



HAL
open science

Validation design I: construction of validation designs via kernel herding

Luc Pronzato, Maria-João Rendas

► **To cite this version:**

Luc Pronzato, Maria-João Rendas. Validation design I: construction of validation designs via kernel herding. 2020. hal-03474805

HAL Id: hal-03474805

<https://hal.science/hal-03474805v1>

Preprint submitted on 10 Dec 2021

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Validation design I: construction of validation designs via kernel herding

Luc Pronzato & Maria-João Rendas

CNRS, Université Côte d’Azur, Laboratoire I3S
Bât. Euclide, Les Algorithmes, 2000 route des lucioles,
06900 Sophia Antipolis cedex, France
{luc.pronzato,rendas}@univ-cotedazur.fr

December 10, 2021

Abstract

We construct validation designs \mathbf{Z}_m aimed at estimating the integrated squared prediction error of a given design \mathbf{X}_n . Our approach is based on the minimization of a maximum mean discrepancy for a particular kernel, conditional on \mathbf{X}_n , so that sequences of nested validation designs can be constructed incrementally by kernel herding. Numerical experiments show that key features for a good validation design are its space-filling properties, in order to fill the holes left by \mathbf{X}_n and properly explore the whole design space, and the suitable weighting of its points, since evaluations far from \mathbf{X}_n tend to overestimate the global error. A dedicated weighting method, based on a particular kernel, is proposed. Numerical simulations with random functions show the superiority the method over more traditional validation based on random designs, low-discrepancy sequences, or leave-one-out cross validation.

keywords validation; design of experiments; computer experiments; discrepancy; space-filling design; greedy algorithm

1 Introduction and motivation

This paper proposes methods to define designs enabling good estimation of the prediction performance of a given non-parametric model, which has been adjusted to a known training dataset. More precisely, we suppose that a design $\mathbf{X}_n = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ with n points in $\mathcal{X} = [0, 1]^d$ has been used to build a predictor of the value of an unknown function f on \mathcal{X} . We denote by $\mathbf{y}_n = [f(\mathbf{x}_1), \dots, f(\mathbf{x}_n)]^\top$ the vector collecting the n evaluations of f at the \mathbf{x}_i and by $\eta_n(\mathbf{x}) = \eta_{[\mathbf{X}_n, \mathbf{y}_n]}(\mathbf{x})$ the corresponding prediction of $f(\mathbf{x})$. The Integrated Squared Error (ISE) over \mathcal{X} is then

$$\text{ISE}(\mathbf{X}_n) = \int_{\mathcal{X}} [\eta_n(\mathbf{x}) - f(\mathbf{x})]^2 \mu(d\mathbf{x}). \quad (1)$$

Above, the measure μ codes the user preferences, penalizing regions of \mathcal{X} which are of particular interest or importance. In the paper we will always consider that μ is the Lebesgue measure on \mathcal{X} , the extension to non-uniform μ requiring only minor modifications. Note that we slightly

abuse notation here, as the dependency of $\text{ISE}(\mathbf{X}_n)$ on \mathbf{X}_n is hidden in $\eta_n(\cdot)$, which is adapted to the training set \mathbf{X}_n . The same shortcut is used throughout the paper.

In practice, the integral in definition (1) is approximated by a discrete sum, which is equivalent to letting μ be a discrete measure with finite support $\mathbf{Z}_m = \{\mathbf{z}_1, \dots, \mathbf{z}_m\} \subset \mathcal{X}$, at which η_n and f are effectively evaluated. The objective of the paper is to propose methods for the construction of *validation designs* \mathbf{Z}_m and investigate the properties of the corresponding estimates of $\text{ISE}(\mathbf{X}_n)$ ¹.

If f were known, we could compute both $\text{ISE}(\mathbf{X}_n)$ and its finite approximation

$$\widehat{\text{ISE}}(\mathbf{Z}_m, \mathbf{X}_n) = \frac{1}{m} \sum_{i=1}^m [\eta_n(\mathbf{z}_i) - f(\mathbf{z}_i)]^2,$$

and directly choose \mathbf{Z}_m to have $\widehat{\text{ISE}}(\mathbf{Z}_m, \mathbf{X}_n) \approx \text{ISE}(\mathbf{X}_n)$ (even if selecting such m points \mathbf{Z}_m would not be an easy task). However, f is unknown and both $\text{ISE}(\mathbf{X}_n)$ and $\widehat{\text{ISE}}(\mathbf{Z}_m, \mathbf{X}_n)$ can only be estimated. To do that, we shall adopt the kriging framework, which we briefly recall in Section 2.

In our study, we consider that the design \mathbf{X}_n is given, making no assumption on how it has been chosen². We are interested in particular in situations where m is not specified in advance and one wishes to construct an increasing sequence of imbedded designs $\mathbf{Z}_k \subset \mathbf{Z}_{k+1} \subset \mathbf{Z}_{k+2} \subset \dots$ such that the \mathbf{Z}_k have increasingly good performance as estimators of $\text{ISE}(\mathbf{X}_n)$ when k increases. As a consequence of the underlying Gaussian framework chosen, the methods proposed for the construction of \mathbf{Z}_m will not depend on the function evaluations \mathbf{y}_n .

The paper is organized as follows. A criterion measuring the quality of a validation design, based on Gaussian process modelling, is introduced in Section 2. In Section 3 we see how the proposed criterion can be optimized by kernel herding, detailing application of the general algorithm to it. The properties of the proposed design construction are investigated numerically in Section 4, exposing two important features: the completed design $\mathbf{X}_m \cup \mathbf{Z}_m$ must be space-filling, the contributions of the individual errors in $\widehat{\text{ISE}}(\mathbf{Z}_m, \mathbf{X}_n)$ must be under-weighted to avoid overestimation of $\text{ISE}(\mathbf{X}_n)$. These findings are confirmed in Section 5 where random test functions are used to illustrate achieved validation performance: the estimation of $\text{ISE}(\mathbf{X}_n)$ is significantly more accurate than with leave-one-out cross validation or uniformly weighted random or space-filling designs, which all seriously overestimate $\text{ISE}(\mathbf{X}_n)$.

2 An ISE-based criterion for validation design

2.1 A Gaussian process model

As mentioned above, f is unknown and we cannot choose \mathbf{Z}_m by minimizing $|\widehat{\text{ISE}}(\mathbf{Z}_m, \mathbf{X}_n) - \text{ISE}(\mathbf{X}_n)|$ directly. Assumptions on the behavior of f must be made. Considering the worst case for f in a given class of functions would be an option. Here we shall follow another, simpler, route and assume that f is the realization of a Gaussian Process (GP), or Gaussian Random Field, \mathcal{F}_x indexed by \mathcal{X} , with given second-order characteristics.

¹Here the integral (1) will be estimated directly by a discrete sum over \mathbf{Z}_m . The situation is different when the validation design \mathbf{Z}_m is used to predict the behavior of the error process $\varepsilon_n(\mathbf{x}) = \eta_n(\mathbf{x}) - f(\mathbf{x})$, in order to estimate $\text{ISE}(\mathbf{X}_n)$ by $\int_{\mathcal{X}} \varepsilon_n^2(x) \mu(dx)$. This alternative construction will be considered in a companion paper.

²Methods similar to those we propose for the construction of \mathbf{Z}_m can also be used to construct \mathbf{X}_n ; see Section 3.1 and the examples in Section 4.

For the sake of simplicity, we suppose that $E\{\mathcal{F}_x\} = 0$ for all $\mathbf{x} \in \mathcal{X}$. Extension to the case of a linearly parameterized mean, with $E\{\mathcal{F}_x\} = \boldsymbol{\beta}^\top \mathbf{h}(\mathbf{x})$ for a vector $\boldsymbol{\beta}$ of unknown parameters and a vector $\mathbf{h}(\mathbf{x}) = [h_1(\mathbf{x}), \dots, h_p(\mathbf{x})]^\top$ of p known functions of \mathbf{x} (including the constant) is possible via some adaptation. We also suppose that $E\{\mathcal{F}_x \mathcal{F}_{x'}\} = K(\mathbf{x}, \mathbf{x}')$, a known covariance function. In practice, K may be known up to a (variance) scaling coefficient σ^2 and parameterized by some parameters $\boldsymbol{\theta}$, setting in particular the correlation lengths and the smoothness of the functions that belong to the Reproducing Hilbert Space (RKHS) \mathcal{H}_K associated with K . Both σ^2 and $\boldsymbol{\theta}$ can be estimated from the data $\mathcal{F}_n = \{\mathbf{X}_n, \mathbf{y}_n\}$, e.g. by maximum likelihood, see for instance Santner et al. (2003), or cross validation; see Bachoc (2013) and Section 5.2 for an example.

The GP assumption defines a prior distribution for f , which can be updated given \mathcal{F}_n into a posterior distribution, with mean $E\{\mathcal{F}_x | \mathcal{F}_n\} = \mathbf{k}_n^\top(\mathbf{x}) \mathbf{K}_n^{-1} \mathbf{y}_n$ and covariance

$$E\{\mathcal{F}_x \mathcal{F}_{x'} | \mathcal{F}_n\} = K_{|n}(\mathbf{x}, \mathbf{x}') = K(\mathbf{x}, \mathbf{x}') - \mathbf{k}_n^\top(\mathbf{x}) \mathbf{K}_n^{-1} \mathbf{k}_n(\mathbf{x}') \geq 0, \quad (2)$$

for any \mathbf{x}, \mathbf{x}' in \mathcal{X} , where

$$\begin{aligned} \mathbf{k}_n(\mathbf{x}) &= [K(\mathbf{x}, \mathbf{x}_1) \dots, K(\mathbf{x}, \mathbf{x}_n)]^\top, \\ \{\mathbf{K}_n\}_{i,j} &= K(\mathbf{x}_i, \mathbf{x}_j), \quad i, j = 1, \dots, n, \end{aligned}$$

the $n \times n$ matrix \mathbf{K}_n being positive definite. The Integrated Mean Squared Error (IMSE)

$$\begin{aligned} \int_{\mathcal{X}} E\left\{[\eta_n(\mathbf{x}) - f(\mathbf{x})]^2 | \mathcal{F}_n\right\} \mu(d\mathbf{x}) &= \int_{\mathcal{X}} E\left\{[\eta_n(\mathbf{x}) - \mathbf{k}_n^\top(\mathbf{x}) \mathbf{K}_n^{-1} \mathbf{y}_n]^2 | \mathcal{F}_n\right\} \mu(d\mathbf{x}) \\ &\quad + \int_{\mathcal{X}} K_{|n}(\mathbf{x}, \mathbf{x}) \mu(d\mathbf{x}) \end{aligned}$$

is minimum when the prediction $\eta_n(\mathbf{x})$ equals the posterior mean $\mathbf{k}_n^\top(\mathbf{x}) \mathbf{K}_n^{-1} \mathbf{y}_n$, which yields

$$\text{IMSE}(\mathbf{X}_n) = \int_{\mathcal{X}} K_{|n}(\mathbf{x}, \mathbf{x}) \mu(d\mathbf{x}). \quad (3)$$

Note that $K_{|n}(\mathbf{x}, \mathbf{x}_i) = 0$ for any design point \mathbf{x}_i and any \mathbf{x} in \mathcal{X} and that η_n interpolates the observations \mathbf{y}_n . The extension to the case where η_n is not a interpolator does not raise particular difficulties, only yielding a kernel $\bar{K}_{|n}$ different from (5), having a slightly more complicated expression where the errors $f(\mathbf{x}_i) - \eta_m(\mathbf{x}_i)$ intervene.

2.2 Validation designs minimizing the expected squared ISE difference

Replacing $f(\mathbf{x})$ by a realization \mathcal{F}_x of the GP model of Section 2.1 in $\text{ISE}(\mathbf{X}_n)$ and $\widehat{\text{ISE}}(\mathbf{Z}_m, \mathbf{X}_n)$, we define

$$\begin{aligned} \bar{\Delta}^2(\mathbf{Z}_m, \mathbf{X}_n) &= E\left\{\left[\text{ISE}(\mathbf{X}_n) - \widehat{\text{ISE}}(\mathbf{Z}_m, \mathbf{X}_n)\right]^2 | \mathcal{F}_n\right\} \\ &= E\left\{\left[\int_{\mathcal{X}} [\mathcal{F}_x - \eta_n(\mathbf{x})]^2 (\zeta_m - \mu)(d\mathbf{x})\right]^2 | \mathcal{F}_n\right\}, \end{aligned}$$

the mean squared error of the ISE estimator $\widehat{\text{ISE}}(\mathbf{Z}_m, \mathbf{X}_n)$, and where ζ_m denotes the discrete measure $\zeta_m = (1/m) \sum_{i=1}^m \delta_{\mathbf{z}_i}$, with $\delta_{\mathbf{z}}$ the delta measure at \mathbf{z} ($\zeta_m - \mu$ is thus a signed measure with total mass 0).

For any positive definite kernel $C(\cdot, \cdot)$ and probability measures ξ and ν , denote by $\gamma_C(\xi, \nu)$ the Maximum Mean Discrepancy (MMD) between ξ and ν , defined by

$$\gamma_C^2(\xi, \nu) = \int_{\mathcal{X}^2} C(\mathbf{x}, \mathbf{x}') (\xi - \nu)(d\mathbf{x})(\xi - \nu)(d\mathbf{x}')$$

see (Sejdinovic et al., 2013, Def. 10). The minimization of $\gamma_C(\xi, \nu)$ with respect to ξ for a given ν can be performed by kernel herding (Section 3), yielding a sequence of finitely supported measures $\xi^{(t)}$ such that $\gamma_C(\xi^{(t)}, \nu) \rightarrow 0$ as $t \rightarrow \infty$.

