

Monograph

Author: Ksenia N. Kyzyurova

Analysis of scientific computer models

Methodology in numerical simulator data analysis

September 21, 2019

This monograph provides methodological foundation for analysis of data from scientific computationally challenging mathematical models. Mathematical topics include statistical modeling, Bayesian inference, stochastic processes, and decision-theoretic model assessment.



Copyright © by Ksenia N. Kyzyurova
All rights reserved

Prerequisites

Science attempts to describe complex natural or engineering phenomena by use of mathematical processes. These computer models are usually either computationally slow and/or too resources demanding for operation of the model and sometimes even for storage of the data produced by the model. This makes the computer model output at any desirable input not available. However, fast approximations to such an output may be obtained, once an emulator of the model (its statistical approximation), using only a handful of input-output data points, is constructed. Construction of the ‘default’ emulator within its objective implementation is outlined below.

Statistical approximation to a computer model

Suppose that a computer model is represented by a smooth function $f(\cdot)$, which takes a d -dimensional input $\zeta \in Z \subseteq \mathcal{R}^d, d \geq 1$ and produces a d^* -dimensional output $\mathbf{f}(\zeta) \subseteq \mathcal{R}^{d^*}, d^* \geq 1$. Both d and d^* are positive integers. In the first two chapters of the present monograph, computer models with one-dimensional output $d^* = 1$ are considered. The problem of construction of statistical approximations to computer models with multivariate output is considered in chapter .

Suppose m inputs to the simulator $\zeta = (\zeta_1, \zeta_2, \dots, \zeta_m)$ are chosen, and the computer model f has been evaluated at those inputs, resulting in outputs $\Xi(\zeta) = (f(\zeta_1), f(\zeta_2), \dots, f(\zeta_m))$. Using these pairs of inputs and outputs $\{\zeta_i, f(\zeta_i)\}_{i=1}^m$ and assuming a Gaussian process prior $\Xi(\cdot) = \mathcal{GP}(\mu(\cdot), \text{cov}(\cdot, \cdot))$ (specified by some mean and covariance function) on the computer model output $f(\cdot)$, we aim at constructing a probabilistic representation of the model at any new input in its domain Z .

A Gaussian process is defined as a stochastic process, whose every finite sample of variables follows (all compatible) multivariate normal distribution. Since any multivariate normal distribution is specified by its mean and covariance matrix, a Gaussian process may be specified by the mean $\mathbf{E}\Xi(\zeta)$ and covariance $\text{Cov}(\Xi(\zeta_k), \Xi(\zeta_\ell))$ for any two inputs ζ_k and ζ_ℓ . For infinite number of random variables, the mean and covariance are convenient to specify parametrically (to satisfy the compatibility requirements). For example, customarily is to specify the mean of the Gaussian process at point ζ as *linear* with two parameters, intercept β_0 and slope β_1

$$\mathbf{E}\Xi(\zeta) = \beta_0 + \beta_1 \zeta_e, \tag{1}$$

where $e \in 1, 2, \dots, d$ is one of the input dimensions.

Equivalently, one may write $\mathbf{E}\Xi(\zeta) = \mathbf{H}\boldsymbol{\beta} = (1 \quad \zeta_e)\boldsymbol{\beta}$, where $\boldsymbol{\beta} = (\beta_0 \quad \beta_1)^\top$ is a vector of linear model regression coefficients. \mathbf{H}^\top is a vector of regression functions evaluated at ζ . In this case, two regression functions are the constant $\psi_0(\zeta) = 1$ and $\psi_1(\zeta) = \zeta_e$. These are the first and second element of the vector \mathbf{H}^\top .

In theoretical development easy to extend the mean specification to any number of regression functions. If the number of regression functions is s , then one must specify each of the functions $h_0(\cdot), h_1(\cdot), \dots, h_{s-1}(\cdot)$; then $\mathbf{H} = (h_0(\zeta) \quad h_1(\zeta) \quad \dots \quad h_{s-1}(\zeta))$.

Specification of covariance for infinite number of random variables is a more complicated matter, since it must be *valid*. Typically covariance between outputs of the computer model at corresponding inputs is specified as the product of variance parameter σ^2 and correlation function $c(\cdot, \cdot)$ evaluated at those inputs, i.e. for any two inputs ζ_k and ζ_ℓ

$$\text{Cov}(\Xi(\zeta_k), \Xi(\zeta_\ell)) = \sigma^2 c(\zeta_k, \zeta_\ell). \tag{2}$$

If an input to the model is one-dimensional, i.e. $d = 1$, a convenient and useful choice for the purpose of emulation of a computer model is squared-exponential correlation function, which is specified by the parameter ω

$$c(\zeta_k, \zeta_\ell) = \exp \left\{ - \exp(\omega) (\zeta_k - \zeta_\ell)^2 \right\}. \tag{3}$$

If an input to the model is multi-dimensional with $d \geq 2$, then convenient is to specify correlation function in the form of a product of correlation functions along each of the d input dimensions, that is $c(\zeta_k, \zeta_\ell) = \prod_{j=1}^d c(\zeta_{kj}, \zeta_{lj})$.

Lemma. The product of correlation functions is a valid correlation function.

Parameters of the mean and covariance form a vector of parameters $\boldsymbol{\theta} = (\boldsymbol{\beta}, \sigma^2, \omega)$ which altogether specify the Gaussian process. While β_0 and β_1 define a linear mean trend, parameter σ^2 is responsible for variability of the process, and parameter ω determines the range scale at which correlations among Gaussian process random variables remain strong.

Function $f(\zeta) = \exp(-\zeta) + \sin(4\zeta)$ in the range $\zeta \in [-1, 1]$ is taken as a prototype of a computer model. Suppose this model has been evaluated at eight input points $\zeta = (\zeta_1, \zeta_2, \dots, \zeta_8)$ which are equadistantly placed between -0.94 and 0.94 resulting in outputs $\mathbf{f}(\zeta) = (f(\zeta_1), f(\zeta_2), \dots, f(\zeta_8))$ respectively.

Suppose for now that the parameters of the Gaussian process θ are known and given. Let these parameters be $\beta_0 = 1.5$, $\beta_1 = 2$, $\sigma^2 = 2$ and $\omega = 2$. Once these are set, the construction of the Gaussian process emulator with parameters of the process being given is possible: assuming that this *training data*¹ $\{\zeta, \mathbf{f}(\zeta)\}$ belong to its to-be constructed approximation Ξ , such that $\Xi(\zeta) = \mathbf{f}(\zeta)$ and that Ξ has the Gaussian process prior, the model for the data is given by its likelihood

$$\Xi(\zeta) \sim \mathcal{N}(\mathbf{H}\beta, \sigma^2 \mathbf{C}_\omega), \quad (4)$$

where \mathbf{C}_ω is the correlation matrix given by a correlation function $c(\cdot, \cdot)$ such that the element of the matrix at row ι and column ι' equals to $c(\zeta_\iota, \zeta_{\iota'})$.

Let ζ' be a vector of input points at which the approximated solution is desired to be found. The joint distribution for inputs ζ and ζ' is multivariate normal

$$\begin{pmatrix} \Xi(\zeta) \\ \Xi(\zeta') \end{pmatrix} \sim \mathcal{N} \left(\begin{pmatrix} \mathbf{H}\beta \\ \mathbf{H}'\beta \end{pmatrix}, \sigma^2 \begin{pmatrix} \mathbf{C}_\omega & c(\zeta, \zeta')_\omega \\ c(\zeta', \zeta)_\omega & \mathbf{C}'_\omega \end{pmatrix} \right). \quad (5)$$

The approximation Ξ at points ζ' conditional on the training data $\{\zeta, \mathbf{f}(\zeta)\}$ is given by the predictive multivariate normal distribution

$$\Xi(\zeta') \mid \Xi(\zeta), \theta \sim \mathcal{N}(\mathbf{H}'\beta + c(\zeta', \zeta)_\omega \mathbf{C}_\omega^{-1}(\Xi(\zeta) - \mathbf{H}\beta), \sigma^2(\mathbf{C}'_\omega - c(\zeta', \zeta)_\omega \mathbf{C}_\omega^{-1} c(\zeta, \zeta')_\omega)). \quad (6)$$

This predictive distribution **defines** the statistical *emulator*, approximation to a computer model f at inputs ζ' . Figure 1 demonstrates the approximation constructed for the simulator f .^a

Theorem. $\Xi(\cdot) \mid \cdot, \Xi(\zeta), \theta$ is the stochastic process.

Proof. This is indeed so by construction. All joint distributions of any sequence of random variables are defined and are compatible. ■

Definition. $\Xi(\cdot) \mid \cdot, \Xi(\zeta), \theta$ is called the Gaussian process stochastic emulator (or approximation) to a computer model f .

$$\Xi(\cdot) \mid \cdot, \Xi(\zeta), \theta \sim \mathcal{GASP}(\mu^*(\cdot), \sigma^{*2}(\cdot, \cdot)), \quad (7)$$

where GASP stands for Gaussian stochastic process, $\mu^*(\cdot)$ is the mean of the predictive process and $\sigma^{*2}(\cdot, \cdot)$ is its variance-covariance, which for any vector of inputs follow the distribution 6.

^aIn literature words “emulator”, “surrogate”, “metamodel” are used interchangeably meaning “a probabilistic statistical approximation to a computer model”.

