

GUILLAUME SAGNOL, HANS-CHRISTIAN HEGE, MARTIN WEISER

# **Using sparse kernels to design computer experiments with tunable precision**

Zuse Institute Berlin  
Takustrasse 7  
D-14195 Berlin-Dahlem

Telefon: 030-84185-0  
Telefax: 030-84185-125

e-mail: [bibliothek@zib.de](mailto:bibliothek@zib.de)  
URL: <http://www.zib.de>

ZIB-Report (Print) ISSN 1438-0064  
ZIB-Report (Internet) ISSN 2192-7782

# Using sparse kernels to design computer experiments with tunable precision

Guillaume Sagnol, Hans-Christian Hege, Martin Weiser

June 3, 2016

## Abstract

Statistical methods to design computer experiments usually rely on a Gaussian process (GP) surrogate model, and typically aim at selecting design points (combinations of algorithmic and model parameters) that minimize the average prediction variance, or maximize the prediction accuracy for the hyperparameters of the GP surrogate. In many applications, experiments have a *tunable precision*, in the sense that one software parameter controls the tradeoff between accuracy and computing time (e.g., mesh size in FEM simulations or number of Monte-Carlo samples). We formulate the problem of allocating a budget of computing time over a finite set of candidate points for the goals mentioned above. This is a continuous optimization problem, which is moreover convex whenever the tradeoff function accuracy vs. computing time is concave. On the other hand, using non-concave weight functions can help to identify sparse designs. In addition, using sparse kernel approximations drastically reduce the cost per iteration of the multiplicative weights updates that can be used to solve this problem.

**Keywords:** Optimal design of computer experiments, Gaussian process, Sparse kernels

## 1 Introduction

We consider a computer code taking an input  $\mathbf{x} \in \mathcal{X}$  (called a *design point*) in a compact set  $\mathcal{X} \subset \mathbb{R}^d$  and a parameter  $\tau$  specifying the time allowed for the computation, and returning an output of the form

$$Y(\mathbf{x}, \tau) = \eta(\mathbf{x}) + \epsilon(\mathbf{x}, \tau), \quad (1)$$

where  $\eta(\cdot)$  is an unknown function and  $\epsilon(\mathbf{x}, \tau)$  represents uncorrelated errors:  $\mathbb{E}[\epsilon(\mathbf{x}, \tau)] = 0$ ,  $\mathbf{u} \neq \mathbf{v} \Rightarrow \mathbb{E}[\epsilon(\mathbf{u}, \tau_u)\epsilon(\mathbf{v}, \tau_v)] = 0$  for all  $\tau_u, \tau_v > 0$ , where  $\mathbb{E}[X]$  stands for the expectation of  $X$ . We assume that the experiments have a *tunable precision*, in the sense that the variance of the error  $\mathbb{V}[\epsilon(\mathbf{x}, \tau)]$  is a decreasing function of the time  $\tau$  spent to compute an approximation of  $\eta(\mathbf{x})$ . Specifically, we assume that there is a *known* parameter  $\sigma_N^2$  (where the subscript  $N$  stands for *noise*) and a *known* differentiable, nondecreasing function  $w : \mathbb{R}_+ \mapsto \mathbb{R}_+$  satisfying  $w(0) = 0$ , such that  $\mathbb{V}[\epsilon(\mathbf{x}, \tau)] = \sigma_N^2 w(\tau)^{-1}$ . Let the experimental design (or simply *the design*) be represented by  $\xi = \{\mathbf{x}_i, \tau_i\}_{i \in \{1, \dots, n\}}$ , where the *given candidate points*  $\mathbf{x}_i \in \mathcal{X}$  are distinct, and  $\tau_i \geq 0$  is the computing time spent on design point  $\mathbf{x}_i$ . Note that  $\tau_i = 0$  means that no computation is carried out at  $\mathbf{x}_i$ , and formally we have  $\mathbb{V}[\epsilon(\mathbf{x}_i, \tau_i)] = +\infty$ . Candidate points  $\mathbf{x}_i$  with  $\tau_i > 0$  are called *support points* and are those which are actually selected for the design.

We assume that a two-stage approach is used, and observations have already been collected during the initial stage from a design  $\xi_{\text{init}} = \{\mathbf{x}_i^0, \tau_i^0\}_{i=1, \dots, n_0}$ , with  $\mathbf{x}_i^0 \neq \mathbf{x}_j$  for all  $i, j$ . The purpose of this article is to develop efficient algorithms for the computation of near-optimal computing times  $\tau_i$ , subject to a constraint on the total computing time allowed for the second stage:  $\sum_{i=1}^n \tau_i = T$ . In other words, we search for a (near-)optimal design within the class of all designs that assign a total computing time of  $T$ , and whose support is a subset of the  $\mathbf{x}_i$ 's. In practice, this two-stage approach can be turned into a sequential one, as follows. Given an optimized design  $\tau^*$ , select one support point  $\mathbf{x}_i$  and compute  $Y(\mathbf{x}_i, \tau_i^*)$ . Then, append  $\{\mathbf{x}_i, \tau_i^*\}$  to  $\xi_{\text{init}}$ , decrement  $T$  by  $\tau_i^*$ , remove  $\mathbf{x}_i$  from the list of candidate points, update the surrogate model for  $\eta(\cdot)$ , compute the next design  $\tau^*$ , and iterate. This procedure can also be generalized to work on a parallel architecture, where several design points can be processed simultaneously.

Our assumption about the existence of an information function  $w()$  is not common. Most authors focus on the search for *exact designs*, i.e.  $w_i = w(\tau_i) \in \{0, 1\}$ , and  $w_i = 1$  indicates that the design point  $\mathbf{x}_i$  belongs to  $\xi$ . We refer the reader to [PM12] for a comprehensive review on exact designs for computer experiments. For the standard linear model, a popular technique is to relax the integer constraint on  $w_i$ , which led to the success story of the *theory of approximate designs* [Páz86, Puk93]. Approximate designs are used most often as a heuristic to find good exact designs, typically by rounding. For computer experiments, however, the total computing time is of much more importance than the number of design points, which motivates to study the tradeoff between tunable accuracy and computing time in more detail. We give two examples:

- In the case of Monte-Carlo simulations, the variance is inversely proportional to the number of samples, and hence  $w(\tau) = \tau$ . In fact,  $\tau$  can take integer values only, but we expect the approximation to be good enough for a large number of simulation runs.
- The standard a priori error estimate for finite element solutions of ansatz order  $p$  and mesh width  $h$  for sufficiently regular stationary elliptic problems in  $\Omega \subset \mathbb{R}^d$  is  $\|u - u_h\|_{H^1(\Omega)} = \mathcal{O}(h^p)$ . The optimal computational complexity is  $\tau = \mathcal{O}(h^{-d})$ , which means  $\tau = \mathcal{O}(\|u - u_h\|_{H^1(\Omega)}^{-d/p})$ , see, e.g., [DW12]. We could thus model the error by a noise of variance  $\mathbb{V}[\epsilon(\mathbf{x}, \tau)] = \mathcal{O}(\tau^{-2p/d})$ , i.e.  $w(\tau) = \mathcal{O}(\tau^{2p/d})$ .

Note, however, that errors at two design points might be correlated, as is usually the case in finite element models. An alternative could be to use a Gaussian process to model the error process. We leave this for future research, and simply assume as an approximation that the  $\epsilon(\mathbf{x}_i, \tau_i)$  are uncorrelated, which might be acceptable if the  $\mathbf{x}_i$ 's are far enough from each other.

Following the kriging methodology, we assume that  $\eta(\mathbf{x}) = f(\mathbf{x})^T \boldsymbol{\beta} + Z(\mathbf{x})$ , where  $f : \mathbb{R}^d \mapsto \mathbb{R}^m$  is a spatial regression function (generally a polynomial),  $\boldsymbol{\beta} \in \mathbb{R}^m$  is an unknown vector of parameters, and  $Z$  is a Gaussian random field with zero mean and *known correlation structure*,

$$\mathbb{E}[Z(\mathbf{x})] = 0, \quad \mathbb{E}[Z(\mathbf{u})Z(\mathbf{v})] = \sigma_Z^2 C(\mathbf{u}, \mathbf{v}),$$

where  $C : \mathbb{R}^d \times \mathbb{R}^d \mapsto \mathbb{R}$  is a positive semidefinite kernel satisfying  $C(\mathbf{u}, \mathbf{u}) = 1$  for all  $\mathbf{u} \in \mathcal{X}$ .