When $\eta_n(\mathbf{x}) = \mathbf{k}_n(\mathbf{x})^\top \mathbf{K}_n^{-1} \mathbf{y}_n$, direct calculation gives

$$\begin{aligned} \bar{\Delta}^2(\mathbf{Z}_m, \mathbf{X}_n) &= \int_{\mathcal{X}^2} \mathbb{E} \{ [\mathcal{F}_x - \eta_n(\mathbf{x})]^2 [\mathcal{F}_{x'} - \eta_n(\mathbf{x}')]^2 | \mathcal{F}_n \} (\zeta_m - \mu)(d\mathbf{x})(\zeta_m - \mu)(d\mathbf{x}') \\ &= \int_{\mathcal{X}^2} \bar{K}_{|n}(\mathbf{x}, \mathbf{x}') (\zeta_m - \mu)(d\mathbf{x})(\zeta_m - \mu)(d\mathbf{x}') = \gamma_{\bar{K}_{|n}}^2(\zeta_m, \mu) \end{aligned} \quad (4)$$

where, for all \mathbf{x}, \mathbf{x}' in \mathcal{X} , we denote

$$\bar{K}_{|n}(\mathbf{x}, \mathbf{x}') = 2 K_{|n}^2(\mathbf{x}, \mathbf{x}') + K_{|n}(\mathbf{x}, \mathbf{x}) K_{|n}(\mathbf{x}', \mathbf{x}'), \quad (5)$$

with $K_{|n}$ defined by (2). Note that $\bar{K}_{|n}$ is positive definite (but not strictly positive definite, see Appendix A). Indeed, the Hadamard product $\mathbf{C}_n^{\circ 2}$ with elements $\{\mathbf{C}_n^{\circ 2}\}_{i,j} = C^2(\mathbf{x}_i, \mathbf{x}_j)$, $i, j = 1, \dots, n$, is positive definite when the matrix \mathbf{C}_n with elements $\{\mathbf{C}_n\}_{i,j} = C(\mathbf{x}_i, \mathbf{x}_j)$ is positive definite. Hence, $K_{|n}^2$ is positive definite since $K_{|n}$ is positive definite, which implies that $\bar{K}_{|n}$ is positive definite. The fact that $\bar{\Delta}^2(\mathbf{Z}_m, \mathbf{X}_n)$ does not depend on \mathbf{y}_n although it relies on conditioning on \mathcal{F}_n is a direct consequence of using a GP model.

3 Kernel herding for validation designs

3.1 A summary of kernel herding

Let C denote a positive definite kernel. For any signed measure ξ on \mathcal{X} , let

$$\mathcal{E}_C(\xi) = \int_{\mathcal{X}^2} C(\mathbf{x}, \mathbf{x}') \xi(d\mathbf{x})\xi(d\mathbf{x}') \geq 0 \quad (6)$$

denote the energy of ξ for C , so that $\gamma_{\bar{K}_{|n}}^2(\zeta_m, \mu) = \mathcal{E}_{\bar{K}_{|n}}(\zeta_m - \mu)$ in (4). A kernel C is called *characteristic* when $\gamma_C(\cdot, \cdot)$ defines a metric on the set of probability measures on \mathcal{X} , implying in particular that, for two probability measures ζ and μ , $\gamma_C(\zeta, \mu) = 0$ if and only if $\zeta = \mu$. The kernel $\bar{K}_{|n}$ is not characteristic, see Appendix A, but we can nevertheless consider the minimization of $\gamma_{\bar{K}_{|n}}^2(\zeta_m, \mu)$.

For any $\alpha \in [0, 1]$, we have $(1-\alpha)\mathcal{E}_C(\xi) + \alpha\mathcal{E}_C(\nu) - \mathcal{E}_C[(1-\alpha)\xi + \alpha\nu] = \alpha(1-\alpha)\mathcal{E}_C(\xi - \nu) \geq 0$, showing that $\mathcal{E}_C(\cdot)$ is convex; see Pronzato and Zhigljavsky (2020), and we can minimize the squared MMD criterion $\gamma_C^2(\xi, \mu) = \mathcal{E}_C(\xi - \mu)$ with respect to ξ by a simple descent algorithm.

Denote by $F_{C,\mu}(\xi; \nu)$ the directional derivative of $\gamma_C^2(\cdot, \mu)$ at ξ in the direction ν ,

$$F_{C,\mu}(\xi; \nu) = \lim_{\alpha \rightarrow 0^+} \frac{\mathcal{E}_C[(1-\alpha)\xi + \alpha\nu - \mu] - \mathcal{E}_C(\xi - \mu)}{\alpha}.$$

Straightforward calculation gives

$$F_{C,\mu}(\xi; \nu) = 2 \left[\int_{\mathcal{X}^2} C(\mathbf{x}, \mathbf{x}') (\nu - \mu)(d\mathbf{x})(\xi - \mu)(d\mathbf{x}') - \mathcal{E}_C(\xi - \mu) \right].$$

In particular, for $\nu = \delta_{\mathbf{x}}$, we get

$$F_{C,\mu}(\xi; \delta_{\mathbf{x}}) = 2 [P_{C,\xi}(\mathbf{x}) - P_{C,\mu}(\mathbf{x}) - \mathcal{E}_C(\xi) + \mathcal{E}_C(\xi, \mu)], \quad (7)$$

where $\mathcal{E}_C(\xi, \nu) = \int_{\mathcal{X}^2} C(\mathbf{x}, \mathbf{x}') \xi(d\mathbf{x})\nu(d\mathbf{x}')$ and

$$P_{C,\xi}(\mathbf{x}) = \int_{\mathcal{X}} C(\mathbf{x}, \mathbf{x}') \xi(d\mathbf{x}')$$

(respectively, $P_{C,\mu}(\mathbf{x}) = \int_{\mathcal{X}} C(\mathbf{x}, \mathbf{x}') \mu(d\mathbf{x}')$) is called the potential of ξ (respectively, of μ), at \mathbf{x} , associated with C . $P_{C,\mu}(\cdot)$ is also called the kernel embedding of μ in the RKHS associated with C (Sejdicinovic et al., 2013, Def. 9). Standard kernel-herding corresponds to the Frank-Wolfe conditional gradient algorithm (Bach et al., 2012), that is, to the vertex-direction method with predefined step-length, commonly used in optimal experimental design since the pioneering work of Wynn (1970) and Fedorov (1972).

The general form of the algorithm, with step length $\alpha_k \in (0, 1)$ at iteration k , is as follows: starting with some probability measure $\zeta^{(k_1)}$ on \mathcal{X} , we take, for all $k \geq k_1$,

$$\zeta^{(k+1)} = (1 - \alpha_k) \zeta^{(k)} + \alpha_k \delta_{\mathbf{z}_{k+1}} \quad (8)$$

where $\mathbf{z}_{k+1} \in \text{Arg min}_{\mathbf{z} \in \mathcal{X}} F_{C,\mu}(\zeta^{(k)}; \delta_{\mathbf{x}})$. From (7), this is equivalent to

$$\mathbf{z}_{k+1} \in \text{Arg min}_{\mathbf{z} \in \mathcal{X}} P_{C,\zeta^{(k)}}(\mathbf{z}) - P_{C,\mu}(\mathbf{z}). \quad (9)$$

If the initial measure $\zeta^{(k_1)}$ is finitely supported on a set $\mathcal{S}^{(k_1)}$, then $\zeta^{(k)}$ remains finitely supported for all k . Selecting the optimal α_k at each iteration corresponds to Fedorov's algorithm (1972) used in optimal design for parametric models. If $\mathcal{S}^{(k_1)}$ has k_1 elements, $\mathcal{S}^{(k_1)} = \mathbf{Z}^{(k_1)} = \{\mathbf{z}_1, \dots, \mathbf{z}_{k_1}\}$, and $\zeta^{(k_1)} = (1/k_1) \sum_{i=1}^{k_1} \delta_{\mathbf{z}_i}$ is uniform on $\mathbf{Z}^{(k_1)}$, by choosing $\alpha_k = 1/(k+1)$ for all $k \geq k_1$ we obtain that $\zeta^{(k)} = (1/k) \sum_{i=1}^k \delta_{\mathbf{z}_i}$ for all k , and

$$P_{C,\zeta^{(k)}}(\mathbf{x}) = \frac{1}{k} \sum_{i=1}^k C(\mathbf{x}, \mathbf{z}_i).$$

In particular, we can take $\alpha_k = 1/(k+1)$, $\zeta^{(1)} = \delta_{\mathbf{z}_1}$ for some $\mathbf{z}_1 \in \mathcal{X}$, and \mathbf{z}_1 can be chosen by maximizing $P_{C,\mu}(\mathbf{z})$. For stationary kernels such that $C(\mathbf{x}, \mathbf{x}')$ only depends on $\mathbf{x}' - \mathbf{x}$, it amounts at taking \mathbf{z}_1 at the center of \mathcal{X} .

We shall denote by $\mathbf{Z}_k = \text{KH}(\mathbf{Z}_{k_1}, C, k)$ the k -point design obtained in this way, after k iterations of kernel herding initialized at \mathbf{Z}_{k_1} containing k_1 elements, with $\alpha_k = 1/(k+k_1)$ for all $k \geq 1$; $\text{KH}(\emptyset, C, k)$ selects \mathbf{z}_1 by maximization of $P_{C,\mu}$.

In practice, the search for \mathbf{z}_{k+1} in (9) is generally made within a finite subset \mathcal{X}_Q of \mathcal{X} , with Q elements. The cost of the determination of \mathbf{z}_{k+1} in (9) is $\mathcal{O}(Q)$ if we compute $C(\mathbf{x}, \mathbf{z}_k)$ for all $\mathbf{x} \in \mathcal{X}_Q$ and update the sum $\sum_{i=1}^{k-1} C(\mathbf{x}, \mathbf{z}_i)$; the cost for k iterations then scales as $\mathcal{O}(kQ)$,

including the initial cost for the computation of $P_{C,\mu}(\mathbf{x})$ for all $\mathbf{x} \in \mathcal{X}_Q$. Another option when μ is approximated by the uniform measure μ_Q on \mathcal{X}_Q and $C(\mathbf{x}, \mathbf{x}')$ only depends on $\|\mathbf{x} - \mathbf{x}'\|$, is to compute in advance the $Q(Q-1)/2$ distances between all pairs of points in \mathcal{X}_Q (feasible only if Q is not too large).

The minimum-norm variant of Bach et al. (2012) replaces $\zeta^{(k)}$ in (9) by the measure having the same support $\mathcal{S}^{(k)}$ but optimal weights, positive and summing to one; these optimal weights are solution of a convex quadratic programming problem. Here we shall consider a simplified version where $\zeta^{(k)}$ is replaced by $\hat{\zeta}^{(k)}$ having weights $\hat{w}_i^{(k)}$ summing to one and such that such that $\mathcal{E}_C(\hat{\zeta}^{(k)} - \mu)$ is minimal. For a measure ζ_k with support $\mathcal{S}^{(k)}$ and weights $\mathbf{w}^{(k)}$, we have

$$\mathcal{E}_C(\zeta_k - \mu) = \mathbf{w}^{(k)\top} \mathbf{C}_k \mathbf{w}^{(k)} - 2 \mathbf{w}^{(k)\top} \mathbf{p}_{C,k}(\mu) + \mathcal{E}_C(\mu), \quad (10)$$

where $\{\mathbf{C}_k\}_{i,j} = C(\mathbf{z}_i, \mathbf{z}_j)$, $i, j = 1 \dots, k$ and $\mathbf{p}_{C,k}(\mu) = [P_{C,\mu}(\mathbf{z}_1), \dots, P_{C,\mu}(\mathbf{z}_k)]^\top$. Its minimization under the constraint $\mathbf{1}_k^\top \mathbf{w}^{(k)} = 1$, with $\mathbf{1}_k$ the k -dimensional vector with all components equal to one, gives the optimal weights

$$\hat{\mathbf{w}}^{(k)} = (\hat{w}_1^{(k)}, \dots, \hat{w}_k^{(k)})^\top = \left(\mathbf{C}_k^{-1} - \frac{\mathbf{C}_k^{-1} \mathbf{1}_k \mathbf{1}_k^\top \mathbf{C}_k^{-1}}{\mathbf{1}_k^\top \mathbf{C}_k^{-1} \mathbf{1}_k} \right) \mathbf{p}_{C,k}(\mu) + \frac{\mathbf{C}_k^{-1} \mathbf{1}_k}{\mathbf{1}_k^\top \mathbf{C}_k^{-1} \mathbf{1}_k}. \quad (11)$$

By construction, $\hat{\zeta}^{(k)}$ minimizes $\mathcal{E}_C(\zeta_k - \mu)$ with respect to measures ζ_k of total mass one supported on $\mathcal{S}^{(k)}$, and one can show (Pronzato and Zhigljavsky, 2020) that its potential $P_{C,\hat{\zeta}^{(k)}}(x)$ satisfies

$$P_{C,\hat{\zeta}^{(k)}}(x) - P_{C,\mu}(x) - \mathcal{E}_C(\hat{\zeta}^{(k)}, \mu) + \mathcal{E}_C(\mu) = 0, \quad \forall x \in \mathcal{S}^{(k)},$$

showing that $P_{C,\hat{\zeta}^{(k)}}(x) - P_{C,\mu}(x)$ is constant on $\mathcal{S}^{(k)}$.

When \mathcal{X} is discretized into \mathcal{X}_Q , the substitution of $\hat{\zeta}^{(k)}$ for $\zeta^{(k)}$ requires the storage of all $C(\mathbf{x}, \mathbf{z}_i)$, $i = 1, \dots, k$, $\mathbf{x} \in \mathcal{X}_Q$, in order to compute $P_{C,\hat{\zeta}^{(k)}}(\mathbf{x}) = \sum_{i=1}^k \hat{w}_i^{(k)} C(\mathbf{x}, \mathbf{z}_i)$ in (9). At iteration k , the computation of $\hat{\mathbf{w}}^{(k)}$ by (11) also induces an additional computational cost of $\mathcal{O}(k^3)$ (reduced to $\mathcal{O}(k^2)$ if rank-one updating is used to compute \mathbf{C}_k^{-1}); $\gamma_C^2(\hat{\zeta}^{(k)}, \mu)$ decreases faster than $\gamma_C^2(\zeta^{(k)}, \mu)$; see Pronzato (2021).