Within objective implementation, suppose that parameter ω is fixed, then the prior for the rest of parameters β and σ^2 for this normal model (4) is

$$p(\beta, \sigma^2) \propto \frac{1}{\sigma^2}. \quad (8)$$

The likelihood (4) for the model parameters $L(\omega, \beta, \sigma^2; \Xi(\zeta)) = p(\Xi(\zeta) \mid \omega, \beta, \sigma^2)$, which is, while omitting all expression which do not contain parameters β, σ^2, ω , is proportional to

$$L(\omega, \beta, \sigma^2; \Xi(\zeta)) \propto |\sigma^2 \mathbf{C}_\omega|^{-\frac{1}{2}} \exp \left\{ -\frac{1}{2} (\Xi(\zeta) - \mathbf{H}\beta)^\top \frac{1}{\sigma^2} \mathbf{C}_\omega^{-1} (\Xi(\zeta) - \mathbf{H}\beta) \right\} = (\sigma^2)^{-\frac{m}{2}} |\mathbf{C}_\omega|^{-\frac{1}{2}} \exp \left\{ -\frac{1}{2} (\Xi(\zeta) - \mathbf{H}\beta)^\top \frac{1}{\sigma^2} \mathbf{C}_\omega^{-1} (\Xi(\zeta) - \mathbf{H}\beta) \right\}. \quad (9)$$

¹Also called *design* data points.

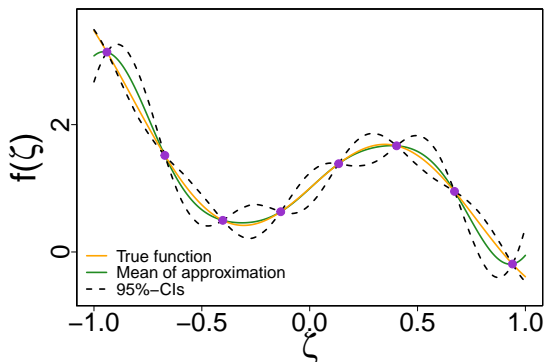


Figure 1: Constructed statistical approximation to the computer model f . The orange curve is the true function; the green — the mean of the emulator; the black dashed lines are those of the central 95%-credible area. The purple points are design points.

Posterior distributions of parameters. Joint posterior distribution of parameters β and σ^2 is convenient to write in conditional form

$$p(\beta, \sigma^2 \mid \Xi(\zeta), \omega) = p(\beta \mid \sigma^2, \Xi(\zeta), \omega) p(\sigma^2 \mid \Xi(\zeta), \omega). \quad (10)$$

These distributions are

$$\beta \mid \Xi(\zeta), \sigma^2, \omega \sim \mathcal{N}_\beta \left((\mathbf{H}^T \mathbf{C}_\omega^{-1} \mathbf{H})^{-1} \mathbf{H}^T \mathbf{C}_\omega^{-1} \Xi(\zeta), \sigma^2 (\mathbf{H}^T \mathbf{C}_\omega^{-1} \mathbf{H})^{-1} \right), \quad (11)$$

$$\sigma^2 \mid \Xi(\zeta), \omega \sim \mathcal{IG}_{\sigma^2} \left(\frac{m-s}{2}, \frac{\Xi(\zeta)^T (\mathbf{C}_\omega^{-1} - \mathbf{C}_\omega^{-1} \mathbf{H} (\mathbf{H}^T \mathbf{C}_\omega^{-1} \mathbf{H})^{-1} \mathbf{H}^T \mathbf{C}_\omega^{-1}) \Xi(\zeta)}{2} \right). \quad (12)$$

Instead of using predefined *ad hoc* values of β and σ^2 , the maximum *a posteriori* estimates obtained *from the data* are readily available from the joint posterior. Modes of posterior distributions 11 and 12 are

$$\hat{\beta} = (\mathbf{H}^T \mathbf{C}_\omega^{-1} \mathbf{H})^{-1} \mathbf{H}^T \mathbf{C}_\omega^{-1} \Xi(\zeta), \quad (13)$$

$$\hat{\sigma}^2 = \frac{\Xi(\zeta)^T (\mathbf{C}_\omega^{-1} - \mathbf{C}_\omega^{-1} \mathbf{H} (\mathbf{H}^T \mathbf{C}_\omega^{-1} \mathbf{H})^{-1} \mathbf{H}^T \mathbf{C}_\omega^{-1}) \Xi(\zeta)}{m-s+2}. \quad (14)$$

Point estimates of these modes if the training data $\{\zeta, \Xi(\zeta)\}$ is plugged in are

$$\hat{\beta} = \begin{pmatrix} 1.8344 \\ -0.1047 \end{pmatrix}, \quad \hat{\sigma}^2 = 11.906. \quad (15)$$

As with respect to parameter ω , in this theoretical exposition, the maximum likelihood estimate of this parameter is used, while parameters β and σ^2 being integrated out. (In practice, the objective implementation considered in [1] and [2] is preferable over the MLE estimate.) That is, the likelihood function $L(\omega; \Xi(\zeta))$ is

$$L(\omega; \Xi(\zeta)) = p(\Xi(\zeta) \mid \omega) = \int \int p(\Xi(\zeta) \mid \omega, \beta, \sigma^2) p(\beta, \sigma^2) d\beta d\sigma^2 \quad (16)$$

$$\propto \int \int (\sigma^2)^{-\frac{m}{2}-1} |\mathbf{C}_\omega|^{-\frac{1}{2}} \exp \left\{ -\frac{1}{2} \Xi(\zeta)^T \frac{1}{\sigma^2} \mathbf{C}_\omega^{-1} \Xi(\zeta) \right\} \exp \left\{ -\frac{1}{2\sigma^2} \beta^T \mathbf{H}^T \mathbf{C}_\omega^{-1} \mathbf{H} \beta \right\} \exp \left\{ -\frac{1}{2\sigma^2} (-2\Xi(\zeta)^T \mathbf{C}_\omega^{-1} \mathbf{H} (\mathbf{H}^T \mathbf{C}_\omega^{-1} \mathbf{H})^{-1} (\mathbf{H}^T \mathbf{C}_\omega^{-1} \mathbf{H}) \beta) \right\} d\beta d\sigma^2 = \quad (17)$$

$$= \int \int (\sigma^2)^{-\frac{m}{2}-1} |\mathbf{C}_\omega|^{-\frac{1}{2}} |\sigma^2 (\mathbf{H}^T \mathbf{C}_\omega^{-1} \mathbf{H})^{-1}|^{\frac{1}{2}} \exp \left\{ -\frac{1}{2\sigma^2} \Xi(\zeta)^T (\mathbf{C}_\omega^{-1} - \mathbf{C}_\omega^{-1} \mathbf{H} (\mathbf{H}^T \mathbf{C}_\omega^{-1} \mathbf{H})^{-1} \mathbf{H}^T \mathbf{C}_\omega^{-1}) \Xi(\zeta) \right\} \mathcal{N}_\beta \left((\mathbf{H}^T \mathbf{C}_\omega^{-1} \mathbf{H})^{-1} \mathbf{H}^T \mathbf{C}_\omega^{-1} \Xi(\zeta), \sigma^2 (\mathbf{H}^T \mathbf{C}_\omega^{-1} \mathbf{H})^{-1} \right) d\beta d\sigma^2$$

$$\propto \int \int |\mathbf{C}_\omega|^{-\frac{1}{2}} |\mathbf{H}^T \mathbf{C}_\omega^{-1} \mathbf{H}|^{-\frac{1}{2}} \mathcal{N}_\beta \left((\mathbf{H}^T \mathbf{C}_\omega^{-1} \mathbf{H})^{-1} \mathbf{H}^T \mathbf{C}_\omega^{-1} \Xi(\zeta), \sigma^2 (\mathbf{H}^T \mathbf{C}_\omega^{-1} \mathbf{H})^{-1} \right) \mathcal{IG}_{\sigma^2} \left(\frac{m-s}{2}, \frac{\Xi(\zeta)^T (\mathbf{C}_\omega^{-1} - \mathbf{C}_\omega^{-1} \mathbf{H} (\mathbf{H}^T \mathbf{C}_\omega^{-1} \mathbf{H})^{-1} \mathbf{H}^T \mathbf{C}_\omega^{-1}) \Xi(\zeta)}{2} \right) d\beta d\sigma^2$$

$$\propto (\Xi(\zeta)^T (\mathbf{C}_\omega^{-1} - \mathbf{C}_\omega^{-1} \mathbf{H} (\mathbf{H}^T \mathbf{C}_\omega^{-1} \mathbf{H})^{-1} \mathbf{H}^T \mathbf{C}_\omega^{-1}) \Xi(\zeta))^{-\frac{m-s}{2}} |\mathbf{C}_\omega|^{-\frac{1}{2}} |\mathbf{H}^T \mathbf{C}_\omega^{-1} \mathbf{H}|^{-\frac{1}{2}},$$

where s is the number of parameters in vector β .

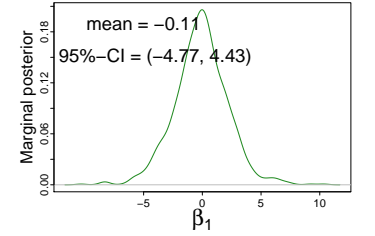
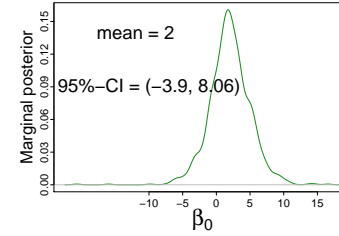
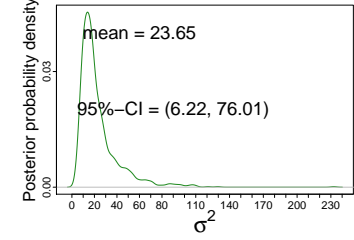


Figure 2: Plots of marginal posterior distributions for parameters σ^2 , β_0 and β_1 of the Gaussian process.

The corresponding log-likelihood function is

$$\begin{aligned} \ell(\omega; \Xi(\zeta)) &= \log(L(\omega; \Xi(\zeta))) = \log(p(\Xi(\zeta) | \omega)) = \text{const} \\ &\quad - \frac{m-s}{2} \log(\Xi(\zeta)^\top (\mathbf{C}_\omega^{-1} - \mathbf{C}_\omega^{-1} \mathbf{H} (\mathbf{H}^\top \mathbf{C}_\omega^{-1} \mathbf{H})^{-1} \mathbf{H}^\top \mathbf{C}_\omega^{-1}) \Xi(\zeta)) \\ &\quad - \frac{1}{2} \log(\det(\mathbf{C}_\omega)) - \frac{1}{2} \log(\det(\mathbf{H}^\top \mathbf{C}_\omega^{-1} \mathbf{H})). \end{aligned} \quad (18)$$

The plot of the log-likelihood is given in Figure 3. The emulator then can be constructed using the *estimated* parameters $\hat{\beta}, \sigma^2, \hat{\omega}$.

