

# Finite-dimensional Gaussian approximation with linear inequality constraints

Andrés F. López-Lopera<sup>1</sup>, François Bachoc<sup>2</sup>, Nicolas Durrande<sup>1,3</sup>, and Olivier Roustant<sup>1</sup>

<sup>1</sup>Mines Saint-Étienne, UMR CNRS 6158, LIMOS, F-42023 Saint-Étienne, France.

<sup>2</sup>Institut de Mathématiques de Toulouse, Université Paul Sabatier - Toulouse III, France.

<sup>3</sup>PROWLER.io, 66-68 Hills Road, Cambridge, CB2 1LA, UK.

## Abstract

Introducing inequality constraints in Gaussian process (GP) models can lead to more realistic uncertainties in learning a great variety of real-world problems. We consider the finite-dimensional Gaussian approach from Maatouk and Bay (2017) which can satisfy inequality conditions everywhere (either boundedness, monotonicity or convexity). Our contributions are threefold. First, we extend their approach in order to deal with general sets of linear inequalities. Second, we explore several Markov Chain Monte Carlo (MCMC) techniques to approximate the posterior distribution. Third, we investigate theoretical and numerical properties of the constrained likelihood for covariance parameter estimation. According to experiments on both artificial and real data, our full framework together with a Hamiltonian Monte Carlo-based sampler provides efficient results on both data fitting and uncertainty quantification.

## 1 Introduction

Gaussian processes (GPs) are one of the most famous non-parametric Bayesian frameworks for modelling stochastic processes. In principle, GP models place prior distributions over function spaces, and prior assumptions (e.g. smoothness, stationarity, sparsity) are encoded in covariance functions (Paciorek and Schervish, 2004; Rasmussen and Williams, 2005; Snelson and Ghahramani, 2006). Because GP provides a well-founded approach to learning, its properties have been explored in many decision tasks in regression (Kriging) and classification problems (Nickisch and Rasmussen, 2008; Rasmussen and Williams, 2005). Computer science, engineering, physics, biology, and neuroscience are some fields where GP models have been applied successfully (Murphy, 2012; Rasmussen and Williams, 2005).

Despite the reliable performance of GPs, they provide less realistic uncertainties when physical systems satisfy inequality constraints (Da Veiga and Marrel, 2012; Golchi et al., 2015; Maatouk and Bay, 2017). Quantifying properly the uncertainties is crucial for understanding real-world phenomena. For example, in nuclear safety criticality assessment, experimental settings typically demand expensive and risky procedures to evaluate neutron productions. Hence, emulators are required to infer these production rates and should assume a priori that the output is positive and usually monotonic with respect to a given set of input parameters. In this sense, to obtain more accurate predictions, both conditions have to be considered in the uncertainty quantification. Other test cases where data exhibit specific inequality constraints are given in computer networking (monotonicity) (Golchi et al., 2015), social system analysis (monotonicity) (Riihimäki and Vehtari, 2010), and econometrics (monotonicity or positivity) (Cousin et al., 2016).

Several studies have shown that including inequality constraints in GP frameworks can lead to more realistic uncertainty quantifications in learning from real data (Da Veiga and Marrel, 2012; Golchi et al., 2015; Riihimäki and Vehtari, 2010). In most of the cases, it is assumed that the inequalities are satisfied on a finite set of input locations. Then, the posterior distribution is approximated given those constrained inputs. To the best of our knowledge, the framework from (Maatouk and Bay, 2017) is the only Gaussian approach proposed in the literature which satisfies specific inequalities everywhere in the input space. There, the GP samples are approximated in the finite-dimensional space of functions such as piecewise linear functions. It is shown in (Bay et al., 2016) that the posterior mode converges to the one provided by thin plate splines. This approach has been applied on several real-data (e.g. econometrics, geostatistics) (Cousin et al., 2016; Maatouk and Bay, 2017), resulting in more realistic uncertainties than unconstrained Kriging.

The framework proposed in (Maatouk and Bay, 2017) still presents some limitations. First, the focus is on either boundedness, monotonicity or convexity conditions. Second, the proposed rejection sampling method for

estimating the posterior (Maatouk and Bay, 2016) results in costly computations when either the order of the finite approximation increases or the inequality constraints become more complex. Third, the proposed leave-one-out (LOO) technique for parameter estimation (Maatouk et al., 2015), restricts the optimal values to be on a finite grid of possible values, and provides the same estimation of correlation parameters as for unconstrained GP parameters. In order to address these limitations, our contributions are threefold. First, we extend the framework to deal with general sets of linear inequality constraints. Second, we evaluate efficient Markov Chain Monte Carlo (MCMC) algorithms that can be used to approximate the posterior distribution. Third, we investigate theoretical and numerical properties of the conditional likelihood for covariance parameter estimation. According to experiments on both artificial and real data, the resulting framework provides efficient results on both data fitting and uncertainty quantification.

This paper is organised as follows. In Section 2, we briefly describe GP modelling with inequality constraints. In Section 3, continuing with the finite-dimensional approach from (Maatouk and Bay, 2017), we propose a general formulation to deal with sets of linear inequalities. In Section 4, we apply several MCMC techniques to approximate the posterior distribution, and we compare their performances with respect to exact Monte Carlo (MC) algorithms. In Section 5, we write the conditional likelihood for the covariance parameter estimation, providing theoretical and empirical properties. In Section 6, we assess our framework in two-dimensional Kriging tasks. Finally, in Section 7, we summarize the conclusions, as well as the potential future works.

## 2 Gaussian process modelling with inequality constraints

### 2.1 Finite-dimensional approximation

Let  $Y$  be a zero-mean Gaussian process (GP) on  $\mathbb{R}$  with covariance function  $k$ . Consider  $x \in \mathcal{D}$  with compact input space  $\mathcal{D} = [0, 1]$ , and a set of knots  $t_1, \dots, t_m \in \mathbb{R}$ . For simplicity, we will consider equally-spaced knots  $t_j = j\Delta_m$  with  $\Delta_m = 1/(m-1)$ , but this assumption can be relaxed. Then, define a finite-dimensional GP, denoted by  $Y_m$ , as the piecewise linear interpolation of  $Y$  at knots  $t_1, \dots, t_m$ :

$$Y_m(x) = \sum_{j=1}^m Y(t_j)\phi_j(x), \quad (1)$$

where  $\phi_1 \dots, \phi_m$  are hat basis functions given by

$$\phi_j(x) := \begin{cases} 1 - \left| \frac{x-t_j}{\Delta_m} \right| & \text{if } \left| \frac{x-t_j}{\Delta_m} \right| \leq 1, \\ 0 & \text{otherwise.} \end{cases} \quad (2)$$

We illustrate the finite-dimensional representation of Equation (1) in Figure 1 for a deterministic function that satisfies two types of inequality constraints: boundedness and monotonicity (non-decreasing). Now, let  $\xi_j := Y(t_j)$  for  $j = 1, \dots, m$ . We aim at computing the distribution of  $Y_m$  conditionally on  $Y_m \in \mathcal{E}$  where  $\mathcal{E}$  is a convex set of functions defined by some inequality constraints. For instance, we may have

$$\mathcal{E} = \mathcal{E}_\kappa := \begin{cases} \{f \in C(\mathcal{D}, \mathbb{R}) \text{ s.t. } \ell \leq f(x) \leq u, \forall x \in \mathcal{D}\} & \text{if } \kappa = 0, \\ \{f \in C(\mathcal{D}, \mathbb{R}) \text{ s.t. } f \text{ is non-decreasing}\} & \text{if } \kappa = 1, \\ \{f \in C(\mathcal{D}, \mathbb{R}) \text{ s.t. } f \text{ is convex}\} & \text{if } \kappa = 2, \end{cases} \quad (3)$$

which corresponds to boundedness, monotonicity, and convexity constraints. The benefit of using hat functions and the finite-dimensional approximation  $Y_m$ , is that satisfying the inequality conditions  $Y_m(x) \in \mathcal{E}$ , for all  $x \in \mathcal{D}$ , is equivalent to satisfying only a finite number of inequality constraints (Maatouk and Bay, 2017). More precisely, for many natural choices of  $\mathcal{E}$ , we have

$$Y_m \in \mathcal{E} \quad \Leftrightarrow \quad \boldsymbol{\xi} \in \mathcal{C}, \quad (4)$$

where  $\mathcal{C}$  is a convex set of  $\mathbb{R}^m$  and  $\boldsymbol{\xi} = [\xi_1, \dots, \xi_m]^\top$ . For instance, for the convex set  $\mathcal{E}_\kappa$  of Equation (3), we have

$$\mathcal{C} = \mathcal{C}_\kappa := \begin{cases} \{\mathbf{c} \in \mathbb{R}^m; \forall j = 1, \dots, m : \ell \leq c_j \leq u\} & \text{if } \kappa = 0, \\ \{\mathbf{c} \in \mathbb{R}^m; \forall j = 2, \dots, m : c_j \geq c_{j-1}\} & \text{if } \kappa = 1, \\ \{\mathbf{c} \in \mathbb{R}^m; \forall j = 3, \dots, m : c_j - c_{j-1} \geq c_{j-1} - c_{j-2}\} & \text{if } \kappa = 2. \end{cases} \quad (5)$$

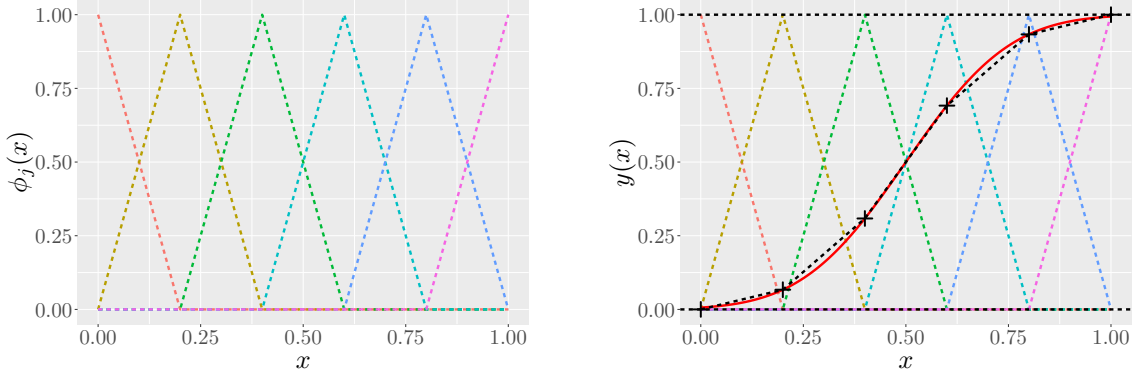


Figure 1: Illustration of the finite-dimensional approximation of Equation (1). (Left) Hat functions  $\phi_j$  for  $j = 1, \dots, 6$ . (Right) Approximation of the function  $y(x) = \Phi\left(\frac{x-0.5}{0.2}\right)$ , where  $\Phi$  is the standard cumulative normal distribution. Solid red and dashed black lines are the function  $y$ , and its finite approximation with six knots given by black crosses, respectively. Horizontal black dashed lines denote the bounds.

## 2.2 Conditioning with interpolation and inequality constraints

Consider the finite dimensional representation of the GP as in Equation (1), given the interpolation and inequality constraints

$$Y_m(x) = \sum_{j=1}^m \xi_j \phi_j(x), \quad \text{s.t.} \quad \begin{cases} Y_m(x_i) = y_i & \text{(interpolation conditions),} \\ Y_m \in \mathcal{E} & \text{(inequality conditions),} \end{cases} \quad (6)$$

where  $x_i \in \mathcal{D}$  and  $y_i \in \mathbb{R}$  for  $i = 1, \dots, n$ . Given a design of experiment (DoE)  $\mathbf{x} = [x_1, \dots, x_n]^\top$ , we have matrixially:

$$\mathbf{Y}_m = [Y_m(x_1), \dots, Y_m(x_n)]^\top = \mathbf{\Phi} \boldsymbol{\xi},$$

where  $\mathbf{\Phi}$  is the  $n \times m$  matrix defined by  $\Phi_{i,j} = \phi_j(x_i)$ . Let  $\mathbf{y} = [y_1, \dots, y_n]^\top$  be a realization of  $\mathbf{Y}_m$  as in Equation (6). From Equation (4), the conditional distribution of  $Y_m$ , under the inequality constraints  $Y_m \in \mathcal{E}$  and interpolation conditions  $Y_m(x_i) = y_i$  for  $i = 1, \dots, n$ , can be obtained from the conditional distribution of  $\boldsymbol{\xi}$  given  $\boldsymbol{\xi} \in \mathcal{C}$  and  $\mathbf{\Phi} \boldsymbol{\xi} = \mathbf{y}$  (see Equation (4)).

Observe that the vector  $\boldsymbol{\xi}$  of the values at the knots is a zero-mean Gaussian vector with covariance matrix  $\mathbf{\Gamma} = (k(t_i, t_j))_{1 \leq i, j \leq m}$ . Then, the distribution of  $\boldsymbol{\xi}$  given both interpolation and inequality conditions is truncated multinormal:

$$\boldsymbol{\xi} \sim \mathcal{N}(\mathbf{0}, \mathbf{\Gamma}) \quad \text{s.t.} \quad \begin{cases} \mathbf{\Phi} \boldsymbol{\xi} = \mathbf{y} & \text{(interpolation conditions),} \\ \boldsymbol{\xi} \in \mathcal{C} & \text{(inequality conditions),} \end{cases} \quad (7)$$

with  $\mathcal{C}$  as in Equation (4). For sampling purposes (see Algorithm 1), we need to compute the posterior mode which is given by the maximum of the probability density function of the posterior, i.e.  $\boldsymbol{\mu}_{\boldsymbol{\xi}}^* = \min\{\boldsymbol{\xi}^\top \mathbf{\Gamma}^{-1} \boldsymbol{\xi} \mid \mathbf{\Phi} \boldsymbol{\xi} = \mathbf{y}, \boldsymbol{\xi} \in \mathcal{C}\}$  (maximum a posteriori, MAP). Notice that  $\boldsymbol{\mu}_{\boldsymbol{\xi}}^*$  converges uniformly to the solution provided by thin plate splines when  $m \rightarrow \infty$  (Bay et al., 2016). More details and theoretical properties are provided in (Bay et al., 2016; Maatouk and Bay, 2017).