We shall denote by $\mathbf{Z}_k = \text{MN}(\mathbf{Z}_{k_1}, C, k)$ the k -point design obtained after k iterations, initialized at \mathbf{Z}_{k_1} ($\text{MN}(\emptyset, C, k)$ chooses \mathbf{z}_1 that maximizes $P_{C,\mu}(\mathbf{z})$). We write $[\mathbf{Z}_k, \hat{\mathbf{w}}^{(k)}] = \text{MN}(\mathbf{Z}_{k_1}, C, k)$ when we are also interested in the weights $\hat{\mathbf{w}}^{(k)}$ given by (11), and for any m -point design \mathbf{Z}_m we denote by $\hat{\mathbf{w}}(\mathbf{Z}_m, C)$ the weights computed by (11).

Example 1. To illustrate the behavior of the algorithms above, we consider a small one-dimensional example with $\mathcal{X} = [0, 1]$ and $C = K_{3/2,\theta}$, the Matérn 3/2 kernel

$$K_{3/2,\theta}(x, x') = (1 + \sqrt{3}\theta |x - x'|) \exp(-\sqrt{3}\theta |x - x'|). \quad (12)$$

The measure μ is approximated by the uniform discrete distribution on \mathcal{X}_Q given by the first $Q = 256$ points of a scrambled Sobol' sequence in \mathcal{X} ; \mathbf{z}_{k+1} in (9) is searched within the same set \mathcal{X}_Q ; we take $\theta = 10$ in $K_{3/2,\theta}$.

The left panel of Figure 1 shows $P_{C,\zeta^{(k)}}(x)$ (black solid line) and $P_{C,\mu}(x)$ (blue dashed line) as functions of $x \in \mathcal{X}$, with $\mathbf{X}_n = \text{KH}(\emptyset, C, n)$; the right panel shows $P_{C,\zeta^{(k)}}(x) - P_{C,\mu}(x)$. The figure is for $n = 4$ and \mathbf{X}_n is indicated by black squares.

Kernel herding is used on the top row: $\zeta^{(k)}$ for $k = 3$ is supported by the points in $\text{KH}(\mathbf{X}_n, C, 3)$ indicated with a red triangle; the next point z_4 chosen by the algorithm — the location of minimum of the right panel — is indicated by the red star. $P_{C,\zeta^{(k)}}(x)$ decreases when x moves away from its closest prediction point x_i or validation point z_i , whereas $P_{C,\mu}(\cdot)$ is a fixed function, independent of the x_i and z_i . The bottom row is for the minimum-norm variant $\text{MN}(\mathbf{X}_n, C, k)$ of kernel herding: at iteration k we replace $\zeta^{(k)}$ by $\hat{\zeta}^{(k)}$ having weights given by (11). The right panel illustrates the property that $P_{C,\hat{\zeta}^{(k)}}(x) - P_{C,\mu}(x)$ is constant on the support $\mathcal{S}^{(k)} = \mathbf{X}_n \cup \{z_1, \dots, z_k\}$ of $\hat{\zeta}^{(k)}$. Note that $P_{C,\hat{\zeta}^{(k)}}(x)$ is closer to $P_{C,\mu}(x)$ than in the first row, indicating a better approximation of μ in the sense of the MMD criterion. \triangleleft

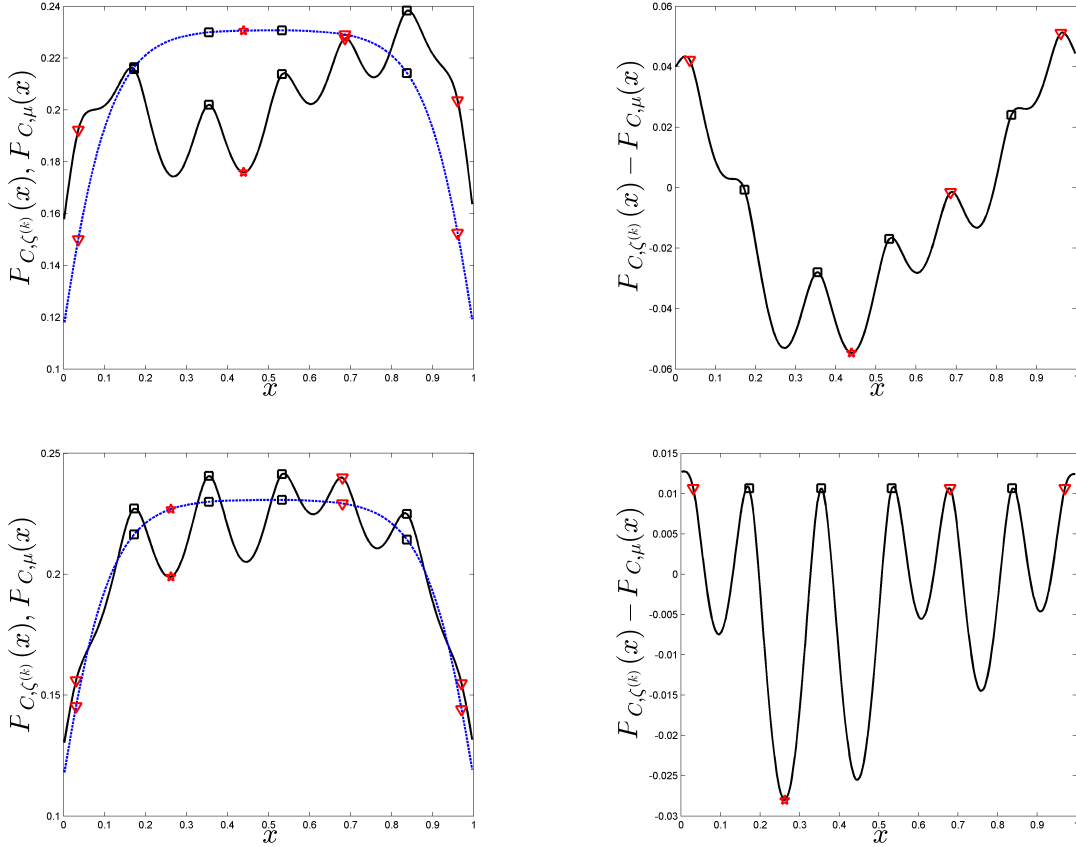


Figure 1: Left: $P_{C,\zeta^{(k)}}(x)$ (black solid line) and $P_{C,\mu}(x)$ (blue dashed line); Right: $P_{C,\zeta^{(k)}}(x) - P_{C,\mu}(x)$. Design points x_i , $i = 1, \dots, n = 4$: black \square ; validation points z_i , $i \leq k = 3$: red ∇ ; z_4 : red \star . Top row: kernel herding; bottom row: MN variant; $C = K_{3/2,\theta}$.

To summarise, we proposed two distinct validation designs in this section:

- The kernel herding solution $\mathbf{Z}_k = \text{KH}(\mathbf{Z}_{k_1}, C, k)$, a design of size k obtained by iteratively minimising $P_{C,\zeta^{(k)}}(\mathbf{z}) - P_{C,\mu}(\mathbf{z})$, see (9), starting from \mathbf{Z}_{k_1} of size k_1 , and with weights

$\alpha_k = 1/(k + k_1)$ for all $k \geq 1$. A variant of this notation will be introduced in the next subsection, where we will introduce $\mathbf{Z}_m = \text{KH}(\mathbf{Z}_{k_1}, C, k, \setminus \mathbf{X}_n)$ to denote the design of size m , with no repeated points and empty intersection with \mathbf{X}_n , obtained after k iterations.

- The minimum norm solution $\mathbf{Z}_k = \text{MN}(\mathbf{Z}_{k_1}, C, k)$, a design of size k obtained after k iterations of (9) initiated at \mathbf{Z}_{k_1} , but using a measure $\zeta^{(k)}$ with the optimal weights (11) in the computation of $P_{C, \zeta^{(k)}}(\mathbf{z})$. As we will see below, this solution is not well defined when $C = K|_n$ or $C = \overline{K}|_n$, and a slightly different definition, dropping the constraint of unitary sum of the weights of $\zeta^{(k)}$, is required. It will be denoted by $\mathbf{Z}_k = \text{MN}_2(\mathbf{Z}_{k_1}, C, k)$.

3.2 Incremental construction of space-filling and validation designs

To apply a kernel-herding algorithm to the minimization of $\gamma_{\overline{K}|_n}(\zeta_m, \mu)$ given by (4) with respect to ζ_m , we simply substitute the conditional kernel $\overline{K}|_n$, given by (5), for the kernel C in (8, 9). The construction has the advantage of being incremental³: it generates a design sequence $\mathbf{z}_1, \mathbf{z}_2, \dots$ which can be interrupted at any design size m .

However, the application of kernel herding to kernels C such that $C(\mathbf{x}, \mathbf{x}_i) = 0$ for all design points \mathbf{x}_i , which occurs when $C = \overline{K}|_n$, requires a specific treatment. Indeed, it may then happen that the two potentials $P_{C, \zeta^{(k)}}(\mathbf{x}_i)$ and $P_{C, \mu}(\mathbf{x}_i)$ used in (9) satisfy $P_{C, \mu}(\mathbf{x}) \leq P_{C, \zeta^{(k)}}(\mathbf{x})$ for all $\mathbf{x} \in \mathcal{X}$, with $P_{C, \mu}(\mathbf{x}) = P_{C, \zeta^{(k)}}(\mathbf{x})$ for $\mathbf{x} \in \mathbf{X}_n$ and the inequality being strict otherwise. In that case, (9) necessarily chooses \mathbf{z}_{k+1} among \mathbf{X}_n . If a $\zeta^{(\ell)}$ has one of its support points \mathbf{z}_j in \mathbf{X}_n , (10) indicates that the associated weight $w_j^{(\ell)}$ does not contribute to $\mathcal{E}_C(\zeta^{(\ell)} - \mu)$. The selection of $\mathbf{z}_{\ell+1}$ among \mathbf{X}_n is thus equivalent to a reduction of the total mass of other points that contribute to $\mathcal{E}_C(\zeta^{(\ell)} - \mu)$.

The possible selection of \mathbf{z}_{k+1} within \mathbf{X}_n has several consequences on kernel herding.

- (i) When it happens that \mathbf{z}_{k+1} is chosen among \mathbf{X}_n at an iteration (9) of standard kernel herding (with uniform weighting), we can nevertheless continue iterations until the number of selected points not in \mathbf{X}_n reaches the desired value m ; we denote by $\mathbf{Z}_m = \text{KH}(\mathbf{Z}_{k_1}, C, k, m, \setminus \mathbf{X}_n)$ the corresponding m -point design. Since the selection is made within a finite set, it may also happen that the same point is selected several times. In that case, we may also impose that \mathbf{Z}_m contains m distinct points and continue iterations until this condition is satisfied; the weights given by the algorithm to the \mathbf{z}_i in \mathbf{Z}_m are then multiple of $1/m'$, with $m' \geq m$ (they are not necessarily all equal to $1/m$).
- (ii) When a support point \mathbf{z}_j of $\zeta^{(k)}$ coincides with a design point \mathbf{x}_i , \mathbf{C}_k is singular and we cannot compute $\hat{\mathbf{w}}^{(k)}$ by (11); that is, the minimum-norm variant of kernel herding cannot be used.
- (iii) When a support point \mathbf{z}_j belongs to \mathbf{X}_n , the optimal weights allocated to the points that do not belong to \mathbf{X}_n are obtained by minimizing (10) with respect to $\mathbf{w}^{(k)}$ *without the constraint* $\mathbf{1}_k^\top \mathbf{w}^{(k)} = 1$.

³However, it does not provide the optimal design for m fixed: the construction of one-shot m -point designs minimizing a MMD criterion is considered for instance in (Pronzato and Zhigljavsky, 2020); we do not develop this aspect here.

To account for the possibility that the algorithm may choose \mathbf{z}_{k+1} in \mathbf{X}_n , we consider a new version of kernel herding where, at iteration k , $\zeta^{(k)}$ is replaced by $\check{\zeta}^{(k)}$ having the same support $\mathcal{S}^{(k)}$ but weights $\check{\mathbf{w}}^{(k)}$ that minimize $\mathcal{E}_C(\zeta^{(k)} - \mu)$ given by (10) with respect to $\mathbf{w}^{(k)}$ without constraints on $\mathbf{w}^{(k)}$, contrarily to $\hat{\mathbf{w}}^{(k)}$, given by (11), which satisfies $\sum_{i=1}^k \hat{w}_i^{(k)} = 1$. Direct calculation gives

$$\check{\mathbf{w}}^{(k)} = (\check{w}_1^{(k)}, \dots, \check{w}_k^{(k)})^\top = \mathbf{C}_k^{-1} \mathbf{p}_{C,k}(\mu). \quad (13)$$

The measures $\zeta^{(k)}$, $\hat{\zeta}^{(k)}$ and $\check{\zeta}^{(k)}$, with respective weights $\mathbf{1}_k$, $\hat{\mathbf{w}}^{(k)}$ and $\check{\mathbf{w}}^{(k)}$, satisfy

$$\mathcal{E}_C(\zeta^{(k)} - \mu) \geq \mathcal{E}_C(\hat{\zeta}^{(k)} - \mu) \geq \mathcal{E}_C(\check{\zeta}^{(k)} - \mu).$$

In general, both $\hat{\mathbf{w}}^{(k)}$ and $\check{\mathbf{w}}^{(k)}$ may have negative components. We shall denote by $\mathbf{Z}_k = \text{MN}_2(\mathbf{Z}_{k_1}, K, k)$ the design obtained after k iterations of this variant of kernel herding, initialized at \mathbf{Z}_{k_1} ($\text{MN}_2(\emptyset, K, k)$ chooses \mathbf{z}_1 that maximizes $P_{K,\mu}$). We write $[\mathbf{Z}_k, \check{\mathbf{w}}^{(k)}] = \text{MN}_2(\mathbf{Z}_{k_1}, K, k)$ when we are also interested in the weights $\check{\mathbf{w}}^{(k)}$ given by (13), and for any m -point design \mathbf{Z}_m , we denote by $\check{\mathbf{w}}(\mathbf{Z}_m, K)$ the weights computed by (13).

