

# Active Emulation of Computer Codes with Gaussian Processes – Application to Remote Sensing

Daniel Heestermans Svendsen<sup>a,\*</sup>, Luca Martino<sup>b</sup>, Gustau Camps-Valls<sup>a</sup>

<sup>a</sup>*Image Processing Lab (IPL), Universitat de València, C/ Cat. José Beltrán, 2. 46980 Paterna, Spain.*

<sup>b</sup>*Dep. Signal Processing, Universidad Rey Juan Carlos (URJC), Camino del Molino 5, 28943 Fuenlabrada, Spain*

## Abstract

Many fields of science and engineering rely on running simulations with complex and computationally expensive models to understand the involved processes in the system of interest. Nevertheless, the high cost involved hamper reliable and exhaustive simulations. Very often such codes incorporate heuristics that ironically make them less tractable and transparent. This paper introduces an active learning methodology for adaptively constructing surrogate models, i.e. *emulators*, of such costly computer codes in a multi-output setting. The proposed technique is sequential and adaptive, and is based on the optimization of a suitable acquisition function. It aims to achieve accurate approximations, model tractability, as well as compact and expressive simulated datasets. In order to achieve this, the proposed Active Multi-Output Gaussian Process Emulator (AMOGAPE) combines the predictive capacity of Gaussian Processes (GPs) with the design of an acquisition function that favors sampling in low density and fluctuating regions of the approximation functions. Comparing different acquisition functions, we illustrate the promising performance of the method for the construction of emulators with toy examples, as well as for a widely used remote sensing transfer code.

*Keywords: Active learning, Gaussian process, emulation, design of experiments, computer code, remote sensing, radiative transfer model*

## 1 Introduction

In many areas of science and engineering, systems are analyzed by running computer code simulations which act as convenient approximations of reality. They allow us to simulate many different systems of interest and characterize the involved processes, such as turbulence or energy transfer, and their interactions and relevance. Depending on the body of literature, they are known as physics-based or mechanistic models, or simply *simulators* [30, 39]. Two important limitations are associated with simulators. The first, and perhaps the most important problem of these computer codes, is their often high computational cost, which hampers reliable and exhaustive simulations. This limits the representativity of the simulations, which in turn makes numerical or statistical inversion a hard problem. Secondly, since computer codes rely on decades of intensive development and parametrizations, they often include heuristics that improve accuracy but ironically make them less mathematically tractable and transparent.

\*Corresponding author, daniel.svendsen@uv.es

**Emulation for forward models.** In the last decade, a field collectively known as *surrogate modeling* or *emulation* has emerged as an efficient alternative: emulators try to mimic costly computer codes with machine learning models. The field of emulation has received attention from subfields of statistical signal processing and machine learning [10, 11, 15, 24, 25]. In order to construct an emulator, we need a simulated dataset which is made by evaluating the computer code in different input points. The problem of choosing these points, for which this paper presents an active learning algorithm [37, 42], is treated in different parts of the statistics and machine learning literature. A non-exhaustive overview is given below.

**Related work.** The problem at hand is closely related to that of Design of Experiments (DOE), where one seeks a set of input values which best allows one to determine the relationship between inputs and outputs. Between the algorithms that will be reviewed in this section, there are key some differences between types of algorithms that it would be beneficial to clarify first:

- *Sequential vs. non-sequential* refers to whether the algorithm needs to know a priori how many input points to choose. Non-sequential or one-shot algorithms need this information, while sequential algorithms can simply run until some time limit or accuracy criteria is met, which is a favourable property. Between the two approaches lie the batch-sequential algorithms.
- *Continuous vs. discrete sampling* refers to whether an algorithm aims to either choose input points in a continuous space or choose among a finite set of points. A large part of the literature deals with the latter problem, and relies on greedy and MCMC algorithms to choose points according to some criterion. Many of the methods proposed in the literature can be easily adapted from continuous to discrete sampling and vice versa.

Among the most popular criteria are maximum entropy [32], maximizing distance to nearest neighbour [14], and minimizing integrated root mean squared error [29]. These criteria have been implemented for construction of GP emulators in both a sequential and batch-sequential way [18]. An interesting approach in this field is the Bayesian Experimental Design (BED) which assumes a probabilistic model of the observed data and defines a so-called utility function based on the posterior of the model parameters. The approach then aims to maximize the mean of the utility. Recent relevant work can be found in [9, 28]. A great deal of DOE methods, even the sequential ones, do not assume the ability to *query* a system, due to the way experiments are carried out.

The field of Active Learning (AL), on the other hand, builds on the premise that we can query a system and thus learn something about it in each iteration [23, 35]. Building an emulator sequentially is a problem that fits directly into this category. The algorithms in the AL literature concerned with GP regression often employ criteria based on predictive variance [31] and entropy [36]. Other algorithms are based on triangulation of the input space [1] or gridding [6], followed by a ranking of each triangle or cell. Greedily searching for candidate points which have maximum distance to their nearest neighbours [41] has also

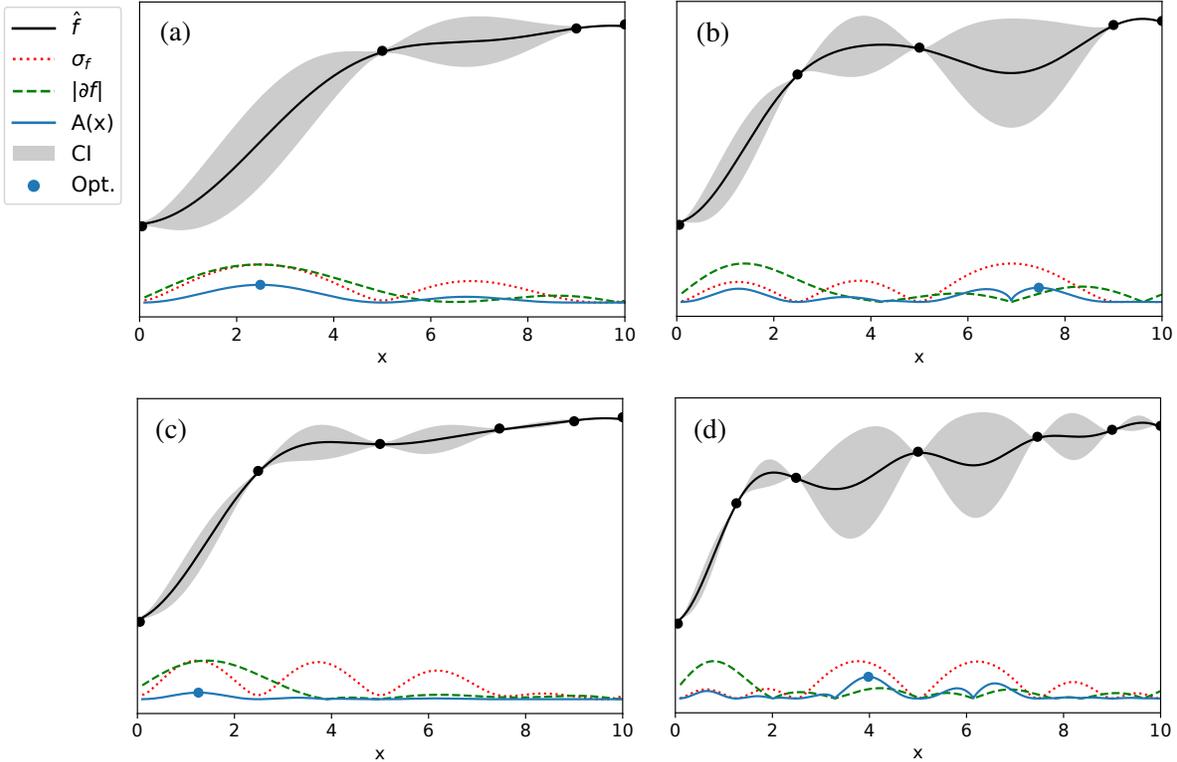


Figure 1: The presented method optimizes the selection of most informative points to approximate an arbitrary multidimensional function iteratively. The example shows the first four iterations in a 1D case. Starting from 4 points, a GP interpolator is built from which some valuable information is derived (the predictive variance -green- and the gradient -red-) and then combined in an acquisition function (blue) that proposes the next point to sample (blue dot). The acquisition function admits many general forms and trades off geometry and diversity terms to account for attractiveness in the sample space.

proven effective. Furthermore, when the input set is comprised of finite discrete values, interesting criteria like mutual information have been employed with success [17].

**Our Contribution: Active emulation as a step forward.** In this paper, we introduce a methodology for developing efficient machine learning emulators of costly physical models based on active emulation. An active learning framework is developed that *sequentially* chooses informative input points, learning about about the underlying function as the algorithm progresses. This active emulation methodology is based on the notion of an *acquisition function* which can be optimized through gradient-based techniques, mirroring approaches in Bayesian Optimization [33]. The goal is to construct an accurate emulator with as few runs of the computer code as possible.

Given a set of initial datapoints, the emulator is built through the online addition of new

nodes<sup>1</sup>, maximizing the acquisition function at each iteration. The acquisition function is constructed to incorporate (a) geometric information of the costly, analytically intractable function  $f$ , and (b) information about the distribution of the current nodes. By using Gaussian processes we can derive both terms analytically, and for multiple outputs at once. The reasoning is that areas of high variability in  $f(\mathbf{x})$  requires the addition of more information, as has also been noted in [6]. In [20] the predictive variance of the gradient norm of a GP is used as a sampling criteria, which is a less straightforward approach than just using the gradient directly as done here. Similarly, regions with a small concentration of nodes requires the introduction of new nodes in order to fill the space (simple exploration, space filling without taking into account the geometrical features of  $f(\mathbf{x})$ ). We show how to define such an acquisition function in a multi-output setting. Figure 1 shows an illustrative example of the building blocks of the active emulation methodology presented here.

The developed methodology of constructing emulators is sequential and searches a continuous input-space, leading to emulators that are *accurate*, so they can be taken as a faithful representation of the physical models and codes, *compact*, and *parsimonious*, as a minimal number of informative points is selected, and *general-purpose* since it is based on properties of Gaussian processes like uncertainty and gradients that can be obtained for any differentiable covariance function. This paper builds on of the preliminary work in [34], extending it in several directions. A general framework is provided before describing some specific implementations, extending the study proposing the use of a range of different acquisition functions. A theoretical demonstration of the utility of a gradient term in active sampling and emulation is also given (see A, for instance). Finally, more thorough experimental results are provided, with more examples and challenging model comparisons, and a more advanced use case.

**Structure of the paper.** The remainder of the paper is organized as follows. We first define active emulation and establish the notation in Section 2. Then, the GP-based active emulation framework is presented in Section 3. The framework defines a general-purpose acquisition function built on optimal search of diversity and uncertainty criteria. Experimental results in synthetic and challenging real problems illustrate the capabilities in Section 4. We will pay special attention to the field of remote sensing, where computer codes, called radiative transfer models (RTMs), are widely used and pose challenges to the design of accurate and compact emulators. Having access to an exhaustive ground truth allows us to analyze performance in terms of convergence and accuracy. We conclude in Section 5 with some remarks and an outline of future work.