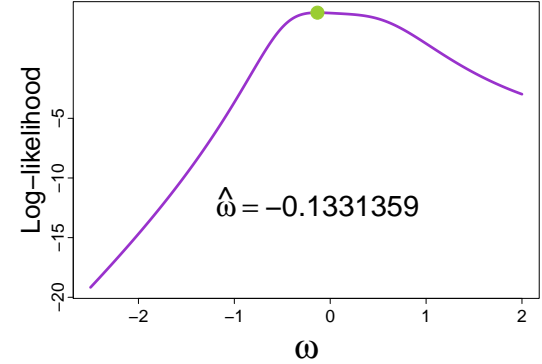


Figure 3: Log-likelihood $\ell(\omega; \Xi(\zeta))$. $\hat{\omega}$ is its maximum.

Accounting for uncertainty in parameters. Incorporation of *uncertainty* in parameters is **defined** as marginalizing over (or integrating out) the parameters. The attempt is to integrate out as many parameters as possible from the corresponding approximation 6. The reason for doing this is two-fold: (a) to formally account for the uncertainty in parameters expressed in their posterior distributions of these, and (b) to reduce the number of parameters which can not be integrated out and, therefore, must be estimated as, for instance, maximum likelihood or maximum *a posteriori*.

Marginalizing over β , the predictive distribution, emulator, becomes:

$$p(\Xi(\zeta') | \Xi(\zeta), \sigma^2, \omega) = \int p(\Xi(\zeta') | \Xi(\zeta), \beta, \sigma^2, \omega) p(\beta | \Xi(\zeta), \sigma^2, \omega) d\beta.$$

β is integrated out following the lemma from [5]. That is, rewrite

$$\Xi(\zeta') | \Xi(\zeta), \beta, \sigma^2, \omega = c(\zeta', \zeta)_\omega \mathbf{C}_\omega^{-1} \Xi(\zeta) + (\mathbf{H}' - c(\zeta', \zeta)_\omega \mathbf{C}_\omega^{-1} \mathbf{H}) \beta + \mathcal{N}(\vec{\mathbf{0}}, \sigma^2 (\mathbf{C}'_\omega - c(\zeta', \zeta)_\omega \mathbf{C}_\omega^{-1} c(\zeta, \zeta')_\omega)) \quad (19)$$

$$\beta | \Xi(\zeta), \sigma^2, \omega = (\mathbf{H}^\top \mathbf{C}_\omega^{-1} \mathbf{H})^{-1} \mathbf{H}^\top \mathbf{C}_\omega^{-1} \Xi(\zeta) + \mathcal{N}(\vec{\mathbf{0}}, \sigma^2 (\mathbf{H}^\top \mathbf{C}_\omega^{-1} \mathbf{H})^{-1}) \quad (20)$$

Summing everything up, we get the following predictive distribution with β marginalized over

$$\Xi(\zeta') | \Xi(\zeta), \sigma^2, \omega \sim \mathcal{N}(\vec{\mu}, \sigma^2 \vec{\Sigma}), \quad (21)$$

where

$$\vec{\mu} = \mathbf{H}' \hat{\beta} + c(\zeta', \zeta)_\omega \mathbf{C}_\omega^{-1} (\Xi(\zeta) - \mathbf{H} \hat{\beta}), \quad (22)$$

$$\vec{\Sigma} = \mathbf{C}'_\omega - c(\zeta', \zeta)_\omega \mathbf{C}_\omega^{-1} c(\zeta, \zeta')_\omega + (\mathbf{H}' - c(\zeta', \zeta)_\omega \mathbf{C}_\omega^{-1} \mathbf{H}) (\mathbf{H}^\top \mathbf{C}_\omega^{-1} \mathbf{H})^{-1} (\mathbf{H}' - c(\zeta', \zeta)_\omega \mathbf{C}_\omega^{-1} \mathbf{H})^\top. \quad (23)$$

By plugging-in estimates of parameters σ^2 and ω we acquire a new emulator — the one which has the uncertainty from parameters β incorporated. The expression for $\vec{\mu}$ contains $\hat{\beta}$ — the mode (and the mean) of the posterior distribution of β given in equation (11). While we may easily plug in the estimate $\hat{\beta}$, this is left so for the convenience and clearance in the expression for $\vec{\mu}$; and to highlight that the mean of the (21) coincides exactly with that of the emulator for which posterior modes of parameters β and σ^2 are used.

Theorem. $\Xi(\cdot) | \cdot, \Xi(\zeta), \sigma^2, \omega$ is the stochastic process.

Corollary. $\Xi(\cdot) | \cdot, \Xi(\zeta), \sigma^2, \omega$ is the Gaussian stochastic process emulator.

Accounting for variability in σ^2 , the following predictive distribution appears

$$p(\Xi(\zeta') | \Xi(\zeta), \omega) = \int p(\Xi(\zeta') | \Xi(\zeta), \sigma^2, \omega) p(\sigma^2 | \Xi(\zeta), \omega) d\sigma^2 = \int \mathcal{N}_{\Xi(\zeta') | \sigma^2, \omega}(\tilde{\boldsymbol{\mu}}, \sigma^2 \tilde{\boldsymbol{\Sigma}}) \mathcal{IG}_{\sigma^2 | \Xi(\zeta), \omega}(\tilde{\alpha}/2, \tilde{\beta}/2) d\sigma^2, \quad (24)$$

where $\tilde{\alpha} = m - s$ and $\tilde{\beta} = \Xi(\zeta)^\top (\mathbf{C}_\omega^{-1} - \mathbf{C}_\omega^{-1} \mathbf{H} (\mathbf{H}^\top \mathbf{C}_\omega^{-1} \mathbf{H})^{-1} \mathbf{H}^\top \mathbf{C}_\omega^{-1}) \Xi(\zeta)$.

$$p(\Xi(\zeta') | \Xi(\zeta), \omega) \propto \int \exp\left(-\frac{(\Xi(\zeta') - \tilde{\boldsymbol{\mu}})^\top \tilde{\boldsymbol{\Sigma}}^{-1} (\Xi(\zeta') - \tilde{\boldsymbol{\mu}})}{2\sigma^2} - \frac{\tilde{\beta}}{2\sigma^2}\right) (\sigma^2)^{-\frac{m}{2} - \frac{\tilde{\alpha}}{2} - 1} d\sigma^2 \propto \left(\tilde{\beta} + (\Xi(\zeta') - \tilde{\boldsymbol{\mu}})^\top \tilde{\boldsymbol{\Sigma}}^{-1} (\Xi(\zeta') - \tilde{\boldsymbol{\mu}})\right)^{-\frac{m+\tilde{\alpha}}{2}} \\ \propto \left(1 + \frac{1}{\tilde{\alpha}} \frac{(\Xi(\zeta') - \tilde{\boldsymbol{\mu}})^\top \tilde{\boldsymbol{\Sigma}}^{-1} (\Xi(\zeta') - \tilde{\boldsymbol{\mu}})}{\tilde{\beta}}\right)^{-\frac{m+\tilde{\alpha}}{2}}. \quad (25)$$

$$\Xi(\zeta') | \Xi(\zeta), \omega \sim \mathcal{T}_n(\tilde{\alpha}, \tilde{\boldsymbol{\mu}}, \gamma^2 \tilde{\boldsymbol{\Sigma}}) \quad (26)$$

is the multivariate (n -dimensional, because n is the length of the vector $\Xi(\zeta')$) Student's t -distribution. Once the estimate of ω is plugged-in, this distribution is an approximation to the computer model f . This distribution is parameterized by *degrees of freedom* $\tilde{\alpha}$, *location* parameter $\tilde{\boldsymbol{\mu}}$ and *shape matrix* $\gamma^2 \tilde{\boldsymbol{\Sigma}}$, where $\gamma^2 = \tilde{\beta}/\tilde{\alpha}$. If the training data are plugged in, then the estimate of $\gamma^2 = 4.6113$.

Theorem. $\Xi(\cdot) | \cdot, \Xi(\zeta), \omega$ is the stochastic process. We call $\Xi(\cdot) | \cdot, \Xi(\zeta), \omega$ - the T -process. This process is not Gaussian, but it can also serve as an emulator of f .

Corollary. A Gaussian process whose mean is $E(\Xi(\cdot) | \cdot, \Xi(\zeta), \omega)$ and covariance is $\text{Cov}(\Xi(\cdot) | \cdot, \Xi(\zeta), \omega)$ is the GASP emulator of function f .

Alternative way to write out the multivariate t -distribution is $\mathcal{T}_{\tilde{\alpha}}(\tilde{\boldsymbol{\mu}}, \gamma^2 \tilde{\boldsymbol{\Sigma}})$ which puts emphasis only on the parameters of this distribution. As with respect to properties of the multivariate Student's t -distribution, the mean of the distribution is its location parameter $\tilde{\boldsymbol{\mu}}$, while variance-covariance equals to $\frac{\tilde{\alpha}}{\tilde{\alpha}-2} \gamma^2 \tilde{\boldsymbol{\Sigma}}$.

The final approximation to the computer model that we have constructed in this chapter is shown in Figure 4. Perhaps, surprisingly, but this emulator, which within objective implementation accounts for uncertainty in all of its parameters except for just one, is a much better emulator than the one we started with, which had somewhat *ad hoc* parameters, with no acknowledging for uncertainty in their estimates. Visually, the approximation is much tighter while closely following the function itself. The first emulator, shown in Figure 1 is not as good, providing somewhat large *uncertainty* by its 95% credible area compared to very small uncertainty given by the last emulator. As we can see, the quality of the approximation highly depends on the choice of the procedure about estimates of parameters of the approximation.