This construction can be interpreted as standard kernel herding applied to a kernel $C(k)$ varying along iterations, given by the conditional kernel $C_{|k}$ at iteration k ,

$$C_{|k}(\mathbf{x}, \mathbf{x}') = C(\mathbf{x}, \mathbf{x}') - \mathbf{c}_k^\top(\mathbf{x}) \mathbf{C}_k^{-1} \mathbf{c}_k(\mathbf{x}'),$$

where $\mathbf{c}_k(\mathbf{x}) = [C(\mathbf{x}, \mathbf{z}_1) \dots, C(\mathbf{x}, \mathbf{z}_k)]^\top$, see (2). Indeed, the potential $P_{C_{|k}, \zeta}(\mathbf{z})$ for a measure ζ on \mathcal{X} is

$$P_{C_{|k}, \zeta}(\mathbf{z}) = P_{C, \zeta}(\mathbf{z}) - \mathbf{c}_k^\top(\mathbf{z}) \mathbf{C}_k^{-1} \mathbf{p}_{C,k}(\zeta), \quad (14)$$

where $\mathbf{p}_{C,k}(\zeta) = [P_{C, \zeta}(\mathbf{z}_1), \dots, P_{C, \zeta}(\mathbf{z}_k)]^\top$. Since $C_{|k}(\mathbf{z}, \mathbf{z}_i) = 0$ for all \mathbf{z}_i , for any ζ_k supported on $\mathbf{Z}_k = \{\mathbf{z}_1, \dots, \mathbf{z}_k\}$ we have, for all $\mathbf{z} \in \mathcal{X}$,

$$P_{C_{|k}, \zeta_k}(\mathbf{x}) - P_{C_{|k}, \mu}(\mathbf{z}) = -P_{C_{|k}, \mu}(\mathbf{z}) = \mathbf{c}_k^\top(\mathbf{z}) \mathbf{C}_k^{-1} \mathbf{p}_{C,k}(\mu) - P_{C, \mu}(\mathbf{z}) \quad (15)$$

$$= P_{C, \check{\zeta}^{(k)}}(\mathbf{z}) - P_{C, \mu}(\mathbf{z}). \quad (16)$$

At iteration k , kernel herding with $C_{|k}$ and the variant with kernel C but optimal weights $\check{\mathbf{w}}^{(k)}$ thus select the same \mathbf{z}_{k+1} that minimizes (16). Note that $P_{C, \check{\zeta}^{(k)}}(\mathbf{z}) - P_{C, \mu}(\mathbf{z}) = 0$ for all \mathbf{z}_i . When we substitute the conditional kernel $K_{|n}$ for C , the variant MN_2 of kernel herding also satisfies the following property, meaning that we do not need to know where the points \mathbf{X}_n are, everything in terms of information being coded in the conditional kernel $K_{|n}$.

Theorem 1. *For any positive definite kernel K , any design \mathbf{X}_n and any $k \geq 1$, there exist choices for \mathbf{z}_{i+1} , $i = 0, \dots, k$ in (9) such that $\text{MN}_2(\mathbf{X}_n, K, k) = \text{MN}_2(\emptyset, K_{|n}, k)$, where $K_{|n}$ is defined by (2).*

Proof. Consider first the case $k = 1$. On the one hand, $\mathbf{z}_1 = \text{MN}_2(\mathbf{X}_n, K, 1)$ minimizes $\mathbf{k}_n^\top(\mathbf{z}) \mathbf{K}_n^{-1} \mathbf{p}_{K,n}(\mu) - P_{K,\mu}(\mathbf{z})$, see (15); on the other hand, $\mathbf{z}'_1 = \text{MN}_2(\emptyset, K_{|n}, 1)$ maximizes $P_{K_{|n}, \mu}(\mathbf{z}) = P_{K, \mu}(\mathbf{z}) - \mathbf{k}_n^\top(\mathbf{z}) \mathbf{K}_n^{-1} \mathbf{p}_{K,n}(\mu)$, see (14). One can therefore choose $\mathbf{z}_1 = \mathbf{z}'_1$.

The identity of the two constructions at any $k > 1$ is a consequence of the conditioning property of GP: at step k , they both use the kernel $K_{|n+k}$. More precisely, \mathbf{z}_{k+1} for the construction of $\text{MN}_2(\mathbf{X}_n, K, k+1)$ minimizes $J(\mathbf{z}) = \mathbf{k}_{n+k}^\top(\mathbf{z})\mathbf{K}_{n+k}^{-1}\mathbf{p}_{K,n+k}(\mu) - P_{K,\mu}(\mathbf{z})$, where

$$\mathbf{k}_{n+k}(\mathbf{z}) = \begin{pmatrix} \mathbf{k}_n(\mathbf{z}) \\ \mathbf{k}_k(\mathbf{z}) \end{pmatrix}, \quad \mathbf{p}_{K,n+k}(\mu) = \begin{pmatrix} \mathbf{p}_{K,n}(\mu) \\ \mathbf{p}_{K,k}(\mu) \end{pmatrix}, \quad \mathbf{K}_{n+k} = \begin{pmatrix} \mathbf{K}_n & \mathbf{K}_{n,k} \\ \mathbf{K}_{k,n} & \mathbf{K}_k \end{pmatrix},$$

while \mathbf{z}'_{k+1} for $\text{MN}_2(\emptyset, K_{|n}, k+1)$ minimizes $J'(\mathbf{z}) = \mathbf{k}_{|n_k}^\top(\mathbf{z})\mathbf{K}_{|n_k}^{-1}\mathbf{p}_{K_{|n},k}(\mu) - P_{K_{|n},\mu}(\mathbf{z})$, where

$$\mathbf{k}_{|n_k}(\mathbf{z}) = \mathbf{k}_k(\mathbf{z}) - \mathbf{K}_{k,n}\mathbf{K}_n^{-1}\mathbf{k}_n(\mathbf{z}), \quad \mathbf{p}_{K_{|n},k}(\mu) = \mathbf{p}_{K,k}(\mu) - \mathbf{K}_{k,n}\mathbf{K}_n^{-1}\mathbf{p}_{K,n}(\mu)$$

and $\mathbf{K}_{|n_k} = \mathbf{K}_k - \mathbf{K}_{k,n}\mathbf{K}_n^{-1}\mathbf{K}_{n,k}$. Direct application of Woodbury identity for matrix inversion and inversion of a block matrix shows that $J(\mathbf{z}) = J'(\mathbf{z})$; we can thus choose the same \mathbf{z}_{k+1} in both constructions in case multiple choices are possible. ■

Example 1 (continued). We consider the same situation as in Example 1 with the same $\mathbf{X}_n = \text{KH}(\emptyset, K, n)$ for $K = K_{3/2,10}$.

Figure 2 illustrates the construction of $\text{KH}(\emptyset, C, m)$ for $C = \overline{K}_{|n}$, with $P_{C,\zeta^{(k)}}(x)$ (black solid line) and $P_{C,\mu}(x)$ (blue dashed line) as functions of $x \in \mathcal{X}$ on the left and $P_{C,\zeta^{(k)}}(x) - P_{C,\mu}(x)$ on the right. \mathbf{X}_n is indicated by black squares; $\zeta^{(k)}$ for $k = 3$ is supported by the points indicated with a red triangle; the next point z_4 chosen by the algorithm corresponds to the red star. $P_{C,\zeta^{(k)}}(x_i) = P_{C,\mu}(x_i)$ for all $x_i \in \mathbf{X}_n$ and large values of potentials are obtained far away from the design points in \mathbf{X}_n only. Note that $P_{C,\mu}(z) < P_{C,\zeta^{(3)}}(z)$ excepted in a small neighborhood around z_4 , and that one of previous points selected by kernel herding (here z_2) coincides with a design point.

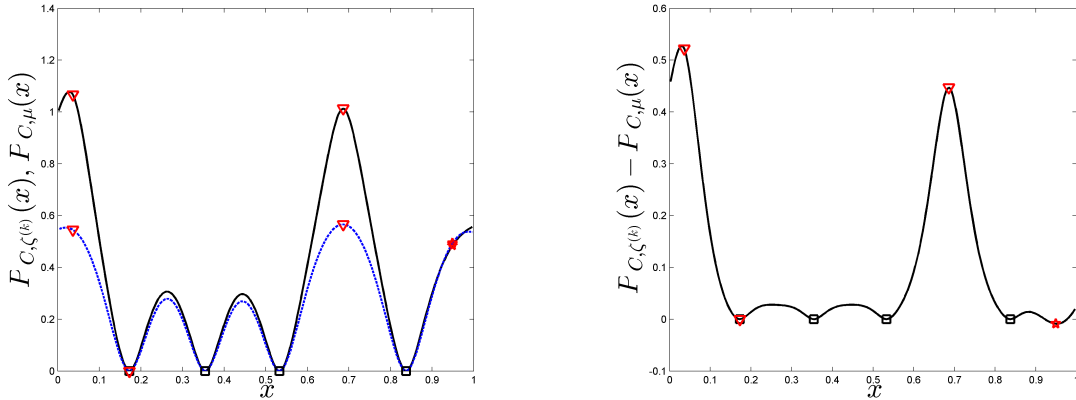


Figure 2: Kernel herding $\text{KH}(\emptyset, C, 4)$ for $C = \overline{K}_{|n}$. Left: $P_{C,\zeta^{(k)}}(x)$ (black solid line) and $P_{C,\mu}(x)$ (blue dashed line); Right: $P_{C,\zeta^{(k)}}(x) - P_{C,\mu}(x)$. Design points $x_i, i = 1, \dots, n = 4$: black \square ; validation points $z_i, i \leq k = 3$: red ∇ ; z_4 : red \star .

We cannot use MN for $C = \overline{K}_{|n}$ since we cannot compute optimal weights through (11); Figure 3 illustrates the construction with the minimum-norm variant MN_2 of kernel herding that uses weights (13). In the first row, we use $C = K$ and at iteration k the support $\mathcal{S}^{(k)}$

equals $\mathbf{X}_n \cup \{z_1, \dots, z_k\}$. The second row corresponds to $C = \overline{K}_{|n}$. Note that in both cases $P_{C, \zeta^{(k)}}(\mathbf{z}_i) = P_{C, \mu}(\mathbf{z}_i)$ for all \mathbf{z}_i in the support $\mathcal{S}^{(k)}$ of $\zeta^{(k)}$. \triangleleft

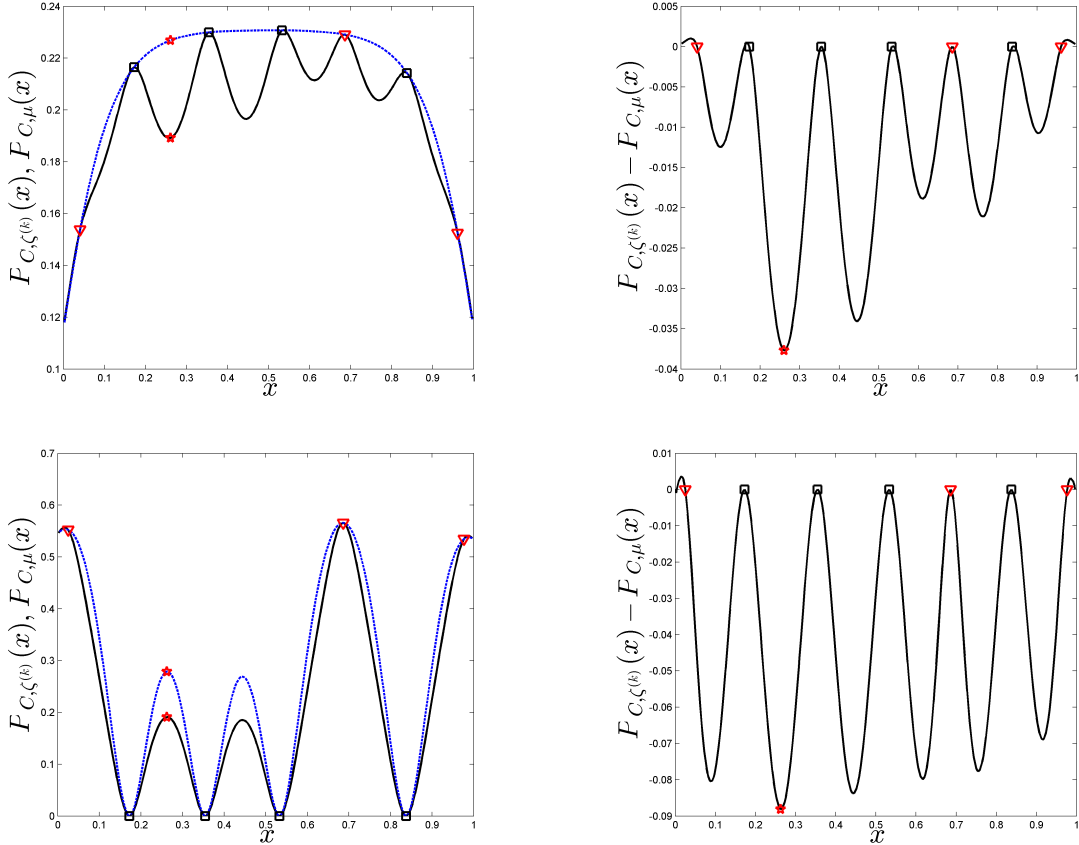


Figure 3: MN₂ variant of kernel herding. Left: $P_{C, \zeta^{(k)}}(x)$ (black solid line) and $P_{C, \mu}(x)$ (blue dashed line); Right: $P_{C, \zeta^{(k)}}(x) - P_{C, \mu}(x)$; $C = K$ (top row) and $C = \overline{K}_{|n}$ (bottom row). Design points x_i , $i = 1, \dots, n = 4$: black \square ; validation points z_i , $i \leq k = 3$: red ∇ ; z_4 : red \star .