## 2 Active Emulation

In this section we describe the generic active emulation (AE) method for a complex system denoted as  $f(\mathbf{x})$ , e.g., an expensive RTM model. We first fix the notation, then present the

---

<sup>1</sup> In the following, the words *node* and *datapoint* will be used interchangeably.

processing scheme. Consider a  $D$ -dimensional bounded input space  $\mathcal{X}$ , i.e.,  $\mathbf{x} \in \mathcal{X} \subset \mathbb{R}^D$ . Furthermore, let  $\mathbf{f}(\mathbf{x}): \mathcal{X} \rightarrow \mathbb{R}^P$  denote a complex system with  $P$  outputs. Finally,  $t \in \mathbb{N}$  denotes the index of the AE algorithm, and  $m_t$  the number of datapoints  $\{\mathbf{x}_k, \mathbf{y}_k\}_{k=1}^{m_t}$  used by the algorithm at iteration  $t$ , where

$$\mathbf{y}_k = \mathbf{f}(\mathbf{x}_k), \quad (1)$$

where  $\mathbf{y}_k = [y_{1,k}, \dots, y_{P,k}]^\top$  and  $k = 1, \dots, m_t$ . Thus, given an input matrix of nodes,  $\mathbf{X}_t = [\mathbf{x}_1, \dots, \mathbf{x}_{m_t}]$  of dimension  $D \times m_t$ , we have a  $P \times m_t$  matrix of outputs,  $\mathbf{Y}_t = [\mathbf{y}_1, \dots, \mathbf{y}_{m_t}]$ . At each iteration  $t$ , given the datapoints  $\{\mathbf{x}_k, \mathbf{y}_k\}_{k=1}^{m_t}$ , the AE method constructs an interpolating function  $\hat{\mathbf{f}}_t(\mathbf{x})$ . Then, an acquisition function  $A_t(\mathbf{x}): \mathbb{R}^D \rightarrow \mathbb{R}$  is built in order to suggest which regions of the space require additional nodes. That is, an optimization step is performed for obtaining the next input  $\mathbf{x}_{m_{t+1}}$ :

$$\mathbf{x}_{m_{t+1}} = \arg \max_{\mathbf{x} \in \mathcal{X}} A_t(\mathbf{x}). \quad (2)$$

The dataset is updated accordingly,  $\mathbf{X}_{t+1} = [\mathbf{X}_t, \mathbf{x}_{m_{t+1}}]$ ,  $\mathbf{Y}_{t+1} = [\mathbf{Y}_t, \mathbf{y}_{m_{t+1}} = \mathbf{f}(\mathbf{x}_{m_{t+1}})]$  adding a new node, and we set  $m_{t+1} = m_t + 1$  and  $t \leftarrow t + 1$ . The procedure is repeated until a stopping condition is met. One possibility is to stop the algorithm when a pre-established maximum number of points  $M$  (determined by the available computational resources) has been included. Theoretically, the user could stop the algorithm when a least a precision error  $\epsilon > 0$  is achieved,  $\|\mathbf{f}(\mathbf{x}) - \hat{\mathbf{f}}_t(\mathbf{x})\| \leq \epsilon$ . However, since  $\mathbf{f}(\mathbf{x})$  is costly and analytically intractable in general<sup>2</sup>, one cannot evaluate and/or approximate the associated error  $\|\mathbf{f}(\mathbf{x}) - \hat{\mathbf{f}}_t(\mathbf{x})\|$ . A practical alternative is to stop the AE method when  $\|\hat{\mathbf{f}}_t(\mathbf{x}) - \hat{\mathbf{f}}_{t-1}(\mathbf{x})\| \leq \epsilon'$  for some  $\epsilon' > 0$ . Figure 1 shows a graphical representation of a generic AE procedure. Table 1 summarizes the main notation of the work. Table 2 shows in details the steps of a generic AE algorithm.

Note that the goal is either to sequentially construct an emulator able to obtain a pre-established error in approximation with the smallest number of nodes possible or, more commonly, the best possible emulator built with a pre-established maximum number of nodes (given some starting points). The constructed emulator will be used for further applications for the interested users, researchers and practitioners. We do not consider time or computational restrictions in the construction stage. Furthermore, our approach is particularly useful when the underlying function is very costly, i.e. when the cost of evaluating this function is significantly greater than the application of one iteration of the proposed algorithm.

## 2.1 Acquisition function

We consider acquisition functions  $A_t(\mathbf{x}): \mathcal{X} \rightarrow \mathbb{R}$  obtained by the multiplication of a *geometry* term  $G_t(\mathbf{x})$  and a *diversity* factor  $D_t(\mathbf{x})$ , i.e. functions of the form:

$$A_t(\mathbf{x}) = [G_t(\mathbf{x})]^{\beta_t} D_t(\mathbf{x}), \quad (3)$$

<sup>2</sup>The system  $\mathbf{y} = \mathbf{f}(\mathbf{x})$  is a black-box mapping, linking the inputs  $\mathbf{x}$  with the outputs  $\mathbf{y}$ . At each new input  $\mathbf{x}'$ , the system returns  $\mathbf{y}' = \mathbf{f}(\mathbf{x}')$ , but does so by way of a computer code which is too complex and slow to lend itself to exhaustive analysis across the input space.

Table 1: Main notation of the work.

$t \in \mathbb{N}$	Iteration index of the active emulator.
$m_t$	Number of data points at the $t$ -th iteration.
$\{\mathbf{x}_k, \mathbf{y}_k\}_{k=1}^{m_t}$	Data points at the $t$ -th iteration.
$\mathbf{x} = [x_1, \dots, x_D]^\top \in \mathcal{X} \subset \mathbb{R}^D$	Input variable.
$\mathbf{y} = [y_1, \dots, y_P]^\top$	Outputs.
$\mathbf{X}_t = [\mathbf{x}_1, \dots, \mathbf{x}_{m_t}]$	$D \times m_t$ input matrix.
$\mathbf{Y}_t = [\mathbf{y}_1, \dots, \mathbf{y}_{m_t}]$	$P \times m_t$ output matrix.
$\mathbf{y} = f(\mathbf{x}) : \mathcal{X} \rightarrow \mathbb{R}^P$	Unknown function/forward model linking $\mathbf{x}$ with $\mathbf{y}$ .
$\hat{\mathbf{y}} = \hat{\mathbf{f}}_t(\mathbf{x}) : \mathcal{X} \rightarrow \mathbb{R}^P$	Interpolator the $t$ -th iteration using $\{\mathbf{x}_k, \mathbf{y}_k\}_{k=1}^{m_t}$ .
$A_t(\mathbf{x}) : \mathcal{X} \rightarrow \mathbb{R}$	Acquisition function at the $t$ -th iteration.
$k(\mathbf{x}, \mathbf{z}) : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$	kernel function.
$\mathbf{K}$	$m_t \times m_t$ kernel matrix.
$\mathbf{k}_x = [k(\mathbf{x}, \mathbf{x}_1), \dots, k(\mathbf{x}, \mathbf{x}_{m_t})]^\top$	$m_t \times 1$ vector.

where  $\beta_t \in [0, 1]$  is a positive non-decreasing function of  $t$ , with  $\lim_{t \rightarrow \infty} \beta_t = 1$ . The function  $G_t(\mathbf{x})$  encodes the geometrical information in  $\mathbf{f}(\mathbf{x})$ , while function  $D_t(\mathbf{x})$  depends on the distribution of the points in the current vector  $\mathbf{X}_t$ . More specifically,  $D_t(\mathbf{x})$  takes greater values around empty areas in  $\mathcal{X}$ , whereas  $D_t(\mathbf{x})$  will be approximately zero close to the nodes and exactly zero at the nodes<sup>3</sup>, i.e.,  $D_t(\mathbf{x}_i) = 0$ , for  $i = 1, \dots, m_t$  and  $\forall t \in \mathbb{N}$ . As a consequence, we have

$$A_t(\mathbf{x}_i) = 0 \quad \forall i, t. \quad (4)$$

Generally, since  $\mathbf{f}(\mathbf{x})$  is analytically intractable, the function  $G_t(\mathbf{x})$  can only be derived from information acquired in advance or by considering the approximation  $\hat{\mathbf{f}}_t(\mathbf{x})$ . The *tempering parameter*,  $\beta_t$ , helps to down-weight the likely less informative estimates of the gradient in the very first iterations. For instance, if  $\beta_t = 0$ , we ignore  $G_t(\mathbf{x})$  and  $A_t(\mathbf{x}) = D_t(\mathbf{x})$ , i.e., only the exploration term is considered. Whereas, if  $\beta_t = 1$ , we have  $A_t(\mathbf{x}) = G_t(\mathbf{x})D_t(\mathbf{x})$ .

## 2.2 Specific implementation

The AE algorithm introduced is completely defined by the choice of the interpolator providing the approximation  $\hat{\mathbf{f}}_t(\mathbf{x})$ , and the functions  $G_t(\mathbf{x})$ ,  $D_t(\mathbf{x})$ , and  $\beta_t$ . Moreover, the initial set of nodes  $\{\mathbf{x}_k, \mathbf{y}_k\}_{k=1}^{m_0}$  and the stopping condition could be considered as additional elements. It is important to note that, in order to choose the interpolating function, we have to take into account the ease of application in high dimensional spaces and the possibility of computing the gradient and other differential geometric measures of  $\hat{\mathbf{f}}_t$  analytically. Different designs of these four elements give rise to different AE techniques. In Section 3, we provide some specific examples of the choice of  $\{\hat{\mathbf{f}}_t, G_t, D_t, \beta_t\}$ .

<sup>3</sup>Note that this is the case only for an interpolator (no output-noise assumed) while for a regressor the value of  $D_t(\mathbf{x})$  will just be very small around already placed nodes.

Table 2: Generic Active Emulator.