Figure 2 shows different Gaussian models for the example of Figure 1. We used a squared exponential (SE) covariance function with parameters  $(\sigma^2 = 1, \theta = 0.2)$ ,<sup>1</sup> and we fixed  $m = 100$ . The posterior distribution was approximated via Hamiltonian Monte Carlo (HMC) (Pakman and Paninski, 2014). From Figures 2(b) and 2(c), we observe that including the inequality constraints in the conditional distribution provides smaller confidence intervals compared to the ones given by the unconstrained GP. However, they do not satisfy both the boundedness and monotonicity conditions exhibited by the function  $y$ . On the other hand, from Figure 2(d), imposing both conditions leads to a more accurate prediction and more realistic confidence intervals. Later in Section 3, we will detail how to obtain the results of Figure 2(d).

<sup>1</sup>SE covariance function:  $k_{\boldsymbol{\theta}}(x - x') = \sigma^2 \exp\left\{-\frac{(x-x')^2}{2\theta^2}\right\}$  with  $\boldsymbol{\theta} = (\sigma^2, \theta)$ .

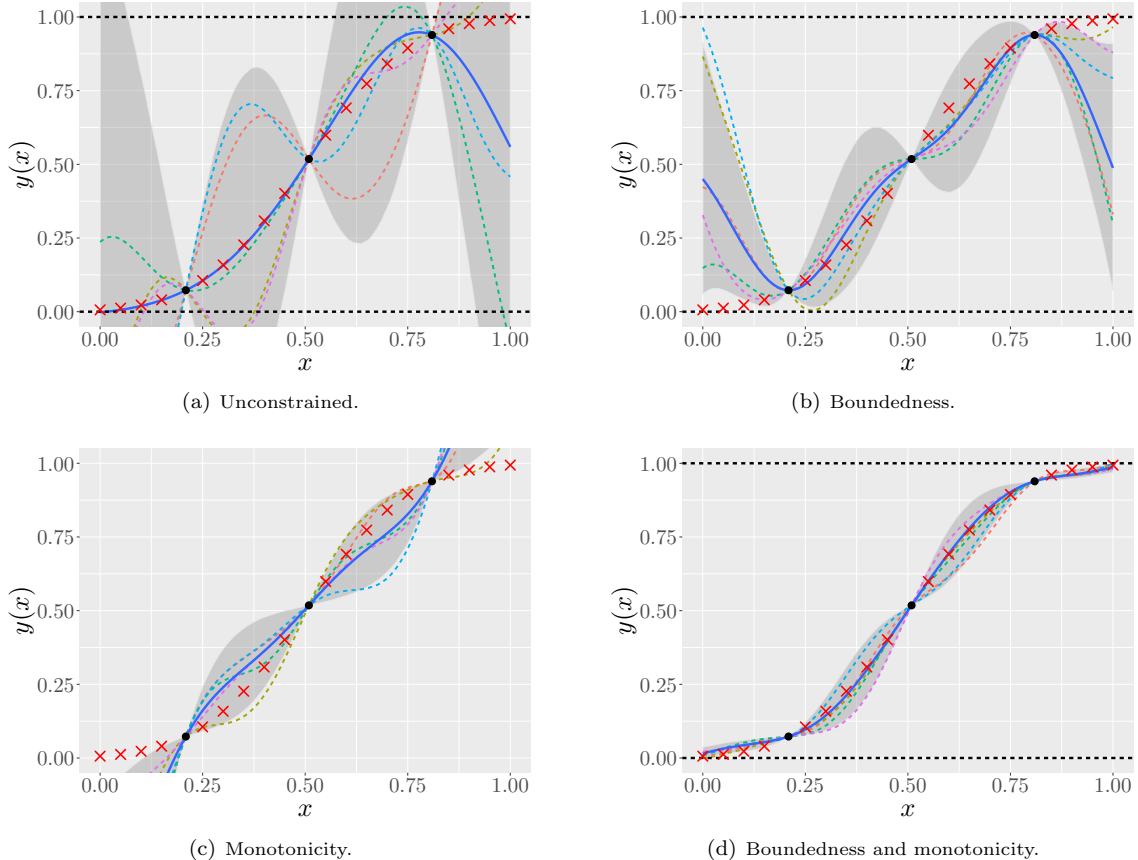


Figure 2: Example of Gaussian models satisfying different types of inequality constraints for interpolating the function  $x \mapsto \Phi\left(\frac{x-0.5}{0.2}\right)$ . Each panel shows: training and test points (black dots and red crosses, respectively), the conditional mean function (blue solid line), the 90% confidence interval (grey region), and conditional realizations (dashed multi colour lines). For boundedness constraints, bounds are showed in black dashed lines.

### 3 Finite-dimensional Gaussian approximation with linear inequality constraints

Now, we consider the case where  $\mathcal{C}$  is composed by a set of  $q$  linear inequalities of the form

$$\mathcal{C} = \left\{ \mathbf{c} \in \mathbb{R}^m; \forall j = 1, \dots, m, \forall k = 1, \dots, q, \lambda_{k,j} \in \mathbb{R} : \ell_k \leq \sum_{j=1}^m \lambda_{k,j} c_j \leq u_k \right\},$$

where the  $\lambda_{k,j}$ 's encode the linear operations, the  $\ell_k$ 's and  $u_k$ 's represent the lower and upper bounds. Notice that the convex sets  $\mathcal{C}_\kappa$  of Equation (5) are particular cases of  $\mathcal{C}$ . Denote  $\mathbf{\Lambda} = (\lambda_{k,j})_{1 \leq k \leq q, 1 \leq j \leq m}$ ,  $\mathbf{l} = (\ell_k)_{1 \leq k \leq q}$ , and  $\mathbf{u} = (u_k)_{1 \leq k \leq q}$ . Hence, Equation (7) is written

$$\boldsymbol{\xi} \sim \mathcal{N}(\mathbf{0}, \mathbf{\Gamma}) \quad \text{s.t.} \quad \begin{cases} \mathbf{\Phi} \boldsymbol{\xi} = \mathbf{y} & \text{(interpolation conditions),} \\ \mathbf{l} \leq \mathbf{\Lambda} \boldsymbol{\xi} \leq \mathbf{u} & \text{(inequality conditions).} \end{cases} \quad (8)$$

We further assume that  $q \geq m$  and that  $\mathbf{\Lambda}$  has rank  $m$ . By the rank-nullity theorem (see e.g. (Meyer, 2000)), it implies that  $\mathbf{\Lambda}$  is injective. In particular a linear system of the form  $\mathbf{\Lambda} \boldsymbol{\xi} = \boldsymbol{\eta}$  admits a unique solution  $\boldsymbol{\xi}$  when  $\boldsymbol{\eta}$  is in the image space of  $\mathbf{\Lambda}$ . This assumption is verified in many practical situations, up to adding inactive constraints. For instance, the monotonicity condition  $\boldsymbol{\xi}_1 \leq \dots \leq \boldsymbol{\xi}_m$ , which involves only  $q = m - 1$  (linearly independent) conditions, can be made compatible by adding the condition  $-\infty \leq \boldsymbol{\xi}_1$  (and/or  $\boldsymbol{\xi}_m \leq \infty$ ).

---

**Algorithm 1** Sampling from the finite-dimensional GP with linear inequality constraints.

---

- 1: **procedure** SAMPLING FROM  $\xi|\{\Phi\xi = \mathbf{y}, l \leq \Lambda\xi \leq \mathbf{u}\}$ , WHERE  $\xi \sim \mathcal{N}(\mathbf{0}, \Gamma)$
  - 2: **Input:**  $\mathbf{y}, \Gamma \in \mathbb{R}^{m \times m}, \Phi \in \mathbb{R}^{n \times m}, \mathcal{C}$ .
  - 3: Compute the conditional mean and covariance of  $\xi|\{\Phi\xi = \mathbf{y}\}$
  - 4:  $\mu = \Gamma\Phi^\top(\Phi\Gamma\Phi^\top)^{-1}\mathbf{y}$ , and
  - 5:  $\Sigma = \Gamma - \Gamma\Phi^\top(\Phi\Gamma\Phi^\top)^{-1}\Phi\Gamma$ .
  - 6: Solve the quadratic problem in  $\mathbb{R}^m$ :  $\mu_\xi^* = \min_{\xi \in \mathbb{R}^m} \{\xi^\top \Gamma^{-1} \xi | \Phi\xi = \mathbf{y}, l \leq \Lambda\xi \leq \mathbf{u}\}$ .
  - 7: Sample from the truncated multinormal distribution
  - 8:  $\Lambda\xi|\{\Phi\xi = \mathbf{y}, l \leq \Lambda\xi \leq \mathbf{u}\} \sim \mathcal{TN}(\Lambda\mu, \Lambda\Sigma\Lambda^\top, l, \mathbf{u})$ .
  - 9: Define  $\eta = \Lambda\xi$ , and solve the linear system to obtain the sample  $\xi$ .
  - 10: **Remark:** use the posterior mode  $\nu_\xi^* = \Lambda\mu_\xi^*$  as a starting state for an MCMC sampler (see Section 4).
- 

We now explain how to sample  $\xi$  from Equation (8). First, we compute the conditional distribution given the interpolation constraints  $\xi|\{\Phi\xi = \mathbf{y}\}$ . Since  $\xi \sim \mathcal{N}(\mathbf{0}, \Gamma)$ , then  $\Phi\xi \sim \mathcal{N}(\mathbf{0}, \Phi\Gamma\Phi^\top)$  and the conditional distribution  $\xi|\{\Phi\xi = \mathbf{y}\}$  is also Gaussian  $\mathcal{N}(\mu, \Sigma)$  (Rasmussen and Williams, 2005), with

$$\mu = \Gamma\Phi^\top[\Phi\Gamma\Phi^\top]^{-1}\mathbf{y}, \quad \text{and} \quad \Sigma = \Gamma - \Gamma\Phi^\top[\Phi\Gamma\Phi^\top]^{-1}\Phi\Gamma. \quad (9)$$

Therefore, we have  $\Lambda\xi|\{\Phi\xi = \mathbf{y}\} \sim \mathcal{N}(\Lambda\mu, \Lambda\Sigma\Lambda^\top)$ . Let  $\mathcal{TN}(\mathbf{m}, \mathbf{C}, \mathbf{a}, \mathbf{b})$  be the truncated multinormal distribution with mean vector  $\mathbf{m}$ , covariance matrix  $\mathbf{C}$ , and bound vectors  $(\mathbf{a}, \mathbf{b})$  such that  $\mathbf{a} \leq \mathbf{b}$ . Thus, the posterior distribution of Equation (8) is obtained from

$$\Lambda\xi|\{\Phi\xi = \mathbf{y}, l \leq \Lambda\xi \leq \mathbf{u}\} \sim \mathcal{TN}(\Lambda\mu, \Lambda\Sigma\Lambda^\top, l, \mathbf{u}). \quad (10)$$

Notice that the inequality conditions are encoded in the posterior mean  $\Lambda\mu$ , the posterior covariance  $\Lambda\Sigma\Lambda^\top$ , and bounds  $(l, \mathbf{u})$ . Finally, the posterior mode is given by  $\nu_\xi^* = \Lambda\mu_\xi^*$  where  $\mu_\xi^*$  is the solution provided in Subsection 2.2. The truncated multinormal of Equation (10) can be approximated using Markov Chain Monte Carlo (MCMC) algorithms. Denoting  $\eta = \Lambda\xi$ , notice that the samples for  $\xi$  can be obtained by using the ones obtained for  $\eta$  if the linear system is solved. Indeed, as mentioned above, we assumed that  $\Lambda$  has rank  $m$ , which implies that the solution of  $\Lambda\xi = \eta$  exists and is unique.

The whole sampling scheme is summarized in Algorithm 1. Now, we illustrate some examples where the proposed framework satisfies different types of inequality conditions. The posterior is approximated via HMC (more details about HMC are given in Section 4).

**Example 1** We continue with the example of Figure 1. As we can fix the structure of the linear inequalities  $(\Lambda, l, \mathbf{u})$ , we can impose both boundedness and monotonicity conditions in the constrained GP. One way to do this is to encode them individually. Let  $l_1 \leq \Lambda_1\xi \leq \mathbf{u}_1$  and  $l_2 \leq \Lambda_2\xi \leq \mathbf{u}_2$  be the sets of conditions to satisfy boundedness and monotonicity constraints, respectively. Then, we can build an extended set of inequalities  $l \leq \Lambda\xi \leq \mathbf{u}$  by stacking the constraints (i.e.  $\Lambda = [\Lambda_1, \Lambda_2]^\top, l = [l_1, l_2]^\top, \mathbf{u} = [u_1, u_2]^\top$ ), so that Algorithm 1 can be used. Notice that one can encode the same information in a reduced set of linear inequalities. Instead of encoding independently the boundedness and monotonicity constraints, which requires  $q = 2m - 1$  inequalities, one can impose boundedness conditions only for the first and last knot, and monotonicity conditions for all the knots except the first one. Due to monotonicity, the intermediate knots will also satisfy the boundaries. In this way, we only need  $q = m + 1$  conditions. In many other cases, the size of specific sets of linear constraints can be reduced. However for general discussions, we will use the full extended set and we will apply efficient samplers to approximate the posterior.

**Example 2** Notice from the previous example that the extension for more than two sets of inequalities is straightforward. Consider for instance  $Q$  different sets of conditions. We can build the posterior from Equation (10) with  $\Lambda = [\Lambda_1, \dots, \Lambda_Q]^\top, l = [l_1, \dots, l_Q]^\top$ , and  $\mathbf{u} = [u_1, \dots, u_Q]^\top$ , and apply Algorithm 1. Figure 3 shows an example with the target function  $y(x) = x^2$ , satisfying three types of inequality constraints: boundedness, monotonicity and convexity. We proposed different models satisfying one or more inequality constraints. We used a SE covariance function with parameters  $(\sigma^2 = 1.0, \theta = 0.2)$ . By imposing the three conditions, we obtain samples that also satisfy the three types of constraints.