Although they follow the same principle of one-step ahead minimization of a convex functional of a measure, the three methods KH, MN and MN₂ rely on quite different functions $P_{C, \zeta^{(k)}}(x) - P_{C, \mu}(x)$ for the selection of support points in (9); see the right columns of Figures 1 to 3. The differences are also important depending on which kernel is used: the original one K , which is stationary in Example 1 above, or $\overline{K}_{|n}$ which accounts for the presence of the n design points in \mathbf{X}_n . Next section contains a numerical comparison of the performances of designs obtained with those different approaches, in particular in terms of $\overline{\Delta}(\mathbf{Z}_m, \mathbf{X}_n)$ given by (4).

4 Properties of validation design constructed by kernel herding

In this section, we investigate and compare the properties of validation designs obtained by minimizing $\gamma_{\overline{K}_{|n}}(\zeta_m, \mu)$ for different choices of n , m and dimension d . In the kernel-herding

algorithm and its variants, we approximate μ by the uniform measure μ_Q on \mathcal{X}_Q given by the first $Q = 2^{12}$ points of a scrambled Sobol' sequence in $\mathcal{X} = [0, 1]^d$; Q is taken small enough to allow the computation of all $Q(Q - 1)/2$ distances between pairs of points in \mathcal{X}_Q . K is the Matérn 3/2 isotropic kernel,

$$K_{3/2,\theta}(\mathbf{x}, \mathbf{x}') = (1 + \sqrt{3}\theta \|\mathbf{x} - \mathbf{x}'\|) \exp(-\sqrt{3}\theta \|\mathbf{x} - \mathbf{x}'\|),$$

with $\theta = n^{1/d}$ and n the size of the prediction design \mathbf{X}_n , given by $\mathbf{X}_n = \text{KH}(\emptyset, K, n)$.

a) Space-filling performance. Although $\bar{\Delta}^2(\mathbf{Z}_m, \mathbf{X}_n) = \gamma_{\bar{K}|n}^2(\zeta_m, \mu)$ given by (4) is not directly related to a space-filling characteristic, below we shall see that kernel herding applied to its minimization may provide designs with attractive space-filling properties.

Figure 4 shows the design $\mathbf{X}_n = \text{KH}(\emptyset, K, n)$ (black squares) and the validation designs $\text{KH}(\mathbf{X}_n, K, m)$ (blue triangles) and $\text{MN}_2(\emptyset, \bar{K}|n, m)$ (red stars) for $n = 50$ and $m = 25$ (left), and $n = 50$, $m = 50$ (right) when $d = 2$ (note that $\text{KH}(\mathbf{X}_n, K, 25) \subset \text{KH}(\mathbf{X}_n, K, 50)$ and $\text{MN}_2(\emptyset, \bar{K}|n, 25) \subset \text{MN}_2(\emptyset, \bar{K}|n, 50)$). For the two values of m considered, $\text{KH}(\mathbf{X}_n, K, m)$ looks more evenly spread in \mathcal{X} than $\text{MN}_2(\emptyset, \bar{K}|n, m)$, even if both designs are well interlaced with \mathbf{X}_n .

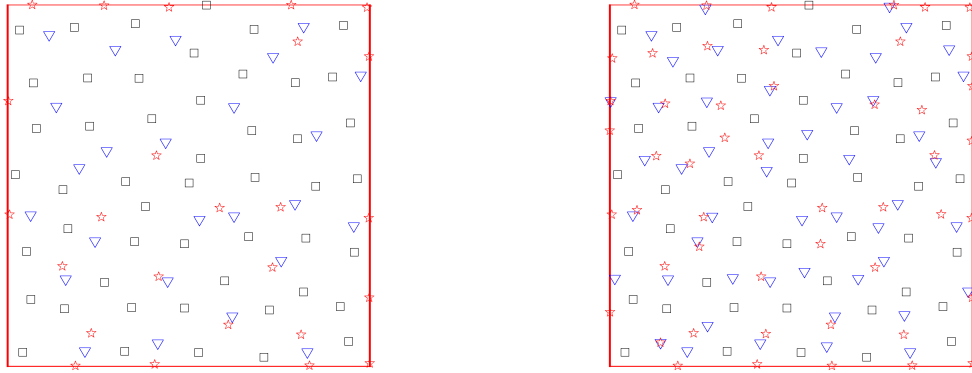


Figure 4: Designs $\mathbf{X}_n = \text{KH}(\emptyset, K, n)$ (black \square), $\text{KH}(\mathbf{X}_n, K, m)$ (blue ∇) and $\text{MN}_2(\emptyset, \bar{K}|n, m)$ (red \star) for $n = 50$ and $m = 25$ (left), $m = 50$ (right).

The quantitative comparison below of the space-filling properties of the different designs considered relies on their covering and packing (or separating) radii, respectively defined by

$$\text{CR}(\mathbf{X}_s) = \max_{\mathbf{x} \in \mathcal{X}} \min_{1 \leq i \leq s} \|\mathbf{x} - \mathbf{x}_i\| \text{ and } \text{PR}(\mathbf{X}_s) = \frac{1}{2} \min_{i \neq j} \|\mathbf{x}_i - \mathbf{x}_j\|$$

when $\mathbf{X}_s = \{\mathbf{x}_1, \dots, \mathbf{x}_s\}$. When $d \leq 4$, the exact value of $\text{CR}(\mathbf{X}_s)$ is calculated by Voronoi tessellation (Pronzato, 2017); when $d > 4$, we under-approximate $\text{CR}(\mathbf{X}_s)$ by $\max_{\mathbf{x} \in \mathcal{X}_{Q'}}$ $\min_{1 \leq i \leq s} \|\mathbf{x} - \mathbf{x}_i\|$, with $\mathcal{X}_{Q'}$ given by the first 2^{19} points of a scrambled Sobol' sequence complemented with a 3^d full factorial design (so that $Q' = 2^{19} + 3^d$).

Figure 5 presents the values of CR and PR (multiplied by $s^{1/d}$ for a design of size s) obtained for $\mathbf{X}_n = \text{KH}(\emptyset, K, n)$ (black squares), $\text{KH}(\mathbf{X}_n, K, m)$ (blue triangles), and $\text{KH}(\emptyset, \bar{K}|n, m)$ (red

circles) and $\text{MN}_2(\emptyset, \overline{K}_{|n}, m)$ (red stars), with $n = m = 50$ on the left column and $n = 200$, $m = 100$ on the right. The magenta diamonds correspond to \mathbf{S}_m given by the first m points of a scrambled Sobol' sequence.

The designs constructed by kernel herding and its variants have good space-filling performance, typically better, and often much better, than Sobol' points \mathbf{S}_m . On the left column $m = n$, and we can directly compare the space-filling performance of \mathbf{Z}_m and \mathbf{X}_n . The good space-filling properties of $\mathbf{X}_n = \text{KH}(\emptyset, K, n)$ tend to deteriorate when considering its continuation $\text{KH}(\mathbf{X}_n, K, m)$. Some other constructions sometimes compare favorably to \mathbf{X}_n in terms of CR, or PR, or even both. It is true in particular for $\text{KH}(\emptyset, \overline{K}_{|n}, k, m, \setminus \mathbf{X}_n)$, see (i) in Section 3.2, for which we continue iterations until number of selected points not in \mathbf{X}_n equals m : it is almost uniformly better than $\text{KH}(\mathbf{X}_n, K, m)$. This opens interesting perspectives in terms of construction of space-filling designs. $\text{MN}_2(\emptyset, \overline{K}_{|n}, m)$ performs significantly worse; its rather poor space-filling properties were already apparent on Figure 4.

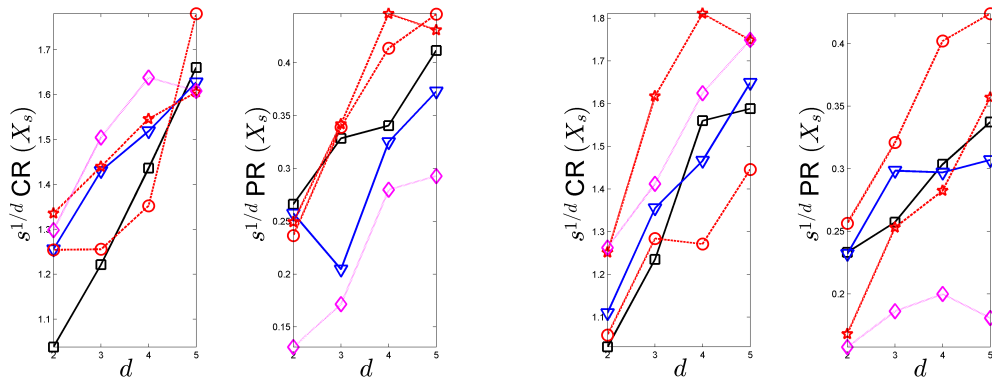


Figure 5: Renormalized values of CR and PR for $\mathbf{X}_n = \text{KH}(\emptyset, K, n)$ (black \square), $\text{KH}(\mathbf{X}_n, K, m)$ (blue ∇), $\text{KH}(\emptyset, \overline{K}_{|n}, k, m, \setminus \mathbf{X}_n)$, see (i) in Section 3.2 (red \circ), and $\text{MN}_2(\emptyset, \overline{K}_{|n}, m)$ (red \star); first m points \mathbf{S}_m of a scrambled Sobol' sequence (magenta \diamond). Left column: $n = m = 50$; right column: $n = 200$, $m = 100$.

b) IMSE and ISE performance. Below we compare the values of $\overline{\Delta}(\mathbf{Z}_m, \mathbf{X}_n)$ given by (4) for different designs \mathbf{Z}_m . For all designs considered, weighted or not, $\overline{\Delta}(\mathbf{Z}_m, \mathbf{X}_n)$ is computed directly from the associated measure ζ_m as $\mathcal{E}_{\overline{K}_{|n}}^{1/2}(\zeta_m - \mu_Q)$, see (4) and (10).

We also consider

$$\Delta(\mathbf{Z}_m, \mathbf{X}_n) = |\widehat{\text{IMSE}}(\mathbf{Z}_m, \mathbf{X}_n) - \text{IMSE}(\mathbf{X}_n)|,$$

where

$$\widehat{\text{IMSE}}(\mathbf{Z}_m, \mathbf{X}_n) = \frac{1}{m} \sum_{i=1}^m \mathbb{E} \left\{ \left[\mathbf{k}_n^\top(\mathbf{z}_i) \mathbf{y}_n - f(\mathbf{z}_i) \right]^2 \right\} = \frac{1}{m} \sum_{i=1}^m K_{|n}(\mathbf{z}_i, \mathbf{z}_i) \quad (17)$$

and $\text{IMSE}(\mathbf{X}_n)$ is given by (3), which is approximated by a discrete sum $\widehat{\text{IMSE}}(\mathcal{X}_{Q''}, \mathbf{X}_n)$, where $\mathcal{X}_{Q''}$ corresponds to $Q'' = 2^{19}$ points of a scrambled Sobol' sequence. When some weights $\mathbf{w}^{(m)}$

are associated with \mathbf{Z}_m , with $\mathbf{w}^{(m)} = \hat{\mathbf{w}}(\mathbf{Z}_m, C)$ or $\mathbf{w}^{(m)} = \check{\mathbf{w}}(\mathbf{Z}_m, C)$ for a kernel C , see (11) and (13), we use

$$\widehat{\text{IMSE}}([\mathbf{Z}_m, \mathbf{w}^{(m)}], \mathbf{X}_n) = \sum_{i=1}^m \{\mathbf{w}^{(m)}\}_i K_{|n}(\mathbf{z}_i, \mathbf{z}_i)$$

in $\Delta(\mathbf{Z}_m, \mathbf{X}_n)$ instead of (17). The minimization of $\bar{\Delta}(\mathbf{Z}_m, \mathbf{X}_n)$ is not equivalent to that of $\Delta(\mathbf{Z}_m, \mathbf{X}_n)$, and we shall see that the designs constructed for the former are not necessarily the most efficient for the latter. Note that the evaluation of $\text{IMSE}(\mathbf{X}_n)$ is much easier than that of $\text{ISE}(\mathbf{X}_n)$; see, e.g., (Gauthier and Pronzato, 2014, 2016, 2017) for the construction of designs \mathbf{X}_n that minimize $\text{IMSE}(\mathbf{X}_n)$.

Performances in terms of $\Delta(\mathbf{Z}_m, \mathbf{X}_n)$ are shown on Figure 6, with Sobol' points \mathbf{S}_m corresponding to magenta diamonds and $\mathbf{Z}_m = \text{KH}(\mathbf{X}_n, K, m)$ to blue triangles down. The designs $\text{KH}(\emptyset, \bar{K}_{|n}, k, m, \setminus \mathbf{X}_n)$ correspond to red circles and $[\mathbf{Z}_m'', \check{\mathbf{w}}^{(m)}] = \text{MN}_2(\emptyset, \bar{K}_{|n}, m)$ to red stars; $m = n = 50$ on the left column, $n = 200$ and $m = 100$ on the right. In $\text{KH}(\emptyset, \bar{K}_{|n}, k, m, \setminus \mathbf{X}_n)$, all points receive the same weight $1/k$ with $k > m$; see (i) in Section 3.2. $\Delta(\mathbf{S}_m, \mathbf{X}_n)$ is much smaller than the values obtained for $\text{KH}(\mathbf{X}_n, K, m)$. This could be anticipated from Figures 2. It is related to the stronger variability of $K_{|n}(\mathbf{z}_i, \mathbf{z}_i)$ for Sobol' points, which are distributed independently of \mathbf{X}_n , than for the designs \mathbf{Z}_m constructed by kernel herding, which tend to fill the holes left by \mathbf{X}_n . For those designs, each \mathbf{z}_i is selected far away from its closest \mathbf{x}_j , all $K_{|n}(\mathbf{z}_i, \mathbf{z}_i)$ tend to be large and $\widehat{\text{IMSE}}(\mathbf{Z}_m, \mathbf{X}_n)$ tends to severely overestimate $\text{IMSE}(\mathbf{X}_n)$. The designs constructed with $\bar{K}_{|n}$ compensate this effect by weight reduction and behave more similarly to Sobol' points: in $\text{KH}(\emptyset, \bar{K}_{|n}, k, m, \setminus \mathbf{X}_n)$ all points receive the same weight $1/k < 1/m$; in $[\mathbf{Z}_m'', \check{\mathbf{w}}^{(m)}] = \text{MN}_2(\emptyset, \bar{K}_{|n}, m)$ the total mass is smaller than one.