<p>1. Set <math>t = 0</math>, select initial points <math>\mathbf{X}_0 = [\mathbf{x}_1, \dots, \mathbf{x}_{m_0}]</math>, and <math>\mathbf{Y}_0 = [\mathbf{y}_1, \dots, \mathbf{y}_{m_0}]</math>, and maximum number of nodes <math>M</math>.</p> <p>2. While <math>m_t &lt; M</math>:</p> <p>(a) Given <math>\mathbf{X}_t = [\mathbf{x}_1, \dots, \mathbf{x}_{m_t}]</math> and <math>\mathbf{Y}_t = [\mathbf{y}_1, \dots, \mathbf{y}_{m_t}]</math>, build function <math>\widehat{\mathbf{f}}_t(\mathbf{x})</math>.</p> <p>(b) Build the acquisition function <math>A_t(\mathbf{x})</math> from <math>\widehat{\mathbf{f}}_t</math>, and obtain the new input</p> $\mathbf{x}_{m_{t+1}} = \arg \max_{\mathbf{x} \in \mathcal{X}} A_t(\mathbf{x}). \quad (5)$ <p>(c) Obtain outputs <math>\mathbf{y}_{m_{t+1}} = \mathbf{f}(\mathbf{x}_{m_{t+1}})</math>.</p> <p>(d) Update <math>\mathbf{X}_{t+1} = [\mathbf{X}_t, \mathbf{x}_{m_{t+1}}]</math>, <math>\mathbf{Y}_{t+1} = [\mathbf{Y}_t, \mathbf{y}_{m_{t+1}}]</math>.</p> <p>(e) Set <math>m_{t+1} = m_t + 1</math> and <math>t \leftarrow t + 1</math>.</p> <p>3. Build the interpolating function <math>\widehat{\mathbf{f}}_t(\mathbf{x})</math>.</p> <p>4. Return final set of optimal nodes <math>\{\mathbf{x}_k, \mathbf{y}_k\}_{k=1}^{m_t}</math> as a Look-up Table (LUT), as well as the gradient and the predictive variance of the predictive model <math>\widehat{\mathbf{f}}_t(\mathbf{x})</math>.</p>
---

### 2.3 Parsimonious sequential approach

It is also important to remark that the active emulation procedure presented in this work is intrinsically a sequential technique. This means that the nodes in  $\mathbf{X}_{t-1}$  are always contained in  $\mathbf{X}_t$ , i.e. the locations of previous nodes are not changed. This solution minimizes the number of evaluations of the complex system  $\mathbf{f}$ . In this sense, the active emulation procedure is a parsimonious sequential technique that applies, at each iteration, all previously obtained information about the underlying function  $\mathbf{f}$ . Namely, all the previous evaluations of  $\mathbf{f}$  are used, and only one additional evaluation of  $\mathbf{f}$  is required at each iteration.

### 2.4 Products of the algorithm

The active emulation procedure proposed in this work is a methodology that delivers: (a) an accurate GP emulator (considering a specific choice of the interpolator) while evaluating the computer code as little as possible, (b) a final set of nodes  $\{\mathbf{x}_k, \mathbf{y}_k\}_{k=1}^{m_t}$  as a Lookup Table (LUT; other interpolation procedures can be applied using the obtained set of points), and (c) useful statistical information about the model  $\mathbf{f}$ , such as predictive variance and gradients of the learned function, which can be further used for model inversion and error propagation analyses.

### 3 Active Multi-Output Gaussian Process Emulator (AMO-GAPE)

An active emulator is completely defined by the choice of the predictive, model  $\widehat{\mathbf{f}}(\mathbf{x})$  and the acquisition function  $A_t(\mathbf{x})$ . In this work, we consider a GP interpolator, as well as the regression formulation [26], which has been successfully used in remote sensing applications recently [8].

#### 3.1 The Gaussian Process Interpolator

For the sake of simplicity, let us first start considering the GP solution for the scalar output case, i.e.,  $P = 1$ . Hence, in this case the vectorial function  $\mathbf{y} = \mathbf{f}(\mathbf{x})$  is a simple function  $y = f(\mathbf{x})$ , and the matrix  $\mathbf{Y}_t = [y_{1,1}, \dots, y_{1,m_t}]$ , becomes a  $1 \times m_t$  vector. Given a generic test input  $\mathbf{x}$ , GPs provide a Gaussian predictive density  $p(y|\mathbf{x}) = \mathcal{N}(y|\mu(\mathbf{x}), \sigma^2(\mathbf{x}))$  with predictive mean  $\mu(\mathbf{x})$  and variance  $\sigma^2(\mathbf{x})$ . The predictive mean gives us the interpolating function and is given by

$$\widehat{f}_t(\mathbf{x}) = \mu_t(\mathbf{x}) = \mathbf{k}_x^\top \mathbf{K}^{-1} \mathbf{Y}_t^\top, \quad (6)$$

where we defined a kernel function  $k(\mathbf{x}, \mathbf{z}): \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ , the corresponding kernel matrix  $[\mathbf{K}]_{ij} := k(\mathbf{x}_i, \mathbf{x}_j)$  of dimension  $m_t \times m_t$  containing all kernel entries, and the kernel vector  $\mathbf{k}_x = [k(\mathbf{x}, \mathbf{x}_1), \dots, k(\mathbf{x}, \mathbf{x}_{m_t})]^\top$  of dimension  $m_t \times 1$ . The interpolating function can be simply expressed as a linear combination of  $\widehat{f}_t(\mathbf{x}) = \mathbf{k}_x^\top \boldsymbol{\alpha} = \sum_{i=1}^{m_t} \alpha_i k(\mathbf{x}, \mathbf{x}_i)$ , where the weights  $\boldsymbol{\alpha} = [\alpha_1, \dots, \alpha_{m_t}]^\top$  are  $\boldsymbol{\alpha} = \mathbf{K}^{-1} \mathbf{Y}_t^\top$ . The GP formulation also provides an expression for the predictive variance

$$\sigma_t^2(\mathbf{x}) = k(\mathbf{x}, \mathbf{x}) - \mathbf{k}_x^\top \mathbf{K}^{-1} \mathbf{k}_x. \quad (7)$$

An example is the exponentiated quadratic kernel function,

$$k(\mathbf{x}, \mathbf{z}) = \exp\left(-\frac{\|\mathbf{x} - \mathbf{z}\|^2}{2\delta^2}\right), \quad (8)$$

where  $\|\cdot\|$  is the  $\ell_2$ -norm, and  $\delta > 0$  is a positive scalar hyper-parameter. Note that the norm of the gradient of the interpolating function  $\widehat{f}_t$  w.r.t. the input data  $\mathbf{x}$  can be easily computed,

$$\text{Gr}_t(\mathbf{x}) = \left\| \nabla_x \widehat{f}_t(\mathbf{x}) \right\| = \left\| \sum_{i=1}^{m_t} \alpha_i \nabla_x k(\mathbf{x}, \mathbf{x}_i) \right\|. \quad (9)$$

The gradient vector of  $k(\mathbf{x}, \mathbf{x}_i)$  with  $\mathbf{x} = [x_1, \dots, x_D]^\top$  and  $\mathbf{x}_i = [x_{1,i}, \dots, x_{D,i}]^\top$ , is

$$\nabla_x k(\mathbf{x}, \mathbf{x}_i) = -\frac{k(\mathbf{x}, \mathbf{x}_i)}{\delta^2} [(x_1 - x_{1,i}), \dots, (x_D - x_{D,i})]^\top, \quad (10)$$

which can be easily computed analytically, and by automatic differentiation software. At this point an intuitive choice of acquisition function of Eq. (3) presents itself. The predictive

variance which describes the uncertainty of the GP prediction, and which largely depends on distance to nearby training points, is a natural choice for the diversity term  $D_t(\mathbf{x}) = \sigma_t^2(\mathbf{x})$ . Furthermore, since the emulator is differentiable, we can use the gradient as a measure of function variation and choose the geometry term as  $G_t(\mathbf{x}) = \text{Gr}_t(\mathbf{x})$ .

### 3.2 Multi-output GP interpolator

Several multi-output GP schemes have been proposed with the aim of exploiting the correlation among the output variables [2–4, 12, 19, 40]. These models are especially well suited for multitask problems where little data is available or for gap filling, which is not the scenario of this work [4]. We do not face such problems in our particular remote sensing application since the RTMs provide all vector components when executed in forward mode. We adopt a simpler yet highly effective approach, simply treating each output independently. For simplicity, we consider an isotopic case where to each input  $\mathbf{x}_k$  we have  $P$  different outputs,  $[y_{1,k}, \dots, y_{P,k}]^\top$ ; see the description of isotopic and heterotopic models in [4]. We also define the  $p$ -th row of the matrix  $\mathbf{Y}_t$  as  $\tilde{\mathbf{y}}_{p,t} = [y_{p,1}, \dots, y_{p,m_t}]$ , with  $p = 1, \dots, P$ , so that  $\mathbf{Y}_t$  is matrix of dimension  $P \times m_t$ . Here, for the sake of simplicity, we apply one GP interpolator for each output independently, i.e.,

$$\hat{\mathbf{f}}_t(\mathbf{x}) = \begin{cases} \hat{f}_{1,t}(\mathbf{x}) = \mathbf{k}_{x,1}^\top \mathbf{K}_1^{-1} \tilde{\mathbf{y}}_{1,t}^\top \\ \vdots \\ \hat{f}_{P,t}(\mathbf{x}) = \mathbf{k}_{x,P}^\top \mathbf{K}_P^{-1} \tilde{\mathbf{y}}_{P,t}^\top \end{cases}, \quad (11)$$

where the vectors  $\mathbf{k}_{x,p}$  have all dimension  $m_t \times 1$  and the matrices  $\mathbf{K}_p$  have dimension  $m_t \times m_t$ . The subindex  $p$  in the kernel vector  $\mathbf{k}_{x,p}$  and the kernel matrix  $\mathbf{K}_p$  denotes the dependence to a different hyper-parameter  $\delta_p$  (we learn one for each output). More generally, we can consider a different kernel for each output which allows for much model flexibility. Hence, for each output, we have a different variance

$$\sigma_{p,t}^2(\mathbf{x}) = k_p(\mathbf{x}, \mathbf{x}) - \mathbf{k}_{x,p}^\top \mathbf{K}_p^{-1} \mathbf{k}_{x,p}. \quad (12)$$

Similarly, we have one gradient norm for each interpolating function  $\text{Gr}_{p,t}(\mathbf{x})$ . It is important to note here that any multi-output GP framework will fit in the AMOGAPE method as long as it provides a differentiable predictive variance and gradient, which for most commonly used kernels is the case.

### 3.3 The acquisition function

Note that  $\sigma_p^2(\mathbf{x}_i) = 0$  for all  $i = 1, \dots, m_t$  and all  $p$ , and each  $\sigma_p^2(\mathbf{x})$  depends on the distance among the support points  $\mathbf{x}_t$ , the chosen kernel function  $k$ , and the value of the corresponding hyper-parameter  $\delta_p$ . For this reason, it is reasonable to consider as diversity term the following function that combines them all:

$$D_t(\mathbf{x}) := \sigma_{1,t}^2(\mathbf{x}) \odot \sigma_{2,t}^2(\mathbf{x}) \odot \sigma_{3,t}^2(\mathbf{x}) \dots \odot \sigma_{P,t}^2(\mathbf{x}), \quad (13)$$

where  $\odot$  represents a generic mathematical operation such as sum (+) or multiplication ( $\times$ ). We wish to use the geometric information term to sample where the norm of the gradient is high and thus define similarly

$$G_t(\mathbf{x}) := \text{Gr}_{1,t}(\mathbf{x}) \odot \text{Gr}_{2,t}(\mathbf{x}) \odot \text{Gr}_{3,t}(\mathbf{x}) \dots \odot \text{Gr}_{P,t}(\mathbf{x}). \tag{14}$$

The intuition behind this choice is that wavy regions of  $\mathbf{f}$  (estimated by  $\hat{\mathbf{f}}_t$ ) require more support points than flat regions. In A we demonstrate the importance of the gradient term using the simple example of a piecewise-constant interpolator. As previously mentioned, we define the acquisition function as

$$A_t(\mathbf{x}) = [G_t(\mathbf{x})]^{\beta_t} D_t(\mathbf{x}). \tag{15}$$

Table 3 shows several combinations that generate different acquisition functions according to the choice of the operator  $\odot$ .