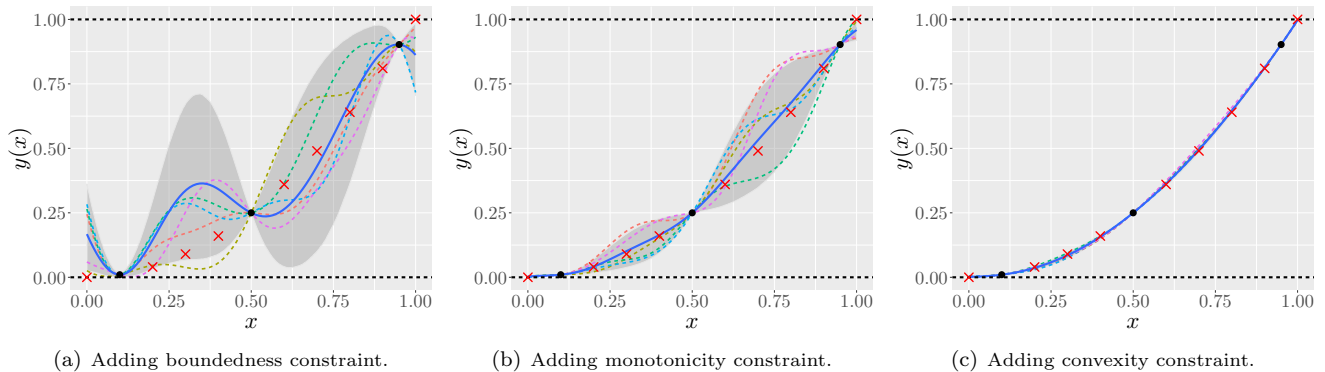


Figure 3: Examples of Gaussian models satisfying one or several types of inequality constraints for interpolating the square function  $x \mapsto x^2$ . Panel description is the same as in Figure 2.

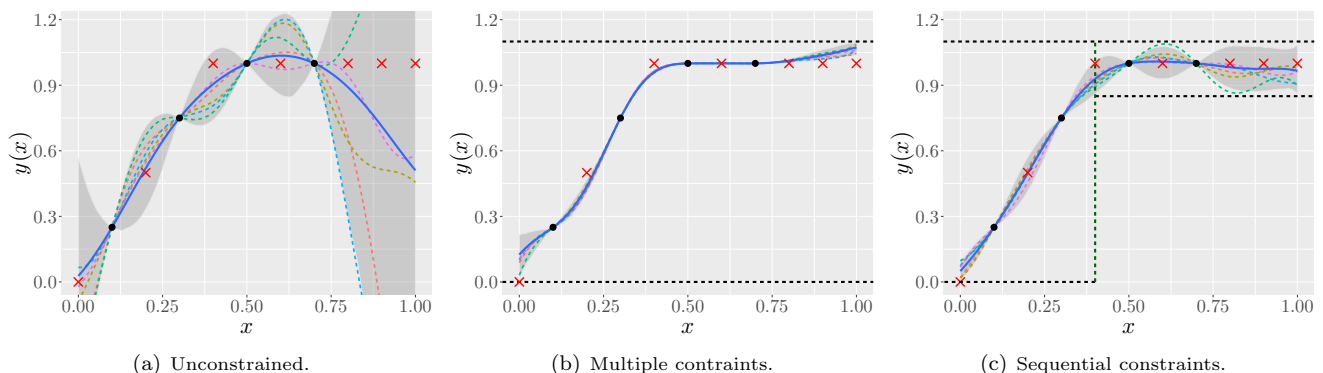


Figure 4: Examples of Gaussian models with different types of constraints for the example 3 from Section 3. (b) Boundedness and monotonicity constraints are imposed. (c) The two non-overlapping intervals are divided by a vertical dashed line at 0.4. In the first interval, boundedness and monotonicity constraints are taken into account. In the second interval, only boundedness is imposed. Panel description is the same as in Figure 2.

**Example 3** Since the bounds  $(\mathbf{l}, \mathbf{u})$  are not forced to be the same everywhere, it is possible to fix specific constraints over non-overlapping intervals. For instance if the interval is partitioned in  $G$  subintervals, we consider the corresponding partition  $\boldsymbol{\xi} = [\xi_1, \dots, \xi_G]^\top$ . Then, we can impose different types of inequality conditions in each group by considering the same structure used in Example 2. Figure 4 shows an example where the function  $y$  satisfies different behaviours in two non-overlapping intervals. The output increases monotonically and peaks at  $y(0.4) = 1.0$ . This kind of profile is met in different applications (e.g. step responses in control theory, protein profiles in molecular biology) (Kocijan, 2016; Murphy, 2012). We trained three models satisfying different conditions. For the case of multiple constraints, we imposed boundedness and monotonicity. For the case of sequential conditions, we divided the profile in two non-overlapping intervals satisfying different types of constraints. We used a SE covariance function with parameters  $(\sigma^2 = 1.0, \theta = 0.2)$ . By imposing sequentially the constraints, we obtain less restricted uncertainties and more accurate models for data fitting.

## 4 Simulating from the posterior distribution

As shown in Equation (10), the posterior distribution  $\boldsymbol{\Lambda}\boldsymbol{\xi} | \{\Phi\boldsymbol{\xi} = \mathbf{y}, \mathbf{l} \leq \boldsymbol{\Lambda}\boldsymbol{\xi} \leq \mathbf{u}\}$  is truncated multinormal. It is supported on  $\mathbb{R}^m$ , where  $m$  is the number of basis functions. Notice that  $m$  should be chosen large enough for better approximations. A Monte Carlo (MC) algorithm based on rejection sampling was proposed in (Maatouk and Bay, 2016) using the posterior mode. This method, called rejection sampling from the mode (RSM), is an exact sampler that provides independent and identically distributed (iid) sample paths. However, the acceptance rate from RSM



decreases when  $m$  gets larger, providing a poor performance for high dimensional spaces. Another MC-based exact sampler was introduced by (Botev, 2017) to deal with truncated multinormals in higher dimensions. It can both simulate multinormals under linear constraints, and estimate the probabilities that these constraints are satisfied, via minimax exponential tilting (ET). As RSM and ET are exact methods, we will use them as gold standards to evaluate the performance of the MCMC techniques that we describe now.

## 4.1 MCMC for truncated multinormal distributions

MCMC approaches assume that samples follow a Markov chain, providing correlated samples but with a higher acceptance rate. Recently, efficient algorithms have been proposed for truncated multinormal distributions such as Gibbs sampling (Taylor and Benjamini, 2017), Metropolis-Hastings (Murphy, 2012), and Hamiltonian Monte Carlo (Pakman and Paninski, 2014). In this section, we apply them to simulate from the posterior distribution of Equation (10).

**Gibbs sampling** Algorithms based on Gibbs sampling are widely used to sample from truncated multinormals due to their easy implementation, and their reliable performances (Brooks et al., 2011; Murphy, 2012). They sample each variable in turn conditionally on the values of the other ones (Murphy, 2012). Therefore, sampling from a truncated multinormal is reduced to sampling sequentially from conditional truncated (univariate) normals. Unlike RSM, there is no rejection step. However, the “single site updating” property may produce strong correlations, requiring to discard intermediate samples (thinning effect). Several studies have proposed efficient algorithms to obtain less correlated sample paths (e.g. collapsed Gibbs sampling, blocked Gibbs sampling) (Murphy, 2012). In this paper, we will use the fast Gibbs sampler proposed in (Taylor and Benjamini, 2017).

**Metropolis Hastings (MH)** MH-based algorithms propose to move all the coordinates at a time in each step to obtain less correlated simulations. Given a proposed state  $\mathbf{x}'$ , we either accept or reject the new state according to a given acceptance rule (Murphy, 2012). If the proposal is accepted, the new state is  $\mathbf{x}'$ , otherwise the new state remains at the previous state  $\mathbf{x}$ . For multinormal distributions, a symmetric Gaussian proposal is commonly used, i.e.  $q(\mathbf{x}'|\mathbf{x}) = \mathcal{N}(\mathbf{x}, \eta \Sigma_{\mathbf{x}'|\mathbf{x}})$  where  $\eta$  is a scale factor. This approach is known as random walk Metropolis algorithm (Murphy, 2012). One can increase the acceptance rate by tuning properly the value of  $\eta$ .

**Hamiltonian Monte Carlo (HMC)** Nowadays, hybrid methods have been subject to great attention from the statistical community due to the inclusion of physical interpretation that may provide useful intuition (Brooks et al., 2011; Neal, 1996). In (Duane et al., 1987), an efficient hybrid approach was introduced using the properties of Hamiltonian dynamics. Later in (Neal, 1996), the hybrid approach from (Duane et al., 1987) was extended to statistical applications, and was introduced formally as Hamiltonian Monte Carlo (HMC). The Hamiltonian dynamics provides distant proposal distributions producing less correlated sample paths without diminishing the acceptance rate. In this paper, we use the HMC-based approach for truncated multinormals introduced in (Pakman and Paninski, 2014).

## 4.2 Results

In Table 1, we evaluate the efficiency of the MC and MCMC approaches described in Subsection 4.1 on the examples from Figure 2. In order to reduce the simulation cost, we used  $m = 30$  hat basis functions. Hence, the problem is to sample a vector of length 30 from a truncated multinormal distribution. We set the tuning hyperparameters such that the effective sample size (ESS) is within the ranges produced by both RSM and ET (grey columns). The ESS is a heuristic used commonly to evaluate the quality of correlated sample paths, and it gives an intuition on how many samples from the path can be considered independent (Gong and Flegal, 2016). A standard ESS is given by  $ESS = n_s / (1 + 2 \sum_{k=1}^{n_s} \rho_k)$  where  $n_s$  is the size of the sample path and  $\rho_k$  is the sample autocorrelation with lag  $k$ . However, the drawback of this indicator is that it accepts negative correlations to evaluate the quality of mean estimators (e.g. for variance reduction). Thus, we suggest an alternative ESS which penalizes both positive and negative correlations:  $ESS = n_s / (1 + 2 |\sum_{k=1}^{n_s} \rho_k|)$ . Notice that it is equal to  $n_s$  for iid sample paths, or smaller otherwise. We compute the ESS indicator for each coordinate of  $\boldsymbol{\xi}$ , i.e.  $ESS_j = ESS(\xi_j^1, \dots, \xi_j^{n_s})$  for the  $j$ -th component of  $\boldsymbol{\xi}$  with  $j = 1, \dots, m$ . We then compute quantiles ( $q_{10\%}, q_{50\%}, q_{90\%}$ ) over the 30 resulting ESS values.

Table 1: Efficiency of MC/MCMC samplers (by rows) in term of ESS-based indicators (by columns). **Samplers:** Rejection Sampling from the Mode (RSM) (Maatouk and Bay, 2016), Exponential Tilting (ET) (Botev, 2017), Gibbs Sampling (Gibbs) (Taylor and Benjamini, 2017), Metropolis-Hasting (MH) (Murphy, 2012), Hamiltonian Monte Carlo (HMC) (Pakman and Paninski, 2014). **Indicators:** effective sample size (ESS):  $\text{ESS} = n/(1 + 2|\sum_{\forall k} \rho_k|)$ , multivariate ESS (mvESS) (Vats et al., 2017), time normalised ESS (TN-ESS) (Lan and Shahbaba, 2016).

Toy Example	Method	CPU Time [s]	ESS [ $\times 10^4$ ]	mvESS	TN-ESS	Hyperparameter
			( $q_{10\%}, q_{50\%}, q_{90\%}$ )	[ $\times 10^4$ ]	[ $\times 10^4 s^{-1}$ ]	
Toy Figure 2(b) (Boundedness)	RSM	99.06	(0.81, 0.89, 0.96)	1.22	0.01	-
	ET	<b>0.44</b>	(0.83, 0.90, 0.99)	1.17	<b>1.88</b>	-
	Gibbs	3.54	(0.81, 0.91, 0.93)	1.16	0.23	thinning = 200
	MH	52.21	(0.77, 0.87, 0.95)	1.21	0.01	$\eta = 1$
	HMC	<b>0.44</b>	(0.72, 0.76, 0.91)	1.26	1.64	-
Toy Figure 2(c) (Monotonicity)	RSM	190.62	(0.86, 0.90, 0.93)	1.21	0.00	-
	ET	0.77	(0.84, 0.91, 0.97)	1.18	1.09	-
	Gibbs	3.04	(0.80, 0.93, 0.99)	1.15	0.26	thinning = 200
	MH	96.64	(0.81, 0.91, 0.99)	1.23	0.01	$\eta = 1$
	HMC	<b>0.33</b>	(0.73, 0.79, 0.88)	1.28	<b>2.22</b>	-
Toy Figure 2(d) (Bounded Monotonicity)	RSM	-	-	-	-	-
	ET	41.16	(0.80, 0.88, 0.95)	1.23	0.02	-
	Gibbs	40.28	(0.28, 0.44, 0.77)	1.09	0.01	thinning = 1000
	MH	-	-	-	-	-
	HMC	<b>12.92</b>	(0.72, 0.85, 0.99)	1.26	<b>0.06</b>	-

To take into account cross-correlations from multivariate MCMC, we also compute the multivariate ESS (mvESS) proposed in (Vats et al., 2017). For mvESS, values higher than  $n_s$  indicate the presence of negative correlations. In our case, we are interested in being around  $n_s$ . The size  $n_s = 10^4$  is chosen to be larger than the minimum ESS required to obtain a proper estimation of the vector  $\xi \in \mathbb{R}^{30}$ :  $\text{minESS}(30) = 8563$  (Gong and Flegal, 2016; Vats et al., 2017). Finally, using the procedure proposed in (Lan and Shahbaba, 2016), we test the efficiency of each method by computing the time normalised ESS (TN-ESS) at  $q_{10\%}$  (worst case) using the CPU time in seconds, i.e.  $\text{TN-ESS} = q_{10\%}(\text{ESS})/(\text{CPU Time})$ .

Table 1 shows the efficiency of MC/MCMC algorithms in terms of ESS indicators. Notice that for the two examples of Figures 2(b) and 2(c), the MC/MCMC techniques tend to produce similar ESS intervals, but RSM and MH are the most expensive procedures due to their high rejection rates. Although the Gibbs sampler requires to discard a large amount of simulations in order to be within reasonable ESS ranges, it also presents accurate results in both efficiency and CPU time. In general, both ET and HMC methods present more efficient results than the other samplers in the first two examples. For more complex constraints as in the example of Figure 2(d), the efficiency is reduced dramatically for all the methods. For example, the acceptance rates of both RSM and MH are so small that sampling was not feasible in a reasonable time. For the other methods, the TN-ESS rates are smaller but HMC still gives a reasonable value (almost three times larger than for ET), concluding that HMC is an efficient sampler for the proposed framework.