Consider now our criterion of interest, $\bar{\Delta}(\mathbf{Z}_m, \mathbf{X}_n)$. The same symbols as above are used to represent the different designs, but two more designs are considered: $[\mathbf{S}_m, \check{\mathbf{w}}(\mathbf{S}_m, \bar{K}_{|n})]$ with magenta plus and $[\text{KH}(\mathbf{X}_n, K, m), \check{\mathbf{w}}(\text{KH}(\mathbf{X}_n, K, m), \bar{K}_{|n})]$ with blue triangles up. We can see that the introduction of weights $\check{\mathbf{w}}(\mathbf{Z}_m)$ has a major effect on the reduction of $\bar{\Delta}(\mathbf{Z}_m, \mathbf{X}_n)$; $\text{KH}(\emptyset, \bar{K}_{|n}, k, m, \setminus \mathbf{X}_n)$ and $[\mathbf{Z}_m'', \check{\mathbf{w}}^{(m)}] = \text{MN}_2(\emptyset, \bar{K}_{|n}, m)$ have very good performance too.

Figure 7 shows the total mass $\sum_{i=1}^m \check{w}_i$ for the designs $[\mathbf{S}_m, \check{\mathbf{w}}(\mathbf{S}_m, \bar{K}_{|n})]$ (magenta plus) $[\mathbf{Z}_m, \check{\mathbf{w}}(\mathbf{Z}_m, \bar{K}_{|n})]$ (blue triangles up), $[\mathbf{Z}_m'', \check{\mathbf{w}}^{(m)}] = \text{MN}_2(\emptyset, \bar{K}_{|n}, m)$ (red stars) and m/k for the design $\text{KH}(\emptyset, \bar{K}_{|n}, k, m, \setminus \mathbf{X}_n)$ (red circles); $m = n = 50$ on the left column, $n = 200$ and $m = 100$ on the right. There is no strict relation between total mass and performance in terms of $\bar{\Delta}(\mathbf{Z}_m, \mathbf{X}_n)$ shown on Figure 6, indicating that it is the interplay between the location of the points and the weighing that matters. Note in particular that $[\mathbf{Z}_m, \check{\mathbf{w}}(\mathbf{Z}_m, \bar{K}_{|n})]$ and $[\mathbf{Z}_m'', \check{\mathbf{w}}^{(m)}] = \text{MN}_2(\emptyset, \bar{K}_{|n}, m)$ have quite different weightings although they have similar values of $\bar{\Delta}(\mathbf{Z}_m, \mathbf{X}_n)$ on Figure 6.

5 Examples of validation design for ISE estimation

5.1 Separable kernels

The substitution of a finite set \mathcal{X}_Q for \mathcal{X} and of the uniform measure on \mathcal{X}_Q for μ yields a drastic simplification of calculations in the evaluation of the MMD $\gamma_{\bar{K}_{|n}}(\zeta_m, \mu)$ and in the algorithmic construction of designs by kernel herding and its variants. However, for large d we

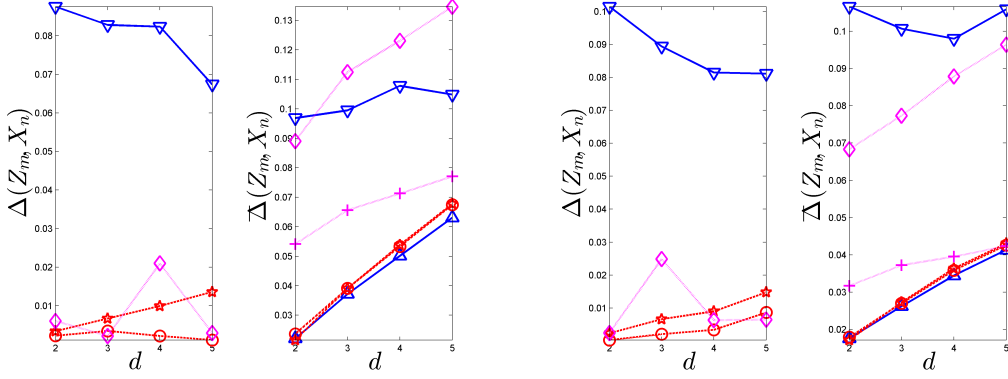


Figure 6: $\Delta(\mathbf{Z}_m, \mathbf{X}_n)$ and $\bar{\Delta}(\mathbf{Z}_m, \mathbf{X}_n)$ for $\mathbf{Z}_m = \text{KH}(\mathbf{X}_n, K, m)$ (blue ∇) and $[\mathbf{Z}_m, \check{\mathbf{w}}(\mathbf{Z}_m, \bar{K}_{|n})]$ (blue \triangle), $\text{KH}(\emptyset, \bar{K}_{|n}, k, m, \setminus \mathbf{X}_n)$ with weights $1/k$ (red \circ) and $[\mathbf{Z}_m'', \check{\mathbf{w}}^{(m)}] = \text{MN}_2(\emptyset, \bar{K}_{|n}, m)$ (red \star); first m points \mathbf{S}_m of a scrambled Sobol' sequence (magenta \diamond), $[\mathbf{S}_m, \check{\mathbf{w}}(\mathbf{S}_m, \bar{K}_{|n})]$ (magenta $+$). Left column: $n = m = 50$; right column: $n = 200, m = 100$.

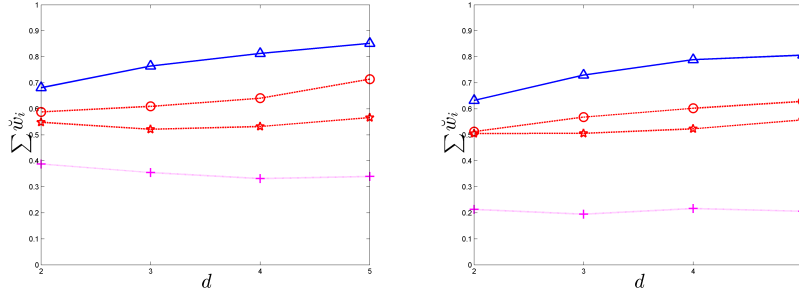


Figure 7: $\sum_{i=1}^m \check{w}_i$ for $[\mathbf{Z}_m, \check{\mathbf{w}}(\mathbf{Z}_m, \bar{K}_{|n})]$ (blue \triangle), $[\mathbf{Z}_m'', \check{\mathbf{w}}^{(m)}] = \text{MN}_2(\emptyset, \bar{K}_{|n}, m)$ (red \star) and $[\mathbf{S}_m, \check{\mathbf{w}}(\mathbf{S}_m, \bar{K}_{|n})]$ (magenta $+$), with \mathbf{S}_m given by the first m points of a scrambled Sobol' sequence. For $\text{KH}(\emptyset, \bar{K}_{|n}, k, m, \setminus \mathbf{X}_n)$ (red \circ), the total mass equals m/k . Left column: $n = m = 50$; right column: $n = 200, m = 100$.

need to take Q very large to make \mathcal{X}_Q dense enough in \mathcal{X} , and another approach is required if we want to maintain a reasonable accuracy.

A bottleneck in the application of kernel herding is the need to calculate $P_{\bar{K}_{|n}, \mu}(\mathbf{z})$ for many \mathbf{z} in order to choose \mathbf{z}_{k+1} in (9). An additional difficulty for the evaluation of $\gamma_{\bar{K}_{|n}}(\zeta_m, \mu)$ is the need to compute $\mathcal{E}_{\bar{K}_{|n}}(\mu)$, see (10). However, when K is a separable (tensor-product) kernel, both $P_{\bar{K}_{|n}, \mu}$ and $\mathcal{E}_{\bar{K}_{|n}}(\mu)$ can be calculated explicitly.

Since μ is uniform on $\mathcal{X} = [0, 1]^2$, we can write $\mu(d\mathbf{x}) = \prod_{i=1}^d \mu_1(dx_i)$ with μ_1 the uniform measure on $[0, 1]$ and $\mathbf{x} = (x_1, \dots, x_d)^\top$. For a separable (or tensor-product) kernel K , such that

$$K(\mathbf{x}, \mathbf{x}') = \prod_{i=1}^d K_i(x_i, x'_i),$$

where $\mathbf{x} = (x_1, \dots, x_d)^\top$ and $\mathbf{x}' = (x'_1, \dots, x'_d)^\top$, we have

$$\mathcal{E}_K(\mu) = \prod_{i=1}^d \mathcal{E}_{K_i}(\mu_1) \text{ and } P_{K,\mu}(\mathbf{x}) = \prod_{i=1}^d \int_{\mathcal{X}_i} K_i(x_i, x'_i) \mu_1(dx'_i) = \prod_{i=1}^d P_{K_i,\mu_1}(x_i).$$

One may refer to Szabó and Sriperumbudur (2018) for connections between positive-definiteness properties of the K_i and those of K . The expressions of $\mathcal{E}_{K_i}(\mu_1)$ and $P_{K_i,\mu_1}(\cdot)$ are available for many kernels K_i ; see Pronzato and Zhigljavsky (2020) and the references therein.

Before deriving the expressions of $P_{\bar{K}|n,\mu}(\mathbf{x})$ and $\mathcal{E}_{\bar{K}|n}(\mu)$, we introduce some notation. Denote by $\bar{\mathbf{\Omega}}_{K,n}$ and $\bar{\mathbf{\Gamma}}_{K,n}$ the $n \times n$ matrices with respective elements

$$\{\bar{\mathbf{\Omega}}_{K,n}\}_{j,k} = \prod_{i=1}^d \beta_{K_i}(x_{j_i}, x_{k_i}) \text{ and } \{\bar{\mathbf{\Gamma}}_{K,n}\}_{j,k} = \prod_{i=1}^d \gamma_{K_i}(x_{j_i}, x_{k_i}),$$

and by $\bar{\omega}_{K,n}(\mathbf{x})$ the vector with j -th component

$$\{\bar{\omega}_{K,n}(\mathbf{x})\}_j = \prod_{i=1}^d \beta_{K_i}(x_{j_i}, x_i),$$

where x_{j_i} (respectively, x_{k_i}) is the i -th component of \mathbf{x}_j (respectively, \mathbf{x}_k), and

$$\begin{aligned} \beta_{K_i}(r, s) &= \int_{\mathcal{X}} K_i(r, t) K_i(s, t) \mu_1(dt), \quad i = 1, \dots, d, \\ \gamma_{K_i}(r, s) &= \int_{\mathcal{X}^2} K_i(r, t) K_i(s, u) K_i(t, u) \mu_1(dt) \mu_1(du), \quad i = 1, \dots, d. \end{aligned}$$

Then, using (5), direct calculation gives

$$\begin{aligned} P_{\bar{K}|n,\mu}(\mathbf{x}) &= 2 P_{K^2,\mu}(\mathbf{x}) - 4 \mathbf{k}_n^\top(\mathbf{x}) \mathbf{K}_n^{-1} \bar{\omega}_{K,n}(\mathbf{x}) + 2 \mathbf{k}_n^\top(\mathbf{x}) \mathbf{K}_n^{-1} \bar{\mathbf{\Omega}}_{K,n} \mathbf{K}_n^{-1} \mathbf{k}_n(\mathbf{x}) \\ &\quad + \left[1 - \mathbf{k}_n^\top(\mathbf{x}) \mathbf{K}_n^{-1} \mathbf{k}_n(\mathbf{x}) \right] \left[1 - \text{trace}(\mathbf{K}_n^{-1} \bar{\mathbf{\Omega}}_{K,n}) \right], \\ \mathcal{E}_{\bar{K}|n}(\mu) &= 2 \mathcal{E}_{K^2}(\mu) - 4 \text{trace}(\mathbf{K}_n^{-1} \bar{\mathbf{\Gamma}}_{K,n}) + 2 \text{trace}[(\mathbf{K}_n^{-1} \bar{\mathbf{\Gamma}}_{K,n})^2] + \left[1 - \text{trace}(\mathbf{K}_n^{-1} \bar{\mathbf{\Omega}}_{K,n}) \right]^2. \end{aligned}$$

The expressions of $P_{K^2,\mu_1}(x)$, $\mathcal{E}_{K^2}(\mu_1)$, $\beta_K(u, v)$ and $\gamma_K(u, v)$, $x, u, v \in [0, 1]$, for μ_1 uniform on $[0, 1]$ and $K_i(x, x')$ a Matérn 3/2 kernel (12) are given in Appendix B, making the expressions of $P_{\bar{K}|n,\mu}(\mathbf{x})$ and $\mathcal{E}_{\bar{K}|n}(\mu)$ available in closed form when $K(\mathbf{x}, \mathbf{x}')$ is the product of uni-dimensional Matérn 3/2 kernels and μ is uniform on $\mathcal{X} = [0, 1]^d$. Similar calculations can be conducted for other kernels.

