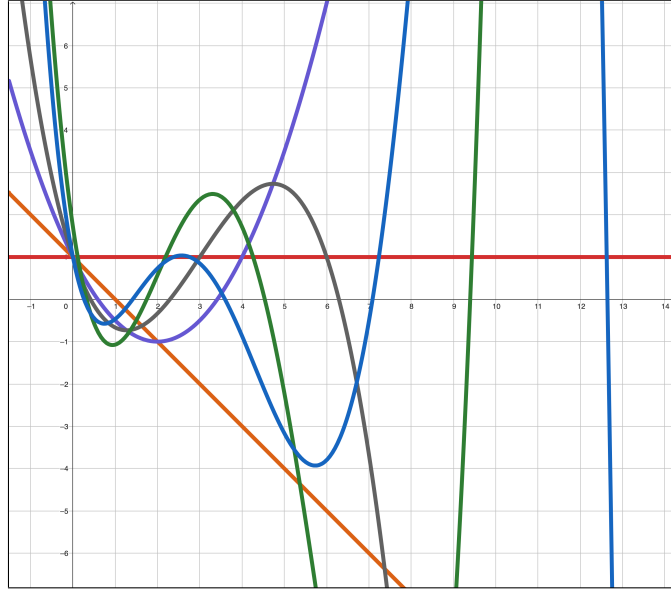
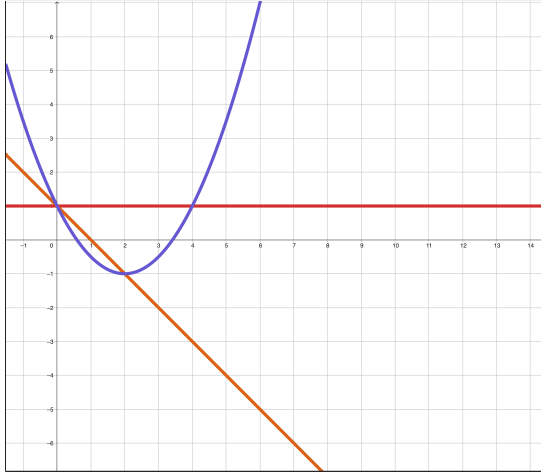


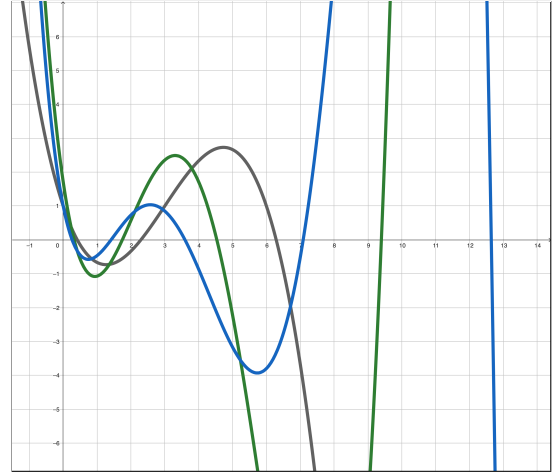
Figure 8: The first six Laguerre polynomials. Uploaded figure. I think this figure is very busy so there is also an option to separate this into two figures.

## References

- [1] C. E. Rasmussen, C. K. Williams, et al., Gaussian processes for machine learning, Vol. 1, Springer, 2006.
- [2] D. G. Krige, A statistical approach to some mine valuation and allied problems on the witwatersrand: By dg krige, Ph.D. thesis, University of the Witwatersrand (1951).
- [3] G. Matheron, Principles of geostatistics, Economic geology 58 (8) (1963) 1246–1266.
- [4] J. Sacks, W. J. Welch, T. J. Mitchell, H. P. Wynn, The design and analysis of computer experiments, Statistical Science 4 (1989).
- [5] T. J. Santner, B. J. Williams, W. I. Notz, B. J. Williams, The design and analysis of computer experiments, Vol. 1, Springer, 2003.



(a) The Laguerre polynomials  $L_0, L_1$  and  $L_2$ .



(b) The Laguerre polynomials  $L_3, L_4$  and  $L_5$ .

Figure 9: The first six Hermite polynomials Uploaded figure.