## 5 Covariance parameter estimation with inequality constraints

### 5.1 Conditional maximum likelihood

Let  $\{k_\theta; \theta \in \Theta\}$ , with  $\Theta \subset \mathbb{R}^p$ , be a parametric family of covariance functions. We assume in this section that the zero-mean GP  $Y$  has covariance function  $k_\theta$  for an unknown  $\theta^* \in \Theta$ . We consider the problem of estimating  $\theta^*$ . Commonly,  $\theta^*$  is estimated by maximising the unconstrained Gaussian likelihood  $p_\theta(\mathbf{Y}_m)$  with respect to  $\theta \in \Theta$  (maximum likelihood, ML), with  $\mathbf{Y}_m = [Y_m(x_1), \dots, Y_m(x_n)]^\top$ . Let  $\mathcal{L}_m(\theta)$  be the log likelihood of  $\theta$

$$\mathcal{L}_m(\theta) = \log p_\theta(\mathbf{Y}_m) = -\frac{1}{2} \log(\det(\mathbf{K}_\theta)) - \frac{1}{2} \mathbf{Y}_m^\top \mathbf{K}_\theta^{-1} \mathbf{Y}_m - \frac{n}{2} \log 2\pi, \quad (11)$$

with  $\mathbf{K}_\theta = \Phi \Gamma_\theta \Phi^\top$  and  $\Gamma_\theta = (k_\theta(t_i, t_j))_{1 \leq i, j \leq m}$ . Then, the ML estimation (MLE) is

$$\hat{\theta}_{\text{MLE}} = \arg \max_{\theta \in \Theta} \mathcal{L}_m(\theta). \quad (12)$$



When we maximise the likelihood of Equation (12), we are looking for a parameter  $\theta$  that improves the ability of our model to explain the data (Rasmussen and Williams, 2005). However, because the unconstrained ML itself does not take into account the constraints  $\xi \in \mathcal{C}$ , the estimated  $\hat{\theta}_{\text{MLE}}$  may produce less realistic models. Here, we suggest to use the constrained likelihood. Let  $p_{\theta}(\mathbf{Y}_m | \xi \in \mathcal{C})$  be the conditional probability density function of  $\mathbf{Y}_m$  given  $\xi \in \mathcal{C}$ , when  $Y$  has covariance function  $k_{\theta}$ . By using Bayes' theorem, the constrained log likelihood  $\mathcal{L}_{\mathcal{C},m}(\theta) = \log p_{\theta}(\mathbf{Y}_m | \xi \in \mathcal{C})$  is

$$\mathcal{L}_{\mathcal{C},m}(\theta) = \log p_{\theta}(\mathbf{Y}_m) + \log P_{\theta}(\xi \in \mathcal{C} | \Phi \xi = \mathbf{Y}_m) - \log P_{\theta}(\xi \in \mathcal{C}), \quad (13)$$

where the first term is the unconstrained log-likelihood, and the last two terms depend on the inequality constraints. Then, the constrained ML (CML) estimator is given by

$$\hat{\theta}_{\text{CMLE}} = \arg \max_{\theta \in \Theta} \mathcal{L}_{\mathcal{C},m}(\theta). \quad (14)$$

Notice that  $P_{\theta}(\xi \in \mathcal{C} | \Phi \xi = \mathbf{Y}_m)$  and  $P_{\theta}(\xi \in \mathcal{C})$  are Gaussian orthant probabilities. As they have no explicit expressions, numerical procedures have been investigated (Botev, 2017; Genz, 1992). Hence, the likelihood evaluation and optimisation of Equations (13) and (14) have to be done numerically.

## 5.2 Simulation study

To assess the performance of the estimator of Equation (14), we simulated sample paths from a zero-mean constrained GP  $Y$  using a Matérn 5/2 covariance function with  $\theta^* = (1, 0.2)$ .<sup>2</sup> We sampled 100 realizations of  $Y$  on  $\mathcal{D} = [0, 1]$  such that  $Y \in [-1, 1]$ . Then, for each realization, we trained a constrained model assuming boundedness conditions with bounds  $[-1, 1]$ . We used 10 training points regularly spaced in  $\mathcal{D}$  and  $m = 50$  hat basis functions. For ML and CML optimisations, we used multistart with ten initial vectors of covariance parameters located on a maximin Latin hypercube DoE with  $\sigma^2 \in [0, 2]$  and  $\theta \in [0.04, 0.40]$ . As the parameters of the Matérn 5/2 covariance function are non-microergodic for one-dimensional input spaces, they cannot be estimated consistently (Zhang, 2004). Therefore, we evaluated the quality of the likelihood estimators using the consistently estimable ratio  $\rho = \sigma^2/\theta^5$ . In Figure 5(a), we show the boxplots of the estimated ratios obtained with the 100 simulations drawn from the GP. Notice that the estimated logged ratios  $\log \hat{\rho}_{\text{MLE}}$  and  $\log \hat{\rho}_{\text{CMLE}}$  are reasonably close to the true value  $\log \rho^* = \log(1^2/0.2^5)$ , but the one using CMLE is slightly better in terms of variance and bias.

We also evaluate the efficiency of the two estimators in terms of prediction accuracy. For each realization, we estimated the covariance parameters  $\theta^*$  by MLE and CMLE. We then simulated the posterior at 50 new regularly spaced locations using the estimated covariance parameter  $\hat{\theta}$ . The conditional sample paths were simulated via HMC. We used the  $Q^2$  and predictive variance adequation (PVA) criteria to assess the quality of predictions over the 50 new values. Denoting by  $n_t$  the number of test points,  $z_1, \dots, z_{n_t}$  and  $\hat{z}_1, \dots, \hat{z}_{n_t}$  the sets of test and predicted observations (respectively), then  $Q^2 = 1 - \sum_{i=1}^{n_t} (\hat{z}_i - z_i)^2 / \sum_{i=1}^{n_t} (\bar{z} - z_i)^2$ , where  $\bar{z}$  is the mean of the test data. Notice that for noise-free observations, the  $Q^2$  indicator is equal to one if the predictors  $\hat{z}_1, \dots, \hat{z}_{n_t}$  are exactly equal to the test data (ideal case), zero if they are equal to the constant prediction  $\bar{z}$ , and negative if they perform worse than  $\bar{z}$ . On the other hand, PVA assesses the quality of predictive variances  $\hat{\sigma}_i^2$ , for all  $i = 1, \dots, n_t$ , using the prediction errors:  $\text{PVA} = |\log(\frac{1}{n_t} \sum_{i=1}^{n_t} (z_i - \hat{z}_i)^2 / \hat{\sigma}_i^2)|$  (Bachoc, 2013). Notice that, if the standardized residuals  $(z_i - \hat{z}_i)/\hat{\sigma}_i$  were an iid sample, then the PVA value would be close to 0 by the law of large numbers. Departure from 0 may indicate either a lack of independence, or a biased estimation of the prediction uncertainty  $\sigma_i$ . In that sense, smaller PVA values may correspond to more reliable confidence intervals.

Figure 5 shows the inferred sample paths for one realization using 5(d)  $\hat{\theta}_{\text{MLE}}$ , 5(e)  $\hat{\theta}_{\text{CMLE}}$ , and 5(f)  $\theta^*$ . We observe that, in the three cases, the models tend to fit properly the test data with accurate confidence intervals. According to Figures 5(b) and 5(c), we see they provide  $Q^2$  and PVA median values close to the ones obtained when the true  $\theta^*$  is used. Although the predictive accuracies obtained using CMLE are slightly better than for MLE in terms of bias, we observe larger variances in the PVA criterion. Since the orthant Gaussian terms from the conditional likelihood of Equation (13) have to be approximated, we believe that this affects the effectiveness of CMLE. Furthermore, existing estimators of orthant Gaussian probabilities present some numerical instabilities

<sup>2</sup>Matérn 5/2 kernel function:  $k_{\theta}(x - x') = \sigma^2 \left(1 + \frac{\sqrt{5}|x-x'|}{\theta} + \frac{5}{3} \frac{(x-x')^2}{\theta^2}\right) \exp\left\{-\frac{\sqrt{5}|x-x'|}{\theta}\right\}$  with  $\theta = (\sigma^2, \theta)$ .

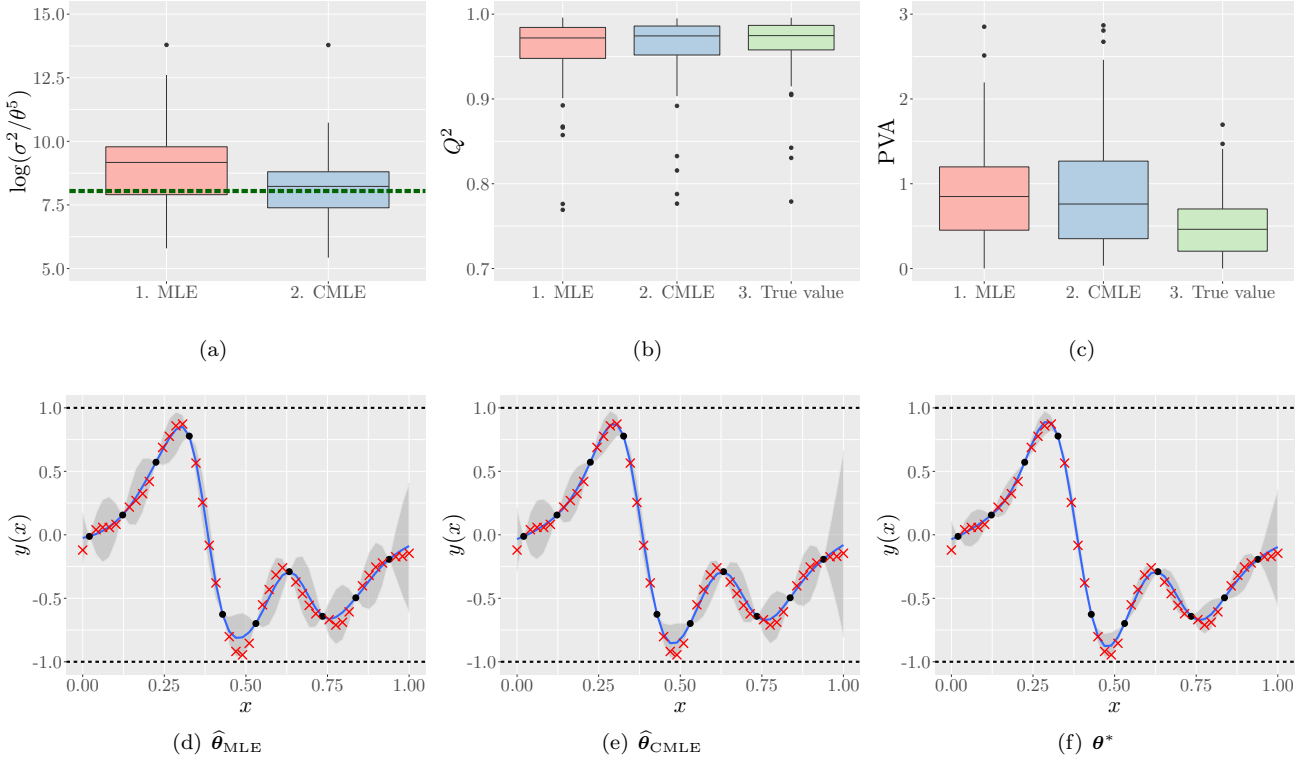


Figure 5: Assessment of the likelihood (ML) and conditional likelihood (CML) estimators for 100 samples drawn from a GP with true parameters  $\theta^* = (1, 0.2)$ , and satisfying the bounds  $[-1, 1]$ . (a) Estimated values of the log-ratio  $\log \rho^* = \log(1^2/0.2^5)$  (dashed green line) using MLE and CML. Predictive accuracies are evaluated using the (b)  $Q^2$  and (c) PVA criteria. Predictions are showed for one sample using (d)  $\hat{\theta}_{\text{MLE}}$ , (e)  $\hat{\theta}_{\text{CMLE}}$ , and (f)  $\theta^*$ . For the predictions, panel description is the same as Figure 2.

limiting the CML optimisation routine and providing suboptimal results. Finally, notice that MLE also provides reliable predictions. This suggests that, if we properly take into account the inequality constraints in the posterior distribution, the unconstrained ML optimisation can be used for practical implementation.

### 5.3 Asymptotic properties

Now, we study the asymptotic properties of likelihood-based estimators for constrained GPs. We consider the fixed-domain asymptotic setting (Stein, 1999), with a dense sequence of observation points in a bounded domain. It should be noted that, when the GP is not constrained, significant contributions have been provided to study the consistency or asymptotic normality of the ML estimator (Du et al., 2009; Loh, 2005; Loh and Lam, 2000; Ying, 1993; Zhang, 2004). In this paper, we show that, loosely speaking, any consistency result for ML with unconstrained GPs, is preserved when adding either boundedness, monotonicity or convexity constraint. Furthermore, this consistency occurs for both the unconditional and conditional likelihood functions.

For  $\kappa \in \{0, 1, 2\}$ , let  $Y$  be a GP with  $C^\kappa$  trajectories on a bounded set  $\mathbb{X} \subset \mathbb{R}^d$ . Let  $\mathcal{E}_\kappa$  be one of the following convex set of functions

$$\mathcal{E}_\kappa = \begin{cases} f : \mathbb{X} \rightarrow \mathbb{R}, f \text{ is } C^0 \text{ and } \forall \mathbf{x} \in \mathbb{X}, \ell \leq f(\mathbf{x}) \leq u & \text{if } \kappa = 0, \\ f : \mathbb{X} \rightarrow \mathbb{R}, f \text{ is } C^1 \text{ and } \forall \mathbf{x} \in \mathbb{X}, \forall i = 1, \dots, d, \frac{\partial}{\partial x_i} f(\mathbf{x}) \geq 0 & \text{if } \kappa = 1, \\ f : \mathbb{X} \rightarrow \mathbb{R}, f \text{ is } C^2 \text{ and } \forall \mathbf{x} \in \mathbb{X}, \frac{\partial^2}{\partial \mathbf{x}^2} f(\mathbf{x}) \text{ is a non-negative definite matrix} & \text{if } \kappa = 2. \end{cases} \quad (15)$$

For the purpose of asymptotic analysis, we do not consider the hat basis functions anymore, and we focus on the GP  $Y$  and the observation vector  $\mathbf{Y}_n = [Y(x_1), \dots, Y(x_n)]^\top$ . We study the (unconstrained) likelihood function based

on  $p_{\theta}(\mathbf{Y}_n)$  and the constrained likelihood function based on  $p_{\theta}(\mathbf{Y}_n|Y \in \mathcal{E}_{\kappa})$ . Notice that these quantities are more challenging to evaluate in practice than for Subsections 5.1 and 5.2, but the purpose is a theoretical analysis.