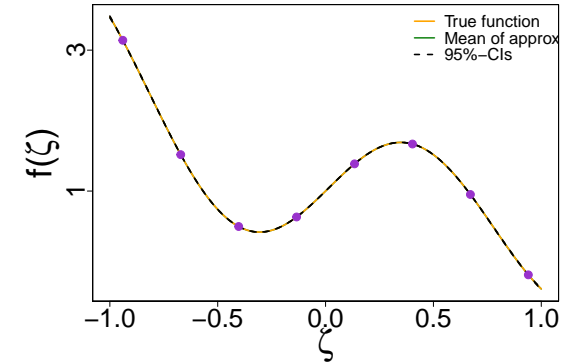


Figure 4: Gaussian process emulator which accounts for uncertainty in parameters β and σ^2 within its objective implementation.

Next chapter discusses formal ways to assessing the quality of a constructed statistical approximation of a computer model (or any other statistical model). While in this exposition we have compared the performance of the emulator as an approximation to the true function visually, in practice, assessing the quality of the constructed approximation is not as trivial, because the number of computer model *testing* data points available for assessment of the emulator is limited.

Chapter 1. Assessment of a statistical model

Protagoras argued that a man is a measure of all things. Within the decision-theoretic approach this is mathematically shown to be indeed so: scoring rules calculated for predictive model evaluation and comparison are SUBJECTIVE. This means that the choice of a scoring rule for model comparison affects the results of the comparison and, therefore, the decision on a model choice. What's more, the scoring rules are hardly interpretable. Recommendation is to, instead, employ three independent frequency estimates of the quality of model predictions: (1) empirical frequency coverage, (2) predictive bias, and (3) uncertainty (variability) in predictions.

In order to assess quality of the predictive performance of an emulator, in addition to training data, one must obtain *testing* data. Since the computer model is computationally slow, a small number of additional runs n has been obtained at inputs ζ^* resulting in outputs $\mathbf{f}(\zeta^*)$. As for the illustrative example $n = 7$ testing inputs are shown in Figure 5. These are chosen right in-between the pairs of the neighbouring training inputs. In this example (1) all testing outputs are captured within the 95% central credible area provided by the emulator, indicating good performance, (2) the bias — discrepancy between true values and the mean (on average or the maximum) of the stochastic approximation is rather small, (3) the average (or maximum) length of credible intervals is short, so that the emulator provides tight approximations for testing points as well as along the whole domain of its inputs.

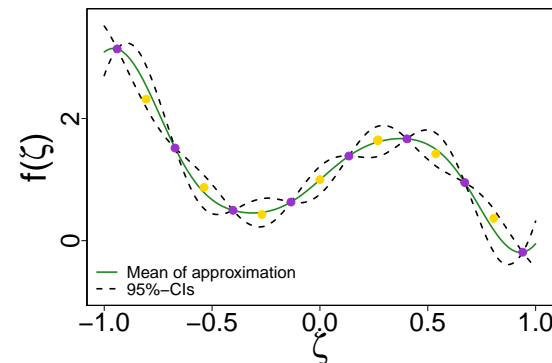


Figure 5: The first emulator constructed in the previous chapter. Training points are those coloured in purple. The emulator is an interpolator at these points. Testing points are coloured in yellow.

Denote the marginal predictive distribution for the i th input ζ_i^* as p_i . Let the mean of this distribution be μ_i and variance of the distribution σ_i^2 . Formally the *frequency estimates* are: *empirical frequency coverage* (EFC), the proportion of true values captured by the corresponding predictive distribution within its central 95%-credible interval²

$$\text{EFC} = \left(\frac{\sum_{i=1}^n I_{f(\zeta_i^*) \in \text{CI}_i}}{n} \right) 100\% \in (0, 100), \quad (27)$$

where CI_i is the 95% credible area; *root-mean-square predictive error* (RMSPE) — the estimate of the bias³ —

$$\text{RMSPE} = \sqrt{\frac{\sum_{i=1}^n (f(\zeta_i^*) - \mu_i)^2}{n}} \in (0, \infty), \quad (28)$$

and the average length of the 95% credible intervals over n points

$$\overline{\text{L}_{\text{CI}}} = \frac{\sum_{i=1}^n \text{CI}_i}{n} \in (0, \infty). \quad (29)$$

For normal predictive distributions the last estimate is defined as $\overline{\text{L}_{\text{CI}}} = \frac{\sum_{i=1}^n 2 \times 1.96 \sigma_i}{n}$. Since $\overline{\text{L}_{\text{CI}}}$ does not have true values involved in its calculation, one may also obtain this estimate over greater number of possible inputs — to investigate the overall emulator's uncertainty in predictions. Useful is to verify that the maximal values among

²In my experience 95% threshold for a credible area has served as a useful nominal value. This choice is subjective. Other nominal values may be employed.

³Alternatively, one may calculate the average absolute mean bias as $\frac{\sum_{i=1}^n |f(\zeta_i^*) - \mu_i|}{n} \in (0, \infty)$.

predictive credible intervals are acceptable, and are not substantially greater than average ones. Empirical values of EFC which correspond to the *nominal* (theoretical) value of 95% are desired to be close to each other. However, one needs to subjectively assess if, say, 80% empirical coverage versus corresponding 95% area is good enough. For other two frequency estimates, it is desirable that RMSPE and \overline{L}_{CI} are small.

Three estimates provide clear and intuitive *interpretation* of what they mean in the evaluation of the predictive performance of a model. Altogether, they *attempt* to provide *objective* assessment because these are *frequency* estimates. These estimates must be considered independently. The attempt to collapsing them into one leads to an *improper* criterion (an improper scoring rule)⁴, that is the rule that favours a wrong forecast rather than the true one. If one emulator (or any statistical model) is substantially better than another one, then comparison with the frequency estimates reveal this fact. If two statistical models are *on par* in their performance, then comparison between the models with any estimate or criterion is challenging.

There have been attempts to find a one single criterion which would make predictive model comparison possible. The methodology of *scoring rules* have been proposed. Turns out, this task is not achievable. Two main properties of why this is so are: scoring rules are subjective; and being evaluated, the scores are hardly interpretable.

Definition. A *scoring rule* S is a function which takes a probabilistic forecast, that is, a random variable ζ whose distribution is F , and a true forecast, random variable ω , producing a real-valued *score* $S(\zeta, \omega)$.

Definition. A scoring rule is said to be *proper* if its maximum is reached at the true random variable, that is $\max(S(\zeta, \omega)) = S(\omega, \omega)$.⁵

Theorem. Every scoring rule is *subjective*, that is no scoring rule is equivalent to any other scoring rule.

Proof. Consider two scoring rules S_1 and S_2 that are distinct (that is functions $S_1(\cdot, \cdot) \neq S_2(\cdot, \cdot)$) and two forecasts $\zeta_1 \neq \zeta_2$ in distribution, and the true variable ω , then $f(S_1(\zeta_1, \omega), S_1(\zeta_2, \omega)) \neq f(S_2(\zeta_1, \omega), S_2(\zeta_2, \omega))$, where $f(x_1, x_2)$ is a function of two arguments used to compare scores, that is, for model comparison $f(x_1, x_2)$ is either the difference $x_1 - x_2$ or the ratio of corresponding scores x_1/x_2 . ■

This theorem holds, even if S_2 is a monotone increasing transformation of S_1 . Therefore, the theorem asserts that conclusion on *how much better* one model is compared to another is not possible to make.

Illustrative example. The illustrative example demonstrates that within a scoring rule the assessment may be inadequate. In particular, a popular logarithmic scoring rule is chosen to demonstrate its deficiencies. Let p denote a predictive density of a forecast, real-valued random variable ζ . Let the random variable ω have a degenerate distribution, that is $P(\omega = v) = 1$. Therefore, v is the true realization. Consider a *logarithmic scoring rule*

$$\log S(\zeta, \omega = v) = \log p(v). \tag{30}$$

Consider a class of normal predictive distributions with mean μ and standard deviation σ . Then for $\zeta \sim \mathcal{N}(\mu, \sigma^2)$ with predictive density p , the logarithmic scoring rule with respect to v is

$$\log S(\zeta, \omega = v) = -\frac{\log(2\pi\sigma^2)}{2} - \frac{(v - \mu)^2}{2\sigma^2}. \tag{31}$$

Assuming that the true value v equals to zero, a plot with contourlines of scores which in this case depends simply on the mean and standard deviation of the random variable ζ is shown in Figure 6. The contourlines show that since there are infinitely many points along a particular contourline, infinitely many distributions exist which have exactly the same value of a log-score. The same value of a score makes these distributions to be considered equally good predictive distributions, but these distributions may be very different from each other.

⁴Predictive model choice criterion (PMCC) is improper.

⁵If in this definition the scoring rule is maximized, then the scoring rule is said to be *positively oriented*. Alternatively, one can choose to minimize a scoring rule, then the rule is said to be *negatively oriented*.

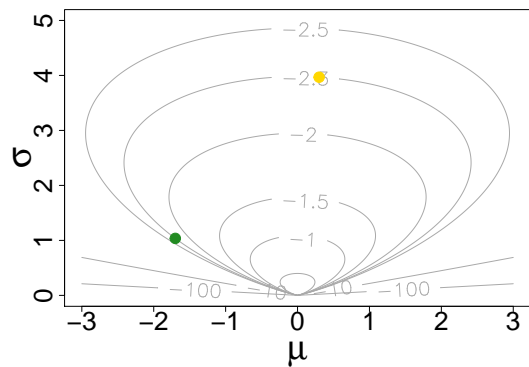


Figure 6: Contourlines of the logarithmic score given for a class of normal distributions characterized by two parameters: mean μ and standard deviation σ . Two coloured points correspond to two distributions.