5.2 Numerical results

We use test functions given by random multivariate polynomials in dimension $d = 2, \dots, 10$, with $n = 100$ and $m = 50$, generated as indicated in Appendix C, the set \mathbb{L} in (19) being constrained by $N = n/2$, $p = 7$ and $p_T = 25$. We take $\alpha = 1/2$ in (20), $\lambda_i = 1/[(i+1)^2 \tau^i]$, where $\tau = \max_{i=1, \dots, d} \sum_{j=1}^d |\{\mathbf{Q}\}_{i,j}|$ (the renormalization by τ^i accounts for the fact that points $\{\mathbf{Q}(\mathbf{x} - \mathbf{1}_d/2) + \mathbf{1}_d/2\}_i$ do not belong to $[0, 1]$).

For each $d = 2, \dots, 10$, we generate $r = 100$ random functions $f^{(j)}$, $j = 1, \dots, r$. For each $f^{(j)}$, \mathbf{X}_n corresponds to the first n points of a scrambled Sobol' sequence, the next m points of the sequence are denoted \mathbf{S}_m and form one of the validation designs considered in the comparison. The second design considered is $\mathbf{Z}_m = \text{KH}(\mathbf{X}_n, K, m)$, constructed by kernel herding with a candidate set \mathcal{X}_Q given by the first $Q = 2^{16}$ points of another scrambled Sobol' sequence. We also consider random designs \mathbf{R}_m made of m points independently uniformly distributed in $[0, 1]^d$. A different design \mathbf{X}_n , \mathbf{S}_m , \mathbf{Z}_m , \mathbf{R}_m and candidate set \mathcal{X}_Q is used for each random $f^{(j)}$ generated, but we omit the index j in the notation. The kernel K is the tensor product of univariate Matérn 3/2 kernels (12). The construction of $\text{KH}(\mathbf{X}_n, K, m)$ by kernel herding and the computation of the weights $\check{\mathbf{w}}_m$ given by (13) exploit the results of Section 5.1.

We set $\theta = n^{1/d}$ in (12) to construct $\text{KH}(\mathbf{X}_n, K, m)$ (it is the space-filling property of \mathbf{Z}_m that matters here), but to estimate $\text{ISE}(\mathbf{X}_n)$ we use $\theta = \theta_n^{(j)}$ estimated by Leave-One-Out Cross Validation (LOO CV) applied to the centered data $\tilde{\mathbf{y}}_n^{(j)} = \mathbf{y}_n^{(j)} - \bar{y}_n^{(j)} \mathbf{1}_n$, with $\bar{y}_n^{(j)} = \mathbf{1}_n^\top \mathbf{y}_n^{(j)} / n$ the empirical mean of $\mathbf{y}_n^{(j)} = (f^{(j)}(\mathbf{x}_1), \dots, f^{(j)}(\mathbf{x}_n))$. Following Dubrule (1983), $\theta_n^{(j)}$ minimizes

$$\widehat{\text{ISE}}_{\text{LOO}}(j) = \frac{1}{n} \sum_{i=1}^n [f^{(j)}(\mathbf{x}_i) - \eta_{n,-i}^{(j)}(\mathbf{x}_i)]^2 = \frac{1}{n} (\tilde{\mathbf{y}}_n^{(j)})^\top \mathbf{K}_n^{-1} \mathbf{D}_n \mathbf{K}_n^{-1} \tilde{\mathbf{y}}_n^{(j)} \quad (18)$$

with respect to $\theta \in \mathbb{R}^+$, where $\eta_{n,-i}^{(j)}(\mathbf{x})$ uses the $n-1$ points in $\mathbf{X}_n \setminus \{\mathbf{x}_i\}$ and \mathbf{D}_n is the diagonal matrix with elements $\{\mathbf{D}_n\}_{i,i} = \{\mathbf{K}_n^{-1}\}_{i,i}^{-2}$; \mathbf{K}_n depends on θ through (12).

The exact value of $\text{ISE}^{(j)} = \text{ISE}^{(j)}(\mathbf{X}_n)$ given by (1) is approximated by a discrete sum, with the uniform measure on \mathcal{X}_Q substituted for μ . For each one of the designs \mathbf{S}_m , \mathbf{Z}_m and \mathbf{R}_m we compute the optimum weights $\check{\mathbf{w}}$ for the kernel $\bar{K}|_n$, and for each $f^{(j)}$ we compute

$$\widehat{\text{ISE}}^{(j)}(\mathbf{Z}_m, \mathbf{X}_n) = \frac{1}{m} \sum_{i=1}^m [f^{(j)}(\mathbf{z}_i) - \eta_n^{(j)}(\mathbf{z}_i)]^2, \quad \widehat{\text{ISE}}^{(j)}(\mathbf{Z}_m, \check{\mathbf{w}}_m, \mathbf{X}_n) = \sum_{i=1}^m \check{w}_i [f^{(j)}(\mathbf{z}_i) - \eta_n^{(j)}(\mathbf{z}_i)]^2,$$

for the unweighted and weighted design, respectively, where $\eta_n^{(j)}(\mathbf{x}) = \mathbf{k}_n^\top(\mathbf{x}) \mathbf{K}_n^{-1} \tilde{\mathbf{y}}_n^{(j)} + \bar{y}_n^{(j)} \mathbf{1}_n$ and $\theta = \theta_n^{(j)}$ in \mathbf{k}_n and \mathbf{K}_n . For each function $f^{(j)}$ and each design, weighted or not, we denote by $\rho^{(j)}$ the relative error.

$$\rho^{(j)} = \frac{\widehat{\text{ISE}}^{(j)} - \text{ISE}^{(j)}}{\text{ISE}^{(j)}}$$

The left panel of Figure 8 presents the empirical means $\widehat{\text{E}}\{|\rho^{(j)}|\} = (1/r) \sum_{j=1}^r |\rho^{(j)}|$ as functions of $d = 2, \dots, 10$ obtained for the different designs considered and for $\widehat{\text{ISE}}_{\text{LOO}}(j)$ given by (18); the right panel shows $\widehat{\text{E}}\{\rho^{(j)}\} = (1/r) \sum_{j=1}^r \rho^{(j)}$. Unsurprisingly, LOO CV strongly overestimates $\text{ISE}(\mathbf{X}_n)$ since (i) each of the n predictions in the summation in (18) uses $n-1$ design points only, and (ii) each \mathbf{x}_i is far from the $n-1$ other design points. The superiority of the weighted design $[\mathbf{Z}_m, \check{\mathbf{w}}_m(\mathbf{Z}_m, \bar{K}|_n)]$ over the other ones is clear on the left panel; in particular weight reduction by $\check{\mathbf{w}}_m(\mathbf{Z}_m, \bar{K}|_n)$ greatly improves the precision of ISE estimation, compare the two curves with triangles. Sobol' and random points behave similarly, with slightly better performance for Sobol' points in the weighted versions. The right panel provides information on the bias on ISE estimation. The three weighted designs (solid lines) underestimate $\text{ISE}(\mathbf{X}_n)$. The unweighted random design \mathbf{R}_m (dashed line with circles) has a small bias, but is of limited

interest due to its large variability, as shown on the left panel. Both \mathbf{S}_m and \mathbf{Z}_m tend to fill the holes left by \mathbf{X}_n and therefore overestimate $\text{ISE}(\mathbf{X}_n)$. Other designs, in particular based on the kernel $\bar{K}_{|n}$ have also been considered, but they perform worse than $[\mathbf{Z}_m, \check{\mathbf{w}}_m(\mathbf{Z}_m, \bar{K}_{|n})]$ and the results are not shown.

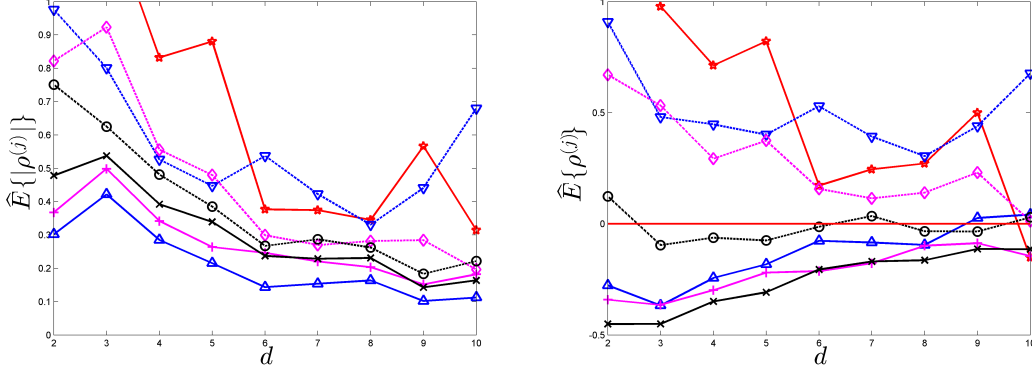


Figure 8: Left: $\hat{E}\{\rho^{(j)}\}$; Right: $\hat{E}\{\rho^{(j)}\}$; for $d = 2, \dots, 10$ and the designs \mathbf{S}_m (magenta \diamond), $[\mathbf{S}_m, \check{\mathbf{w}}_m(\mathbf{S}_m, \bar{K}_{|n})]$ (magenta $+$), $\mathbf{Z}_m = \text{KH}(\mathbf{X}_n, K, m)$ (blue ∇), $[\mathbf{Z}_m, \check{\mathbf{w}}_m(\mathbf{Z}_m, \bar{K}_{|n})]$ (blue \triangle), \mathbf{R}_m (black \circ), $[\mathbf{R}_m, \check{\mathbf{w}}_m(\mathbf{R}_m, \bar{K}_{|n})]$ (black \times); LOO CV (red \star); 100 repetitions, $n = 100$, $m = 50$.

6 Conclusions

The construction of a validation design \mathbf{Z}_m aimed at estimating $\text{ISE}(\mathbf{X}_n)$ for a given \mathbf{X}_n can be casted as the choice of a design minimizing a maximum mean discrepancy for a particular kernel, conditional of \mathbf{X}_n . A sequence of nested validation designs can be obtained by incremental construction via kernel herding. Numerical experiments indicate that the most important characteristics of a good validation design are its space-filling properties (it should populate the holes left by \mathbf{X}_n to properly explore the design space) and the weighting of its points (since evaluations far from the design points tend to overestimate the global error). What one would expect is that some combination of both is needed: if the validation points would sample the error well, no weighting would be needed; if the points were space-filling, and the design very regular, a constant weight smaller than one would be almost optimal. In fact these factors play in antagonistic directions and some compromise is needed. A dedicated weighting method, based on a particular kernel, conditional on \mathbf{X}_n , has been proposed. Numerical simulations with random functions show the effectiveness of this weight reduction when it is applied to random or usual low-discrepancy designs. Performances are still better when the weight reduction is associated with a space-filling design that minimizes a kernel discrepancy: they are significantly better than those obtained with leave-one-out cross validation, which strongly overestimates $\text{ISE}(\mathbf{X}_n)$.

Appendix A: Characteristic kernels

A characteristic kernel C defines a metric on the set of probability measures on \mathcal{X} . It is called *Integrally Strictly Positive Definite* (ISPD) when $\mathcal{E}_C(\nu) > 0$ for any nonzero signed

measure ν on \mathcal{X} , see (6), and *Conditionally Integrally Strictly Positive Definite* (CISPD) when $\mathcal{E}_C(\nu) > 0$ for all nonzero signed measures ν on \mathcal{X} with total mass $\nu(\mathcal{X}) = 0$. An ISPD kernel is CISPD; a bounded ISPD kernel is SPD and defines an RKHS. When C is uniformly bounded, it is characteristic if and only if it is CISPD; see (Sriperumbudur et al., 2010, Lemma 8). For instance, the isotropic squared exponential, Matérn and generalized multiquadric kernels are ISPD.

The kernel $\overline{K}_{|n}$ considered in this paper is positive definite but not strictly positive definite (and thus not ISPD). Indeed, $K_{|n}(\mathbf{x}, \mathbf{x}_i) = \overline{K}_{|n}(\mathbf{x}, \mathbf{x}_i) = 0$ for all \mathbf{x} and all \mathbf{x}_i , $i = 1, \dots, n$, implying that $\gamma_{\overline{K}_{|n}}(\zeta, \xi) = 0$ for any measures ζ and ξ supported on \mathbf{X}_n . Since μ is uniform and not supported on \mathbf{X}_n , it nevertheless makes sense to minimize $\gamma_{\overline{K}_{|n}}(\zeta_m, \mu)$. The investigation of conditions under which $\gamma_{\overline{K}_{|n}}(\zeta, \mu) = 0$ would imply $\zeta = \mu$, exploiting for instance the notion of universal kernel (Sriperumbudur et al., 2011), is beyond the scope of this paper and we simply mention the following two points, concerning respectively $K_{|n}$ and $\overline{K}_{|n}$.