In Proposition 5.1, we prove that if ML is consistent, when considering the (unconditional) distribution of  $Y$ , then it remains consistent when conditioning to  $Y \in \mathcal{E}_{\kappa}$ . In Proposition 5.2, we prove that, under mild conditions, implying the consistency of the ML estimator with the (unconditional) distribution of  $Y$ , the CML remains consistent when adding the constraint  $Y \in \mathcal{E}_{\kappa}$ . The proofs of Propositions 5.1 and 5.2 require supplementary conditions and lemmas, which are given in Appendix A.

**Proposition 5.1.** *Let  $Y$  be a zero-mean GP on a bounded set  $\mathbb{X} \subset \mathbb{R}^d$  with covariance function  $k$  satisfying Condition A.1. Let  $\Theta$  be a compact set on  $(0, \infty)^{d+1}$ . Let  $k_{\theta}$  be the covariance function of  $x \rightarrow \sigma Y(\theta_1 x_1, \dots, \theta_d x_d)$  for  $\theta = (\sigma^2, \theta_1, \dots, \theta_d) \in \Theta$ . Let  $\theta^* = (1, \dots, 1)$ . Remark that  $k = k_{\theta^*}$  and assume that  $\theta^* \in \Theta$ . Let  $(\mathbf{x}_i)_{i \in \mathbb{N}}$  be a dense sequence in  $\mathbb{X}$ . Let  $\mathbf{Y}_n = [Y(x_1), \dots, Y(x_n)]^{\top}$ . Let*

$$\mathcal{L}_n(\theta) = -\frac{1}{2} \log(\det(\mathbf{R}_{\theta})) - \frac{1}{2} \mathbf{Y}_n^{\top} \mathbf{R}_{\theta}^{-1} \mathbf{Y}_n - \frac{n}{2} \log 2\pi,$$

with  $\mathbf{R}_{\theta} = (k_{\theta}(x_i, x_j))_{1 \leq i, j \leq n}$ . Let  $\hat{\theta} \in \arg \max_{\theta \in \Theta} \mathcal{L}_n(\theta)$ . Assume that  $\forall \varepsilon > 0$ ,

$$P(\|\hat{\theta} - \theta^*\| \geq \varepsilon) \xrightarrow{n \rightarrow \infty} 0.$$

Let  $\kappa \in \{0, 1, 2\}$ . Let  $\mathcal{E}_{\kappa}$  be as in Equation (15). Then, we have  $P(Y \in \mathcal{E}_{\kappa}) > 0$  from Lemmas A.3 to A.5, and thus

$$P(\|\hat{\theta} - \theta^*\| \geq \varepsilon \mid Y \in \mathcal{E}_{\kappa}) \xrightarrow{n \rightarrow \infty} 0.$$

*Proof.* We have

$$P(\|\hat{\theta} - \theta^*\| \geq \varepsilon \mid Y \in \mathcal{E}_{\kappa}) = \frac{P(\|\hat{\theta} - \theta^*\| \geq \varepsilon, Y \in \mathcal{E}_{\kappa})}{P(Y \in \mathcal{E}_{\kappa})} \leq \frac{P(\|\hat{\theta} - \theta^*\| \geq \varepsilon)}{P(Y \in \mathcal{E}_{\kappa})}.$$

Since  $P(Y \in \mathcal{E}_{\kappa}) > 0$  is fixed, and  $P(\|\hat{\theta} - \theta^*\| \geq \varepsilon) \xrightarrow{n \rightarrow \infty} 0$ , the result follows.  $\square$

**Proposition 5.2.** *We use the same notations and assumptions as in Proposition 5.1. Let  $\kappa \in \{0, 1, 2\}$  be fixed. Let  $P_{\theta}$  be the distribution of  $Y$  with covariance function  $k_{\theta}$ . Let*

$$\mathcal{L}_{C,n}(\theta) = \mathcal{L}_n(\theta) + \log P_{\theta}(Y \in \mathcal{E}_{\kappa} \mid \mathbf{Y}_n) - \log P_{\theta}(Y \in \mathcal{E}_{\kappa}).$$

Assume that  $\forall \varepsilon > 0$  and  $\forall M < \infty$ ,

$$P\left(\sup_{\|\theta - \theta^*\| \geq \varepsilon} (\mathcal{L}_n(\theta) - \mathcal{L}_n(\theta^*)) \geq -M\right) \xrightarrow{n \rightarrow \infty} 0.$$

Then,

$$P\left(\sup_{\|\theta - \theta^*\| \geq \varepsilon} (\mathcal{L}_{C,n}(\theta) - \mathcal{L}_{C,n}(\theta^*)) \geq -M \mid Y \in \mathcal{E}_{\kappa}\right) \xrightarrow{n \rightarrow \infty} 0.$$

Consequently

$$\operatorname{argmax}_{\theta \in \Theta} \mathcal{L}_n(\theta) \xrightarrow[n \rightarrow \infty]{P} \theta^*, \quad \text{and} \quad \operatorname{argmax}_{\theta \in \Theta} \mathcal{L}_{C,n}(\theta) \xrightarrow[n \rightarrow \infty]{P|Y \in \mathcal{E}_{\kappa}} \theta^*,$$

where  $\xrightarrow[n \rightarrow \infty]{P}$  denotes the convergence in probability under the distribution of  $Y$ , and  $\xrightarrow[n \rightarrow \infty]{P|Y \in \mathcal{E}_{\kappa}}$  denotes the convergence in probability under the distribution of  $Y$  given  $Y \in \mathcal{E}_{\kappa}$ .

*Proof.* We have from Lemmas A.1 and A.2 that  $\forall \varepsilon > 0$

$$P\{\log(P_{\theta^*}(Y \in \mathcal{E}_{\kappa} \mid \mathbf{Y}_n)) \geq \log(1 - \varepsilon) \mid Y \in \mathcal{E}_{\kappa}\} \xrightarrow{n \rightarrow \infty} 1.$$

Hence  $\forall \delta > 0$

$$P\left\{\sup_{\|\theta - \theta^*\| \geq \varepsilon} \log(P_{\theta}(Y \in \mathcal{E}_{\kappa} \mid \mathbf{Y}_n)) - \log(P_{\theta^*}(Y \in \mathcal{E}_{\kappa} \mid \mathbf{Y}_n)) \geq \delta \mid Y \in \mathcal{E}_{\kappa}\right\} \xrightarrow{n \rightarrow \infty} 0.$$

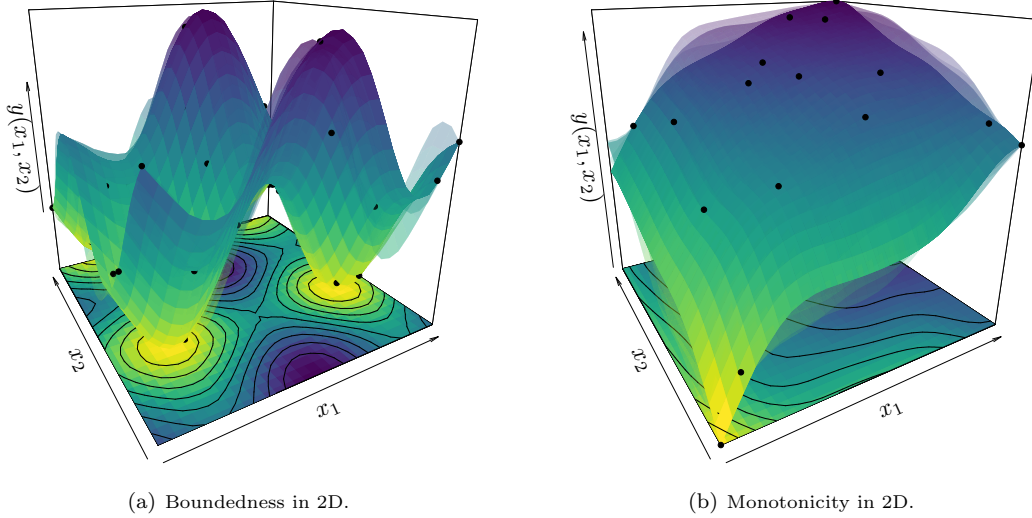


Figure 6: Examples of 2D Gaussian models with different types of constraints for interpolating the toy examples from Subsection 6.1. Each panel shows: training points (black dots), the conditional function (solid surface), and some conditional realizations (light surfaces).

Also, from Lemma A.6, there exists  $\Delta > 0$  so that we have

$$\inf_{\|\theta - \theta^*\| \geq \varepsilon} P_\theta(Y \in \mathcal{E}_\kappa) \geq \Delta > 0,$$

so that

$$\sup_{\|\theta - \theta^*\| \geq \varepsilon} -\log(P_\theta(Y \in \mathcal{E}_\kappa)) + \log(P_{\theta^*}(Y \in \mathcal{E}_\kappa)) \leq -\log(\Delta) < \infty.$$

Hence, the proposition follows. □

## 6 Extension to multidimensional input spaces

### 6.1 2D dimensional case

The finite-dimensional Gaussian representation of Section 3 can be extended to  $d$  dimensional input spaces by tensorisation. For readability, we focus on the case  $d = 2$  with  $\mathcal{D} = [0, 1]^2$  and  $m_1 \times m_2$  knots located on a regular grid. Then, the finite approximation is given by

$$Y_{m_1, m_2}(x_1, x_2) := \sum_{j_1=1}^{m_1} \sum_{j_2=1}^{m_2} \xi_{j_2, j_1} \phi_{j_1}^1(x_1) \phi_{j_2}^2(x_2), \text{ s.t. } \begin{cases} Y_{m_1, m_2}(x_1^i, x_2^i) = y_i, & (i = 1, \dots, n) \\ \xi_{j_2, j_1} \in \mathcal{C}, \end{cases} \quad (16)$$

where  $\xi_{j_2, j_1} = Y(t_{j_1}, t_{j_2})$  and  $(x_1^1, x_2^1), \dots, (x_1^n, x_2^n)$  constitute a DoE. If we follow a similar procedure as in Section 3, we observe that  $\boldsymbol{\xi} = [\xi_{1,1}, \dots, \xi_{1,m_1}, \dots, \xi_{m_2,1}, \dots, \xi_{m_2,m_1}]^\top$  is a zero-mean Gaussian vector with covariance matrix  $\boldsymbol{\Gamma}$  as in Equation (8). Notice that each row of the new matrix  $\boldsymbol{\Phi}$  is given by

$$\boldsymbol{\Phi}_{i,\cdot} = [\phi_1^1(x_1^i) \phi_1^2(x_2^i) \quad \dots \quad \phi_{m_1}^1(x_1^i) \phi_1^2(x_2^i) \quad \dots \quad \phi_1^1(x_1^i) \phi_{m_2}^2(x_2^i) \quad \dots \quad \phi_{m_1}^1(x_1^i) \phi_{m_2}^2(x_2^i)],$$

for all  $i = 1, \dots, n$ . Finally, the posterior distribution of Equation (10) can be computed, and the routine follows Algorithm 1. The module from Equation (16) was also used in (Maatouk and Bay, 2017) for monotonicity constraints in multidimensional cases.

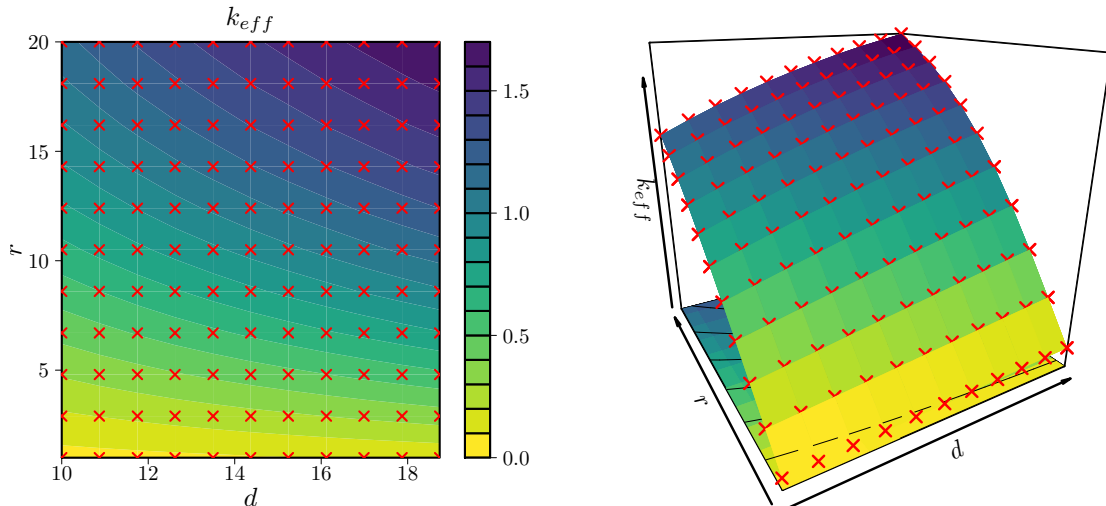


Figure 7: Nuclear criticality safety assessments: Godiva’s dataset. (Left) 2D visualization of the  $k_{eff}$  values measured over a regular grid. (Right) 3D visualization of the  $k_{eff}$  data.

Figure 6 shows two examples where boundedness or monotonicity are exhibited. We used a 2D SE covariance function with parameters ( $\sigma^2 = 1.0$ ,  $\theta_1 = 0.2$ ,  $\theta_2 = 0.2$ ).<sup>3</sup> The training points were generated with a maximin Latin hypercube DoE over  $[0, 1]^2$ . The functions are: 6(a)  $y(x_1, x_2) = -\frac{1}{2}[\sin(9x_1) - \cos(9x_2)]$ , and 6(b)  $y(x_1, x_2) = \arctan(5x_1) + \arctan(x_2)$ .