As example, consider two points on the plot: the green point corresponds to the distribution $\mathcal{N}(\mu = -1.7, \sigma = 1.038363)$ and orange point corresponds to $\mathcal{N}(\mu = 0.3, \sigma = 3.967767)$. Indeed, the distributions are very different: the first distribution has a large bias — the mean of the distribution is far from the true value of zero. The orange distribution has bias more than 5 times smaller of that of the green one. However, the orange distribution has standard deviation more than three times greater than that of the green one, which translates to variance being more than 14 times greater. Yet, two distributions are assigned exactly the same value of a logarithmic score.

Consider another example of two similar forecasts which could have resulted from two other models, namely $\mathcal{N}(\mu = 0.9, \sigma = 1)$ and $\mathcal{N}(\mu = 1, \sigma = 0.9)$. Providing virtually the same information on the predictive distribution in practice, the scores, as compared to the true value of zero, are -1.323939 and -1.430862 respectively, thus differentiating between the two distributions much more than one desires.

These examples demonstrate that logarithmic score can not distinguish between very different distributions, assigning the same scores to considerably different shapes; and at the same time assigning much more different value to very similar distributions. It may be demonstrated that other scores provide similar but different contourplots for the class of normal distributions, behaving in analogous way as have been observed with the logarithmic scoring rule.

Chapter 2. Computer model with multivariate output

Computer models often produce multivariate output for every single run of the model. There have been attempts to account for correlation among outputs in the construction of a Gaussian process emulator of such a model with the goal of achieving a more accurate emulator. Both, theoretical evidence and simulations are presented here which demonstrate that multivariate emulator does not lead to “better” (that is, more accurate, precise or less uncertain) emulation results compared to independent modeling of each component of the output.

Suppose a computer model produces bivariate output, which may be described as two smooth functions $f_1(\cdot)$ and $f_2(\cdot)$. Suppose the simulator data is obtained at n d -dimensional points $\zeta = (\zeta_1, \dots, \zeta_n)$ from the input space $\mathcal{X} = \mathcal{X}_1 \times \dots \times \mathcal{X}_d$. In other words, at the i th input ζ_i , the output of a simulator is a two-dimensional vector $(f_1(\zeta_i), f_2(\zeta_i))^T$.

The simplest emulator of such bivariate output of a computer model may be constructed choosing independent emulators of each output. That is, if ι th function is assumed to come from a Gaussian process characterized by parameters $\theta_\iota = (\beta_\iota, \sigma_\iota^2, \omega_\iota)$

$$\Phi_\iota(\cdot) \mid \theta_\iota \sim \mathcal{GP}(\mu_\iota(\cdot), \text{cov}_\iota(\cdot, \cdot)). \quad (32)$$

A computer model bivariate output from the simulator by $\Phi_1(\zeta) = (f_1(\zeta_1), \dots, f_1(\zeta_n))$ and $\Phi_2(\zeta) = (f_2(\zeta_1), \dots, f_2(\zeta_n))$.

$$\Phi_\iota(\cdot) \mid \cdot, \theta_\iota, \Phi_\iota(\zeta) \sim \mathcal{GASP}(\mu_\iota^*(\cdot), \sigma_\iota^{*2}(\cdot, \cdot)). \quad (33)$$

The resulting independent emulators are the posterior predictive distributions which are independent normal distributions $\mathcal{N}(\mu_\iota^*(\zeta'), \sigma_\iota^{*2}(\zeta', \zeta'))$, $\iota = 1, 2$.

We want to compare such emulators to an emulator of a joint bivariate Gaussian process $(f_1(\cdot), f_2(\cdot))^T \mid \zeta$ in terms of their predictive distributions. This is especially of interest if $f_1(\cdot)$ and $f_2(\cdot)$ are suspected to be highly correlated. The answer if joint modeling is beneficial compared to independent modeling of multivariate output for emulation purpose is no.

Suppose $f_2(\cdot) = \lambda f_1(\cdot) + \gamma$, where λ and γ are unknown constants. This is a situation of two perfectly correlated functions, for which one might expect to achieve the most value from multivariate modeling. (Example of such functions is given in Figure 7.) The following proposition establishes this is false.

Proposition. Suppose the bivariate simulator f_1 and f_2 is observed at the input points ζ and that $f_2 = \lambda f_1 + \gamma$. GASP approximations that would result from emulating each output independently are given by (33). Let $c_1(\cdot, \cdot)$ and $c_2(\cdot, \cdot)$ have the same functional form, but with possibly different parameters ω_1 and ω_2 . Let each Gaussian process mean be linear with the same vector of regression covariates $\mathbf{h}(\cdot)^T$, such that $h_0(\cdot) = 1$ defining an intercept, i.e., $\mu_i(\cdot) = h(\cdot)\beta^{(i)}$. Then the maximum likelihood estimates of the model parameters have the following properties:

The estimates of the correlation parameters are the same, i.e. $\hat{\omega}_1 = \hat{\omega}_2$.

The estimates of the variances satisfy $\hat{\sigma}_2 = \lambda \hat{\sigma}_1$.

The estimates of the intercepts satisfy $\hat{\beta}_0^{(2)} = \lambda \hat{\beta}_0^{(1)} + \gamma$ and, for the other regression coefficients, $\hat{\beta}_j^{(2)} = \lambda \hat{\beta}_j^{(1)}$.

Corollary. Let the predictive distribution of the i 's output at a new point ζ' be denoted as

$$\Phi_i(\zeta') \sim \mathcal{N}(\mu_i^*(\zeta'), \sigma_i^*(\zeta', \zeta')). \quad (34)$$

At a new point ζ' the predictive means and variances of the two independent emulators $\Phi_1(\cdot)$ and $\Phi_2(\cdot)$ (with maximum likelihood estimates of parameters) are related as

$$\mu_2^*(\zeta') = \gamma + \lambda \mu_1^*(\zeta'), \quad (35)$$

$$\sigma_2^{*2}(\zeta', \zeta') = \lambda^2 \sigma_1^{*2}(\zeta', \zeta'). \quad (36)$$

Theorem. Let the joint predictive distribution of the two functions be specified conditionally, i.e.,

$$p(\Phi_2(\zeta'), \Phi_1(\zeta') \mid \zeta') = p(\Phi_2(\zeta') \mid \Phi_1(\zeta'), \zeta') p(\Phi_1(\zeta') \mid \zeta'). \quad (37)$$

Then marginal emulators $\Phi_1(\cdot)$ and $\Phi_2(\cdot)$ of each output $f_1(\cdot)$ and $f_2(\cdot)$ coincide in either their independent or joint modeling if maximum likelihood or maximum a posteriori estimates of parameters of the processes are employed.

Therefore, no benefit is due to construction of a multivariate emulator of a multivariate output computer model compared to independent emulation of each output. Proofs of the mathematical statements in this chapter are given in my thesis.

Chapter 3. Approximation to a system of computer models

Direct approximation of a system which consists of several computer models is difficult for computational and logistical reasons. The methodology of the linked Gaussian approximation has been demonstrated to be a successful alternative. This is outlined in this chapter.

Consider two computer models, f_1 and f_2 , whose inputs and outputs are real-valued; and for which the corresponding Gaussian process approximations (introduced in chapter) may be constructed. Let Ξ and Υ be the corresponding emulators. The emulator $\Xi(\cdot)$ at any new input, given pairs of model runs $\{\zeta, \Xi(\zeta) = \mathbf{f}_1(\zeta)\}$ and

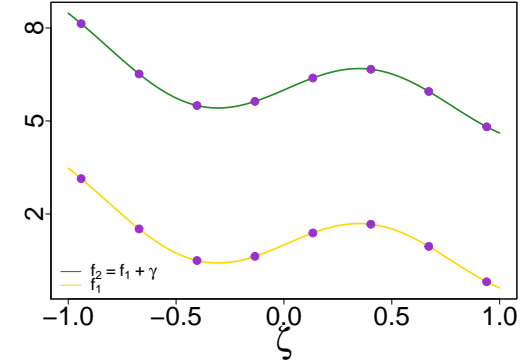


Figure 7: Two curves represent two output functions $f_1(\zeta)$ and $f_2(\zeta)$ of a computer model. The circles correspond to the design input-output points that are used to construct the joint emulator.

this emulator's vector of parameters θ_{Ξ} , is $\Xi(\cdot) \mid \cdot, \Xi(\zeta), \theta_{\Xi} \sim \mathcal{GASP}(\mu_{\Xi}^*(\cdot), \sigma_{\Xi}^{*2}(\cdot, \cdot))$. Likewise, the emulator $\Upsilon(\cdot)$ of the model f_2 , given $\{\kappa, \Upsilon(\kappa) = \mathbf{f}_2(\kappa)\}$ and θ_{Υ} , is $\Upsilon(\cdot) \mid \cdot, \Upsilon(\kappa), \theta_{\Upsilon} \sim \mathcal{GASP}(\mu_{\Upsilon}^*(\cdot), \sigma_{\Upsilon}^{*2}(\cdot, \cdot))$. The parameters θ_{Ξ} and θ_{Υ} here are assumed known. In practice, they are estimated.

Having introduced the framework for an emulator of a single computer model I consider a question how to approximate a system of two computer models $f_1 \circ f_2$ defined as the composite model if only the two submodels $f_1(\cdot)$ and $f_2(\cdot)$ may be independently computed at a few data points, while the composition $f_1 \circ f_2(\cdot)$ can not be observed. For any new input u to the system of computer models $f_1 \circ f_2$, the approximation is defined as $\Xi \circ \Upsilon$, which is denoted as Φ , i.e. $p(\Phi(u) \mid \Xi(\zeta), \Upsilon(\kappa), \theta_{\Xi}, \theta_{\Upsilon}, u) = \int p(\Xi(\Upsilon(u)) \mid \Xi(\zeta), \Upsilon(u), \theta_{\Upsilon})p(\Upsilon(u) \mid \Upsilon(\kappa), \theta_{\Xi})d\Upsilon(u)$.