- (i) Suppose that K is ISPD. For any signed measure ξ on \mathcal{X} and $\mathbf{w} \in \mathbb{R}^n$, define $\xi^{[\mathbf{w}]} = \xi + \sum_{i=1}^n w_i \delta_{\mathbf{x}_i}$. Then, using the notation of Section 3.1, $\mathcal{E}_K(\xi^{[\mathbf{w}]} - \mu) = \mathcal{E}_K(\xi - \mu) + \mathbf{w}^\top \mathbf{K}_n \mathbf{w} + 2 \mathbf{w}^\top \mathbf{p}_{K,n}(\xi - \mu) \geq 0$, with equality if and only if $\xi^{[\mathbf{w}]} = \mu$ since K is ISPD. Direct calculation gives $\min_{\mathbf{w}} \mathcal{E}_K(\xi^{[\mathbf{w}]} - \mu) = \mathcal{E}_{K_{|n}}(\xi - \mu)$, and therefore $\mathcal{E}_{K_{|n}}(\xi - \mu) \geq 0$ with equality if and only if $\xi^{[\mathbf{w}]} = \mu$. As μ has no discrete component, $\xi^{[\mathbf{w}]} = \mu$ implies $\mathbf{w} = \mathbf{0}$, and we get $\xi = \mu$.
- (ii) Suppose now that K is ISPD and continuous on \mathcal{X} and consider its Mercer decomposition: $K(\mathbf{x}, \mathbf{x}') = \sum_{i \geq 1} \lambda_i \varphi_i(\mathbf{x}) \varphi_i(\mathbf{x}')$, $\lambda_i > 0$. It yields the following decomposition for K^2 : $K^2(\mathbf{x}, \mathbf{x}') = \sum_{i,j \geq 1} \lambda_i \lambda_j \varphi_i(\mathbf{x}) \varphi_j(\mathbf{x}) \varphi_i(\mathbf{x}') \varphi_j(\mathbf{x}')$, and $\mathcal{E}_{K^2}(\nu) = 0$ for some signed measure ν on \mathcal{X} implies that $\int_{\mathcal{X}} \varphi_i(\mathbf{x}) \varphi_j(\mathbf{x}) \nu(d\mathbf{x}) = 0$ for all i and j . When the constant 1 belongs to the RKHS \mathcal{H}_K associated with K , there exists constants α_i , $i \geq 1$, such that $\sum_{i \geq 1} \alpha_i \varphi_i(\mathbf{x}) = 1$ for all $\mathbf{x} \in \mathcal{X}$, and $\mathcal{E}_{K^2}(\nu) = 0$ implies that $\int_{\mathcal{X}} \varphi_j(\mathbf{x}) \nu(d\mathbf{x}) = 0$ for all $j \geq 1$. Therefore, $\mathcal{E}_K(\nu) = 0$, and $\nu = 0$ since K is ISPD. However, this argumentation cannot be combined with (i) above to show that $\mathcal{E}_{K_{|n}}(\xi - \mu) = 0$ implies that $\xi = \mu$ since $1 \notin \mathcal{H}_{K_{|n}}$: indeed, $f(\mathbf{x}_i) = 0$ for any $f \in \mathcal{H}_{K_{|n}}$ and any $\mathbf{x}_i \in \mathbf{X}_n$.

Appendix B: expressions of $P_{K_i^2, \mu_1}(x)$, $\mathcal{E}_{K_i^2}(\mu_1)$, $\beta_{K_i}(u, v)$ and $\gamma_{K_i}(u, v)$ for Matérn 3/2 kernel and μ_1 uniform on $[0, 1]$

When $K_i(x, x') = K_{3/2, \theta/\sqrt{3}}(x, x')$ given in (12), we have (Ginsbourger et al., 2014)

$$\begin{aligned} \mathcal{E}_{K_i}(\mu_1) &= \frac{2}{\theta^2} [(\theta + 3)e^{-\theta} + 2\theta - 3], \\ P_{K_i, \mu_1}(x) &= S_\theta(x) + S_\theta(1 - x), \text{ with } S_\theta(x) = \frac{1}{\theta} [2 - (2 + \theta x)e^{-\theta x}], \quad x \in [0, 1]. \end{aligned}$$

Straightforward but lengthy calculation gives

$$\begin{aligned} \mathcal{E}_{K_i^2}(\mu_1) &= \frac{1}{4\theta^2} [(2\theta^2 + 8\theta + 9)e^{-2\theta} + 10\theta - 9], \\ P_{K_i^2, \mu_1}(x) &= T_\theta(x) + T_\theta(1 - x), \text{ with } T_\theta(x) = \frac{1}{4\theta} [5 - (5 + 6\theta x + 2\theta^2 x^2)e^{-2\theta x}], \quad x \in [0, 1]. \end{aligned}$$

Also, $\beta_{K_i}(u, v) = B_\theta(u, v) - C_\theta(u, v) - C_\theta(1 - u, 1 - v)$, $u, v \in [0, 1]$, with

$$\begin{aligned} B_\theta(u, v) &= \frac{e^{-\theta|u-v|}}{6\theta} [15(1 + \theta|u - v|) + 6\theta^2|u - v|^2 + \theta^3|u - v|^3] , \\ C_\theta(u, v) &= \frac{e^{-\theta(u+v)}}{4\theta} [5 + 3\theta(u + v) + 2\theta^2uv] , \end{aligned}$$

and $\gamma_{K_i}(u, v) = G_\theta(u, 1 - v) + G_\theta(v, 1 - u) - H_\theta(u, v) - H_\theta(1 - u, 1 - v) + I_\theta(u, v)$, $u, v \in [0, 1]$, with

$$\begin{aligned} G_\theta(u, v) &= \frac{e^{-\theta(1+u+v)}}{16\theta^2} \{21 + \theta[9 + 13(u + v)] + \theta^2[6(u + v) + 8uv] + 4\theta^3uv\} , \\ H_\theta(u, v) &= \frac{e^{-\theta(u+v)}}{24\theta^2} \{126 + 96\theta(u + v) + 24\theta^2(u + v)^2 + 3\theta^3(u + v)^3 \\ &\quad + \theta^2uv[24 + 6\theta(u + v) + 2\theta^2(u^2 + v^2)]\} , \\ I_\theta(u, v) &= \frac{e^{-\theta|u-v|}}{120\theta^2} \{945 + 945\theta|u - v| + 420\theta^2|u - v|^2 + 105\theta^3|u - v|^3 + 15\theta^4|u - v|^4 \\ &\quad + \theta^5|u - v|^5\} . \end{aligned}$$

Appendix C: random polynomials

Consider the family of Legendre polynomials, orthonormal for the uniform measure μ_1 on $\mathcal{X}_1 = [0, 1]$:

$$\begin{aligned} P_0(x) &= 1 \\ P_1(x) &= \sqrt{3}(2x - 1) \\ P_2(x) &= \sqrt{5}(6x^2 - 6x + 1) \\ P_3(x) &= \sqrt{7}(20x^3 - 30x^2 + 12x - 1) \\ P_4(x) &= 3(70x^4 - 140x^3 + 90x^2 - 20x + 1) \\ &\vdots \end{aligned}$$

satisfying $\int_0^1 P_i(x)P_j(x)dx = \delta_{i,j}$ (the Kronecker delta). To each P_i we associate a $\lambda_i \in \mathbb{R}^+$, with $\lambda_0 = 1$ and $\lambda_i > \lambda_{i+1}$ for all i . A reasonable choice is $\lambda_i = 1/(i + 1)^\gamma$ for some $\gamma > 0$. Denote by \mathbb{L} a subset of \mathbb{N}^d containing multi-indices $\underline{\ell} = \{\ell_1, \dots, \ell_d\}$, with each $\ell_i \in \mathbb{N}$ pointing to a polynomial P_{ℓ_i} . The multivariate polynomials we consider have the form

$$P(\mathbf{x}) = \sum_{\underline{\ell} \in \mathbb{L}} \beta_{\underline{\ell}} \Psi_{\underline{\ell}}(\mathbf{x}) , \quad (19)$$

where $\Psi_{\underline{\ell}}(\mathbf{x}) = \prod_{i=1}^d P_{\ell_i}(x_i)$ and the $\beta_{\underline{\ell}}$ are independent normal variables $\mathcal{N}(0, \Lambda_{\underline{\ell}})$ with $\Lambda_{\underline{\ell}} = \prod_{i=1}^d \lambda_{\ell_i}$. If we only constrain the maximum degree p in each variable, that is, if we consider all $\underline{\ell}$ with $\ell_i \leq p$ for all i , then \mathbb{L} contains $(p + 1)^d$ elements; if we constrain the total degree p_T of $P(\mathbf{x})$, \mathbb{L} has $\binom{p_T+d}{d}$ elements. In both cases, the evaluation of f quickly becomes very costly when d , p or p_T increase. For that reason, we shall set a constraint on the number of elements of \mathbb{L} and only retain the largest $\Lambda_{\underline{\ell}}$; that is, we use

$$\mathbb{L}_N = \{\underline{\ell}_1, \dots, \underline{\ell}_M \in \mathbb{N}^d, \text{ with } M \text{ the smaller integer } \geq N \text{ such that } \Lambda_{\underline{\ell}_M} < \Lambda_{\underline{\ell}_{M+1}}\} ;$$

see Pronzato (2019) for implementation details.

To avoid favouring too much the use of separable kernels, we apply a random linear transformation to \mathbf{x} before computing f and set $f(\mathbf{x}) = P[\mathbf{Q}(\mathbf{x} - \mathbf{1}_d/2) + \mathbf{1}_d/2]$, with

$$\mathbf{Q} = \alpha \mathbf{Q}_R(d) + (1 - \alpha) \mathbf{I}_d, \quad (20)$$

where $\alpha \in [0, 1]$, and $\mathbf{Q}_R(d)$ is a random rotation matrix in the orthogonal group $\mathcal{O}(d - 1)$. To generate a random matrix $\mathbf{Q}_R(d)$ uniformly distributed in $\mathcal{O}(d - 1)$ we proceed as follows, see Diaconis and Shahshahani (1987, Lemma 3.1). For $d = 2$, we take

$$\mathbf{Q}_R(d) = \begin{pmatrix} \cos \theta & \sin \theta \\ a \sin \theta & a \cos \theta \end{pmatrix},$$

with θ uniformly distributed in $[0, 2\pi]$ and $a = \pm 1$ with probability 1/2. For larger d , we construct $\mathbf{Q}_R(d)$ recursively as

$$\mathbf{Q}_R(d) = \left(\mathbf{I}_d - 2 \frac{[\mathbf{e}_1 - \mathbf{u}(d)][\mathbf{e}_1 - \mathbf{u}(d)]^\top}{\|\mathbf{e}_1 - \mathbf{u}(d)\|^2} \right) \begin{pmatrix} 1 & 0 & \cdots & 0 \\ 0 & & & \\ \vdots & & \mathbf{Q}_R(d-1) & \\ 0 & & & \end{pmatrix},$$

with $\mathbf{e}_1 = (1, 0, \dots, 0)^\top$ and $\mathbf{u}(d)$ uniformly distributed on the d dimensional unit sphere (for instance, we can take $\mathbf{u}(d) = \mathbf{v}/\|\mathbf{v}\|$ with \mathbf{v} having the standard normal distribution $\mathcal{N}(\mathbf{0}_d, \mathbf{I}_d)$).

Acknowledgments

This work was partly supported by project INDEX (INcremental Design of EXperiments) ANR-18-CE91-0007 of the French National Research Agency (ANR).

References

- Bach, F., Lacoste-Julien, S., and Obozinski, G. (2012). On the equivalence between herding and conditional gradient algorithms. In *Proc. 29th Annual International Conference on Machine Learning*, pages 1355–1362.
- Bachoc, F. (2013). Cross validation and maximum likelihood estimations of hyper-parameters of Gaussian processes with model misspecification. *Comput. Statist. Data Anal.*, 66:55–69.
- Diaconis, P. and Shahshahani, M. (1987). The subgroup algorithm for generating uniform random variables. *Probability in the Engineering and Informational Sciences*, 1(1):15–32.
- Dubrulle, O. (1983). Cross validation of kriging in a unique neighborhood. *Journal of the International Association for Mathematical Geology*, 15(6):687–699.
- Fedorov, V. (1972). *Theory of Optimal Experiments*. Academic Press, New York.
- Gauthier, B. and Pronzato, L. (2014). Spectral approximation of the IMSE criterion for optimal designs in kernel-based interpolation models. *SIAM/ASA J. Uncertainty Quantification*, 2:805–825. DOI 10.1137/130928534.

- Gauthier, B. and Pronzato, L. (2016). Approximation of IMSE-optimal designs via quadrature rules and spectral decomposition. *Communications in Statistics – Simulation and Computation*, 45(5):1600–1612.
- Gauthier, B. and Pronzato, L. (2017). Convex relaxation for IMSE optimal design in random field models. *Computational Statistics and Data Analysis*, 113:375–394.
- Ginsbourger, D., Roustant, O., Schuhmacher, D., Durrande, N., and Lenz, N. (2014). On ANOVA decompositions of kernels and Gaussian random field paths. *preprint arXiv:1409.6008*.
- Pronzato, L. (2017). Minimax and maximin space-filling designs: some properties and methods for construction. *Journal de la Société Française de Statistique*, 158(1):7–36.
- Pronzato, L. (2019). Sensitivity analysis via Karhunen-Loève expansion of a random field model: estimation of Sobol’ indices and experimental design. *Reliability Engineering and System Safety*, 187:93–109. hal-01545604v2.
- Pronzato, L. (2021). Performance analysis of greedy algorithms for minimising a maximum mean discrepancy. *hal-03114891, arXiv:2101.07564*.
- Pronzato, L. and Zhigljavsky, A. (2020). Bayesian quadrature, energy minimization and space-filling design. *SIAM/ASA J. Uncertainty Quantification*, 8(3):959–1011.
- Santner, T., Williams, B., and Notz, W. (2003). *The Design and Analysis of Computer Experiments*. Springer, Heidelberg.
- Sejdinovic, S., Sriperumbudur, B., Gretton, A., and Fukumizu, K. (2013). Equivalence of distance-based and RKHS-based statistics in hypothesis testing. *The Annals of Statistics*, 41(5):2263–2291.
- Sriperumbudur, B., Fukumizu, K., and Lanckriet, G. (2011). Universality, characteristic kernels and RKHS embedding of measures. *Journal of Machine Learning Research*, 12:2389–2410.
- Sriperumbudur, B., Gretton, A., Fukumizu, K., Schölkopf, B., and Lanckriet, G. (2010). Hilbert space embeddings and metrics on probability measures. *Journal of Machine Learning Research*, 11:1517–1561.
- Szabó, Z. and Sriperumbudur, B. (2018). Characteristic and universal tensor product kernels. *Journal of Machine Learning Research*, 18:1–29.
- Wynn, H. (1970). The sequential generation of D -optimum experimental designs. *Annals of Math. Stat.*, 41:1655–1664.