## 6.2 2D application: nuclear criticality safety

For assessing the stability of neutron production in nuclear reactors, safety criteria based on the effective neutron multiplication factor  $k_{eff}$  are commonly used (IAE, 2014; Fernex et al., 2005). This factor is defined as the ratio of the total number of neutrons produced by a fission chain reaction to the total number of neutrons lost by absorption and leakage. Besides the geometry and composition of fissile materials (e.g. mass, density),  $k_{eff}$  is sensitive to other types of parameters like the structure materials characteristics (e.g. concrete), and the presence of specific materials (e.g. moderators). Since the optimal control of an individual parameter or a combination of them can lead to safe conditions, the understanding of their influence in criticality safety assessment is crucial.

In this section, we applied the proposed framework to a dataset provided by the “Institut de Radioprotection et de Sûreté Nucléaire” (IRSN), France. The  $k_{eff}$  factor was obtained from a nuclear reactor called “Lady Godiva device” originally situated at Los Alamos National Laboratory (LANL), New Mexico, U.S., where uranium materials were managed. Two input parameters of the uranium sphere are considered: its radius  $r$  and its density  $d$ . The dataset contains 121 observations in a regular grid of  $11 \times 11$  locations (see Figure 7). Notice that, on the domain considered for the input variables,  $k_{eff}$  increases as the geometry and density of the uranium sphere increase.

We trained different Gaussian models whether the inequality constraints are considered or not. For all the models, we normalised the input space to be in  $[0, 1]^2$ . We used the same 2D SE covariance functions as for the example from Figure 6. For the unconstrained model, we used multistart for the ML optimisation with six initial vectors of covariance parameters  $\Theta = (\sigma^2, \theta_1, \theta_2)$  with  $\sigma^2 \in [0.2, 1]$  and  $\theta_1, \theta_2 \in [0.1, 0.9]$ . For the constrained models, since the  $k_{eff}$  factor indicates the production rate of neutron population, the output of the constrained processes has to be positive. Taking also into account the non-decreasing behaviour, we also consider the monotonicity constraints. We estimated the covariance parameters by MLE or CMLE. We trained both unconstrained and constrained models with a fixed maximin Latin hypercube DoE at eight locations extracted from the normalised grid. We used the remaining data to asses the quality of prediction tasks.

<sup>3</sup>2D SE covariance function:  $k_{\theta}(\mathbf{x} - \mathbf{x}') = \sigma^2 \exp \left\{ -\frac{(x_1 - x'_1)^2}{2\theta_1^2} - \frac{(x_2 - x'_2)^2}{2\theta_2^2} \right\}$  with  $\theta = (\sigma^2, \theta_1, \theta_2)$ .

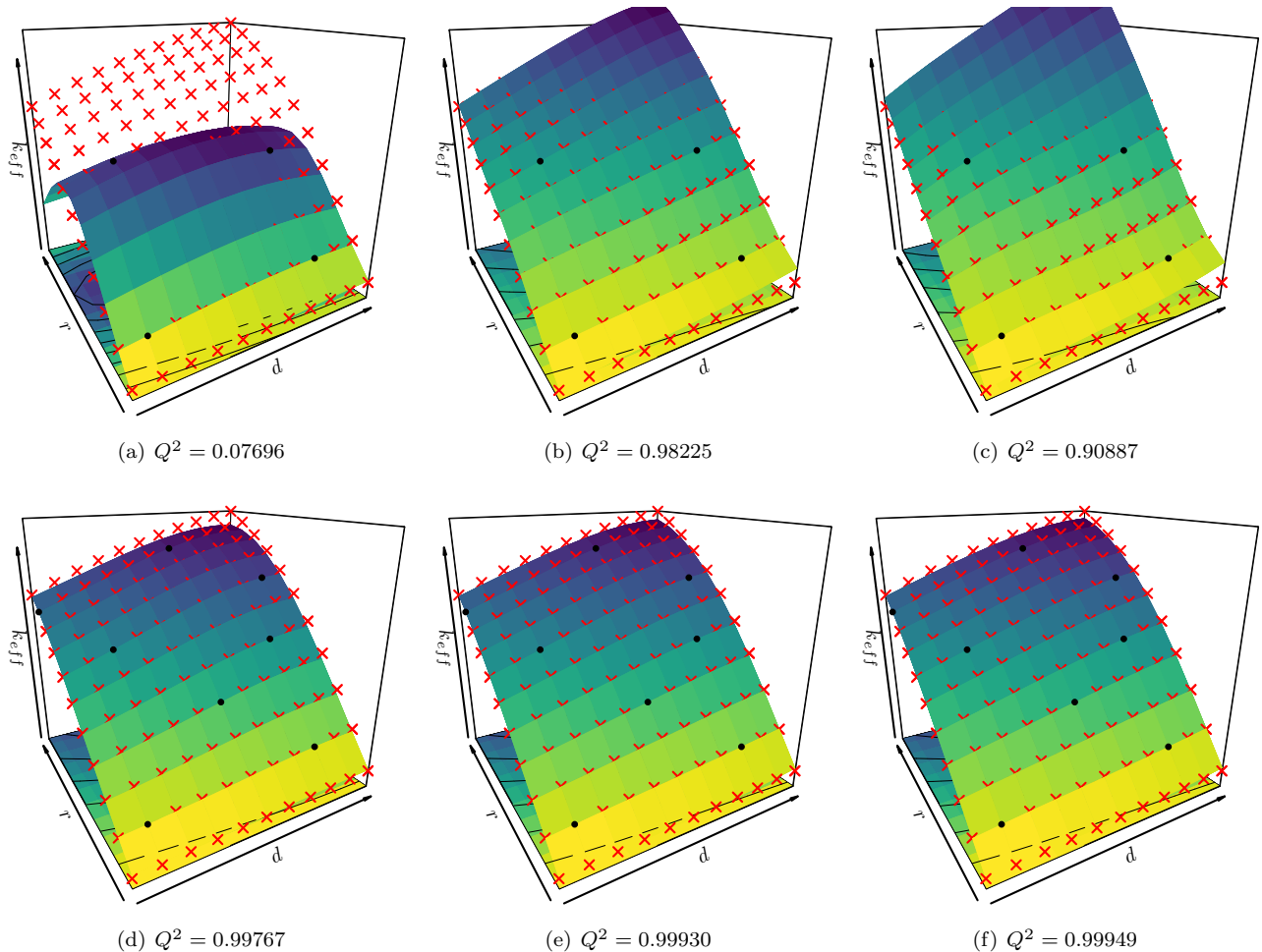


Figure 8: 2D Gaussian models for interpolating the Godiva’s dataset. Unconstrained GP models are trained using MLE (left column) and using (a) four, or (d) eight training points from the proposed maximin Latin hypercube DoE. Constrained GP models are trained either using MLE (middle column) and CMLE (right column). Each panel shows: training and test points (black dots and red crosses), the conditional mean function (solid surface), and the  $Q^2$  criterion (subcaptions).

Figure 8 shows the performance of the proposed models using four or eight points from the proposed fixed DoE. For the unconstrained models, we observe that the quality of the predictions depends strongly on both the amount of training data and their distribution in the input space. Notice from Figure 8(a) that if only few training points are available, predictions are poor and they do not satisfy positive and non-decreasing behaviours. In Figure 8(d), we observe that if the training data are large enough and cover the input space, the unconstrained model behaves well and provides reliable predictions. On the other hand, we observe that the constrained models produce accurate prediction results also when the training set is small.

Because the prediction accuracy depends on the training set, we repeated the procedure with twenty different random Latin hypercube DoEs using several values of  $n$ . We used the  $Q^2$  and PVA criteria to evaluate the quality of the predictions (see Subsection 5.2). Figure 9 shows that the constrained models often outperform the unconstrained ones. Notice that although the  $Q^2$  results obtained by the unconstrained model are comparable with the constrained ones when the number of training points is large enough, we observe, according to the PVA criterion, that the constrained models provide more reliable confidence intervals. This means that, if we consider both positivity and monotonicity conditions to take into account the physics of the  $k_{eff}$  factor, we can obtain more informative and robust models. Furthermore, we observe that the unconstrained MLE achieves a good tradeoff between prediction accuracy/reliability and computational cost.



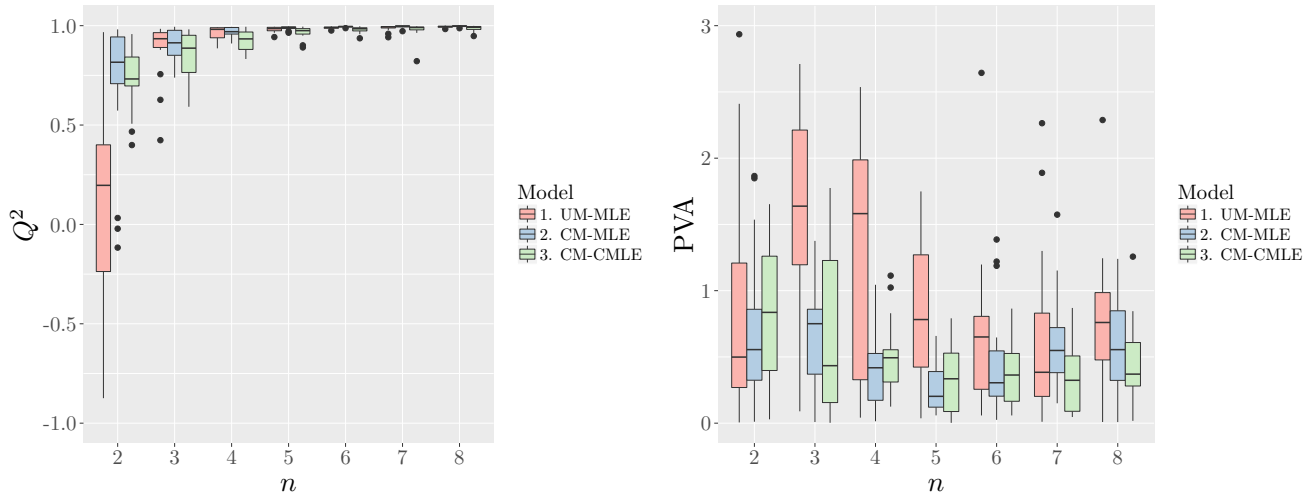


Figure 9: Assessment of the Gaussian models for interpolating the dataset from Figure 7 using different number of training points  $n$  and using twenty different random Latin hypercube designs. Predictive accuracy is evaluated using the (left)  $Q^2$  and (right) PVA criteria. Results are showed for the unconstrained model (UM) using MLE (red), and the constrained models (CM) using either MLE (blue) or CMLE (green).

## 7 Conclusions

Continuing the approach proposed in (Maatouk and Bay, 2017), we have introduced a full Gaussian-based framework to satisfy linear sets of inequality constraints. The proposed finite-dimensional approach takes into account the inequalities for both data interpolation and covariance parameter estimation. Because the posterior distribution is expressed as a truncated multinormal distribution, we compared different MCMC methods as well as exact sampling methods. According to experiments, we concluded that the Hamiltonian Monte Carlo-based sampler adapts to our needs. For parameter estimation, we suggested the constrained likelihood which takes into account the inequality constraints. We showed that, loosely speaking, any consistency result for ML with unconstrained GPs, is preserved when adding boundedness, monotonicity and convexity constraints. Furthermore, this consistency occurs for both the unconditional and conditional likelihood functions.

We tested our model in both synthetic and real-world data in 1D or 2D. According to the experimental results under different types of inequalities, the proposed framework fits properly the observations and provides realistic confidence intervals. Our approach is also flexible enough to satisfy multiple inequality conditions and to deal with specific types of constraints which are sequentially activated. Finally, as we showed in the 2D nuclear criticality safety assessment, the proposed framework provides reliable predictions on both data prediction and uncertainty quantification satisfying the inequality constraints exhibited by the neutron population (positivity and monotonicity conditions). We also observe that the unconstrained MLE achieves a good tradeoff between prediction accuracy/reliability and computational cost.

The framework presented in this paper can be improved in different ways. First, the precision of the results depends on the number of knots  $m$  used in the finite approximation. For higher values of  $m$ , the interpolation is better but more expensive. In this sense, we could consider the optimal location of the knots over the input space instead of using regular grids. This potentially allows to reduce the computational cost of the full framework. Second, as we discussed for 2D input spaces, the model can be generalized to  $d$  dimensional problems. However, due to its tensor structure, its practical application could be time-consuming. Hence, there is a need to find an extension of the model to higher dimensions that can be applied in real-world problems. Finally, the estimation of the orthonormal Gaussian probabilities can be improved in order to exploit the advantages of the constrained likelihood.

## A Conditional maximum likelihood: asymptotic properties

We detail in this appendix section the conditions and lemmas we used in Propositions 5.1 and 5.2 from Subsection 5.3. We use the same notations and assumptions as in this subsection.

**Condition A.1.** Let  $\mathbf{x}, \mathbf{x}' \in \mathbb{X}$ . For a fixed  $\kappa \in \{0, 1, 2\}$ , assume one of the following conditions:

- If  $\kappa = 0$ . Assume that  $Y$  has continuous trajectories. Let  $k$  be the covariance function of  $Y$ . Let

$$d_k(\mathbf{x}, \mathbf{x}') = \sqrt{k(\mathbf{x}, \mathbf{x}) + k(\mathbf{x}', \mathbf{x}') - 2k(\mathbf{x}, \mathbf{x}')}.$$

Let  $N(\mathbb{X}, d_k, \rho)$  be the minimum number of balls with radius  $\rho$  (w.r.t. the distance  $d_k$ ), required to cover  $\mathbb{X}$ . Assume that

$$\int_0^\infty \sqrt{\log(N(\mathbb{X}, d_k, \rho))} d\rho < \infty.$$

Assume also that the Fourier transform  $\widehat{k}$  of  $k$  satisfies

$$\exists P < \infty \quad \text{so that as } \|w\| \rightarrow \infty, \quad \widehat{k}(w) \|w\|^P \rightarrow \infty.$$

- If  $\kappa = 1$ . Assume that  $Y$  has  $C^1$  trajectories. Let  $k_i^{[1]}$  be the covariance function of  $\frac{\partial}{\partial x_i} Y$ . Let  $d_{k_i^{[1]}}$  and  $N(\mathbb{X}, d_{k_i^{[1]}}, \rho)$  be defined as  $d_k$  and  $N(\mathbb{X}, d_k, \rho)$  for  $\kappa = 0$ . Assume that

$$\int_0^\infty \sqrt{\log(N(\mathbb{X}, d_{k_i^{[1]}}, \rho))} d\rho < \infty, \quad \forall i = 1, \dots, d.$$

Assume also that the Fourier transform  $\widehat{k}_i^{[1]}$  of  $k_i^{[1]}$  satisfies the same conditions as for  $\kappa = 0$ , for  $i = 1, \dots, d$ .