Denote by  $Y(\xi) = [Y(\mathbf{x}_1), \dots, Y(\mathbf{x}_n)]^T$  the vector of observations associated with the design  $\xi$ . We have

$$Y(\xi) \sim \mathcal{N}(\mathbf{F}^T \boldsymbol{\beta}, \sigma_Z^2 \mathbf{C} + \boldsymbol{\Delta}), \quad (2)$$

where  $\mathbf{F} = [f(\mathbf{x}_1), \dots, f(\mathbf{x}_n)] \in \mathbb{R}^{m \times n}$ ,  $\{\mathbf{C}\}_{ij} := C(\mathbf{x}_i, \mathbf{x}_j)$ , and  $\boldsymbol{\Delta} = \sigma_N^2 \text{Diag}(w(\tau))^{-1}$ . We also define the *signal-to-noise ratio*  $\gamma = (\sigma_Z/\sigma_N)^2$  and  $\boldsymbol{\Sigma} = \mathbf{C} + \text{Diag}(\gamma w(\tau))^{-1}$ , so that  $\mathbb{V}[Y(\xi)] = \sigma_Z^2 \boldsymbol{\Sigma}$ . The best linear unbiased estimator (BLUE) of  $\boldsymbol{\beta}$  and its variance are given by

$$\hat{\boldsymbol{\beta}} = (\mathbf{F} \boldsymbol{\Sigma}^{-1} \mathbf{F}^T)^{-1} \mathbf{F} \boldsymbol{\Sigma}^{-1} Y(\xi), \quad \mathbb{V}[\hat{\boldsymbol{\beta}}] = \sigma_Z^2 (\mathbf{F} \boldsymbol{\Sigma}^{-1} \mathbf{F}^T)^{-1}.$$

Then, the best linear unbiased predictor (BLUP) of the unknown function  $\eta$  at  $\mathbf{x} \in \mathcal{X}$  based on the observations  $Y(\xi)$  is given by

$$\hat{\eta}(\mathbf{x}|\xi) = f(\mathbf{x})^T \hat{\boldsymbol{\beta}} + c(\mathbf{x})^T \Sigma^{-1} (Y(\xi) - \mathbf{F}^T \hat{\boldsymbol{\beta}}),$$

where  $\{c(\mathbf{x})\}_i := C(\mathbf{x}, \mathbf{x}_i)$  is the vector of cross-covariances between  $\mathbf{x}$  and the design points, and the mean-squared prediction error (MSPE) is

$$\begin{aligned} \rho(\mathbf{x}) &:= \mathbb{E}[(\hat{\eta}(\mathbf{x}|\xi) - \eta(\mathbf{x}))^2] \\ &= \sigma_Z^2 \left\{ 1 - c(\mathbf{x})^T \Sigma^{-1} c(\mathbf{x}) + (f(\mathbf{x}) - \mathbf{F} \Sigma^{-1} c(\mathbf{x}))^T (\mathbf{F} \Sigma^{-1} \mathbf{F}^T)^{-1} (f(\mathbf{x}) - \mathbf{F} \Sigma^{-1} c(\mathbf{x})) \right\}. \end{aligned}$$

The above expression reduces to  $\rho(\mathbf{x}) = \sigma_Z^2 (1 - c(\mathbf{x})^T \Sigma^{-1} c(\mathbf{x}))$  when the trend parameter  $\boldsymbol{\beta}$  is known. Note that  $\rho(\mathbf{x})$  depend on the design  $\xi$  through  $\Sigma$ . A standard approach is to choose  $\xi$  so as to minimize the integrated mean squared error (IMSE):

$$\text{IMSE}(\xi) := \int_{\mathcal{X}} \rho(\mathbf{x}) d\mu(\mathbf{x}).$$

The IMSE criterion depends on a measure  $\mu$  on  $\mathcal{X}$ , which can be used to weigh the interest of the experimenter for knowing the value of  $\eta$  at  $\mathbf{x}$ . E.g., if the goal is to minimize  $\eta(\mathbf{x})$  over  $\mathcal{X}$ , or to estimate the probability that  $\eta(\mathbf{x})$  lies below some threshold,  $\mu$  should weigh regions of  $\mathcal{X}$  such as to balance the exploration/exploitation tradeoff; see, e.g., [JSW98, BGL<sup>+</sup>12].

It was shown (for the standard case where the variance of  $\epsilon$  is not a function of  $\tau$ ) in [Fed96, FF97] that model (1) can be approximated arbitrarily well by a Bayesian linear model of the form

$$Y(\mathbf{x}) \simeq [f(\mathbf{x})^T, g(\mathbf{x})^T] \begin{bmatrix} \boldsymbol{\beta} \\ \boldsymbol{\alpha} \end{bmatrix} + \epsilon(\mathbf{x}), \quad (3)$$

where  $\boldsymbol{\alpha}$  is a random regression parameter with prior  $\boldsymbol{\alpha} \sim \mathcal{N}(0, \sigma_Z^2 \mathbf{I}_s)$ ,  $\mathbf{I}_s$  is the  $s \times s$ -identity matrix, and the function  $g : \mathbb{R}^d \mapsto \mathbb{R}^s$  can be obtained by truncating the Mercer's expansion of the kernel  $C(\cdot, \cdot)$ . In our case, recall that observations have already been collected during an initial stage at  $\xi_{\text{init}} = \{\mathbf{x}_i^0, \tau_i^0\}_{i=1, \dots, n_0}$ ,  $\mathbb{V}[\epsilon(\mathbf{x}, \tau)] = \sigma_N^2 w(\tau)^{-1}$  and the noise is uncorrelated. Then, by using standard results from the literature on Bayesian designs, see, e.g., [Pil91], one obtains the following approximation of the Kriging variance for the design  $\xi = \{\mathbf{x}_i, \tau_i\}_{i=1, \dots, n}$ :

$$\tilde{\rho}(\mathbf{x}) \simeq \sigma_Z^2 h(\mathbf{x})^T \mathbf{M}(\xi)^{-1} h(\mathbf{x}),$$

where  $h(\mathbf{x}) = [f(\mathbf{x})^T, g(\mathbf{x})^T]^T$  and  $\mathbf{M}(\xi)$  is the (scaled) Fisher information matrix for  $(\boldsymbol{\beta}, \boldsymbol{\alpha})$ :

$$\mathbf{M}(\xi) := \sum_{i=1}^n \gamma w(\tau_i) \begin{bmatrix} f(\mathbf{x}_i) \\ g(\mathbf{x}_i) \end{bmatrix} \begin{bmatrix} f(\mathbf{x}_i) \\ g(\mathbf{x}_i) \end{bmatrix}^T + \underbrace{\sum_{i=1}^{n_0} \gamma w(\tau_i) \begin{bmatrix} f(\mathbf{x}_i^0) \\ g(\mathbf{x}_i^0) \end{bmatrix} \begin{bmatrix} f(\mathbf{x}_i^0) \\ g(\mathbf{x}_i^0) \end{bmatrix}^T}_{\boldsymbol{\Gamma}} + \begin{bmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_s \end{bmatrix}. \quad (4)$$

Further, the IMSE criterion can be approximated by a criterion of Bayesian A-optimality:

$$\text{IMSE}(\xi) = \int_{\mathcal{X}} \tilde{\rho}(\mathbf{x}) d\mu(\mathbf{x}) = \text{trace } \mathbf{M}(\xi)^{-1} \mathbf{L}$$

with a coefficient matrix  $\mathbf{L} = \sigma_Z^2 \int_{\mathcal{X}} h(\mathbf{x}) h(\mathbf{x})^T d\mu(\mathbf{x})$ . We restrict our attention to the situation in which  $\boldsymbol{\beta}$  is estimable from the observations collected with the initial design  $\xi_{\text{init}}$ , which ensures that  $\boldsymbol{\Gamma}$  is positive definite, and  $\mathbf{M}(\xi)$  is invertible for all designs.