**Optimization approaches.** The maximization of Eq. (15) can be performed by using different optimization algorithms, e.g., gradient ascent or simulated annealing. As can be seen in the simple 1-D example of Fig. 1, the acquisition function has many local optima. Thus, while it is useful to have access to the gradient of Eq. 15, we find that it is important to incorporate stochasticity in the optimization. This can be done, for example by performing a number of random searches and then performing gradient ascent, initialized at the best candidate point.

**Tempering of the geometric information.** The parameter  $\beta_t \in [0, 1]$  indicates how the acquisition function should “trust” the provided geometric information and must be a non-decreasing function of  $t$ . Indeed, recall that the geometric information is given by analyzing the interpolating function  $\hat{\mathbf{f}}$ , instead of the complex system  $\mathbf{f}$ , since it is analytically intractable. Clearly,  $\beta_t$  must be an increasing function with respect to  $t$ , since at each iteration the interpolating function  $\hat{\mathbf{f}}$  is improved and becomes step-by-step more reliable. One possible choice is  $\beta_t = 1 - \exp(-\gamma t)$ , where  $\gamma \geq 0$  is a positive scalar established by the user or, alternatively,  $\beta_t = 1 - \frac{1}{t}$ , for instance.

### 3.4 From interpolation to regression

So far we have described the emulation as an interpolation problem since RTMs are deterministic models: running the code multiple times will always return identical answers. Hence, we have assumed an observation equation of type  $y = f(\mathbf{x})$ <sup>4</sup>. However, in some cases, it is preferable to consider an observation equation of type  $y = f(\mathbf{x}) + \epsilon$  where  $\epsilon \sim \mathcal{N}(0, v^2)$  represents a Gaussian noise perturbation with zero mean and variance  $v^2$ . There are three main

---

<sup>4</sup>In this section, we assumed only one output in the equation, just for the sake of simplicity. Clearly, the same considerations are valid for the multi-output case.

Table 3: Acquisition functions for a multi-output emulator and their shorthand notation used in Section 4.

$A_t(\mathbf{x})$	Shorthand
$\sum_{p=1}^P \sigma_{p,t}^2(\mathbf{x})$	$\Sigma D$
$\prod_{p=1}^P \sigma_{p,t}^2(\mathbf{x})$	$\Pi D$
$\sum_{p=1}^P \sigma_{p,t}^2(\mathbf{x}) \sum_{p=1}^P \text{Gr}_{p,t}(\mathbf{x})$	$\Sigma D \times \Sigma G$
$\sum_{p=1}^P \sigma_{p,t}^2(\mathbf{x}) \prod_{p=1}^P \text{Gr}_{p,t}(\mathbf{x})$	$\Sigma D \times \Pi G$
$\prod_{p=1}^P \sigma_{p,t}^2(\mathbf{x}) \sum_{p=1}^P \text{Gr}_{p,t}(\mathbf{x})$	$\Pi D \times \Sigma G$
$\prod_{p=1}^P \sigma_{p,t}^2(\mathbf{x}) \prod_{p=1}^P \text{Gr}_{p,t}(\mathbf{x})$	$\Pi D \times \Pi G$

reasons, both theoretical and practical, for considering noisy outputs: (a) the system to emulate actually contains stochastic elements (i.e., it is not a completely deterministic system), (b) to increase the prediction power of the emulator function  $\widehat{f}_t(\mathbf{x})$  providing more flexibility to the GP model, and (c) in order to avoid numerical problems, increasing also the stability of the computation. This last point is due to the fact that the noise variance  $v^2$  plays the role of a regularization term which is added to the diagonal of the kernel matrix (also called a *nugget* in kriging literature). Indeed, when noisy outputs are assumed and by denoting the  $m_t \times m_t$  identity matrix as  $\mathbf{I}$ , then the GP regression equations become

$$\widehat{f}_t(\mathbf{x}) = \mathbf{k}_x^\top (\mathbf{K} + v^2 \mathbf{I})^{-1} \mathbf{Y}_t^\top, \quad (16)$$

$$\sigma_t^2(\mathbf{x}) = v^2 + k(\mathbf{x}, \mathbf{x}) - \mathbf{k}_x^\top (\mathbf{K} + v^2 \mathbf{I})^{-1} \mathbf{k}_x. \quad (17)$$

Note that, if we set again  $D_t(x) = \sigma_t^2(\mathbf{x})$  (with  $\sigma_t^2(\mathbf{x})$  defined above), then  $A_t(\mathbf{x})$  does not fulfil Eq. (4). However,  $A_t(x)$  still takes greater values far from nodes  $\mathbf{x}_i$ , and smaller values close to points  $\mathbf{x}_i$ . If the application strictly requires that the condition in Eq. (4) must be satisfied, then we can simply define  $D_t(x) = k(\mathbf{x}, \mathbf{x}) - \mathbf{k}_x^\top \mathbf{K}^{-1} \mathbf{k}_x$ , i.e.,  $\sigma^2(\mathbf{x})$  without the noise term. With this definition, we have again  $A_t(\mathbf{x}_i) = 0$ . This means that the noise term is only used in the GP equations and not for the construction of the acquisition function. Finally, note that, if a regressor is applied instead of an interpolator, then two hyperparameters must be tuned,  $\delta$  and  $v$ , instead of just only  $\delta$ , assuming the kernel in Eq. (8). The user might also wish to decide a value of  $v^2$  in advance instead of learning it, using it as a regularization term in order to guarantee the numerical stability of the method. Hyperparameter tuning can be performed with standard Cross Validation (CV) procedures, or maximizing the marginal likelihood function by gradient ascent or other optimization techniques [16, 21]. In the interpolation case or when  $v^2$  is decided in advance by the user, another interesting approach is to find the maximum value of bandwidth  $\delta$  which still allows the numeric inversion of the matrix  $\mathbf{K}$  (imposing a upper bound for its condition number).

## 4 Experimental results

This section presents experimental results of our AE framework in synthetic and real (Earth-observation) systems. The AMOGAPE<sup>5</sup> method is compared to standard algorithms in the literature, namely random exploration/sampling and most notably Sobol’s sampling [5] and the Latin Hypercube Sampling (LHS) method [22]. Algorithms are compared in terms of accuracy and convergence rates in problems of different input and output dimensionality. The real experiments involve a widely used code that models the relation between vegetation parameters and the corresponding reflectance signal.

### 4.1 Toy Experiment 1: Example of unidimensional multi-output emulation

We consider a multi-output toy example with scalar inputs  $x \in \mathbb{R}$  where we can easily compare the achieved approximation  $\hat{\mathbf{f}}_t(x) = [\hat{f}_1(x), \hat{f}_2(x)]$  with the underlying function  $\mathbf{f}(x) = [f_1(x), f_2(x)]$ . In this way, we can exactly check the true accuracy of the obtained approximation using different schemes. For the sake of simplicity, we consider the following multi-output mapping

$$\mathbf{f}(x) = [\log(x), 0.5 \log(3x)], \quad x \in (0, 10], \quad (18)$$

then  $D = 1$  and  $P = 2$  (two outputs). Even in this simple scenario, the procedure used for selecting new points is relevant. We start with  $m_0 = 4$  support points,  $\mathbf{X}_0 = [0.1, 3.4, 6.7, 10]$ , apply an independent GP per output, and for AMOGAPE we use the acquisition function denoted as  $\Pi D \times \Pi G$  in Table 3 with the tempering function  $\beta_t = 1 - \frac{1}{t}$ . We also set  $\nu^2 = 0.02$  as a regularization term, in order to avoid numerical issues.

#### 4.1.1 Comparison among sequential methods

It is important to remark that all the active emulators presented in this work are intrinsically sequential techniques. This means that the nodes in  $\mathbf{X}_{t-1}$  are always contained in  $\mathbf{X}_t$ , i.e., the previous configuration of points is always kept. Therefore, for a fair comparison we have to consider other sequential algorithms. We add to  $\mathbf{X}_t$  sequentially 20 additional points, using different sampling strategies: AMOGAPE, uniform points randomly generated in  $(0, 10]$ , a sequential Sobol sequence, and a sequential version of the Latin Hypercube Sampling procedure (Seq-LHS). Seq-LHS simply generates 20 nodes following the LHS procedure and then adds one to  $\mathbf{X}_t$  at each iteration (without replacement). Note that, at each run, the results can vary even for the deterministic procedure due to the optimization of the hyperparameters. We use simulated annealing, which is a stochastic optimization technique [16, 21], both for hyperparameter and acquisition function optimization. We average all the results over 500 independent runs. For model comparison, we compute the root mean square error (RMSE) between  $\hat{\mathbf{f}}_t(x)$  and  $\mathbf{f}(x)$  at each iteration, and show the evolution of the (averaged) RMSE

---

<sup>5</sup>Code available at <https://github.com/dhsvendsen/AMOGAPE>

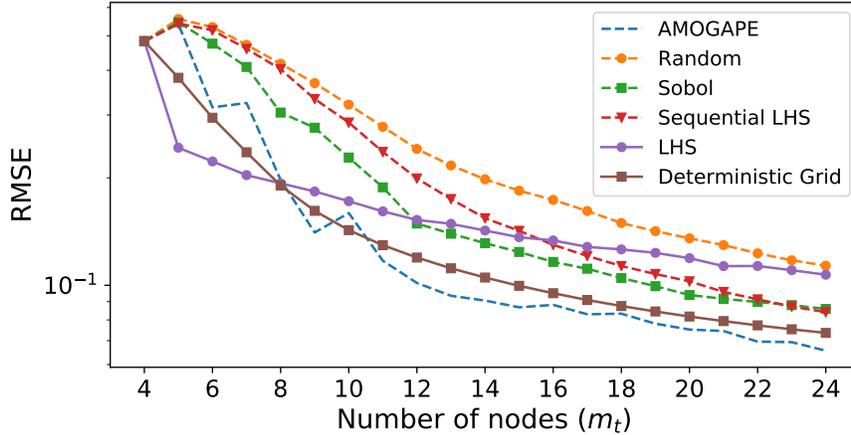


Figure 2: RMSE (in log-scale) between  $f(x)$  and  $\hat{f}_t(x)$  versus the number of nodes  $m_t$ , that is  $m_t = t + 4$  in this example ( $D = 1$  and  $P = 2$ ). Sequential methods, which are more comparable as they utilize  $m_T$  evaluations of  $f(x)$ , are shown with dashed lines. The comparison with two non-sequential methods, using  $\sum_{t=1}^T m_t = \frac{m_T^2 + m_T}{2}$  evaluations of  $f(x)$ , are shown with solid lines.

versus the number of support points  $m_t$  (that is  $m_t = t + m_0$ ) in Figure 2. We can observe that the AMOGAPE scheme outperforms the other methods, providing the smallest RMSEs between  $f(x)$  and  $\hat{f}_t(x)$ .