- If  $\kappa = 2$ . Assume that  $Y$  has  $C^2$  trajectories. Let  $k_{i,j}^{[2]}$  be the covariance function of  $\frac{\partial^2}{\partial x_i \partial x_j} Y$ . Let  $d_{k_{i,j}^{[2]}}$  and  $N(\mathbb{X}, d_{k_{i,j}^{[2]}}, \rho)$  be defined as  $d_k$  and  $N(\mathbb{X}, d_k, \rho)$  for  $\kappa = 0$ . Assume that

$$\int_0^\infty \sqrt{\log(N(\mathbb{X}, d_{k_{i,j}^{[2]}}, \rho))} d\rho < \infty, \quad \forall i, j = 1, \dots, d.$$

Assume also that the Fourier transform  $\widehat{k}_{i,j}^{[2]}$  of  $k_{i,j}^{[2]}$  satisfies the same conditions as for  $\kappa = 0$ , for  $i, j = 1, \dots, d$ .

**Lemma A.1.** Let  $0 \leq \ell < u \leq \infty$ . Let

$$P_{n,\ell,u}(\mathbf{Y}_n) = P_{\theta^*}(Y \in \mathcal{E}_0 | \mathbf{Y}_n).$$

Then,  $\forall \varepsilon \geq 0$ , we have

$$P(P_{n,\ell,u}(\mathbf{Y}_n) \leq 1 - \varepsilon | Y \in \mathcal{E}_0) \xrightarrow{n \rightarrow \infty} 0.$$

*Proof.* From Lemma A.3 we have  $P(Y \in \mathcal{E}_0) > 0$ . Hence, it is sufficient to show

$$P(P_{n,\ell,u}(\mathbf{Y}_n) \leq 1 - \varepsilon, Y \in \mathcal{E}_0) \xrightarrow{n \rightarrow \infty} 0.$$

The term  $P_{n,\ell,u}(\mathbf{Y}_n)$ , being a conditional expectation, is a martingale with respect to the  $\sigma$ -algebra generated by  $Y(\mathbf{x}_1), \dots, Y(\mathbf{x}_n)$ . Furthermore,  $0 \leq P_{n,\ell,u}(\mathbf{Y}_n) \leq 1$ . Hence

$$P_{n,\ell,u}(\mathbf{Y}_n) \xrightarrow[n \rightarrow \infty]{a.s.} P(Y \in \mathcal{E}_0 | \mathcal{F}_\infty),$$

where  $\mathcal{F}_\infty$  is the  $\sigma$ -algebra generated by  $[Y(\mathbf{x}_i)]_{i \in \mathbb{N}}$  using Theorem 6.2.3 from (Kallenberg, 2002). Let  $\mu_n$  and  $k_n$  be the mean and the covariance function (respectively) of  $Y$  given  $\mathbf{Y}_n$ . From proposition 2.8 in (Bect et al., 2016), the conditional distribution of  $Y$  given  $\mathcal{F}_\infty$  is the distribution of a GP with mean function  $\mu_\infty$  and covariance function

$k_\infty$ . Furthermore, a.s.,  $\mu_n$  and  $k_n$  converge uniformly to  $\mu_\infty$  and  $k_\infty$ , respectively. Hence we can show simply that, because  $(x_i)_{i \in \mathbb{N}}$  is dense in  $\mathbb{X}$ , we have a.s.  $\mu_\infty = Y$  and  $k_\infty$  is the zero function. Hence a.s. if  $Y \in \mathcal{E}_0$  holds, then

$$P(Y \in \mathcal{E}_0 \mid \mathcal{F}_\infty) = 1, \quad \text{so that} \quad P_{n,\ell,u}(\mathbf{Y}_n) \xrightarrow{n \rightarrow \infty} 1.$$

Hence by the dominated convergence theorem

$$P(P_{n,\ell,u}(\mathbf{Y}_n) \leq 1 - \varepsilon, Y \in \mathcal{E}_0) \xrightarrow{n \rightarrow \infty} 0.$$

□

**Lemma A.2.** *Let  $\kappa = \{1, 2\}$ . Let*

$$P_n(\mathbf{Y}_n) = P_{\theta^*}(Y \in \mathcal{E}_\kappa \mid \mathbf{Y}_n).$$

*Then,  $\forall \varepsilon > 0$ , we have*

$$P(P_n(\mathbf{Y}_n) \leq 1 - \varepsilon \mid Y \in \mathcal{E}_\kappa) \xrightarrow{n \rightarrow \infty} 0.$$

*Proof.* The proof is the same as that of Lemma A.1. In particular, we remark that  $\mathbf{1}_{Y \in \mathcal{E}_\kappa}$  is a measurable random variable, as  $Y$  has  $C^\kappa$  trajectories. □

**Lemma A.3.** *Let  $\kappa = 0$ . Assume that Condition A.1 is satisfied. Then*

$$P(Y \in \mathcal{E}_0) > 0, \quad \text{for} \quad -\infty \leq \ell < u \leq \infty.$$

*Proof.* We first prove that for all  $\delta > 0$

$$P(\forall \mathbf{x} \in \mathbb{X} : |Y(\mathbf{x})| \leq \delta) > 0.$$

This result is true and appears implicitly in the literature about small ball estimates for GP (Li and Linde, 1999). We nevertheless provide a proof of it for self-consistency. Let  $(\mathbf{v}_i)_{i \in \mathbb{N}}$  be a dense sequence in  $\mathbb{X}$ . Let  $\mathbf{Y}_v = [Y(\mathbf{v}_1), \dots, Y(\mathbf{v}_n)]^\top$ . Let  $\mu_n$  and  $k_n$  be the mean and the covariance function of  $Y$  given  $\mathbf{Y}_v$ . Then we let

$$d_{k_n}^2(\mathbf{x}, \mathbf{x}') = \text{var} \{ (Y(\mathbf{x}) - Y(\mathbf{x}')) \mid \mathcal{F}_n \},$$

where  $\mathcal{F}_n = \sigma(Y(\mathbf{v}_1), \dots, Y(\mathbf{v}_n))$ . Thus

$$d_{k_n}^2(\mathbf{x}, \mathbf{x}') = \mathbb{E} \{ \text{var} \{ (Y(\mathbf{x}) - Y(\mathbf{x}')) \mid \mathcal{F}_n \} \} \leq \text{var} \{ (Y(\mathbf{x}) - Y(\mathbf{x}')) \} = d_k^2(\mathbf{x}, \mathbf{x}'),$$

from the law of total variance. Hence  $N(\mathbb{X}, d_{k_n}, \rho) \leq N(\mathbb{X}, d_k, \rho) \forall \rho$ . Also, from Theorem 2.10 in (Azaïs and Wschebor, 2009) (together with an union bound and using that  $\max_{\mathbf{x} \in \mathbb{X}} Y(\mathbf{x})$  and  $\max_{\mathbf{x} \in \mathbb{X}} [-Y(\mathbf{x})]$  have the same law) we have, with  $C$  an universal constant,

$$\begin{aligned} \mathbb{E} \left\{ \max_{\mathbf{x} \in \mathbb{X}} |Y(\mathbf{x}) - \mu_n(\mathbf{x})| \right\} &\leq C \int_0^\infty \sqrt{\log(N(\mathbb{X}, d_{k_n}, \rho))} d\rho \\ &= C \int_0^{2\sqrt{\sup_{\mathbf{x} \in \mathbb{X}} k_n(\mathbf{x}, \mathbf{x})}} \sqrt{\log(N(\mathbb{X}, d_{k_n}, \rho))} d\rho \\ &\leq C \int_0^{2\sqrt{\sup_{\mathbf{x} \in \mathbb{X}} k_n(\mathbf{x}, \mathbf{x})}} \sqrt{\log(N(\mathbb{X}, d_k, \rho))} d\rho. \end{aligned}$$

This last integral goes to 0 as  $n \rightarrow \infty$  because  $\sup_{\mathbf{x} \in \mathbb{X}} k_n(\mathbf{x}, \mathbf{x}) \rightarrow 0$  (see the proof of Lemma A.1), and because of Condition A.1. Hence  $\max_{\mathbf{x} \in \mathbb{X}} |Y(\mathbf{x}) - \mu_n(\mathbf{x})|$  goes to 0 in probability. Furthermore,  $\mathcal{P} = P(\forall \mathbf{x} \in \mathbb{X}, -\delta \leq Y(\mathbf{x}) \leq \delta)$  satisfies

$$\begin{aligned} \mathcal{P} &\geq P \left( \forall \mathbf{x} \in \mathbb{X}, -\frac{\delta}{2} \leq \mu_n(\mathbf{x}) \leq \frac{\delta}{2}, -\frac{\delta}{2} \leq Y(\mathbf{x}) - \mu_n(\mathbf{x}) \leq \frac{\delta}{2} \right) \\ &= P \left( \forall \mathbf{x} \in \mathbb{X}, -\frac{\delta}{2} \leq \mu_n(\mathbf{x}) \leq \frac{\delta}{2} \right) P \left( \forall \mathbf{x} \in \mathbb{X}, -\frac{\delta}{2} \leq Y(\mathbf{x}) - \mu_n(\mathbf{x}) \leq \frac{\delta}{2} \right), \end{aligned}$$

since the distribution of  $Y - \mu_n$  does not depend on  $\mathbf{Y}_v$ . We now fix  $n \in \mathbb{N}$  for which the second probability is non-zero (the existence is guaranteed from above). Then, the first probability is non-zero by continuity since, when  $\mathbf{Y}_v = \mathbf{0}$ , then  $\mu_n$  is the zero function. Hence we have

$$P(\forall \mathbf{x} \in \mathbb{X} : |Y(\mathbf{x})| \leq \delta) > 0.$$

Let  $f$  be a  $C^\infty$  function on  $\mathbb{R}^d$ , square integrable, satisfying

$$\forall \mathbf{x} \in \mathbb{X}, \quad \ell + \delta \leq f(\mathbf{x}) \leq u - \delta,$$

for  $\delta > 0$ . ( $f$  exists for  $\delta > 0$  small enough, and can be taken for instance as  $f(\mathbf{x}) = \exp\{-\tau\|\mathbf{x} - \mathbf{x}_0\|^2\} \left[\frac{u+\ell}{2}\right]$  with  $\tau > 0$  small enough, and for any  $\mathbf{x}_0 \in \mathbb{X}$ ). Let  $Z$  be a GP with covariance function  $k$  and mean function  $f$ . Then, from what we have shown before, we have

$$P(\forall \mathbf{x} \in \mathbb{X} : |Z(\mathbf{x}) - f(\mathbf{x})| \leq \delta) > 0,$$

so that

$$P(\forall \mathbf{x} \in \mathbb{X} : \ell \leq Z(\mathbf{x}) \leq u) > 0.$$

From (Yadrenko, 1983) (p.138), as discussed by (Stein, 1999) (p.121), the Gaussian measures of  $Y$  and  $Z$  are equivalent. Thus

$$P(Y \in \mathcal{E}_0) = P(\forall \mathbf{x} \in \mathbb{X} : \ell \leq Y(\mathbf{x}) \leq u) > 0.$$

□

**Lemma A.4.** *Let  $\kappa = 1$ . Assume that Condition A.1 is satisfied. Then*

$$P(Y \in \mathcal{E}_1) > 0.$$

*Proof.* We first prove that for all  $\delta > 0$

$$P\left(\forall i = 1, \dots, d, \forall \mathbf{x} \in \mathbb{X} : \left| \frac{\partial}{\partial x_i} Y(\mathbf{x}) \right| \leq \delta\right) > 0.$$

We let  $(\mathbf{v}_i)_{i \in \mathbb{N}}$  and  $\mathbf{Y}_v$  be defined as in the proof of Lemma A.3. Then, as in this proof we can show that for all  $i = 1, \dots, d$

$$\max_{\mathbf{x} \in \mathbb{X}} \left| \frac{\partial}{\partial x_i} Y(\mathbf{x}) - \mathbb{E} \left\{ \frac{\partial}{\partial x_i} Y(\mathbf{x}) \middle| \mathbf{Y}_v \right\} \right| \xrightarrow[n \rightarrow \infty]{P} 0.$$

Furthermore,  $\mathcal{P} = P\left(\forall i = 1, \dots, d, \forall \mathbf{x} \in \mathbb{X} : \left| \frac{\partial}{\partial x_i} Y(\mathbf{x}) \right| \leq \delta\right)$  satisfies

$$\begin{aligned} \mathcal{P} &\geq P\left(\forall i = 1, \dots, d, \forall \mathbf{x} \in \mathbb{X}, -\frac{\delta}{2} \leq \mathbb{E} \left\{ \frac{\partial}{\partial x_i} Y(\mathbf{x}) \middle| \mathbf{Y}_v \right\} \leq \frac{\delta}{2}, \right. \\ &\quad \left. \forall i = 1, \dots, d, \forall \mathbf{x} \in \mathbb{X}, -\frac{\delta}{2} \leq \frac{\partial}{\partial x_i} Y(\mathbf{x}) - \mathbb{E} \left\{ \frac{\partial}{\partial x_i} Y(\mathbf{x}) \middle| \mathbf{Y}_v \right\} \leq \frac{\delta}{2} \right) \\ &= P\left(\forall i = 1, \dots, d, \forall \mathbf{x} \in \mathbb{X}, -\frac{\delta}{2} \leq \mathbb{E} \left\{ \frac{\partial}{\partial x_i} Y(\mathbf{x}) \middle| \mathbf{Y}_v \right\} \leq \frac{\delta}{2} \right) \\ &\quad \times P\left(\forall i = 1, \dots, d, \forall \mathbf{x} \in \mathbb{X}, -\frac{\delta}{2} \leq \frac{\partial}{\partial x_i} Y(\mathbf{x}) - \mathbb{E} \left\{ \frac{\partial}{\partial x_i} Y(\mathbf{x}) \middle| \mathbf{Y}_v \right\} \leq \frac{\delta}{2} \right). \end{aligned}$$