This technique was recently used by [SP10, GP16], who compute approximate designs by using standard algorithms of Bayesian A-optimality, and use rounding heuristics to find exact designs. One

disadvantage is that it requires the knowledge of a Mercer's expansion of the kernel. To tackle this problem, a polar spectral approximation of the kernel has been used [SP10, SP15], but it is not clear whether this can be generalized for parameter spaces of dimension  $d > 2$ . In [GP16] it is assumed that  $\mu$  has a finite support containing the candidate points  $\mathbf{x}_i$ , so the computation of the  $g(\mathbf{x}_i)$  reduce to a standard matrix eigenproblem. In Section 3 we establish a link between this approach and the class of SOR kernels commonly used in machine learning.

In this article, we focus on two classes of sparse kernels, which are commonly used in Machine Learning, and for which there is a simple, finite Mercer's expansion. We will show in Section 3 that using sparse kernels with  $s$  inducing points reduce the cost per iteration of the multiplicative weights update algorithm from  $\mathcal{O}(n^3)$  to  $\mathcal{O}(ns^2)$ , when the goal is to compute the weights of a design minimizing the IMSE over  $n$  *predefined candidate points*. This reduction is crucial for computer experiments with parameter dimension  $d \geq 4$ , where a very large number  $n$  of candidate points is required to fill the space of parameters. We show further in Section 4 that the same complexity reduction can be achieved for the search of optimal designs for the prediction of hyperparameters of the kernel. However, we point out that the optimization problems we consider are, in general, not convex. Hence, global optimality cannot be guaranteed, but one can implement standard strategies to try and escape local optima. One exception are the sparse IMSE-optimal design problems studied in Section 3, which are convex when the information function  $w(\cdot)$  is concave. Finally, some numerical experiments illustrate our method in Section 5.

The second goal of this article, covered in the next section, is to extend the theory of approximate optimal designs to the situation in which the weights  $w_i$  depend on the true design parameters  $\tau_i$  via an information function  $w(\cdot)$ .

## 2 Approximate designs in presence of an information function

For a design  $\xi = \{\mathbf{x}_i, \tau_i\}_i$ , define  $\mathbf{M}(\xi) := \sum_{i=1}^n w(\mathbf{x}_i, \tau_i) h(\mathbf{x}_i) h(\mathbf{x}_i)^T + \mathbf{\Gamma}$ . Throughout this section, we assume that we are given a set  $X = \{\mathbf{x}_1, \dots, \mathbf{x}_n\} \subseteq \mathcal{X}$  of candidate points, and we consider a design problem of the form

$$\min_{\boldsymbol{\tau} \in \Delta_T} \Phi(\xi_{\boldsymbol{\tau}}) := \text{trace } \mathbf{M}(\xi_{\boldsymbol{\tau}})^{-1} \mathbf{L}, \quad (5)$$

where  $\xi_{\boldsymbol{\tau}}$  represents the design  $\{\mathbf{x}_i, \tau_i\}_{i=1, \dots, n}$ . The matrix  $\mathbf{L}$  is positive semidefinite, and the computing times are constrained in the set  $\Delta_T := \{\boldsymbol{\tau} \in \mathbb{R}_+^n : \sum_{i=1}^n \tau_i = T\}$ . For the sake of generality, the information function  $w$  is allowed to depend on the design point  $\mathbf{x}_i$ . For all  $\mathbf{x} \in \mathcal{X}$ , the function  $w(\mathbf{x}, \cdot)$  is assumed to be continuously differentiable and nondecreasing on  $\mathbb{R}_+$ . We restrict our attention to the case of a positive definite  $\mathbf{\Gamma}$  for the sake of simplicity (so that  $\mathbf{M}(\xi_{\boldsymbol{\tau}})$  is invertible for all  $\boldsymbol{\tau} \in \Delta_T$ ), but we stress that the results presented here can be extended to the case of a positive semidefinite  $\mathbf{\Gamma}$ .

It is well known that  $\Phi(\xi_{\boldsymbol{\tau}})$  is a convex function of the vector of design weights  $\mathbf{w} = [w(\mathbf{x}_1, \tau_1), \dots, w(\mathbf{x}_n, \tau_n)]^T$ . Using the fact that  $\Phi(\xi_{\boldsymbol{\tau}})$  is a nonincreasing function of  $w(\mathbf{x}_i, \tau_i)$ , standard composition theorems yield the following:

**Proposition 2.1.** *If the information function  $\tau \mapsto w(\mathbf{x}, \tau)$  is nondecreasing and concave for all  $\mathbf{x} \in \mathcal{X}$ , then the function  $\boldsymbol{\tau} \mapsto \Phi(\xi_{\boldsymbol{\tau}})$  is convex over  $\Delta_T$ .*

The success of the theory of approximate designs is largely due to equivalence theorems such as the Kiefer-Wolfowitz theorem [KW60], that give simple means to check the optimality of a design.

First we state the Karush-Kuhn-Tucker (KKT) necessary optimality conditions for a vector of design weights  $\boldsymbol{\tau}$  over  $\Delta_T$ , which are valid even if the  $w(\mathbf{x}_i, \cdot)$  are not concave:

**Proposition 2.2.** *Let  $\mathbf{L} = \mathbf{K}\mathbf{K}^T$  and  $\mathbf{x}_1, \dots, \mathbf{x}_n$  be given candidate points in  $\mathcal{X}$ . For all  $i \in \{1, \dots, n\}$ , define*

$$d_i(\boldsymbol{\tau}) = \frac{\partial w(\mathbf{x}_i, \tau_i)}{\partial \tau_i} \|h(\mathbf{x}_i)^T \mathbf{M}(\xi_{\boldsymbol{\tau}})^{-1} \mathbf{K}\|^2.$$

If  $\boldsymbol{\tau}$  is a local minimizer of  $\Phi(\xi_{\boldsymbol{\tau}})$  over  $\Delta_T$ , then we have:  $\forall i \in \{1, \dots, n\}$ ,  $d_i(\boldsymbol{\tau}) \leq \frac{1}{T} \sum_{k=1}^n \tau_k d_k(\boldsymbol{\tau})$ . Moreover, the inequality becomes an equality for support points of  $\xi_{\boldsymbol{\tau}}$ , i.e. for all  $i$  such that  $\tau_i > 0$ .

*Proof.* Observe that  $\frac{\partial \Phi(\xi_{\boldsymbol{\tau}})}{\partial \tau_i} = -\frac{\partial w(\mathbf{x}_i, \tau_i)}{\partial \tau_i} \text{trace } \mathbf{M}(\xi_{\boldsymbol{\tau}})^{-1} h(\mathbf{x}) h(\mathbf{x})^T \mathbf{M}(\xi_{\boldsymbol{\tau}})^{-1} \mathbf{L} = -d_i(\boldsymbol{\tau})$ . Then, the dual feasibility and complementary slackness KKT-conditions of the optimization problem  $\min\{\Phi(\xi_{\boldsymbol{\tau}}) : \forall i \in \{1, \dots, n\}, \tau_i \geq 0, \sum_i \tau_i = T\}$  can be expressed as follows:

$$\exists \lambda \geq 0 : \forall i \in \{1, \dots, n\}, \quad \left( (\tau_i = 0 \text{ and } d_i(\boldsymbol{\tau}) \leq \lambda) \text{ or } (\tau_i \geq 0 \text{ and } d_i(\boldsymbol{\tau}) = \lambda) \right).$$

Moreover, the Lagrange multiplier  $\lambda$  must satisfy  $T\lambda = \sum_i \tau_i \lambda = \sum_i \tau_i d_i(\boldsymbol{\tau})$ . Substituting the value of  $\lambda$  in the KKT conditions yields the proposition.  $\square$

If the information functions are concave, we obtain a much stronger result. For the next theorem we temporarily drop the assumption that the  $\mathbf{x}_i$ 's are given. We characterize optimal designs over the set  $\Xi = \left\{ \xi = \{\mathbf{x}_i, \tau_i\}_{i=1, \dots, n} : n \in \mathbb{N}, \forall i \in \{1, \dots, n\}, \mathbf{x}_i \in \mathcal{X}, \tau_i \geq 0, \sum_i \tau_i = T \right\}$  of all designs with support points in  $\mathcal{X}$ :