Definition. The variable $\xi = \Phi(u) \mid \Xi(\zeta), \Upsilon(\kappa), \theta_{\Xi}, \theta_{\Upsilon}, u$ is called the linked emulator.

Theorem. Linked Gaussian emulator (or Gaussian approximation to a collection of random variables $\Phi_t, t \in T$, where T is an index set) is a stochastic process.

Proof. The linked Gaussian emulator is a Gaussian stochastic process by construction: the mean, variances and correlation structure among *all* variables of this process may be specified. Corresponding formulae rely on the laws of total mean, total variance and total covariance respectively.⁶

Illustrative example. Two functions $f_2(x) = \cos(2x)$ in the domain $x \in [-3, 3]$ and $f_1(z) = \cos(z/2)$ in the domain $z \in [-2, 5]$ are test functions. Model f_2 has been evaluated at 8 training inputs; f_1 — at 5 training data points. The simulators, their designs and emulators are shown in 8. The resulting linked Gaussian approximation to $f_1 \circ f_2$ is shown in 9.

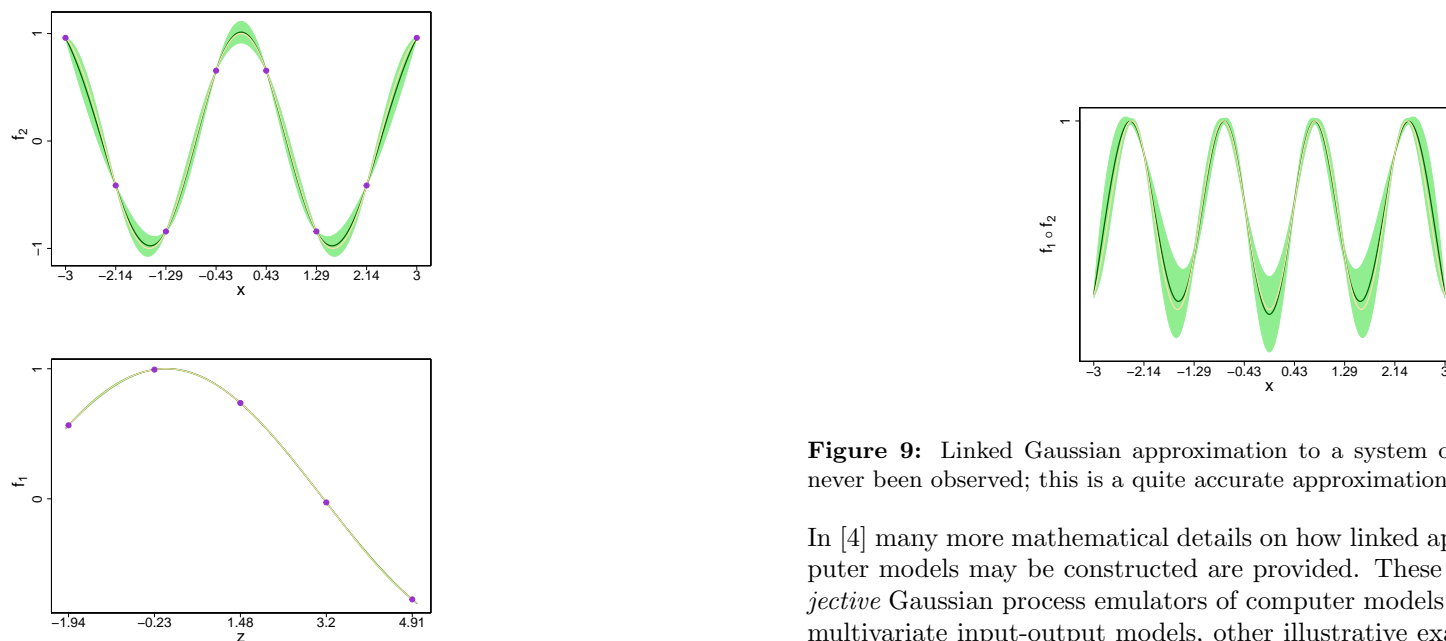


Figure 8: Two test functions: $f_1(x)$ and $f_2(x)$ along with their Gaussian approximations. The dark green lines are the means of emulators, the yellow lines are the true functions. 95% central credible area is shown with the green shaded regions. The purple circles are at the design points.

Figure 9: Linked Gaussian approximation to a system of simulators which has never been observed; this is a quite accurate approximation to the system.

In [4] many more mathematical details on how linked approximations of computer models may be constructed are provided. These include linking of *objective* Gaussian process emulators of computer models, analysis of linking of multivariate input-output models, other illustrative examples, and examples of doing so with more realistic computer models of volcano pyroclastic flows and volcano ash plume dispersion.

⁶See details in [3]

Chapter 4. Calibration of a computer model

If observations corresponding to the output of the model are collected, a theoretical model may be assessed on the agreement to the collected data. Moreover, given a computer model and collected data, one might inquire which values of inputs to the model could have generated the collected data; thus, performing calibration of a model.

Calibration is, therefore, analogous to finding an inverse image of function $f : Z \rightarrow Y$, given a set of n values $\mathbf{y} = (y_1, \dots, y_n) \in Y$, that is $f^{-1}(\mathbf{y}) = \{\zeta = (\zeta_1, \dots, \zeta_n) \in Z : f(\zeta_i) = y_i \ \forall i = 1, \dots, n\}$. However, this is indeed only an analogy, since in practice calibration involves (a) a computationally intensive computer model or a system of such models, and (b) noisy collected data which does not arise from the computer model itself but is obtained in either an experiment or observed in nature. Calibration framework which accounts for these two levels of complexity is presented in this study.

Suppose that n experimental data points $\mathbf{y} = (y_1, \dots, y_n)^T$ have been collected; each i th point $y_i, i = 1, \dots, n$ has been obtained under some input conditions $\zeta_i^* \in Z$. These input conditions are assumed to belong to the same input space $\zeta_i^* \in Z$ as the computer model inputs.

Each i th input ζ_i^* is comprised of a set of control variables $\zeta_i^{*,c} \in \mathcal{R}^{d_1}$ (say, first d_1 inputs) and a set of calibration variables $\zeta_i^{*,p} \in \mathcal{R}^{d_2}$, the set of the rest $d_2 = d - d_1$ variables; i.e. $\zeta^* = \{\zeta^{*,c}, \zeta^{*,p}\}$. A subset of control variables $\zeta^{*,c}$ is set and known. A subset of $\zeta^{*,p}$ is unknown. However, the unknown conditions are assumed to be the same for all data points, i.e. $\zeta_1^{*,p} = \zeta_2^{*,p} = \dots = \zeta_n^{*,p} = \alpha$.

Calibration aims at identifying such values of parameters α , that computer model output (that is of its approximation) $\Xi(\zeta^*) = (\Xi(\zeta_1^{*,c}, \alpha), \dots, \Xi(\zeta_n^{*,c}, \alpha))^T$ obtained at inputs ζ^* with parameters α plugged in matches closely experimental data \mathbf{y} .

Conditional on the computer model data $\{\zeta, \Xi(\zeta)\}$, within the emulation framework, computer model output $\Xi(\zeta^*) = \Xi(\zeta^{*,c}, \alpha)$ is given by the predictive distribution

$$\Xi(\zeta^{*,c}, \alpha) \mid \zeta^{*,c}, \alpha, \zeta, \Xi(\zeta) = (\Xi(\zeta_1^{*,c}, \alpha), \dots, \Xi(\zeta_n^{*,c}, \alpha))^T \mid \zeta^{*,c}, \alpha, \zeta, \Xi(\zeta). \quad (38)$$

This yields the likelihood $L(\alpha; \mathbf{y}) = p(\mathbf{y} \mid \zeta^{*,c}, \alpha, \zeta, \Xi(\zeta))$ for to-be calibrated parameters α , which is given by the distribution of experimental data $\mathbf{y} = (y_1, \dots, y_n)^T$ observed at ζ^*

$$L(\alpha; \mathbf{y}) = p(\Xi(\zeta_1^{*,c}, \alpha) = y_1, \dots, \Xi(\zeta_n^{*,c}, \alpha) = y_n \mid \zeta^{*,c}, \alpha, \zeta, \Xi(\zeta)). \quad (39)$$

To complete specification of the model, one must specify the prior distribution for unknown parameters α . In the calibration setting prior distribution $p(\alpha)$ on parameters α is, though, possibly vague, but an expert-elicited prior.

The solution to the calibration problem is the posterior distribution $p(\alpha \mid \mathbf{y}, \zeta^{*,c}, \zeta, \Xi(\zeta))$. For convenience the short notation is $p(\alpha \mid \mathbf{y})$ which implicitly assumes conditioning on known controlled settings $\zeta^{*,c}$ and computer model data, that is, pairs of observations $\{\zeta, \Xi(\zeta)\}$. Posterior distribution is then obtained

$$p(\alpha \mid \mathbf{y}) = \frac{L(\alpha; \mathbf{y})p(\alpha)}{\int L(\alpha; \mathbf{y})p(\alpha) d\alpha}. \quad (40)$$

This distribution is well-defined since the prior $p(\alpha)$ is defined on closed space of parameters, and, therefore, the posterior is proper. Thus, the numerical solution (40) to the calibration problem always exists.

Qualitative assessment of calibration results. Ideally, posterior $p(\alpha \mid \mathbf{y})$ identifies a small set of most probable parameters α such that computer model data evaluated at these values matches closely collected data \mathbf{y} . In the unfortunate scenario, one may find that even being evaluated at the most probable values of α , computer model output may not match the collected data, thus, resulting in no agreement between the theory and the data. Such a case would indicate a crucial problem with either a theoretical computer model or a prior on calibration parameters $p(\alpha)$, or both; demanding for more inquiry into domain knowledge.