Notice that the last equality holds because the distribution of the process  $\mathbf{x} \rightarrow \frac{\partial}{\partial x_i} Y(\mathbf{x}) - \mathbb{E} \left\{ \frac{\partial}{\partial x_i} Y(\mathbf{x}) \middle| \mathbf{Y}_v \right\}$  does not depend on  $\mathbf{Y}_v$ . We now fix  $n \in \mathbb{N}$  so that the second probability is non-zero (the existence is guaranteed from above). Then, the first probability is non-zero by continuity since, when  $\mathbf{Y}_v = \mathbf{0}$ , then for all  $i = 1, \dots, d$ ,  $\mathbb{E} \left\{ \frac{\partial}{\partial x_i} Y \middle| \mathbf{Y}_v \right\}$  is the zero function. Hence, we have obtained

$$P\left(\forall i = 1, \dots, d, \forall \mathbf{x} \in \mathbb{X} : \left| \frac{\partial}{\partial x_i} Y(\mathbf{x}) \right| \leq \delta\right).$$

We now conclude the proof in the same way as for Lemma A.3. We consider the mean function

$$f(\mathbf{x}) = \left[ \sum_{i=1}^d x_i \right] \exp\{-\tau \|\mathbf{x} - \mathbf{x}_0\|^2\},$$

with  $\mathbf{x}_0 \in \mathbb{X}$  and  $\tau > 0$ . For  $\tau$  small enough,  $f$  is  $C^\infty$ , square integrable, and satisfies

$$\forall i = 1, \dots, d, \quad \forall \mathbf{x} \in \mathbb{X}, \quad \frac{\partial}{\partial x_i} f(\mathbf{x}) \geq \frac{1}{2}.$$

Then, we conclude the proof as in the proof of Lemma A.3.  $\square$

**Lemma A.5.** *Let  $\kappa = 2$ . Assume that Condition A.1 is satisfied. Then,*

$$P(Y \in \mathcal{E}_2) > 0.$$

*Proof.* We first prove that for all  $\delta > 0$

$$P\left(\forall i, j = 1, \dots, d, \forall \mathbf{x} \in \mathbb{X} : \left| \frac{\partial^2}{\partial x_i \partial x_j} Y(\mathbf{x}) \right| \leq \delta\right) > 0.$$

This is done in a similar way as for showing  $P\left(\forall i = 1, \dots, d, \forall \mathbf{x} \in \mathbb{X} : \left| \frac{\partial}{\partial x_i} Y(\mathbf{x}) \right| \leq \delta\right) > 0$  in the proof of Lemma A.4. We then conclude similarly as the rest of the proof this Lemma. In particular, we consider the mean function

$$f(\mathbf{x}) = \left[ \sum_{i=1}^d x_i^2 \right] \exp\{-\tau \|\mathbf{x} - \mathbf{x}_0\|^2\},$$

with  $\mathbf{x}_0 \in \mathbb{X}$  and  $\tau > 0$ . Let  $\lambda_{\inf}(M)$  be the smallest eigenvalue of a symmetric matrix  $M$ . Then, for  $\tau$  small enough,  $f$  is  $C^\infty$ , square integrable, and satisfies

$$\forall \mathbf{x} \in \mathbb{X}, \quad \lambda_{\inf}\left(\frac{\partial^2}{\partial \mathbf{x}^2} f(\mathbf{x})\right) \geq 1.$$

$\square$

**Lemma A.6.** *Let  $\kappa \in \{0, 1, 2\}$ . Assume that Condition A.1 holds. Let  $Y_\theta$  be the GP defined by*

$$Y_\theta(t) = \sigma Y(\theta_1 t_1, \dots, \theta_d t_d).$$

*Let  $P_\theta^\kappa = P(Y_\theta \in \mathcal{E}_\kappa)$  (see Equation (15)). Then,*

$$\inf_{\theta \in \Theta} P_\theta^\kappa > 0.$$

*Proof.* We do the proof for  $\kappa = 2$ . The proof for  $\kappa = 0, 1$  is similar. Let  $\varepsilon > 0$  and let

$$P_{\theta, \varepsilon}^\kappa = \mathbb{E} \left\{ \mathbb{1}_{I(Y_\theta) \geq \varepsilon} + \frac{I(Y_\theta)}{\varepsilon} \mathbb{1}_{0 \leq I(Y_\theta) \leq \varepsilon} \right\},$$

with  $I(Y_\theta) = \inf_{\mathbf{x} \in \mathbb{X}} \lambda_{\inf}\left(\frac{\partial^2}{\partial \mathbf{x}^2} Y_\theta(\mathbf{x})\right)$ . We have  $P_{\theta, \varepsilon}^\kappa \leq P_\theta^\kappa$  for all  $\varepsilon > 0$ . With the proof of Lemma A.5, we also obtain for  $\varepsilon > 0$  small enough

$$\forall \theta \in \Theta \quad P_{\theta, \varepsilon}^\kappa > 0.$$

Hence, the proof is concluded, by compactity, if we show that  $\theta \rightarrow P_{\theta, \varepsilon}^\kappa$  is a continuous function on  $\Theta$ . Let us show this. Let  $\theta = (\sigma_1^2, \theta_1, \dots, \theta_d) \in (0, \infty)^{d+1}$  and  $\theta_n = (\sigma_n^2, \theta_{n1}, \dots, \theta_{nd}) \rightarrow \theta$ . We have

$$\frac{\partial^2}{\partial x_i \partial x_j} Y_{\theta_n}(\mathbf{x}) = \sigma_n \left( (\theta_n)_i (\theta_n)_j \frac{\partial^2}{\partial x_i \partial x_j} Y(\theta_{n1} x_1, \dots, \theta_{nd} x_d) \right).$$

Hence, because  $Y$  is  $C^2$ , we have a.s.

$$\sup_{\mathbf{x} \in \mathbb{X}} \left\| \frac{\partial^2}{\partial \mathbf{x}^2} Y_{\theta_n}(\mathbf{x}) - \frac{\partial^2}{\partial \mathbf{x}^2} Y_{\theta}(\mathbf{x}) \right\| \xrightarrow{n \rightarrow \infty} 0,$$

for any matrix norm  $\|\cdot\|$ . Hence also since  $Y$  is  $C^2$ , we can show, a.s.

$$\left( \inf_{\mathbf{x} \in \mathbb{X}} \lambda_{\inf} \left( \frac{\partial^2}{\partial \mathbf{x}^2} Y_{\theta_n}(\mathbf{x}) \right) - \inf_{\mathbf{x} \in \mathbb{X}} \lambda_{\inf} \left( \frac{\partial^2}{\partial \mathbf{x}^2} Y_{\theta}(\mathbf{x}) \right) \right) \xrightarrow{n \rightarrow \infty} 0.$$

Hence, we conclude by dominated convergence observing that  $t \rightarrow (\mathbf{1}_{t \geq \varepsilon} + \frac{t}{\varepsilon} \mathbf{1}_{0 \leq t \leq \varepsilon})$  is a continuous function on  $\mathbb{R}$ .  $\square$

## Acknowledgement

This research was conducted within the frame of the Chair in Applied Mathematics OQUAIDO, gathering partners in technological research (BRGM, CEA, IFPEN, IRSN, Safran, Storengy) and academia (CNRS, Ecole Centrale de Lyon, Mines Saint-Etienne, University of Grenoble, University of Nice, University of Toulouse) around advanced methods for Computer Experiments. We thank Yann Richet (IRSN) for providing the nuclear criticality safety data.

## References

- (2014). *Criticality Safety in the Handling of Fissile Material: Specific Safety Guide*. International Atomic Energy Agency (IAEA). IAEA Safety Standards Series No. SSG-27.
- Azaïs, J.-M. and Wschebor, M. (2009). *Level sets and extrema of random processes and fields*. John Wiley & Sons.
- Bachoc, F. (2013). Cross validation and maximum likelihood estimations of hyper-parameters of Gaussian processes with model misspecification. *Computational Statistics & Data Analysis*, 66:55 – 69.
- Bay, X., Grammont, L., and Maatouk, H. (2016). Generalization of the Kimeldorf-Wahba correspondence for constrained interpolation. *Electronic journal of statistics*, 10(1):1580–1595.
- Bect, J., Bachoc, F., and Ginsbourger, D. (2016). A supermartingale approach to Gaussian process based sequential design of experiments. *ArXiv e-prints*.
- Botev, Z. I. (2017). The normal law under linear restrictions: simulation and estimation via minimax tilting. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 79(1):125–148.
- Brooks, S., Gelman, A., Jones, G., and Meng, X. (2011). *Handbook of Markov Chain Monte Carlo*. Chapman & Hall/CRC Handbooks of Modern Statistical Methods. CRC Press.
- Cousin, A., Maatouk, H., and Rullière, D. (2016). Kriging of financial term-structures. *European Journal of Operational Research*, 255(2):631–648.
- Da Veiga, S. and Marrel, A. (2012). Gaussian process modeling with inequality constraints. *Annales de la faculté des sciences de Toulouse Mathématiques*, 21(3):529–555.
- Du, J., Zhang, H., and Mandrekar, V. S. (2009). Fixed-domain asymptotic properties of tapered maximum likelihood estimators. *Ann. Statist.*, 37(6A):3330–3361.
- Duane, S., Kennedy, A. D., Pendleton, B. J., and Roweth, D. (1987). Hybrid Monte Carlo. *Physics letters B*, 195(2):216–222.
- Fernex, F., Heulers, L., Jacquet, O., Miss, J., and Richet, Y. (2005). The MORET 4B Monte Carlo code - New features to treat complex criticality systems. In *M&C international conference on mathematics and computation supercomputing, reactor physics and nuclear and biological application, Avignon, France*, volume 59.



- Genz, A. (1992). Numerical computation of multivariate normal probabilities. *Journal of Computational and Graphical Statistics*, 1:141–150.
- Golchi, S., Bingham, D. R., Chipman, H., and Campbell, D. A. (2015). Monotone emulation of computer experiments. *SIAM/ASA Journal on Uncertainty Quantification*, 3(1):370–392.
- Gong, L. and Flegal, J. M. (2016). A practical sequential stopping rule for high-dimensional Markov Chain Monte Carlo. *Journal of Computational and Graphical Statistics*, 25(3):684–700.
- Kallenberg, O. (2002). *Foundations of Modern Probability*. Probability and its Applications. Springer New York.
- Kocijan, J. (2016). *Modelling and Control of Dynamic Systems Using Gaussian Process Models*. Advances in Industrial Control. Springer International Publishing, 1 edition.
- Lan, S. and Shahbaba, B. (2016). *Sampling Constrained Probability Distributions Using Spherical Augmentation*, pages 25–71. Springer International Publishing, Cham.
- Li, W. V. and Linde, W. (1999). Approximation, metric entropy and small ball estimates for Gaussian measures. *Ann. Probab.*, 27(3):1556–1578.
- Loh, W. L. (2005). Fixed-domain asymptotics for a subclass of Matérn-type Gaussian random fields. *Ann. Statist.*, 33(5):2344–2394.
- Loh, W. L. and Lam, T. K. (2000). Estimating structured correlation matrices in smooth Gaussian random field models. *Ann. Statist.*, 28(3):880–904.
- Maatouk, H. and Bay, X. (2016). *A New Rejection Sampling Method for Truncated Multivariate Gaussian Random Variables Restricted to Convex Sets*, pages 521–530. Springer International Publishing, Cham.
- Maatouk, H. and Bay, X. (2017). Gaussian process emulators for computer experiments with inequality constraints. *Mathematical Geosciences*, 49(5):557–582.
- Maatouk, H., Roustant, O., and Richet, Y. (2015). Cross-validation estimations of hyper-parameters of Gaussian processes with inequality constraints. *Procedia Environmental Sciences*, 27:38–44.
- Meyer, C. D. (2000). *Matrix Analysis and Applied Linear Algebra*. Society for Industrial and Applied Mathematics, Philadelphia, PA, USA.
- Murphy, K. P. (2012). *Machine Learning: A Probabilistic Perspective (Adaptive Computation And Machine Learning Series)*. The MIT Press.
- Neal, R. M. (1996). *Bayesian Learning for Neural Networks*. Springer-Verlag New York, Inc., Secaucus, NJ, USA.
- Nickisch, H. and Rasmussen, C. (2008). Approximations for binary Gaussian process classification. *Journal of Machine Learning Research*, 9:2035–2078.
- Paciorek, C. J. and Schervish, M. J. (2004). Nonstationary covariance functions for Gaussian process regression. In *In Proc. of the Conf. on Neural Information Processing Systems*. MIT Press.
- Pakman, A. and Paninski, L. (2014). Exact Hamiltonian Monte Carlo for truncated multivariate Gaussians. *Journal of Computational and Graphical Statistics*, 23(2):518–542.
- Rasmussen, C. E. and Williams, C. K. I. (2005). *Gaussian Processes for Machine Learning (Adaptive Computation and Machine Learning)*. The MIT Press.
- Riihimäki, J. and Vehtari, A. (2010). Gaussian processes with monotonicity information. In *Journal of Machine Learning Research: Workshop and Conference Proceedings*, volume 9, pages 645–652.
- Snelson, E. and Ghahramani, Z. (2006). Sparse Gaussian processes using pseudo-inputs. In Weiss, Y., Schölkopf, P. B., and Platt, J. C., editors, *Advances in Neural Information Processing Systems 18*, pages 1257–1264. MIT Press.

- Stein, M. (1999). *Interpolation of Spatial Data: Some Theory for Kriging*. Springer Series in Statistics. Springer New York.
- Taylor, J. and Benjamini, Y. (2017). RestrictedMVN: multivariate normal restricted by affine constraints. <https://cran.r-project.org/web/packages/restrictedMVN/index.html>. [Online; 02-Feb-2017].
- Vats, D., Flegal, J. M., and Jones, G. L. (2017). Multivariate output analysis for Markov Chain Monte Carlo. *ArXiv e-prints*.
- Yadrenko, M. I. (1983). *Spectral theory of random fields*. Translation series in mathematics and engineering. Optimization Software, New York, NY. Transl. from the Russian.
- Ying, Z. (1993). Maximum likelihood estimation of parameters under a spatial sampling scheme. *Ann. Statist.*, 21(3):1567–1590.
- Zhang, H. (2004). Inconsistent estimation and asymptotically equal interpolations in model-based geostatistics. *Journal of the American Statistical Association*, 99(465):250–261.