#### 4.1.2 Comparison with non-sequential methods

In order to provide an exhaustive numerical analysis we also compare AMOGAPE with non-sequential techniques where the input matrix  $\mathbf{X}_t$  can be completely different from  $\mathbf{X}_{t-1}$  (whereas, in AMOGAPE, the nodes in  $\mathbf{X}_{t-1}$  are all always contained in  $\mathbf{X}_t$ ). This approach would not be used in practice, but serves as an interesting comparison of AMOGAPE with one-shot space-filling algorithms. More specifically, we consider:

- Deterministic grid: at each step, we consider an equal-spaced set of points (deterministically chosen). Thus, at each step, all the points in the previous timestep  $\mathbf{X}_{t-1}$  are not considered but replaced by new nodes.
- Standard LHS: also in this case, at each iteration all the previous points are changed.

Clearly, these two schemes evaluate the underlying function in  $m_t$  new nodes at each iteration and are therefore more costly than AMOGAPE. The total number of evaluations of  $f(x)$  for AMOGAPE is  $m_T$  whereas, for the non-sequential schemes above is  $\sum_{t=1}^T m_t = (m_T^2 + m_T)/2$ . However, even in this unfair comparison for our method, Figure 2 shows that AMOGAPE is able to provide the smallest error when more than 12 new points are incorporated. This illustrates that the gradient term encoded in the AMOGAPE adds useful information to the active learning scheme.

## 4.2 Toy Experiment 2: Example of bidimensional multi-output emulation

In this section, we extend the previous example to consider multi-input and multi-output problems, i.e.  $D = P = 2$ . More specifically, we consider

$$\mathbf{f}(\mathbf{x}) = [\log(|\mathbf{x}|), 0.5 \log(3|\mathbf{x}|)], \quad \mathbf{x} \in (0, 10] \times (0, 10]. \quad (19)$$

We start with  $m_0 = 25$  starting nodes in the input matrix,  $\mathbf{X}_0 = [\mathbf{x}_1, \dots, \mathbf{x}_{m_0=25}]$  where  $\mathbf{x}_i = [x_{i,1}, x_{i,2}]^\top$ , with  $i = 1, \dots, 25$ , distributed as shown in Figure 3(a) with black circles. In order to evaluate the approximation RMSE obtained with the emulators, we consider a thin grid in the square  $(0, 10] \times (0, 10]$  (with step 0.3). The starting nodes in the input matrix,  $\mathbf{X}_0$  (black points), and the thin test grid (green dots) are shown in Fig. 3(a). We apply again one independent GP for each output and, as in the previous example. For AMOGAPE, we apply the acquisition function denoted as  $\Pi D \times \Pi G$  in Table 3 and we use again the tempering function  $\beta_t = 1 - \frac{1}{t}$  and set  $v^2 = 0.02$  as a regularization term, only for avoiding numerical issues. We compare different sampling strategies : AMOGAPE, a sequential Sobol sequence, and sequential LHS. We add 30 additional points to  $\mathbf{X}_0$  in the first two sequential approaches. In LHS all the previous points change at each iteration. The results (averaged over 500 independent runs) are shown in Fig. 3(b), which show a considerable gain in accuracy and convergence rates by the presented algorithms.

The distributions in input-space of the final 55 nodes - 25 on a grid and 30 subsequently chosen with a sampling algorithm - are shown as a Kernel Density Estimation (KDE) plot for each of the methods in Figure. 4. We can observe that AMOGAPE adds points in the border and in a left-bottom corner where the gradient is comparatively high. It makes sense that these points are deemed the most useful since the initial 25 nodes points are well-located. The sampling method using the Sobol algorithm and the sequential LHS algorithms incorporate new points that fill out the input-space but do not pay particular attention to the behaviour of the underlying function.

## 4.3 Application to remote sensing: Emulating a radiative transfer model

Our method is assessed for the emulation of the leaf-canopy PROSAIL RTM, which is the most widely used RTM over the last two decades in remote sensing studies [13]. PROSAIL simulates reflectance as a function of:

1. *Leaf optical properties*, given by the mesophyll structural parameter (N), leaf chlorophyll (Chl), dry matter (Cm), water (Cw), carotenoid (Car) and brown pigment (Cbr) contents.
2. *Canopy level characteristics*, determined by leaf area index (LAI), the average leaf angle inclination (ALA) and the hot-spot parameter (Hotspot). System geometry is described by the solar zenith angle ( $\theta_s$ ), view zenith angle ( $\theta_v$ ), and the relative azimuth angle between both angles ( $\Delta\Theta$ ).

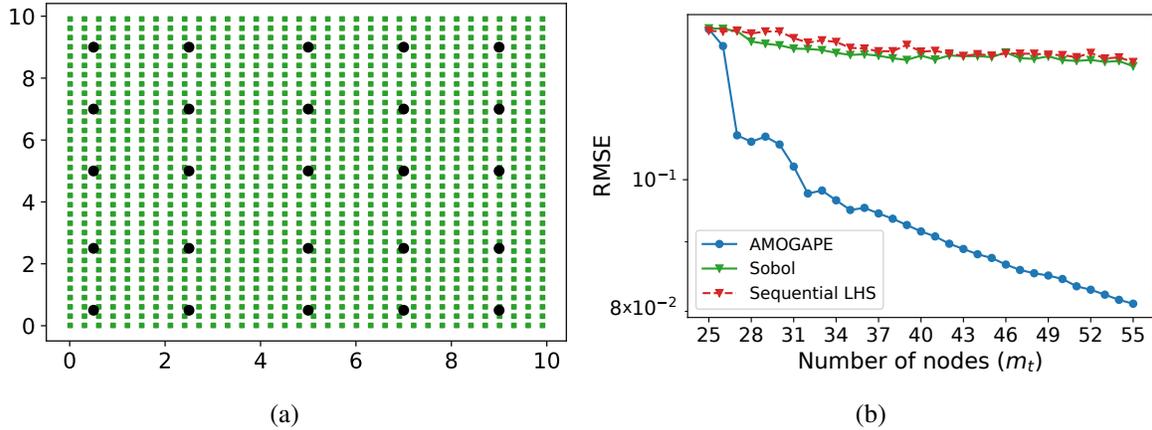


Figure 3: **(a)** The starting points in the input matrix  $\mathbf{X}_0$  are shown with black points, whereas the points in the thin test grid are depicted with green squares. **(b)** RMSE (in log-scale) between  $f(\mathbf{x})$  and  $\hat{f}_t(\mathbf{x})$  versus the number of the number of support points  $m_t$ , that is  $m_t = \dots$  in this example ( $D = 2$  and  $P = 2$ ). Note that the number of initial points is  $m_0 = 25$ .

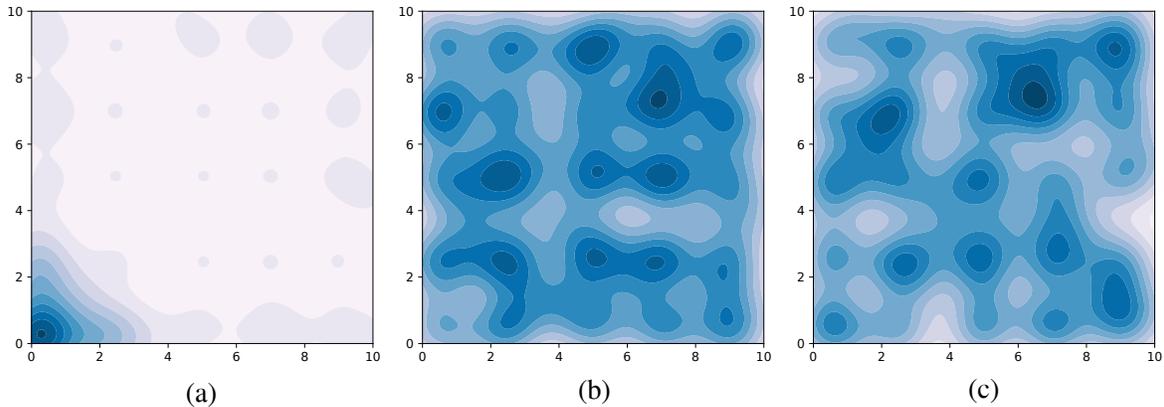


Figure 4: Kernel Density Estimation plot of the final configuration of points in the emulation of Eq. (19), showing the **(a)** AMOGAPE, **(b)** Sobol and **(c)** Sequential LHS algorithms respectively.

We consider PROSAIL for simulating [Landsat-8](#) spectra, a satellite sensor widely used for land cover applications in general and vegetation monitoring in particular. Therefore, the generated, eventually optimized, look-up tables are used for inversion and thus retrieve vegetation parameters with the Landsat-8 satellite imagery. This leaves us with an output-dimension of  $P = 9$  for our problem, i.e. the number of spectral bands of the satellite. Now, depending on the parameters of interest the input dimensionality  $D$  may vary.

### 4.3.1 Sampling a 2-dimensional space for PROSAIL emulation

In this experiment, we chose the most important variables at leaf and canopy-level respectively, namely Chl and LAI, and kept the rest fixed. Table 4 shows the values for the remain-

ing parameters which are set for simulation of rice crops [7]. When generating look-up tables with RTMs it is common practice to use expert knowledge to determine distributions over the biophysical parameters which constitute the RTM input [38]. The desired amount of samples are then drawn, and the model is evaluated in each of these points. A commonly used distribution is the truncated Gaussian  $\mathcal{N}_{\mathcal{T}}(\mathbf{x}|\mu, \sigma, \min, \max)$ . Indeed, the truncated Gaussians  $\mathcal{N}_{\mathcal{T}}(\text{Chl}|45, 30, 20, 90)$  and  $\mathcal{N}_{\mathcal{T}}(\text{LAI}|3.5, 4.5, 0, 10)$  for Chl and LAI, respectively, have proven effective for crop reflectance modeling [7]. We denote their joint distribution, which has no covariance between the variables, as  $\mathcal{N}_{\mathcal{T}}(\text{Chl}, \text{LAI})$ .

Table 4: Characteristics of the simulation used in the PROSAIL model.

<i>Leaf level</i>	N	Cm	Cw	Car	Cbr
	1.5	0.01 $\mu\text{g}/\text{cm}^2$	0.01 $\mu\text{g}/\text{cm}^2$	8 $\text{g}/\text{cm}^2$	0
<i>Canopy level</i>	ALA	Hotspot	$\theta_s$	$\theta_v$	$\Delta\Theta$
	Spherical	0.01	30°	10°	0

In summary, we are emulating a function  $f(\mathbf{x})$  where  $\mathbf{x} = [\text{Chl}, \text{LAI}]$  mapping from an input space of dimension  $D = 2$  to the output space of dimension  $P = 9$ . The search space is restricted to physically meaningful values of  $\text{Chl} \in [20; 90] \mu\text{g}/\text{cm}^2$  and  $\text{LAI} \in [0; 10]$ . In order to gain insight into the relative importance of the Diversity and Geometric terms, an array of different acquisition functions shown in Table 3 are applied. The AMOGAPE sampling schemes are compared with sampling randomly from  $\mathcal{N}_{\mathcal{T}}(\text{Chl}, \text{LAI})$ . This distribution encodes knowledge about the physically feasible region to sample in, which is also encoded in the AMOGAPE, simply by multiplying the truncated density function  $\psi(\text{Chl}, \text{LAI})$  onto the acquisition functions. We set  $\beta_t = 1 \forall t$  in order to simplify the experiments.

