DESIGN OF EXPERIMENTS VIA MULTI-FIDELITY SURROGATES AND STATISTICAL SENSITIVITY MEASURES

David J. Gillcrist, ¹ Negin Alemazkoor, ² Yanlai Chen, ¹ & Mazdak Tootkaboni^{3,*}

¹Department of Mathematics, University of Massachusetts Dartmouth, North Dartmouth, Massachusetts 02747, USA

²Department of Civil and Environmental Engineering, University of Virginia, Charlottesville, Virginia 22903, USA

³Department of Civil and Environmental Engineering, University of Massachusetts Dartmouth, North Dartmouth, Massachusetts 02747, USA

*Address all correspondence to: Mazdak Tootkaboni, Department of Civil and Environmental Engineering Center for Scientific Computing and Data Science Research (CSCDR), University of Massachusetts Dartmouth, North Dartmouth, MA 02747, USA, E-mail: mtootkaboni@umassd.edu

Parameter estimation and optimal experimental design problems have been widely studied across science and engineering. The two are inextricably linked, with optimally designed experiments leading to better-estimated parameters. This link becomes even more crucial when available experiments produce minimal data due to practical constraints of limited experimental budgets. This work presents a novel framework that allows for the identification of optimal experimental arrangement, from a finite set of possibilities, for precise parameter estimation. The proposed framework relies on two pillars. First, we use multi-fidelity modeling to create reliable surrogates that relate unknown parameters to a measurable quantity of interest for a multitude of available choices defined through a set of candidate control vectors. Secondly, we quantify the "estimation potential" of an arrangement from the set of control vectors through the examination of statistical sensitivity measures calculated from the constructed surrogates. The measures of sensitivity are defined using analysis of variance as well as directional statistics. Two numerical examples are provided, where we demonstrate how the correlation between the estimation potential and the frequency of precise parameter estimation can inform the choice of optimal arrangement.

KEY WORDS: parameter estimation, optimal design of experiments, multi-fidelity surrogates, directional statistics, variance-based sensitivity analysis, multi-variate polynomial-basis conversion

1. INTRODUCTION

Many applications in science and engineering rely on parameterized models. It is often the case that the parameters in these models cannot be estimated through direct measurements. Approaches that permit the estimation of the unknown parameters of a model based on

measurements from other observable quantities are therefore needed. The literature is rife with examples of parameter estimation or model calibration (also known as inverse analysis). These include the estimation of structural and load parameters using data from static diagnostic load testing (Rózsás et al., 2022), soil hydraulic properties using groundwater pressure and soil weight data (Romano and Santini, 1999), diffusion coefficients of a large body of water via dye concentration experiments (Gasca-Ortiz et al., 2021), and groundwater contaminant properties via data collected from tracer experiments (Kowalsky et al., 2012; Lile et al., 1997; Rainwater et al., 1987; Zhang et al., 2002). Source localization problems are also often formulated as inverse analysis problems. Examples include acoustic localization of moving vehicles (Sheng and Hu, 2004; Yao et al., 2002) and animal species (Ali et al., 2007) as well as the localization of radiation sources (Rao et al., 2008; Wu et al., 2019).

Similarly, many science and engineering problems rely on optimal design of experiments to accomplish satisfactory parameter estimation, as is the case in building energy-efficient walls (Jumabekova and Berger, 2023; Jumabekova et al., 2020), in biological cell signaling (Bandara et al., 2009), and in detector placement for soil moisture sensing (Wu et al., 2012). Pronzato (2008) discusses the connection between optimal design of experiments and parameter estimation, emphasizing that optimally designed experiments lead to optimally precise estimations. To our knowledge, however, parameter estimation and the design of experiments with minimal data—below what is advised for contemporary estimation methods—is a topic that has not been directly addressed in the literature.

In this work, we propose a novel framework that addresses the optimal design of experiments restrained by minimal data in a deterministic setting. We view optimal design in a general context, encompassing parameter estimation problems, where the goal is to identify the best experimental program, as well as localization problems, where the goal is to select the best arrangement for detector sensors. The selection of the optimal experimental setting is formulated as making the best choice(s) from a finite collection of available options. The proposed framework relies on (i) multi-fidelity surrogates that provide compact relationships between measurable/observable quantities of interest and the set of unknown parameters at manageable computational cost, and (ii) statistical measures of sensitivity that systematically guide the selection process either through analysis of variance or the directional statistics of the sensitivities (in the sense of derivatives) with respect to unknown parameters. While we lay out all components of the proposed framework, topics such as computational complexity will not be discussed and are a subject of our future research.

The organization of this paper is as follows. In Section 2 we introduce the objective of the paper as making the optimal choice from a finite number of experimental arrangements for the purpose of estimating unknown parameters. Section 3 introduces the proposed framework that draws upon multi-fidelity surrogate modeling in conjunction with the application of different statistical sensitivity measures. This includes a brief overview of multi-fidelity surrogate modeling using polynomial chaos expansions and Kaczmarz updating in Section 3.1, some critical definitions such as that of vector-defined quantity of interest and performance function in Section 3.2, and various strategies to quantify the sensitivity of the surrogates and the estimation potential of an experimental arrangement in Section 3.3, which also includes a presentation of variance-based sensitivity measures and those based on directional statistics of the derivative. Section 4 contains the results from two numerical examples: the optimal choice of borehole experiments to determine properties of deep underground aquifers and the optimal placement of detectors in a river to pinpoint the location of contaminant source. Section 5, finally, contains the concluding remarks.