The following posterior corresponding to the calibration of α is

$$p(\alpha \mid \mathbf{y}) \propto \mathcal{N}(\mathbf{y} \mid \mu_f^*(\mathbf{t}, \alpha), \sigma_f^{*2}((\mathbf{t}, \alpha), (\mathbf{t}, \alpha))) p(\alpha). \quad (41)$$

Probabilistic inversion (41) involves numerical computation of the inverse of the covariance matrix $\sigma_f^{*2}((\mathbf{t}, \alpha), (\mathbf{t}, \alpha))$ from the log-likelihood for different values of α . This matrix may provide strong correlation between points, such that its correct numerical inversion becomes a challenging problem even if the number of experimental

data points is small. Numerical errors coming from the inversion of the matrix in log-densities may preclude obtaining the correct solution. In order to overcome this computational problem the information about the correlation structure in the likelihood may be omitted, effectively resulting in the following distribution (42) as an approximation to the true posterior (41)

$$p(\boldsymbol{\alpha} | \mathbf{y}) \propto \mathcal{N}(\mu_f^*(\mathbf{t}, \boldsymbol{\alpha}), \text{diag}(\sigma_f^{*2}((\mathbf{t}, \boldsymbol{\alpha})))) p(\boldsymbol{\alpha}). \quad (42)$$

This distribution does not take information about GASP sample paths into account but only uses the information from the *marginals* of the corresponding likelihood. Computing this approximated solution is trivial.

Proposition. Let the GASP correlation function be dependent of and decaying with distance along its inputs, then the further apart from each other experimental data inputs are, the more their corresponding cross-correlations are approaching zero. Therefore, the closer approximation (42) is to the true posterior (41).

Proof. Following exposition of the GASP, correlation between two inputs ζ_1 and ζ_2 is $c(\zeta_1, \zeta_2) = \prod_{j=1}^d c(\zeta_{1j}, \zeta_{2j})$. For the j th coordinate correlations depend only on the distance $c(\zeta_{1j}, \zeta_{2j}) = c(|\zeta_{1j} - \zeta_{2j}|)$, whose properties are $c(\zeta_{1j}, \zeta_{2j}) \rightarrow 1$ and $c(\zeta_{1j}, \zeta_{2j}) \rightarrow 0$ as $|\zeta_{1j} - \zeta_{2j}| \rightarrow \infty$. The last property proves the proposition.

Calibration with respect to noisy collected data. In practice collected data is not perfect and is observed with noise ϵ . That is, true values of \mathbf{y} are not known but may be estimated from the collected data. Collected data is denoted as \mathbf{y}^E . In the previous section the ideal scenario of the absence of any noise in the collected data was assumed, resulting in that $\mathbf{y}^E = \mathbf{y}$. In practice, true data \mathbf{y} is not known, but posterior estimate on the true values of \mathbf{y} given the collected data, i.e. the distribution $p(\mathbf{y} | \mathbf{y}^E)$ may be obtained.

The more information about the error in the data \mathbf{y}^E is available, the better the estimate of posterior $p(\mathbf{y} | \mathbf{y}^E)$ is going to be. The common knowledge is that the larger the noise in the data, the more replicates is required to be taken in order to estimate and reduce the noise. Ultimately, the quality of collected data matters most.

Having obtained posterior on the true values $\mathbf{y} | \mathbf{y}^E$, calibration becomes

$$p(\boldsymbol{\alpha} | \mathbf{y}^E) = \int_{\mathbf{y}} p(\boldsymbol{\alpha} | \mathbf{y}) p(\mathbf{y} | \mathbf{y}^E) d\mathbf{y}, \quad (43)$$

where $p(\boldsymbol{\alpha} | \mathbf{y})$ is given by equation (40).⁷

This form of the posterior reflects that one may be unsure about the truth, true values \mathbf{y} . The more sure one is about the truth, the closer calibration framework (43) is to (40). Calibration (43) also states that one is able to calibrate only to the extent that the observed data allow us to do it. This acknowledges the possibility that calibration task with corrupt data (e.g., data with large noise or biased data or data with limited information on how the data has been obtained) does not allow to perform calibration.

Chapter 5. Censoring for a computer model with zero-inflated output

Computer model of a volcano pyroclastic flow, given a set of initial conditions, produces an output, maximum height of the flow at thousands of geographical locations. The output is non-negative and often results in exact zero, thus, indicating the absence of a flow, and resulting in that the zero-height value of an output has a non-zero probability to occur, as opposed to all other simulator output values. In order to account for these features of the output, the methodology of a censored GASP approximation to such a computer model is employed. Customarily employed, usual GASP, does not allow to construct an approximation which would incorporate these features of the simulator.

⁷Computationally approximation to $p(\boldsymbol{\alpha} | \mathbf{y}^E)$ may be obtained through sampling from the distribution $p(\mathbf{y} | \mathbf{y}^E)$. For each of the r independent samples $\mathbf{y}_{smp_i} \sim p(\mathbf{y} | \mathbf{y}^E)$, $i = 1, \dots, r$, the distribution $\boldsymbol{\alpha} | \mathbf{y}_{smp_i}$ may be obtained. All combined (to integrate out $\mathbf{y} | \mathbf{y}^E$), marginal $p(\boldsymbol{\alpha} | \mathbf{y}^E)$ is the posterior estimate of the unknown parameter $\boldsymbol{\alpha}$. This approach is easily parallelizable, since $p(\boldsymbol{\alpha} | \mathbf{y}_{smp_i})$ is independent from $p(\boldsymbol{\alpha} | \mathbf{y}_{smp_j})$ for all independent samples \mathbf{y}_{smp_i} , $i = 1, \dots, r$.

Dealing with non-zero probability of a zero-inflated output is inconvenient, and, therefore, it's tempting to ignore this important information this way or another. Two possibilities to doing this exist: first, ignoring zero-output data from the construction of the emulator of the model; second, setting an assumption that the output is a smooth function over the whole range of its values. Both assumptions are unappealing, because either important information from the zero-output data is completely ignored, or the assumption of a smooth function which a priori is known not to hold true is employed. These are the reasons to construct a censored GASP approximation which incorporates appropriate assumptions and its construction involves all possible data.

Let the computer model f be evaluated at n_1 inputs $\zeta^O = (\zeta_1^O, \dots, \zeta_{n_1}^O)$ resulting in respective outputs $f(\zeta_1^O), \dots, f(\zeta_{n_1}^O)$. In addition to this simulator data, at n_2 inputs ζ^C simulator $f(\cdot)$ is known to produce values in the range (a, b) but values of the simulator $f(\cdot)$ at inputs $\zeta^C = \{\zeta_i^C\}_{i=1}^{n_2}$ are not known and are not available. That is, it is known only that $a < f(\zeta_i^C) < b$ for each i th input in a set ζ^C . Originally this work has been motivated by construction of the emulator for the model of the height of a pyroclastic flow, which demands that $a = 0$ and $b = \infty$. Other possibilities for choices of a and b are not considered in this work. Therefore, the output from the computer model itself called $f_a(\cdot) = \max(a, f(\cdot))$, is a censored function. In other words, for each single input ζ_i the computer model output $f_a(\zeta_i)$

$$f_a(\zeta_i) = \begin{cases} f(\zeta_i), & \text{if } f(\zeta_i) > a \\ a, & \text{otherwise} \end{cases} . \quad (44)$$

Let the latent function $f(\cdot)$ be assigned a Gaussian process prior, i.e. $\Xi(\cdot) \sim \mathcal{GP}(\mu(\cdot), \sigma^2 c(\cdot, \cdot))$. Then, conditional on computer model data $\Xi(\zeta^O)$ and $\Xi(\zeta^C)$ the GASP emulator of f is $\Xi(\cdot) \mid \cdot, \Xi(\zeta^O), \Xi(\zeta^C), \theta \sim \mathcal{GASP}(\mu^*(\cdot), \sigma^{*2}(\cdot, \cdot))$, the predictive distribution $\Xi(\zeta') \mid \Xi(\zeta^O), \Xi(\zeta^C)$ at a new point ζ' .

In this work outputs $\Xi(\zeta^C)$ are not available. It is only known that $\Xi(\zeta_i^C) < a$ for each i th input ζ_i^C . The following predictive distribution $\Xi_a(\zeta') \mid \Xi(\zeta^O), \Xi(\zeta^C) < \mathbf{a}$ is of interest to us. Here \mathbf{a} is an m -dimensional vector of ones (of length of a vector ζ^C) times constant a .

For a collection of m d -dimensional inputs $\zeta = \{\zeta_1, \dots, \zeta_m\}$ (given parameters of the GASP θ) statistical approximation to the outputs at inputs ζ is $\Xi_a(\zeta) = \{\Xi_a(\zeta_1), \dots, \Xi_a(\zeta_m)\}$ has a mixture distribution. This distribution consists of three parts: (a) an absolutely continuous distribution with respect to the reference measure that is the sum of a unit point-mass at the m -dimensional vector $(a, \dots, a)^T$, (b) Lebesgue measure on $\mathbb{R}_{>a}^m$ and (c) mixed type measures on $\mathbb{R}_{>a}^{|\bar{S}|} \times a^{|\bar{S}|}$, where S is a subset of indices of m variables, i.e. $S \subseteq 1 : m = \{1, \dots, m\}$ and $\bar{S} = \{1, \dots, m\} \setminus S$ is a complementary subset of S .

The approximation $\Xi(\cdot) \mid \cdot, \Xi(\zeta^O), \Xi(\zeta^C), \theta$ induces another stochastic process $\Xi_a(\cdot)$ — approximation to the simulator $f_a(\cdot)$. Stochastic process $\Xi_a(\cdot)$ is specified by a real-valued threshold a and a triple of a mean function $\mu(\cdot)$, variance σ^2 and correlation function $c(\cdot, \cdot)$.