- [6] F. Bachoc, Cross validation and maximum likelihood estimations of hyper-parameters of gaussian processes with model misspecification, *Computational Statistics & Data Analysis* 66 (2013) 55–69.
- [7] R. Schobi, B. Sudret, J. Wiart, Polynomial-chaos-based kriging, *International Journal for Uncertainty Quantification* 5 (2) (2015).
- [8] G. Blatman, B. Sudret, Adaptive sparse polynomial chaos expansion based on least angle regression, *Journal of computational Physics* 230 (6) (2011) 2345–2367.
- [9] D. Xiu, G. E. Karniadakis, The wiener–askey polynomial chaos for stochastic differential equations, *SIAM journal on scientific computing* 24 (2) (2002) 619–644.
- [10] G. H. Golub, J. H. Welsch, Calculation of gauss quadrature rules, *Mathematics of computation* 23 (106) (1969) 221–230.
- [11] S. A. Smolyak, Quadrature and interpolation formulas for tensor products of certain classes of functions, in: *Doklady Akademii Nauk*, Vol. 148, Russian Academy of Sciences, 1963, pp. 1042–1045.
- [12] M. Berveiller, B. Sudret, M. Lemaire, Stochastic finite element: a non intrusive approach by regression, *European Journal of Computational Mechanics/Revue Européenne de Mécanique Numérique* 15 (1-3) (2006) 81–92.
- [13] N. Luthen, S. Marelli, B. Sudret, Sparse polynomial chaos expansions: Literature survey and benchmark, *SIAM/ASA Journal on Uncertainty Quantification* 9 (2) (2021) 593–649.
- [14] L. Novák, On distribution-based global sensitivity analysis by polynomial chaos expansion, *Computers & Structures* 267 (2022) 106808.

## Nomenclature

### Functions



$P(\cdot)$	Probability measure
$f_X(\cdot)$	Probability density function (PDF) of a random variable $X$
$f_{\mathbf{X}}(\cdot)$	Joint probability density function of the random vector $\mathbf{X}$
$F_X(\cdot)$	Cumulative distribution function (CDF) of a random variable $X$
$F_{\mathbf{X}}(\cdot)$	Joint Cumulative Distribution Function of the random vector $\mathbf{X}$
$R_{XX}(\cdot)$	Autocorrelation function of the random process $X(\cdot)$ This is the previously used autocorrelation function notation
$\mathcal{R}(\cdot)$	Autocorrelation function of the random process $X(\cdot)$ This notation was introduced in this module. Which is preferred?
$\mathbb{E}[\cdot]$	Expected value of a random variable. Also denoted $\mu_X \triangleq E[X]$
$\mathbf{F}$	Matrix of basis functions
$\mathbf{f}(\cdot)$	Vector of basis functions
$\mathbf{R}$	Correlation matrix of points in the training set
$\mathbf{r}(\cdot)$	Vector of correlations
$\delta_{jk}$	Kronecker delta function
$\Gamma(\cdot)$	The Gamma function
$\lceil \cdot \rceil$	Ceiling function
$\mathcal{F}(\cdot)$	Regression model
$\mathcal{K}_\nu$	Modified Bessel function of the second kind
$\mathcal{L}(\cdot)$	Likelihood function
$\mathcal{Y}(\cdot)$	Approximation of a stochastic process
$B(\cdot)$	Beta function
$\Psi_{\boldsymbol{\alpha}}(\cdot)$	Multivariate orthogonal polynomials
$\text{diag}(\cdot)$	The diagonalizing function
$H_n(\cdot)$	Hermite polynomial of order $n$
$I(\cdot)$	Identity function
$K(\cdot)$	Kernel function
$L_n(\cdot)$	Laguerre polynomial
$L_n^{(\alpha)}(\cdot)$	Generalized Laguerre polynomial
$P_n^{(\alpha, \beta)}(\cdot)$	Jacobi polynomial

${}_pF_q$  Hypergeometric series

### Operators

$\langle \cdot, \cdot \rangle$  Inner product

### Variables

$\beta$  Vector of regression coefficients

$\mathbf{X}$  A random vector in  $\mathbb{R}^n$

$\alpha$  Multi-index

$\Sigma_n$  Covariance matrix **Check if this matches other notation**

$\theta$  Parameters of the model (or function)  $f$

$\epsilon$  Truncation error **Check if epsilon is used for error elsewhere**

$\mathcal{A}^{M,p}$  Multi-index set

$\mathcal{C}^0$  The set of continuous functions

$\mathcal{D}$  Set of training data

$\sigma_X$  Standard deviation of the random variable  $X$

$\xi$  Germ

$X$  A random variable

$\Omega$  Sample space

$\omega$  A random event from the sample space  $\Omega$