Evaluation of which sampling method leads to the best emulator is done by computing the test approximation error on a test-set of 5000 points, sampled from the above-mentioned truncated Gaussian distributions. We initialize with 30 points drawn from  $\mathcal{N}_{\mathcal{T}}(\text{Chl}, \text{LAI})$ . The multi-output RMSE for the  $M = 5000$  test points over the  $P = 9$  single-output GP emulators is computed as follows

$$\text{RMSE} = \sqrt{\frac{1}{M} \sum_{i=1}^M \frac{1}{P} \sum_{p=1}^P (y_{p,i} - \hat{y}_{p,i})^2}. \quad (20)$$

The results are averaged over 15 runs. In order to speed up the experiment, hyperparameter and acquisition function optimization are performed through an initial random search of  $10^D$  points, followed by gradient ascent. Results are shown in Fig. 5. We see that it is possible to perform better using the AMOGAPE approach on our test-set than by sampling randomly from  $\mathcal{N}_{\mathcal{T}}(\text{Chl}, \text{LAI})$ . It is interesting to note that methods using  $\Sigma\text{D} \times \Sigma\text{G}$  and  $\Sigma\text{D}$  perform similarly, implying that the  $\Sigma\text{D}$  term is governing the acquisition function. Similarly, methods using  $\Pi\text{D}$  and  $\Pi\text{D} \times \Sigma\text{G}$  perform equally well, showing that  $\Pi\text{D}$  is the most influential term. The acquisition function  $\Sigma\text{D} \times \Pi\text{G}$ , which penalizes a zero-gradient in any of the output dimension, relies too much on geometric information and performs the worst. It seems that

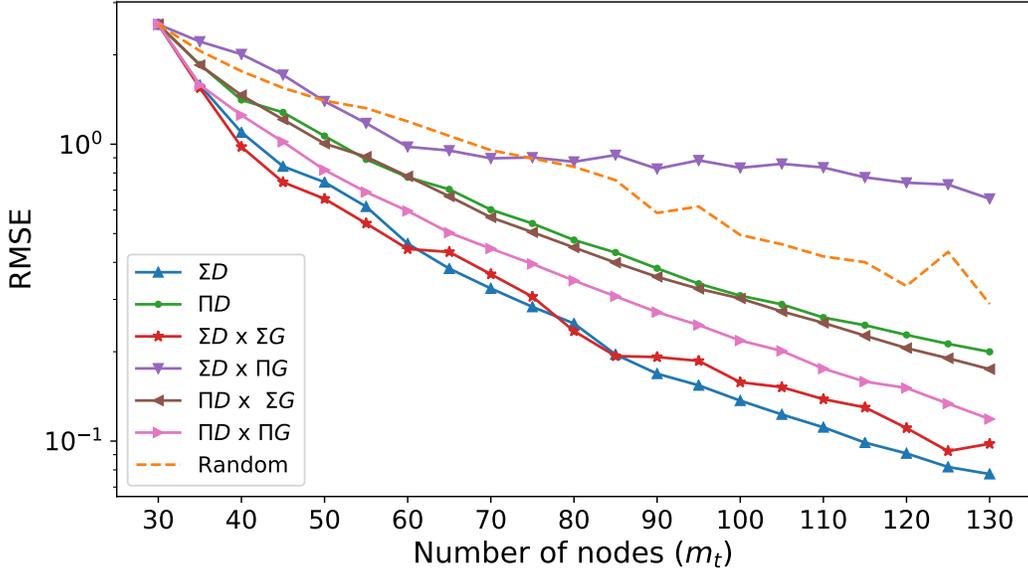


Figure 5: Function approximation errors by different acquisition functions, cf. Tab. 3, and for different number of selected nodes  $m_t$  in a bidimensional PROSAIL problem. Only the best performing acquisition functions are compared here to random sampling.

the information source which is included in product form governs  $A(\mathbf{x})$ . It seems that the  $\Pi D \times \Pi G$  method manages to strike a balance between the two sources of information. All in all, the best performing methods are  $\Sigma D$  and  $\Sigma D \times \Sigma G$ . This hints at the idea that the product form is too restrictive, i.e. considering a point uninteresting if the predictive variance is close to zero in only one of the output-dimensions.

### 4.3.2 Sampling a 3-dimensional space for PROSAIL emulation

We conduct a similar experiment, including now another crucial biophysical parameter in the search space, namely dry matter content (Cm), which is an important parameter to monitor key properties and processes in vegetation and the wider ecosystem. The associated truncated Gaussian used is  $\mathcal{N}_{\mathcal{T}}(\text{Cm} | 0.005, 0.005, 0.003, 0.011)$ . We use a test set of 50000 points generated from the joint truncated Gaussian  $\mathcal{N}_{\mathcal{T}}(\text{Chl}, \text{LAI}, \text{Cm})$ .

We saw earlier that the acquisition function which performs the best was also the most simple, namely  $\Sigma D$ . The acquisition function  $\Pi D \times \Pi G$ , being formulated only in product form, manages not to be dominated by either term and is interesting because it is very selective: It discourages a gradient or predictive variance which is close to zero in any output dimension. For these reasons, along with computational burden, the aforementioned acquisition functions are used for the 3-dimensional experiment. The average results after running the experiment 10 times are shown in Fig. 6. Again, we see that the 1) the two variants of AMOGAPE acquisition functions outperform random sampling, 2) that the acquisition functions behave quite similarly, and 3) that simple acquisitions perform as well as more

complicated ones. Note however, that using a different tempering function than  $\beta_t = 1 \forall t$  would likely make performances diverge.

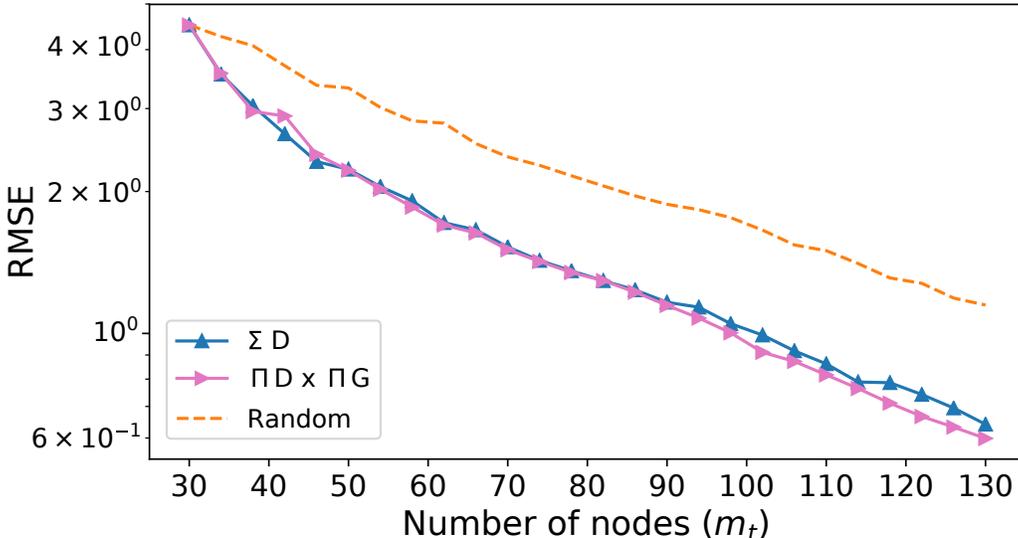


Figure 6: Function approximation errors by different acquisition functions and for different number of selected nodes  $m_t$  in the three-dimensional PROSAIL problem.

## 5 Conclusions

We introduced a simple framework for active construction of emulators for costly physical models used in Earth observation. The proposed framework does not only provide an effective approximating function, but also a compact LUT and some very useful by-products for practitioners, namely confidence intervals for the estimates and information about the gradients.

The methodology iteratively incorporates new sample points that meet both diversity and geometry criteria, thus sampling in low-density and more ‘complex’ regions. This is accomplished by building an acquisition function that takes into account the predictive variance and the norm of the gradient of the GP function used for emulation. The combination of the geometric and diversity sampling criteria was possible because both the GP predictive variance and the gradient of the GP predictive mean are analytic expressions.

We illustrated the promising capabilities of the method through emulation of a popular radiative transfer model. Comparison to established methods in the literature illustrated the favourable performance of the proposed methods. The proposed family of criteria for defining the acquisition functions in emulation allows smart sampling of the input space thus leading to compact and expressive look-up-tables, which can be readily used for model inversion in either statistical or numerical frameworks. The proposed methodology is very general and modular. Alternative acquisition functions, kernel functions and quality measures adapted to

the problem are interesting pathways to explore. In our future work we plan to explore the use of Matérn kernels when function smoothness is not a strict (or even realistic) requirement in a given RTM. Besides, other quality measures other than RMSE could be more interesting for evaluating emulator quality. The information content of the added samples in each iteration by computing maximum differential of entropies in similar ways to the approach in [27].

We anticipate adoption of these methods in the Earth sciences and also in unrelated disciplines where process-based models are widely adopted as well, from econometrics to industry or health sciences. Our future work is centered around speeding up other more complex codes, such as the atmosphere MODTRAN model, as well as to extend the framework to deal with dynamic models. The framework introduced here constitutes the first step towards the ambitious goal of large scale active statistical models that learn Physics models.

## Acknowledgments

This work is supported by the European Research Council (ERC) under the ERC-CoG-2014 SEDAL Consolidator grant (grant agreement 647423).

## References

- [1] A. Ajdari and H. Mahlooji. An adaptive exploration-exploitation algorithm for constructing metamodels in random simulation using a novel sequential experimental design. *Communications in Statistics-Simulation and Computation*, 43(5):947–968, 2014.
- [2] M. Alvarez, D. Luengo, and N. Lawrence. Linear latent force models using Gaussian processes. *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, 35(11):2693–2705, 2013.
- [3] M. A. Alvarez, D. Luengo, M. K. Titsias, and N. D. Lawrence. Efficient multioutput gaussian processes through variational inducing kernels. In *International Conference on Artificial Intelligence and Statistics*, pages 25–32, 2010.
- [4] M. A. Alvarez, L. Rosasco, and N. D. Lawrence. Kernels for vector-valued functions: A review. *Foundations and Trends in Machine Learning*, 4(3):195, 2012.
- [5] P. Bratley and B. Fox. Algorithm 659: Implementing Sobol’s Quasirandom Sequence Generator. *ACM Transactions on Mathematical Software (TOMS)*, 14(1):88–100, 1988.
- [6] D. Busby. Hierarchical adaptive experimental design for gaussian process emulators. *Reliability Engineering & System Safety*, 94(7):1183–1193, 2009.
- [7] M. Campos-Taberner, F. García-Haro, G. Camps-Valls, G. Grau-Muedra, F. Nutini, A. Crema, and M. Boschetti. Multitemporal and multiresolution leaf area index retrieval for operational local rice crop monitoring. *Remote Sensing of Environment*, 187:102 – 118, 2016.