Denote $P(1 : m)$ as a powerset of $1 : m$ and cardinality $|S| = \mathbf{card} S$, then vector $\Xi_a(\zeta)$ has the following joint distribution

$$p(\Xi_a(\zeta)) = \begin{cases} \int_{-\infty}^a \dots \int_{-\infty}^a p(\Xi(\zeta)) \prod_{i=1}^m d\Xi(\zeta_i), & \Xi_a(\zeta_i) = a, i = \overline{1, m}, \\ \forall S \in P(1 : m) \setminus \{\emptyset, 1 : m\} & \Xi_a(\zeta_i) > a, i \in \bar{S} \\ p(\Xi(\zeta)_{\bar{S}}) \int \dots \int_{\Delta} p(\Xi(\zeta)_S \mid \Xi(\zeta)_{\bar{S}}) d\Xi(\zeta)_{\bar{S}}, & \Xi_a(\zeta_i) = a, i \in S \\ p(\Xi(\zeta)) & \Xi_a(\zeta_i) > a, i = \overline{1, m}, \end{cases} \quad (45)$$

where $\Delta = (-\infty, a)^{\mathbf{card} S}$ and $\overline{1, m} = 1, \dots, m$. Here $\Xi(\zeta)$ has a multivariate normal distribution $\mathcal{N}(\mu^*(\zeta), \sigma^{*2}(\zeta, \zeta))$ and $\Xi(\zeta)_S = \{\Xi(\zeta_i)\}_{i \in S}$ is a realization with variables S included in vector $\Xi(\zeta)$.

Theorem. $\Xi_a(\cdot)$ is a stochastic process.

Joint multivariate distribution of the latent emulator $\Xi(\zeta)$ of the simulator f at a set of design points $\zeta = (\zeta^O, \zeta^C)$ (design points ζ^O with corresponding uncensored computer model output and design points ζ^C with corresponding censored output) is given as the following set of two distributions

$$\begin{aligned} \Xi(\zeta^O) &\sim \mathcal{N}\left(\mu(\zeta^O), \sigma^2 \mathbf{C}_{z^O}\right), \\ \Xi(\zeta^C) \mid \Xi(\zeta^O) &\sim \mathcal{N}\left(\mu^*(\zeta^C), \sigma^{*2}(\zeta^C)\right). \end{aligned}$$

Corresponding joint multivariate distribution of the emulator $\Xi_a(\zeta)$ of the simulator output $\mathbf{f}_a(\zeta)$ at design input points ζ is given by the two distributions $\Xi(\zeta^O)$ and $\Xi_a(\zeta^C) \mid \Xi(\zeta^O) = \Xi(\zeta^C) \mid \Xi(\zeta^O), \Xi(\zeta^C) < \mathbf{a}$. Namely

$$\begin{aligned}\Xi(\zeta^O) &\sim \mathcal{N}\left(\boldsymbol{\mu}(\zeta^O), \boldsymbol{\sigma}^2 \mathbf{C}_{z^o}\right), \\ \Xi_a(\zeta^C) \mid \Xi(\zeta^O) &\sim \mathcal{TN}_{(-\infty, \mathbf{a})}\left(\boldsymbol{\mu}^*(\zeta^C), \boldsymbol{\sigma}^{*2}(\zeta^C, \zeta^C)\right).\end{aligned}$$

Predictive distribution at a new input to a computer model. In order to obtain predictive distribution of a simulator $\Xi_a(\zeta)$ at a new input ζ' , consider the statistical approximation to a latent function $f(\cdot)$. Predictive distribution given by the latent emulator $\Xi(\cdot)$ at a new input ζ' , conditional on evaluations of the computer model $\Xi(\zeta)$ at design input points ζ , is

$$\Xi(\zeta') \mid \Xi(\zeta) \sim \mathcal{N}(\boldsymbol{\mu}^*(\zeta'), \boldsymbol{\sigma}^{*2}(\zeta', \zeta')). \quad (46)$$

Let $\Xi(\zeta) = (\Xi(\zeta^O), \Xi(\zeta^C))$ denote a vector of evaluations of the model at inputs ζ^O with uncensored outputs and inputs ζ^C with censored outputs. Computer model outputs $\Xi(\zeta^O)$ is given to us, but $\Xi(\zeta^C)$ are unknown. However, it is known that

$$\Xi(\zeta^C) \mid \Xi(\zeta^O), \Xi(\zeta^C) < \mathbf{a} \sim \mathcal{TN}_{(-\infty, \mathbf{a})}\left(\boldsymbol{\mu}^*(\zeta^C), \boldsymbol{\sigma}^{*2}(\zeta^C, \zeta^C)\right). \quad (47)$$

Distribution of interest is not closed-form, but a numerical approximation may be obtained. Computationally, in order to account for the censored observations ζ^C , samples from this truncated normal distribution (47) of $\Xi(\zeta^C) \mid \Xi(\zeta^O), \Xi(\zeta^C) < \mathbf{a}$ may be obtained. Say, k samples from the distribution of vector $\{\Xi(\zeta^C)_i\}_{i=1}^k$ are available. For each sample a vector $\Xi(\zeta)_i = (\Xi(\zeta^O), \Xi(\zeta^C)_i)$ is used to sample from (46), i.e. distribution $\Xi(\zeta') \mid \Xi(\zeta)_i$ which is represented by

$$\Xi(\zeta') \mid \Xi(\zeta)_i \sim \mathcal{N}(\boldsymbol{\mu}^*(\zeta'), \boldsymbol{\sigma}^{*2}(\zeta', \zeta')). \quad (48)$$

Finally, latent marginal distribution $\Xi(\zeta') \mid \Xi(\zeta) = \Xi(\zeta') \mid \Xi(\zeta^O), \Xi_a(\zeta^C) = \Xi(\zeta') \mid \Xi(\zeta^O), \Xi(\zeta^C) < \mathbf{a}$ is

$$p(\Xi(\zeta') \mid \Xi(\zeta^O), \Xi(\zeta^C) < \mathbf{a}) = \int \cdots \int p(\Xi(\zeta') \mid \Xi(\zeta^O), \Xi(\zeta^C)) p(\Xi(\zeta^C) \mid \Xi(\zeta^O), \Xi(\zeta^C) < \mathbf{a}) d\Xi(\zeta^C). \quad (49)$$

Predictive distribution of the emulator $\Xi_a(\cdot)$ at a new input to the computer model ζ' consists of two parts: a point mass at a and a Lebesgue measure on $\mathbb{R}_{>a}$. Namely, the distribution is

$$\Xi_a(\zeta') \mid \Xi(\zeta^O), \Xi(\zeta^C) < \mathbf{a} = \begin{cases} \Xi(\zeta') \mid \Xi(\zeta^O), \Xi(\zeta^C) < \mathbf{a}, & \Xi_a(\zeta') > a \\ \int_{-\infty}^a p(\Xi(\zeta') \mid \Xi(\zeta^O), \Xi(\zeta^C) < \mathbf{a}) d\Xi(\zeta'), & \Xi_a(\zeta') = a. \end{cases} \quad (50)$$

Theorem. $\Xi_a(\cdot) \mid \Xi(\zeta^O), \Xi(\zeta^C) < \mathbf{a}$ is a stochastic process.

Construction of an adequate emulator is essential for subsequent proper use of such an emulator for decision and policy making. If the purpose is not only to construct an emulator as an approximation to the simulator, but to use this emulator for providing adequate uncertainty estimates and estimation of a probability of an event, e.g., a natural hazard from a volcano pyroclastic flow, then it is preferred to use a valid approximation, which in this case is given by a censored GASP.

One illusive limitation of the proposed framework is that a constructed latent process may have to perform extrapolation in a region with no positive output values. However, the ‘‘censoring’’ operation performs conditioning of a latent process on the zero-output. This is because of this conditioning that censored emulator still performs an interpolation rather than extrapolation.

Chapter 6. Design of experiments for a large-scale computer model

A simulator is defined as `LARGE-SCALE` if the number of inputs is such that construction of its emulator (which involves optimization over its parameters) is prohibitively time-consuming. In order to facilitate the exploration of such a simulator, useful is to divide inputs \mathbf{x} into two groups, that is $\mathbf{x} = (\boldsymbol{\kappa}, \boldsymbol{\lambda})$. First, design over the range of $\boldsymbol{\kappa}$, choosing, say, m points $\{\boldsymbol{\kappa}_i\}_{i=1}^m$. For each $\boldsymbol{\kappa}_i$ develop a Gaussian process emulator over the rest of inputs $\boldsymbol{\lambda}$. Second, for each set of fixed inputs $\boldsymbol{\lambda}$, an emulator over $\boldsymbol{\kappa}$ conditional on fixed $\boldsymbol{\lambda}$ is constructed.

This methodology may be used for facilitating parameter estimation and fast emulation of a model with many inputs. Depending on the purpose and implementation of this methodology in practice, but the Gaussian process over the entire input space may be lost, although useful approximations are constructed.

My biography

This monograph is a follow-up after my dissertation “On Uncertainty Quantification for Systems of Computer Models” with which I completed philosophy doctorate (PhD) in Statistical Science at Duke University, USA in 2017. In my dissertation I have developed and analyzed a fully probabilistic Bayesian framework for testing theoretical scientific models, often realized as scientific computer models, with respect to experimental data.

Bibliography

- [1] James Berger, 2001, Objective Bayesian analysis for spatially correlated data, *Journal of the American Statistical Association* **96**(456), 1361-1374.
- [2] Mengyang Gu, 2018, Robust Gaussian stochastic process emulation, *The Annals of Statistics* **46**, 3038-3066.
- [3] Ksenia Kyzyurova, 2017, On Uncertainty Quantification for Systems of Computer Models, *Philosophy doctorate thesis*, Duke University.
- [4] Ksenia Kyzyurova, 2018, Coupling computer models through linking their statistical emulators, *SIAM/ASA Journal on Uncertainty Quantification* **6**(3), 1151-1171.
- [5] Dennis Lindley and Adrian Smith, 1972, Bayes estimates for the linear model, *Journal of the Royal Statistical Society: Series B (Methodological)* **34**(1),1-18.