**Theorem 2.3.** Let  $\mathbf{L} = \mathbf{K}\mathbf{K}^T$ , and assume that the condition of Proposition 2.1 is satisfied. For all  $\mathbf{x} \in \mathcal{X}$  define

$$d(\xi; \mathbf{x}) = \frac{\partial w(\mathbf{x}, \tau(\mathbf{x}))}{\partial \tau} \|h(\mathbf{x})^T \mathbf{M}(\xi)^{-1} \mathbf{K}\|^2,$$

where  $\tau(\mathbf{x})$  is the computing time spent on the design point  $\mathbf{x}$  (i.e.,  $\tau(\mathbf{x}) = 0$  if  $\mathbf{x} \notin \text{supp}(\xi)$  and  $\tau(\mathbf{x}) = \tau_i$  if  $\mathbf{x} = \mathbf{x}_i \in \text{supp}(\xi)$ ). Then,  $\xi^* = \{\mathbf{x}_i, \tau_i\}_{i=1, \dots, n}$  minimizes  $\Phi$  over  $\Xi$  if and only if

$$\forall \mathbf{x} \in \mathcal{X}, \quad d(\xi^*; \mathbf{x}) \leq \frac{1}{T} \sum_{i=1}^n \tau_i d(\xi^*, \mathbf{x}_i).$$

Moreover the above inequality becomes an equality for all support points of  $\xi^*$ .

*Proof.* First note that the *only if* part of the theorem is a simple consequence of Proposition 2.2. For the *if* part, assume that the condition of the theorem holds. It implies

$$\forall \boldsymbol{\tau}' \in \mathbb{R}_+^n \quad \text{such that} \quad \sum_i \tau'_i = T, \quad \sum_{i=1}^n \tau'_i d(\xi^*, \mathbf{x}_i) \leq \sum_{i=1}^n \tau_i d(\xi^*, \mathbf{x}_i).$$

Using the fact that  $\frac{\partial \Phi(\xi^*)}{\partial \tau_i} = -d(\xi^*, \mathbf{x}_i)$ , this can be rewritten as:

$$\forall \xi' = \{\mathbf{x}_i, \tau'_i\} \in \Xi, \quad \left. \frac{\partial \Phi((1-\alpha)\xi^* + \alpha\xi')}{\partial \alpha} \right|_{\alpha=0} \geq 0.$$

Finally, consider an arbitrary design  $\xi' \in \Xi$ , and define the function  $\psi : \alpha \mapsto \Phi((1-\alpha)\xi^* + \alpha\xi')$ . By proposition 2.1,  $\psi$  is convex on  $[0, 1]$ , and we know that  $\psi'(0) \geq 0$ . So we have  $\Phi(\xi') = \psi(1) \geq \psi(0) + \psi'(0)(1-0) \geq \psi(0) = \Phi(\xi^*)$ , which proves the optimality of  $\xi^*$ .  $\square$

We next adapt the multiplicative weights update algorithm of Titterton [STT78] for Problem (5). The multiplicative algorithm was originally presented in the general setting in which a function  $f$  must be minimized over a unit simplex  $\{\mathbf{w} \in \mathbb{R}_+^n : \sum_i w_i = 1\}$ , so it can be adapted in a straightforward manner to the case of a design problem with information functions, i.e.,  $w_i = w(\mathbf{x}_i, \tau_i)$ . Given an exponent  $q > 0$  and an initial vector  $\boldsymbol{\tau}^{(0)} > \mathbf{0}$ , the iterations are:

$$\forall i \in \{1, \dots, n\}, \quad \tau_i^{(k+1)} \leftarrow T \frac{\tau_i^{(k)} d_i(\boldsymbol{\tau}^{(k)})^q}{\sum_{j=1}^n \tau_j^{(k)} d_j(\boldsymbol{\tau}^{(k)})^q}. \quad (6)$$

In its standard version, that is, when  $w(\mathbf{x}_i, \tau_i) = \tau_i$ , this algorithm converges monotonically towards an  $A$ -optimal design when  $q = \frac{1}{2}$ . The process can be accelerated by pruning candidate points with a sufficiently low weight, which ensures that they do not belong to the support of any optimal design [Pro13]. Convergence for a variety of optimality criteria was shown in [Yul10].

Consider now the general setting of a function  $f$  to be minimized over  $\Delta_T$ , with  $d_i(\boldsymbol{\tau}) = \frac{\partial f(\boldsymbol{\tau})}{\partial \tau_i}$ . If the iterations (6) converge, then the limit point  $\boldsymbol{\tau}^*$  must satisfy the necessary condition of Proposition 2.2, under some mild conditions [GM92]. In practice, we experienced numerical convergence of the above algorithm towards *local minima* of  $f : \boldsymbol{\tau} \mapsto \Phi(\xi_{\boldsymbol{\tau}})$  when  $q$  is well chosen, even in the cases where the information functions  $w(\mathbf{x}_i, \cdot)$  are not concave.

### 3 Sparse covariance kernels

Here we consider two classes of sparse kernel functions commonly used in machine learning for Gaussian process regression with a large number  $n$  of samples. These approximations rely on a small set of *inducing points*,  $\{\mathbf{u}_1, \dots, \mathbf{u}_s\} \subset \mathcal{X}$ , and assume that the covariance  $\text{cov}(Z(\mathbf{x}), Z(\mathbf{y}))$  of the process between the points  $\mathbf{x}$  and  $\mathbf{y} \in \mathcal{X}$  only depends on the covariances between the  $\mathbf{u}_i$ ,  $\mathbf{x}$  and the  $\mathbf{u}_i$ , and  $\mathbf{y}$  and the  $\mathbf{u}_i$ . This reduces the complexity of training a Gaussian process on a dataset with  $n$  samples from  $\mathcal{O}(n^3)$  to  $\mathcal{O}(ns^2)$ . We refer to [QCR05] for a comprehensive review.

**SOR-kernels.** The *Subset of Regressors* (SOR) approximation consists in replacing the correlation function  $C(\mathbf{x}, \mathbf{y})$  by a low-rank kernel,

$$C_{\text{SOR}}(\mathbf{x}, \mathbf{y}) = \mathbf{c}_u(\mathbf{x})^T \mathbf{K}_{u,u}^{-1} \mathbf{c}_u(\mathbf{y}),$$

where  $\{\mathbf{K}_{u,u}\}_{i,j} = C(\mathbf{u}_i, \mathbf{u}_j)$  is the  $s \times s$  matrix of correlations between the  $Z(\mathbf{u}_i)$ , and  $\{\mathbf{c}_u(\mathbf{x})\}_i = C(\mathbf{u}_i, \mathbf{x})$  is the  $s$ -dimensional vector of correlations between  $Z(\mathbf{x})$  and the  $Z(\mathbf{u}_i)$ . Hence, if we let  $\mathbf{J}_u$  be any matrix satisfying  $\mathbf{J}_u \mathbf{J}_u^T = \mathbf{K}_{u,u}^{-1}$ , then the function  $g : \mathbf{x} \mapsto \mathbf{J}_u^T \mathbf{c}_u(\mathbf{x})$  satisfies

$$\forall (\mathbf{x}, \mathbf{y}) \in \mathcal{X} \times \mathcal{X}, \quad C_{\text{SOR}}(\mathbf{x}, \mathbf{y}) = g(\mathbf{x})^T g(\mathbf{y}).$$

Hence, for a SOR-kernel the observation model (1) is equivalent to model (3). Indeed,  $Z(\xi) = [Z(\mathbf{x}_1), \dots, Z(\mathbf{x}_n)]^T \sim \mathcal{N}(0, \sigma_Z^2 \mathbf{C})$  has the same distribution as  $[g(\mathbf{x}_1), \dots, g(\mathbf{x}_n)]^T \boldsymbol{\alpha}$ , where  $\boldsymbol{\alpha} \sim \mathcal{N}(0, \sigma_Z^2 \mathbf{I}_s)$ . To put it in other words, if we assume that the true kernel is  $C_{\text{SOR}}$ , then model (3) is exact, so that  $\text{IMSE}(\xi_{\boldsymbol{\tau}}) = \text{IMSE}(\xi_{\boldsymbol{\tau}}) = \text{trace} \mathbf{M}(\xi_{\boldsymbol{\tau}})^{-1} \mathbf{L}$ , and we can use the multiplicative weights update (6) to compute an optimal design. The complexity of computing  $\text{IMSE}(\xi_{\boldsymbol{\tau}})$  and its gradient  $[d_1(\boldsymbol{\tau}), \dots, d_n(\boldsymbol{\tau})]^T$  is  $\mathcal{O}(n(s+m)^2)$ , which is  $\mathcal{O}(ns^2)$  because the dimension  $m$  of the regression parameter  $\boldsymbol{\beta}$  is a small constant. So the cost of one iteration (6) is  $\mathcal{O}(ns^2)$ . In contrast, for a full kernel the computation involves  $\boldsymbol{\Sigma}^{-1}$  and takes  $\mathcal{O}(n^3)$  operations.