- [8] G. Camps-Valls, J. Verrelst, J. Muñoz Mari, V. Laparra, F. Mateo-Jiménez, and J. Gomez-Dans. A survey on gaussian processes for earth observation data analysis. *IEEE Geoscience and Remote Sensing Magazine*, 2(6), June 2016.
- [9] C. S. Gillespie and R. J. Boys. Efficient construction of Bayes optimal designs for stochastic process models. *arXiv:1803.04254*, pages 1–21, 2018.
- [10] D. Gorissen, I. Couckuyt, P. Demeester, T. Dhaene, and K. Crombecq. A surrogate modeling and adaptive sampling toolbox for computer based design. *Journal of Machine Learning Research*, 11(Jul):2051–2055, 2010.
- [11] D. Gorissen, L. De Tommasi, K. Crombecq, and T. Dhaene. Sequential modeling of a low noise amplifier with neural networks and active learning. *Neural Computing and Applications*, 18(5):485–494, 2009.
- [12] R. K. S. Hankin et al. Introducing multivator: a multivariate emulator. *Journal of Statistical Software*, 46(8):1–20, 2012.
- [13] S. Jacquemoud, W. Verhoef, F. Baret, C. Bacour, P. Zarco-Tejada, G. Asner, C. François, and S. Ustin. PROSPECT+ SAIL models: A review of use for vegetation characterization. *Remote sensing of environment*, 113:S56–S66, 2009.
- [14] M. E. Johnson, L. M. Moore, and D. Ylvisaker. Minimax and maximin distance designs. *Journal of statistical planning and inference*, 26(2):131–148, 1990.
- [15] M. Kennedy and A. O’Hagan. Bayesian calibration of computer models. *Journal of the Royal Statistical Society. Series B: Statistical Methodology*, 63(3):425–450, 2001.
- [16] S. K. Kirkpatrick, C. D. G. Jr., and M. P. Vecchi. Optimization by simulated annealing. *Science*, 220(4598):671–680, May 1983.
- [17] A. Krause, A. Singh, and C. Guestrin. Near-optimal sensor placements in gaussian processes: Theory, efficient algorithms and empirical studies. *Journal of Machine Learning Research*, 9(Feb):235–284, 2008.
- [18] J. L. Loepky, L. M. Moore, and B. J. Williams. Batch sequential designs for computer experiments. *Journal of Statistical Planning and Inference*, 140(6):1452–1464, 2010.
- [19] D. Luengo-Garcia, M. Campos-Taberner, and G. Camps-Valls. Latent force models for earth observation time series prediction. In *2016 IEEE International Workshop on Machine Learning for Signal Processing (MLSP 2016)*, Salerno, Italy, September 2016.
- [20] S. Marmin, D. Ginsbourger, J. Baccou, and J. Liandrat. Warped gaussian processes and derivative-based sequential designs for functions with heterogeneous variations. *SIAM/ASA Journal on Uncertainty Quantification*, 6(3):991–1018, 2018.

- [21] L. Martino, V. Elvira, D. Luengo, J. Corander, and F. Louzada. Orthogonal parallel MCMC methods for sampling and optimization. *Digital Signal Processing*, 58:64–84, 2016.
- [22] M. D. McKay, R. J. Beckman, and W. J. Conover. Comparison of three methods for selecting values of input variables in the analysis of output from a computer code. *Technometrics*, 21(2):239–245, 1979.
- [23] P. Mitra, B. U. Shankar, and S. K. Pal. Segmentation of multispectral remote sensing images using active support vector machines. *Pattern recognition letters*, 25(9):1067–1074, 2004.
- [24] J. Oakley. *Bayesian uncertainty analysis for complex computer codes*. PhD thesis, University of Sheffield, 1999.
- [25] A. O’Hagan. Bayesian analysis of computer code outputs: A tutorial. *Reliability Engineering and System Safety*, 91(10-11):1290–1300, 2006.
- [26] C. E. Rasmussen and C. K. I. Williams. *Gaussian Processes for Machine Learning*. The MIT Press, New York, 2006.
- [27] P. Ruiz, J. Mateos, G. Camps-Valls, R. Molina, and A. K. Katsaggelos. Bayesian active remote sensing image classification. *IEEE Transactions on Geoscience and Remote Sensing*, 52(4):2186–2196, 2013.
- [28] C. M. Ryan, C. C. Drovandi, and A. N. Pettitt. Optimal Bayesian experimental design for models with intractable likelihoods using indirect inference applied to biological process models. *Bayesian Analysis*, 11(3):857–883, 2016.
- [29] J. Sacks, S. B. Schiller, and W. J. Welch. Designs for computer experiments. *Technometrics*, 31(1):41–47, 1989.
- [30] T. Santner, B. Williams, and W. Notz. *The design and analysis of computer experiments*. Springer Verlag, 2003.
- [31] S. Seo, M. Wallat, T. Graepel, and K. Obermayer. Gaussian process regression: Active data selection and test point rejection. In *Mustererkennung 2000*, pages 27–34. Springer, 2000.
- [32] M. C. Shewry and H. P. Wynn. Maximum entropy sampling. *Journal of applied statistics*, 14(2):165–170, 1987.
- [33] J. Snoek, H. Larochelle, and R. P. Adams. Practical bayesian optimization of machine learning algorithms. In *Advances in neural information processing systems*, pages 2951–2959, 2012.

- [34] D. H. Svendsen, L. Martino, J. Vicent, and G. Camps-Valls. Multioutput automatic emulator for radiative transfer models. In *IGARSS 2018-2018 IEEE International Geoscience and Remote Sensing Symposium*, pages 4019–4022. IEEE, 2018.
- [35] D. Tuia, J. Muñoz-Marí, and G. Camps-Valls. Remote sensing image segmentation by active queries. *Pattern Recognition*, 45(6):2180–2192, 2012.
- [36] H. Wang and J. Li. Adaptive gaussian process approximation for bayesian inference with expensive likelihood functions. *Neural computation*, 30(11):3072–3094, 2018.
- [37] Z. Wang, S. Yan, and C. Zhang. Active learning with adaptive regularization. *Pattern Recognition*, 44(10):2375 – 2383, 2011.
- [38] M. Weiss and F. Baret. S2ToolBox Level 2 products: LAI, FAPAR, FCOVER, 2016.
- [39] B. Wescott. *Every Computer Performance Book: How to Avoid and Solve Performance Problems on The Computers You Work With*. CreateSpace Independent Publishing Platform, USA, 1st edition, 2013.
- [40] A. G. Wilson, D. A. Knowles, and Z. Ghahramani. Gaussian process regression networks. In J. Langford and J. Pineau, editors, *Proceedings of the 29th International Conference on Machine Learning (ICML)*, Edinburgh, June 2012. Omnipress.
- [41] D. Wu, C.-T. Lin, and J. Huang. Active learning for regression using greedy sampling. *Information Sciences*, 474:90–105, 2019.
- [42] Y. Yang and M. Loog. A benchmark and comparison of active learning for logistic regression. *Pattern Recognition*, 83:401 – 415, 2018.

## A Importance of a gradient term in the acquisition function

Let us consider the problem of approximating the function  $y = f(\mathbf{x})$ ,  $\mathbf{x} \in \mathcal{D} \subseteq \mathbb{R}^L$ , where  $\mathcal{D} = [a_1, b_1] \times [a_2, b_2] \times \dots \times [a_L, b_L]$ . For the sake of simplicity, we consider one dimensional problems, i.e.  $L = 1$ , with  $\mathcal{D}$  bounded  $a_i, b_i < \infty$ . Moreover, let us consider a set of nodes  $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_M\} \in \mathcal{D}$  and the corresponding values of the function  $y_m = f(\mathbf{x}_m)$ . We use the pairs  $\{\mathbf{x}_m, y_m\}_{m=1}^M$  to perform an interpolation. Given the set of nodes  $\{\mathbf{x}_m\}_{m=1}^M$ , we denote a piecewise constant interpolation (PCI) of the function  $f(\mathbf{x})$  as

$$\hat{y} = \phi(\mathbf{x}|\mathbf{x}_{1:M}) = \phi(\mathbf{x}). \quad (21)$$

In order to measure the discrepancy between function and emulator we introduce a cost function  $C(f, \phi) = C(f(\mathbf{x}), \phi(\mathbf{x})) \geq 0$ . The equality  $C(f, \phi) = 0$  must hold only if  $\phi(\mathbf{x}) = f(\mathbf{x})$ . Here, we consider the  $L_p$  family of cost functions

$$C_p(f, \phi) = \|f(\mathbf{x}) - \phi(\mathbf{x})\|_p = \left( \int_{\mathcal{D}} |f(\mathbf{x}) - \phi(\mathbf{x})|^p d\mathbf{x} \right)^{\frac{1}{p}}. \quad (22)$$

Note that  $C_\infty(f, \phi) = \lim_{p \rightarrow \infty} C_p(f, \phi) = \max_{\mathbf{x} \in \mathcal{D}} |f(\mathbf{x}) - \phi(\mathbf{x})|$ .

### One node ( $M = 1$ ), infinity norm cost functions $p = \infty$

Let us consider  $y = f(x)$ ,  $x \in [a, b] \subseteq \mathbb{R}$ . For simplicity, we assume  $f(x)$  to be strictly monotonic, more specifically increasing. We consider a piecewise constant approximation with  $M = 1$  point  $x_1$  within  $x_1$ , i.e.

$$\phi(x) = \begin{cases} f(a) & x \leq x_1 \\ f(x_1) & x > x_1 \end{cases} \quad (23)$$

Let us consider the  $L_\infty$  distance (i.e.,  $p = \infty$ ),

$$\begin{aligned} C_\infty(x_1) &= \max_{x \in [a, b]} |f(x) - \phi_0(x)| = \max_{x \in [a, b]} [|f(x_1) - f(a)|, |f(x_1) - f(b)|], \\ &= \max_{x \in [a, b]} [f(x_1) - f(a), f(b) - f(x_1)] \end{aligned} \quad (24)$$

The problem consists in finding the optimal node  $x_1^*$  such that  $x_1^* = \arg \min C_\infty(x_1)$ . This optimal point  $x_1^*$  will then satisfy the condition

$$f(x_1^*) - f(a) = f(b) - f(x_1^*). \quad (25)$$

This is because  $c_a \equiv |f(x_1) - f(a)|$  will decrease as  $c_b \equiv |f(x_1) - f(b)|$  increases, and vice versa, due to the monotonicity of  $f(x)$ . Since we are taking the max between the two, the lowest value that  $C_\infty$  can take is the point where they are equal  $c_a = c_b$  (see Fig. 7), as any divergence from that would lead to one of the terms being higher.