There is a vast literature on the selection of inducing points  $\mathbf{u}_i$  to approximate a kernel  $C(\cdot, \cdot)$  by a SOR kernel [Tit09, CBFH15, WN15]. For example, [WN15] uses a regular grid to exploit the Kronecker structure of  $\mathbf{K}_{u,u}$  when  $C$  is a product of one-dimensional kernels, and to speed-up the computations by using fast Fourier transforms. Note that if the points  $\mathbf{u}_1, \dots, \mathbf{u}_s$  are sampled randomly and independently from the probability measure  $\mu$ , the approximation  $C(\mathbf{x}, \mathbf{y}) \simeq C_{\text{SOR}}(\mathbf{x}, \mathbf{y})$  can be regarded as an expansion of the form  $C(\mathbf{x}, \mathbf{y}) = \sum_{i=1}^s \lambda_i \phi_i(\mathbf{x}) \phi_i(\mathbf{y})$ , where the  $\lambda_i$  and  $\phi_i(\cdot)$  are the solutions of the Nyström approximation of the eigenproblem

$$\int_{\mathcal{X}} C(\mathbf{x}, \mathbf{y}) \phi(\mathbf{y}) d\mu(\mathbf{y}) = \lambda \phi(\mathbf{x}).$$



This approximation consists in replacing the integral by  $\frac{1}{s} \sum_{i=1}^s C(\mathbf{x}, \mathbf{u}_i) \phi(\mathbf{u}_i)$ , and reduces the infinite-dimensional eigenproblem to a standard  $s \times s$ -matrix eigenproblem [WS01]. If we choose  $\mathbf{J}_u = \mathbf{U} \mathbf{\Lambda}^{-\frac{1}{2}}$ , where  $\mathbf{U} \mathbf{\Lambda} \mathbf{U}^T$  is a spectral decomposition of  $\mathbf{K}_{u,u}$ , this is equivalent to the approach of [GP14, GP16], where  $\mu$  is approximated by a discrete measure  $\hat{\mu}$  supported by the  $\mathbf{u}_i$ 's (i.e., the IMSE is approximated by a quadrature). In [GP14], the authors further suggest to choose the candidates  $\mathbf{x}_i$  in the support of  $\hat{\mu}$ , i.e.  $n \leq s$ . Then,  $g$  must only be evaluated at the  $\mathbf{u}_i$ 's, and the vectors  $g(\mathbf{u}_i)$  ( $i = 1, \dots, s$ ) are the columns of  $\mathbf{\Lambda}^{\frac{1}{2}} \mathbf{U}^T$ , hence they are orthogonal. This contrasts with our approach, where we generate a large number of candidate points in order to fill the design space, but use a small number of inducing points for the sake of computation ( $n \gg s$ ).

**FITC-kernels.** The FITC approximation (*Fully Independent Training Conditional*) is very similar to SOR, but a diagonal noise is added to the kernel, to ensure that  $C_{\text{FITC}}(\mathbf{x}, \mathbf{x}) = C(\mathbf{x}, \mathbf{x}) = 1$ :

$$C_{\text{FITC}}(\mathbf{x}_i, \mathbf{x}_j) = C_{\text{SOR}}(\mathbf{x}_i, \mathbf{x}_j) + (1 - C_{\text{SOR}}(\mathbf{x}_i, \mathbf{x}_i)) \delta_{ij},$$

If we define as before the function  $g : \mathbf{x} \mapsto \mathbf{J}_u^T c_u(\mathbf{x})$ , where  $\mathbf{J}_u \mathbf{J}_u^T = \mathbf{K}_{u,u}^{-1}$ , we obtain:

$$C_{\text{FITC}}(\mathbf{x}_i, \mathbf{x}_j) = g(\mathbf{x}_i)^T g(\mathbf{x}_j) + (1 - \|g(\mathbf{x}_i)\|^2) \delta_{ij}.$$

It follows that for a FITC kernel, the observation model is equivalent to

$$Y(\mathbf{x}, \tau) = [f(\mathbf{x})^T, g(\mathbf{x})^T] \begin{bmatrix} \boldsymbol{\beta} \\ \boldsymbol{\alpha} \end{bmatrix} + \nu(\mathbf{x}) + \epsilon(\mathbf{x}, \tau), \quad (7)$$

where  $\boldsymbol{\alpha}$  is a random regression parameter with prior  $\boldsymbol{\alpha} \sim \mathcal{N}(0, \sigma_Z^2 \mathbf{I}_k)$ , and  $\nu(\mathbf{x})$  is an unbiased and uncorrelated noise, which is heteroschedastic with  $\mathbb{V}[\nu(\mathbf{x})] = \sigma_Z^2 (1 - \|g(\mathbf{x})\|^2)$ . Under this model, the Fisher information matrix for  $(\boldsymbol{\beta}, \boldsymbol{\alpha})$  becomes (up to a scaling factor  $\sigma_Z^2$ ):

$$\mathbf{M}(\xi) := \sum_{i=1}^n w(\mathbf{x}_i, \tau_i) h(\mathbf{x}_i) h(\mathbf{x}_i)^T + \mathbf{\Gamma}, \quad \text{where } w(\mathbf{x}, \tau) := \frac{\gamma w(\tau)}{1 + (1 - \|g(\mathbf{x})\|^2) \gamma w(\tau)}. \quad (8)$$

Note that  $\mathbf{M}(\xi)$  has the form of the Fisher information matrix of the problem studied in Section 2. Moreover, elementary calculus shows that  $\tau \mapsto w(\mathbf{x}, \tau)$  is concave if  $\tau \mapsto w(\tau)$  is concave. In this situation, Proposition 2.1 shows that (5) is a convex optimization problem for FITC kernels.

## 4 Optimal designs for the estimation of kernel hyperparameters

Until now, we have assumed that the kernel function  $C(\cdot, \cdot)$  was known. In practice however, the kernel depends on a set of hyperparameters  $\boldsymbol{\theta} \in \mathbb{R}^p$ , which must be estimated by maximum likelihood from the set of observations  $Y(\xi)$ . Recall that  $Y(\xi) \sim \mathcal{N}(\mathbf{F}^T \boldsymbol{\beta}, \sigma_Z^2 \boldsymbol{\Sigma}_\theta)$ , where  $\boldsymbol{\Sigma}_\theta = \mathbf{C}_\theta + \mathbf{D}^{-1}$ ,  $\mathbf{D} = \text{Diag}(\gamma w(\boldsymbol{\tau}))$ , and we have inserted the symbol  $\theta$  as subscript to stress the dependency on the hyperparameters. Then, the  $p \times p$  Fisher information matrix for the vector of parameters  $\boldsymbol{\theta}$  can be derived from standard formulas:

$$\{\mathbf{M}_\theta(\xi)\}_{ij} = \frac{1}{2} \text{trace } \boldsymbol{\Sigma}_\theta^{-1} \frac{\partial \mathbf{C}_\theta}{\partial \theta_i} \boldsymbol{\Sigma}_\theta^{-1} \frac{\partial \mathbf{C}_\theta}{\partial \theta_j}. \quad (9)$$

Given a current estimate of  $\boldsymbol{\theta}$ , we propose to search a design  $\xi$  maximizing the criterion of  $D$ -optimality,  $\log \det \mathbf{M}_\theta(\xi)$ . Here, note that we assume that  $\boldsymbol{\beta}$  and  $\sigma_Z^2$  are known. We refer to [PM12] for a review on approaches to deal with a total Fisher information matrix for the set of parameters  $(\boldsymbol{\beta}, \sigma_Z^2, \boldsymbol{\theta})$ . In particular, Müller and Stehlík proposed a compound criterion with a weighing factor that balances the goals of estimating  $\boldsymbol{\beta}$  and estimating  $\boldsymbol{\theta}$  [MS10].

We want to optimize the computing times  $\tau_i$  associated with a large number of candidate points  $\mathbf{x}_i$ . This is a hard optimization problem, since here the  $D$ -criterion is not convex with respect to  $\boldsymbol{\tau}$ . Nevertheless we propose to use the multiplicative update iterations (6), where  $d_i(\boldsymbol{\tau}) := \frac{\partial \log \det \mathbf{M}_\theta(\boldsymbol{\xi}_\boldsymbol{\tau})}{\partial \tau_i}$ , in order to identify good local optima. However, if  $n$  is very large, the computation of  $\mathbf{M}_\theta(\boldsymbol{\xi}_\boldsymbol{\tau})$  and  $d_i(\boldsymbol{\tau})$  is extremely time-consuming. It involves the inversion of the  $n \times n$  matrix  $\boldsymbol{\Sigma}_\theta$ , and many  $n \times n$  matrix-matrix multiplications.

In this section we show that this computational burden can be reduced if sparse kernels are used. From now on, we assume that  $\mathbf{C}_\theta = \mathbf{G}_\theta \mathbf{G}_\theta^T$ , where  $\mathbf{G}_\theta = [g_\theta(\mathbf{x}_1), \dots, g_\theta(\mathbf{x}_n)]^T \in \mathbb{R}^{n \times s}$ . As in previous section, the function  $g_\theta$  is defined by  $g_\theta(\mathbf{x}) = \mathbf{J}_u^T \mathbf{c}_u(\mathbf{x})$ , where  $\mathbf{J}_u \mathbf{J}_u^T = \mathbf{K}_{u,u}^{-1}$ . From now on we set  $\mathbf{J}_u$  to the Cholesky factor of  $\mathbf{K}_{u,u}^{-1}$ , because this choice yields compact formulas.

First note that the low-rank decomposition makes it possible to use the Woodbury matrix identity:

$$\boldsymbol{\Sigma}_\theta^{-1} = (\mathbf{C}_\theta + \mathbf{D}^{-1})^{-1} = \mathbf{D} - \mathbf{D} \mathbf{G}_\theta (\mathbf{I}_s + \mathbf{G}_\theta^T \mathbf{D} \mathbf{G}_\theta)^{-1} \mathbf{G}_\theta^T \mathbf{D}. \quad (10)$$

Then, we also need to compute derivatives of  $g_\theta$  with respect to  $\theta$ . This is possible thanks to the following lemma:

**Lemma 4.1.** *Define the function  $\Phi$  which returns the lower triangle and half the diagonal of a square matrix:*

$$\forall \mathbf{M} \in \mathbb{R}^{n \times n}, \quad \{\Phi(\mathbf{M})\}_{ij} = \begin{cases} M_{ij} & \text{if } i > j \\ \frac{1}{2} M_{ij} & \text{if } i = j \\ 0 & \text{if } i < j. \end{cases}$$

Then, we have:  $\forall \mathbf{x} \in \mathcal{X}$ ,

$$\frac{\partial g_\theta(\mathbf{x})}{\partial \theta_i} = \mathbf{J}_u^T \frac{\partial \mathbf{c}_u(\mathbf{x})}{\partial \theta_i} - \Phi \left( \mathbf{J}_u^T \frac{\partial \mathbf{K}_{u,u}}{\partial \theta_i} \mathbf{J}_u \right)^T \mathbf{J}_u^T \mathbf{c}_u(\mathbf{x}).$$

*Proof.* A formula for the derivative of the Cholesky decomposition  $\mathbf{X} = \mathbf{J} \mathbf{J}^T$  can be found in [Sär13, Theorem A.1], and can be proved by implicit differentiation:

$$\frac{\partial \mathbf{J}}{\partial \theta} = \mathbf{J} \Phi(\mathbf{J}^{-1} \frac{\partial \mathbf{X}}{\partial \theta} \mathbf{J}^{-T}).$$

The formula of the lemma can now be obtained, by applying standard formulas for the differentiation of products and matrix inverse.  $\square$

We can use this lemma to compute the matrices  $\mathbf{G}_i := \frac{\partial \mathbf{G}_\theta}{\partial \theta_i}$ . Now, we also define  $\mathbf{G}_0 := \mathbf{G}_\theta$  to simplify the notation. Substituting  $\frac{\partial \mathbf{C}_\theta}{\partial \theta_i} = \mathbf{G}_i \mathbf{G}_0^T + \mathbf{G}_0 \mathbf{G}_i^T$  and (10) into (9) yields an expression for  $\{\mathbf{M}_\theta(\boldsymbol{\xi})\}_{ij}$  that depends on  $\mathbf{G}_0, \mathbf{G}_1, \dots, \mathbf{G}_p$ . After some simplifications, we obtain

$$\{\mathbf{M}_\theta(\boldsymbol{\xi})\}_{ij} = \text{trace } \mathbf{A}_{0i} \mathbf{A}_{0j} + \mathbf{A}_{00} \mathbf{A}_{ij},$$

where for all  $k, l \in \{0, \dots, p\}$ ,  $\mathbf{A}_{kl} := \mathbf{B}_{kl} - \mathbf{B}_{k0}(\mathbf{I}_s + \mathbf{B}_{00})^{-1} \mathbf{B}_{0l}$  and  $\mathbf{B}_{kl} := \mathbf{G}_k^T \mathbf{D} \mathbf{G}_l$ . From these expressions, it is easy to see that  $\mathbf{M}_\theta(\boldsymbol{\xi})$  can be computed in  $\mathcal{O}(ns^2)$ , which is a great improvement compared to  $\mathcal{O}(n^3)$  for a full kernel.

Similarly, we can compute  $\nabla_\tau \{\mathbf{M}_\theta(\boldsymbol{\xi})\}_{ij} = \left[ \frac{\partial \{\mathbf{M}_\theta(\boldsymbol{\xi})\}_{ij}}{\partial \tau_1}, \dots, \frac{\partial \{\mathbf{M}_\theta(\boldsymbol{\xi})\}_{ij}}{\partial \tau_n} \right]^T$  in  $\mathcal{O}(ns^2)$ . For all  $k \in \{0, \dots, p\}$ , define  $\mathbf{P}_k := \mathbf{G}_k - \mathbf{G}_0(\mathbf{I}_s + \mathbf{B}_{00})^{-1} \mathbf{B}_{0k}$ . Then, we can show that (details omitted):

$$\nabla_\tau \{\mathbf{M}_\theta(\boldsymbol{\xi})\}_{ij} = \gamma \text{Diag}(w'(\boldsymbol{\tau})) \text{Diag} \left( \mathbf{P}_i \mathbf{A}_{0j} \mathbf{P}_0^T + \mathbf{P}_j \mathbf{A}_{0i} \mathbf{P}_0^T + \mathbf{P}_0 \mathbf{A}_{ij} \mathbf{P}_0^T + \mathbf{P}_j \mathbf{A}_{00} \mathbf{P}_i^T \right).$$

Finally, the gradient of the criterion is obtained by  $d_i(\boldsymbol{\xi}) = \text{trace } \mathbf{M}_\theta(\boldsymbol{\xi})^{-1} \frac{\partial \mathbf{M}_\theta(\boldsymbol{\xi})}{\partial \tau_i}$ . Hence, we have shown that the gradient of the criterion can be computed in  $\mathcal{O}(ns^2)$  for a sparse kernel with  $s$  inducing points. In contrast, for a full kernel one requires  $\mathcal{O}(n^3)$  operations.