Using Eq. (25), since we are assuming that  $f$  is monotonic (thus invertible), we can also write

$$f(x_1^*) = \frac{f(b) + f(a)}{2}, \quad (26)$$

and, since we have assumed that  $f(x)$  is monotonic, thus invertible, we have

$$x_1^* = f^{-1} \left( \frac{f(b) + f(a)}{2} \right). \quad (27)$$

Figure 7 illustrates the above reasoning. Note that, if  $f(x)$  is non-linear,  $x_1^* \neq \frac{a+b}{2}$  (as a space filling/Latin hypercube strategy might suggest). The expression of  $x_1^*$  is an extended mean, which takes into account information regarding the non-linearity  $f(x)$ .

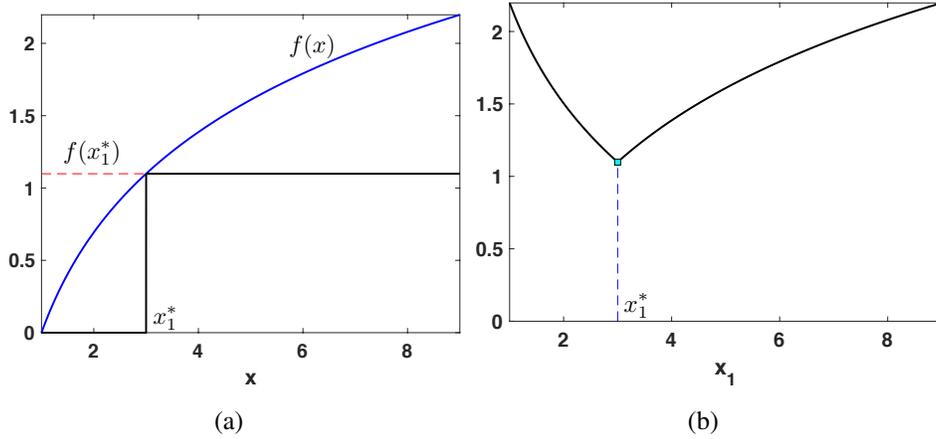


Figure 7: **(a)** Optimal piecewise constant approximation  $\phi(x)$  with  $M = 1$  node. **(b)** The cost function  $C_\infty(x_1)$  and its minimum at  $x_1^*$ .

### Generic number of nodes ( $M > 1$ )

Let us assume now that we may place  $M$  nodes  $x_1, x_2, \dots, x_M$  in order to achieve an optimal emulator of  $f$  with respect to the  $C_\infty$  norm. The optimal nodes  $x_1^*, x_2^*, \dots, x_M^*$  will then satisfy the condition

$$f(x_1^*) - f(a) = f(x_2^*) - f(x_1^*) = \dots = f(x_{M-1}^*) - f(x_{M-2}^*) = f(b) - f(x_{M-1}^*). \quad (28)$$

The point  $(x_1^*, x_2^*, \dots, x_M^*)$  is a minimum for  $C_\infty(x_1, x_2, \dots, x_M)$  and is unique. In order to see this, let us define  $d_1 = f(x_1) - f(a)$ ,  $d_2 = f(x_2) - f(x_1)$ ,  $\dots$ ,  $d_m = f(x_m) - f(x_{m-1})$ ,  $\dots$ ,  $d_M = f(x_M) - f(x_{M-1})$ , and  $d_{M+1} = f(b) - f(x_M)$ . With this definition we reach the minimum for  $C_\infty(x_1, x_2, \dots, x_M)$  when all distances  $\{d_m\}_{m=1}^M$  are equal to  $f(b) - f(a)$  divided by  $M + 1$ :

$$d_1^* = d_2^* = \dots = d_{M+1}^* = \frac{f(b) - f(a)}{M + 1} \equiv c_{\text{MIN}} \quad (29)$$

This can be seen from the fact that the distances satisfy  $d_m \geq 0$  for  $m = 1, 2, \dots, M$  due to the monotonicity of  $f$ , and they sum to a constant

$$\sum_{m=1}^{M+1} d_m = f(b) - f(a) \quad (30)$$

Thus if one  $d_m$  decreases, one or all other have to increase. This implies that any configuration of  $x_1, x_2, \dots, x_M$  resulting in  $d_j < c_{\text{MIN}}$  for some  $j$ , will lead  $d_k > c_{\text{MIN}}$  for one or more  $k$ . Therefore, the configuration of corresponding to (29) is optimal.

Following the above logic, the optimal locations of  $M$  nodes,  $\{x_1, \dots, x_M\}$ , can be obtained by:

1. Dividing with a uniform grid formed by  $M$  points the interval  $[f(a), f(b)]$  (image of  $[a, b]$ ),

$$y_m = f(a) + m \frac{f(b) - f(a)}{M + 1}, \quad m = 1, \dots, M, T \quad (31)$$

2. Finding the  $x_m$  such that  $f(x_m) = y_m$ , i.e., since we assume that  $f(x)$  is invertible,

$$x_m = f^{-1}(y_m), \quad m = 1, \dots, M. \quad (32)$$

See Fig. 8 for an example with  $M = 2$  nodes.

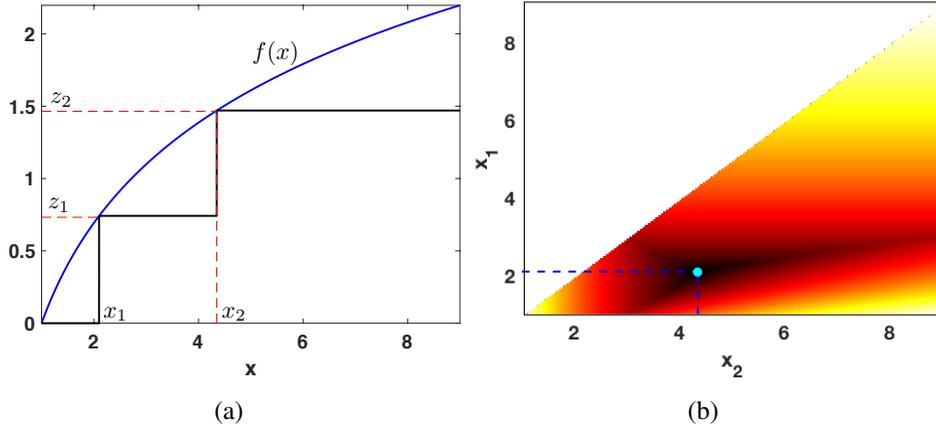


Figure 8: **(a)** Optimal piecewise constant approximation  $\phi_0(x)$  with  $M = 2$  nodes. **(b)** The cost function  $C_\infty(x_1, x_2)$  and its minimum at  $(x_1^*, x_2^*)$ . Note that  $C_\infty(x_1, x_2)$  is defined within the simplex such that  $x_1 \leq x_2$ , by definition; it can be also considered that  $C_\infty(x_1, x_2) = C_\infty(x_2, x_1)$ .

## Distributions of nodes

We have seen that the auxiliary points  $y_m$  are obtained using a uniform grid in the interval  $[f(a), f(b)]$  (image of  $[a, b]$ ). Therefore,  $y_1, \dots, y_M$  is a quasi-Monte Carlo sequence distributed uniformly in  $[f(a), f(b)]$ , i.e.,

$$y_m \sim \mathcal{U}([f(a), f(b)]), \quad m = 1, \dots, M, \quad (33)$$

Following Eq. (32), we can find the distribution of the nodes  $x_m$  since they are obtained by transforming the points  $y_m$  through the function  $f^{-1}(\cdot)$ . Hence, following the expression of the transformation of a random variable, we have

$$x_m \sim p_X(x) = p_Y(f(x)) \left| \frac{df}{dx} \right| \quad (34)$$

$$\propto \left| \frac{df}{dx} \right|, \quad m = 1, \dots, M, \quad (35)$$

Therefore, the set of nodes  $x_1, \dots, x_M$  is a quasi-Monte Carlo sequence with density  $p_X(x) \propto \left| \frac{df}{dx} \right|$  and if  $f$  is increasing, we can write  $p_X(x) \propto \frac{df}{dx}$ . See Fig. 9 for an illustration of this. For higher input dimension than 1 we have

$$\mathbf{x}_m \sim p_X(\mathbf{x}) \propto |\nabla f(\mathbf{x})|. \quad (36)$$

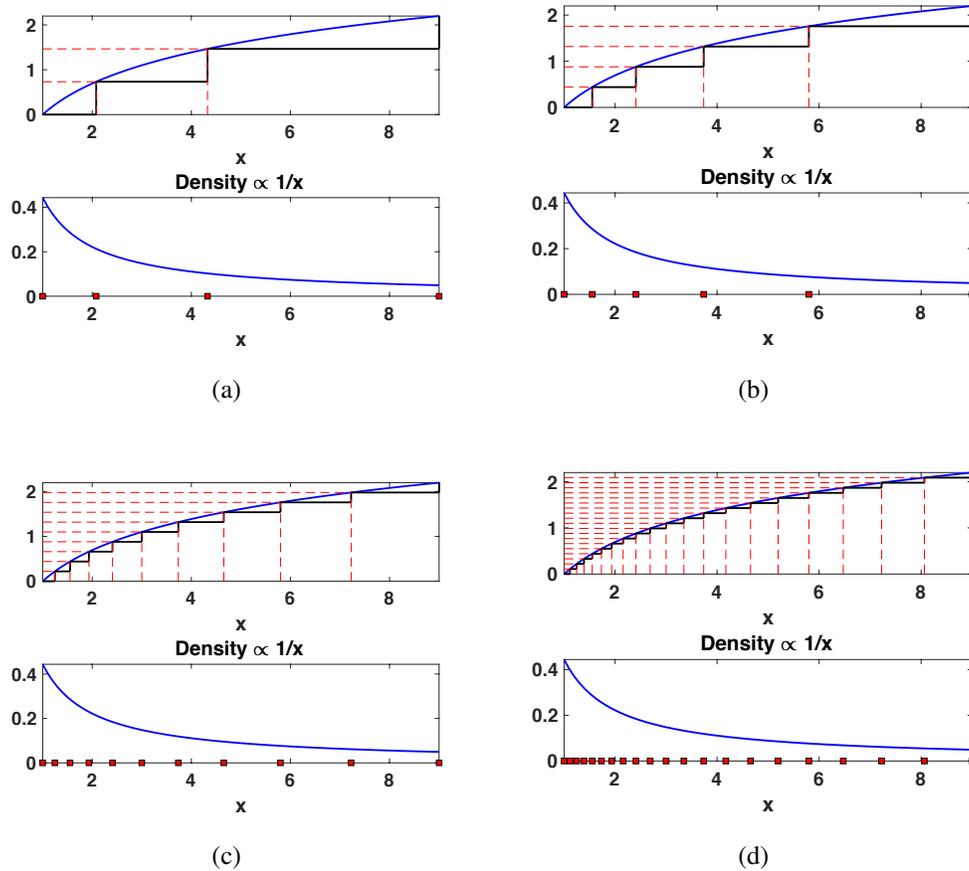


Figure 9: Illustration of function and optimal location of nodes (top) and density proportional to gradient (bottom). For (a)-(d) the number of nodes are 2, 4, 10 and 20 respectively.