## 5 Numerical Experiments

We consider the Ishigami-like function  $\eta$  to illustrate the effect of the information functions:

$$\forall \mathbf{x} \in \mathcal{X} = [0, 1] \times [0, 1], \quad \eta(\mathbf{x}) = 1.1 \sin(\pi(2x_1 - 1)) + 7 \sin^2(\pi(2x_2 - 1)).$$

First, 8 noisy observations of the function  $\eta(\mathbf{x})$  are taken, with  $\sigma_N^2 = 0.05$ , and  $\tau_{init} = \frac{1}{8}$  at each of the 8 locations indicated by yellow dots in Figure 1. These initial values are used to estimate, by maximum likelihood,  $\sigma_Z^2$  and the hyperparameters  $(\ell_1, \ell_2)$  of the Gaussian kernel

$$C(\mathbf{x}, \mathbf{y}) = e^{-\frac{1}{2} \left[ \left( \frac{x_1 - y_1}{\ell_1} \right)^2 + \left( \frac{x_2 - y_2}{\ell_2} \right)^2 \right]}$$

The plots of Figure 1 show some designs for the distribution of  $T = 1$  additional hour of computing time over a regular grid of  $n = 31^2 = 961$  candidate points. The size of the red dots indicate the time to spend on a design point, and the color in the background indicates the prior Kriging variance (after the initial 8 observations; blue: small variance, red: high variance), according to the considered covariance model: In Plot (a) and (b), we respectively used the SOR and FITC Kernel associated with  $C(\cdot, \cdot)$ , for  $s = 12$  inducing points (marked with black squares), that were generated in a space-filling fashion with a Sobol sequence; the number  $s = 12$  is rather small, on purpose, to illustrate the effect of sparsity. Plots (c)-(g) rely on a SOR kernel with  $s = 60$  inducing points (not marked for the sake of visibility); with that many inducing points, the relative errors between  $\text{IMSE}(\xi)$  and  $\tilde{\text{IMSE}}(\xi)$  were in the order of 0,1% for the designs we computed. Also, different information functions were used, cf. Plot (h).

The plots (a),(b),(c),(e),(g) show (near-)optimal weights  $\tau_i$  for the IMSE criterion at the specified  $n$  locations of the  $\mathbf{x}_i$ 's, while (d),(f) are nearly  $D$ -optimal weights for the estimation of  $\boldsymbol{\theta} = (\ell_1, \ell_2)$ . For all computations, the matrix  $\mathbf{L} = \sigma_Z^2 \int_{\mathcal{X}} h(\mathbf{x}) h(\mathbf{x})^T d\mu(\mathbf{x})$  was computed with a Monte-Carlo method with  $N = 10^5$  samples, with  $\mu$  the Lebesgue measure over  $[0, 1]^2$ . The stopping criterion for the multiplicative update iterations was

$$\max_{i=1, \dots, n} d_i(\boldsymbol{\tau}) \leq \frac{1}{T} \sum_{k=1}^n \tau_k d_k(\boldsymbol{\tau}) + \varepsilon, \quad (11)$$

where  $\varepsilon = 10^{-9}$ . Note that the design weights plotted in (a)-(c) and (g) are provably optimal (up to the tolerance  $\varepsilon$ ), because the considered optimization problems are convex. This is not the case for the designs shown in Plots (d), (e) and (f). Here, the multiplicative update algorithm is likely to fall in local optima, so we performed several restarts and kept the best design.

Plots (a) and (b) show the effect of using a sparse kernel, and are to be compared with Plot (c), which can be considered as the optimal design for the full kernel  $C$  when  $w(\tau) = \tau$ . Observe that the kriging variance tends to be underestimated with the SOR kernel (a), while it is overestimated with the FITC kernel (b). As a consequence of the (strict) concavity of  $w(\mathbf{x}, \tau)$ , the FITC design is more spread out than the SOR design. Also, the FITC design has a slightly better efficiency than the SOR design (73% vs. 69%, cf. Formula (12)).

The effect of the information function can be seen by comparing the second column (standard case  $w(\tau) = \tau$ ) to the third and fourth columns. In the second column, the information function is convex near 0, so that we need some minimal amount of computing time to get some information; see blue curve in plot (f). As a consequence, the IMSE-optimal design for this information function are very sparse, a feature that can be very valuable for the experimenter. In contrast, a concave information function is used in the third column (green curve in plot (f)). This incentivizes designs with many design points spread out over regions with a high variance.

Next, we show some results in higher dimensional spaces to illustrate the importance of using sparse approximations of the kernel. We report results for tests in dimension  $d = 4$  and  $d = 7$ . In each case,

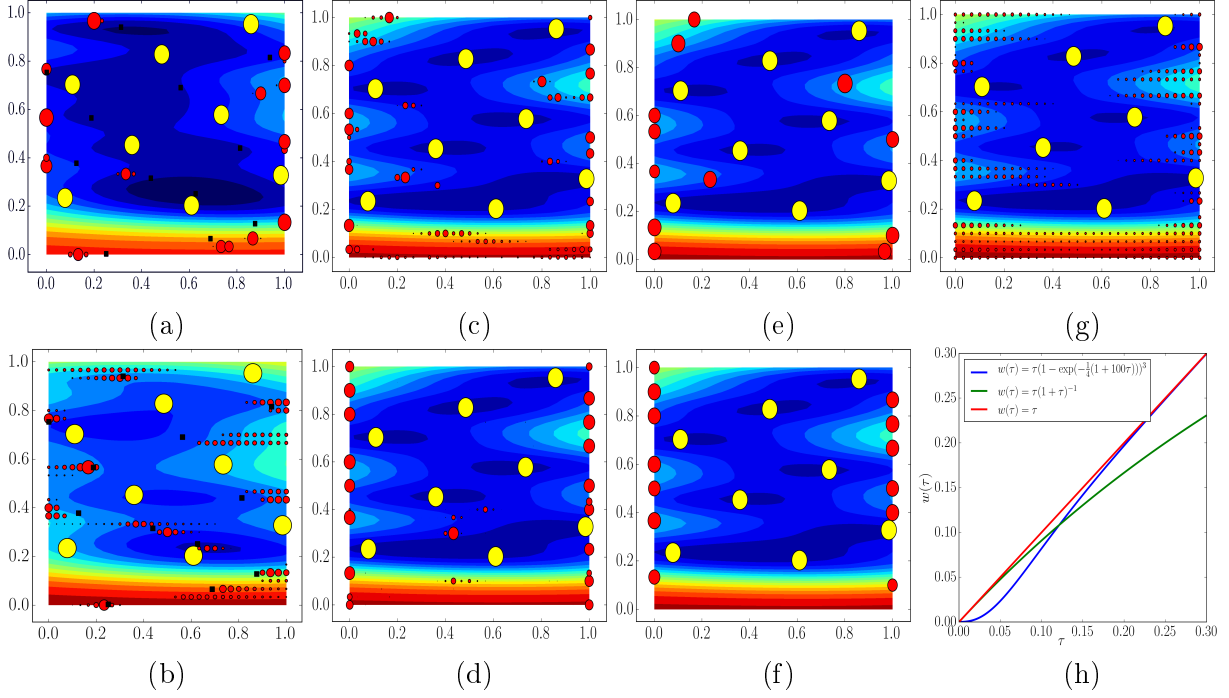


Figure 1: Near-optimal design weights for the test function. Kernel:  $C_{SOR}$  with  $s = 12$  inducing points in (a);  $C_{FITC}$  with  $s = 12$  inducing points in (b), and  $C_{SOR}$  with  $s = 60$  inducing points in (c)-(g). Optimality criterion: IMSE in (a),(b),(c),(e),(g) and D-criterion for  $\theta = (\ell_1, \ell_2)$  in (d),(f). Information function:  $w(\tau) = \tau$  in (a)-(d);  $w(\tau) = \tau(1 - \exp(-\frac{1}{4}(1 + 100\tau)))^3$  in (e),(f);  $w(\tau) = \tau(1 + \tau)^{-1}$  in (g). These information functions are plotted in (h).

we considered 10 instances, corresponding to different functions  $\eta(\cdot)$ ; These functions were dummy rational functions, in which we have selected the coefficients at random.

Our experiments used  $T_{\text{init}} = 1$  hour of computing time distributed uniformly over  $n_{\text{init}} = 50$  initial observations, and aimed at distributing  $T = 1$  additional hour of computing time over  $n = 1500$  randomly generated candidate points for the problems in dimension 4, and  $n = 5300$  points for the problem in dimension 7. The function  $w(\cdot)$  was set to the identity:  $w(\tau) = \tau$ .

For each problem, we have computed the *true* optimal design  $\xi^*$  (within the subset of designs supported over the given candidate points), by using the multiplicative update iterations (6) with a formula for the derivative of the true criterion:  $d_i(\xi) = \frac{\partial \text{IMSE}(\xi)}{\partial \tau_i}$ . The efficiency of a design was evaluated by the following formula:

$$\text{efficiency}(\xi) = \frac{\text{IMSE}(\xi)^{-1} - \text{IMSE}(\xi_{\text{init}})^{-1}}{\text{IMSE}(\xi^*)^{-1} - \text{IMSE}(\xi_{\text{init}})^{-1}}. \quad (12)$$

Here,  $\xi_{\text{init}}$  denotes the initial design supported by the  $n_{\text{init}}$  initial observation points, so the numerator expresses the gain of information provided by  $\xi$ , compared to the situation where no additional measurement is performed. Figure 2 shows the efficiency of designs computed by using a SOR approximation of the kernel, for 10 instances with  $d = 4$  (left) and  $d = 7$  (right). In both cases, we observe an excellent efficiency when  $s \geq 100$ , and even for  $s \geq 70$  for the instances in a 4-dimensional space. In terms of computing time, the speed-up was on the order of x200 on average for  $s = 70$  and  $d = 4$ ,

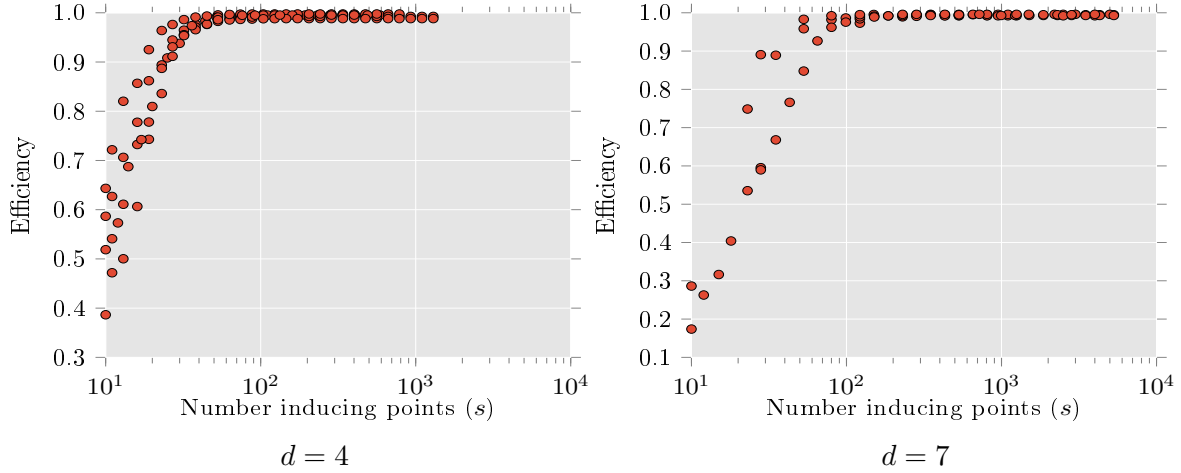


Figure 2: Efficiency of IMSE-optimal design, for an approximation of the kernel based on a SOR kernel with  $s$  inducing points.

and even of x350 for  $s = 100$  and  $d = 7$ . For the latter instances, the computations took more than 36 hours on a Linux PC with 8 cores at 3.60 GHz, while with the SOR kernel a solution was found within 6 minutes (with a tolerance parameter  $\varepsilon$  set to  $10^{-4}$  in the stopping criterion (11)).

**Acknowledgments.** The authors want to thank the two anonymous reviewers for their valuable comments and suggestions that significantly improved the presentation.

## References

- [BGL<sup>+</sup>12] J. Bect, D. Ginsbourger, L. Li, V. Picheny, and E. Vazquez. Sequential design of computer experiments for the estimation of a probability of failure. *Statistics and Computing*, 22(3):773–793, 2012.
- [CBFH15] Y. Cao, M.A. Brubaker, D.J. Fleet, and A. Hertzmann. Efficient optimization for sparse gaussian process regression. *IEEE transactions on Pattern Analysis and Machine Intelligence*, 37(12):2415–2427, 2015.
- [DW12] P. Deuffhard and M. Weiser. *Adaptive numerical solution of PDEs*. Walter de Gruyter, 2012.
- [Fed96] V. Fedorov. Design of spatial experiments: Model fitting and prediction. *Handbook of Statistics*, 13:515–553, 1996.
- [FF97] V.V. Fedorov and D. Flanagan. Optimal monitoring network design based on Mercer’s expansion of covariance kernel. *Journal of Combinatorics, Information and System Sciences*, 23:237–250, 1997.
- [GM92] N. Gaffke and R. Mathar. On a class of algorithms from experimental design theory. *Optimization*, 24(1-2):91–126, 1992.
- [GP14] B. Gauthier and L. Pronzato. Spectral approximation of the IMSE criterion for optimal designs in kernel-based interpolation models. *SIAM/ASA Journal on Uncertainty Quantification*, 2(1):805–825, 2014.

- [GP16] B. Gauthier and L. Pronzato. Optimal design for prediction in random field models via covariance kernel expansions. In *Proceedings of the 11th Workshop on Model-Oriented Data Analysis and Optimum Designs (mODa'11)*, 2016.
- [JSW98] D.R. Jones, M. Schonlau, and W.J. Welch. Efficient global optimization of expensive black-box functions. *Journal of Global optimization*, 13(4):455–492, 1998.
- [KW60] J. Kiefer and J. Wolfowitz. The equivalence of two extremum problems. *Canadian Journal of Mathematics*, 12:363–366, 1960.
- [MS10] W.G. Müller and M. Stehlík. Compound optimal spatial designs. *Environmetrics*, 21(3-4):354–364, 2010.
- [Páz86] A. Pázman. *Foundations of optimum experimental design*. D. Reidel Dordrecht, 1986.
- [Pil91] J. Pilz. *Bayesian estimation and experimental design in linear regression models*, volume 212. John Wiley & Sons Inc, 1991.
- [PM12] L. Pronzato and W.G. Müller. Design of computer experiments: space filling and beyond. *Statistics and Computing*, 22(3):681–701, 2012.
- [Pro13] L. Pronzato. A delimitation of the support of optimal designs for Kiefer’s  $\phi_p$ -class of criteria. *Statistics & Probability Letters*, 83(12):2721–2728, 2013.
- [Puk93] F. Pukelsheim. *Optimal Design of Experiments*. Wiley, 1993.
- [QCR05] J. Quiñonero-Candela and C.E. Rasmussen. A unifying view of sparse approximate gaussian process regression. *The Journal of Machine Learning Research*, 6:1939–1959, 2005.
- [Sär13] S. Särkkä. *Bayesian filtering and smoothing*, volume 3. Cambridge University Press, 2013.
- [SP10] G. Spöck and J. Pilz. Spatial sampling design and covariance-robust minimax prediction based on convex design ideas. *Stochastic Environmental Research and Risk Assessment*, 24(3):463–482, 2010.
- [SP15] G. Spöck and J. Pilz. Incorporating covariance estimation uncertainty in spatial sampling design for prediction with trans-gaussian random fields. *Frontiers in Environmental Science*, 3(39), 2015.
- [STT78] S.D. Silvey, D.M. Titterton, and B. Torsney. An algorithm for optimal designs on a finite design space. *Communications in Statistics - Theory and Methods*, 7(14):1379–1389, 1978.
- [Tit09] M.K. Titsias. Variational learning of inducing variables in sparse gaussian processes. In *International Conference on Artificial Intelligence and Statistics*, pages 567–574, 2009.
- [WN15] A.G. Wilson and H. Nickisch. Kernel interpolation for scalable structured gaussian processes (KISS-GP). In *Proceedings of The 32nd International Conference on Machine Learning*, volume 37 of *JMLR: W&CP*, pages 1775–1784, Lille, France, 2015.
- [WS01] C. Williams and M. Seeger. Using the Nyström method to speed up kernel machines. In *Proceedings of the 14th Annual Conference on Neural Information Processing Systems*, number EPFL-CONF-161322, pages 682–688, 2001.
- [Yu10] Y. Yu. Monotonic convergence of a general algorithm for computing optimal designs. *The Annals of Statistics*, 38(3):1593–1606, 2010.