2. PROBLEM STATEMENT

Let Q be a mathematical model, representing a particular physical phenomenon, that maps a vector of control parameters, referred to henceforth as the control vector \vec{x} , a vector of fixed and known parameters $\vec{\mu}_D$, and a vector of unknown parameters $\vec{\mu}$, to a measurable quantity z, that is:

$$z = Q(\vec{x}; \vec{\mu}_D, \vec{\mu}) \tag{2.1}$$

By an experimental design, we mean the selection of one or more \vec{x} s from a space of candidate control vectors that have the potential to allow for the identification of the unknown parameters (identification exercise, hereafter) from observed values of z. An optimal experimental design is an arrangement of points in the space of candidate control vectors that best serves the identification exercise with acceptable accuracy. We note that the cardinality of the arrangement set is determined by practical constraints and experimental budget. This could mean very limited data, potentially on the order of the number of unknown parameters. In what follows, we lay out the details of an approach that systemically seeks such an optimal arrangement under such restrictive observation/measurement plans. The proposed approach rests on two pillars: (i) multifidelity surrogates that efficiently emulate the mapping Q for the set of candidate control vectors, (ii) statistical measures of sensitivity of the surrogates with respect to the unknown parameters. These measures of sensitivity are used to determine the arrangements most likely to estimate the unknown parameters with good precision in a deterministic setting and with no regard to measurement noise. That is, the estimates of the unknown parameters are simply calculated by minimizing an L^2 -cost functional constructed from the data associated with the selected arrangement.

3. PROPOSED FRAMEWORK

3.1 Multi-Fidelity Surrogate Modeling

An integral part of the proposed framework in this work is surrogate modeling. For the mathematical model Q, the surrogate model serves as an emulator when a quantity of interest (QoI) from the model \tilde{z}_0 —corresponding to a particular realization of the control vector \vec{x}_0 ,

$$\tilde{z}_0(\vec{\mu}) = G(Q(\vec{x}_0; \vec{\mu}_D, \vec{\mu}))$$
 (3.1)

where G is a strictly monotonic function in Eq. (3.1)—needs to be computed many times for different variations of the unknown parameters.

Research in uncertainty quantification (UQ) has heavily utilized surrogate modeling via polynomial chaos expansions (PCEs) (Ghanem and Spanos, 2003; Xiu, 2010; Xiu and Karniadakis, 2002). Here, we utilize one of the most common ways to build PCE surrogates, nonintrusively, through oversampling and least-square regression. This has been done, typically, via Monte Carlo (MC) sampling, which suffers from a poor, sublinear convergence rate of the root mean squared error in estimates of expectations, requiring oversampling of 1.5–3 times the number of basis functions in PCE (Shin and Xiu, 2016). Recent advances in UQ, especially those concerning multi-fidelity modeling (Alemazkoor et al., 2022; Le Gratiet and Cannamela, 2015; Narayan et al., 2014; Parussini et al., 2017; Piazzola et al., 2020; Song et al., 2019; Zhu et al., 2014), however, have allowed for more efficient construction of PCE surrogates. In the multi-fidelity approach to surrogate modeling, computational budget is initially invested in computing a large

number of outputs from a low-fidelity model with inputs that are sampled from the probability density functions of the corresponding unknown parameters, $\vec{\mu}$. The generated data set can be analyzed to select the inputs that contribute minimally to the L^2 -error between a calculated low-fidelity surrogate and the low-fidelity model. Alemazkoor et al. showed that this can be accomplished systematically via a Kaczmarz updating scheme, where the size of the initial pool can be reduced down to the cardinality of the basis set in PCE (Alemazkoor et al., 2022). A new surrogate is then created using the high-fidelity model and points with the highest yield in the parameter space that comprises the reduced pool.

3.1.1 Polynomial Chaos Expansion

The set of independent unknown parameters for the model Q is modeled as a d-dimensional random vector $\mathbf{M}=(M_1,M_2,\ldots,M_d)$. Let M_i be the ith component of this random vector with support I_{M_i} and probability distribution $\eta_i:I_{M_i}\to\mathbb{R}^+$. The support of \mathbf{M} is given by $I_{\mathbf{M}}:= \mathop{\textstyle \times}_{i=1}^d I_{M_i}$, where $\mathop{\textstyle \times}$ represents the tensor product, and the probability distribution of \mathbf{M} is given by $\eta(\mathbf{M}):=\prod_{i=1}^d \eta_i(M_i)$. Without loss of generality, we assume components of \mathbf{M} can be mapped to a set of elementary and independent random variables collected in a random vector $\mathbf{\Xi}$ via a one-to-one mapping $F:\mathcal{D}\to I_{\mathbf{M}}$ with $\mathcal{D}\equiv I_{\mathbf{\Xi}}$ is the tensor product of the supports of the elementary random variables. We then define $u_0(\mathbf{\Xi})\in L^2$ to be the QoI over the domain \mathcal{D} and write it as

$$u_0(\Xi) = \tilde{z}_0(F(\Xi)) = \sum_{\boldsymbol{s}^{(j)} \in \mathbb{N}_a^d} c_{\boldsymbol{s}^{(j)}} \psi_{\boldsymbol{s}^{(j)}}(\Xi)$$
(3.2)

with $\{\psi_{\boldsymbol{s}^{(j)}}\}_{\boldsymbol{s}^{(j)}\in\mathbb{N}_0^d}:\mathcal{D}\to\mathbb{R}$ the set of multivariate orthogonal polynomial basis satisfying the condition

$$\int_{I_{\Xi}} \psi_{s(i)}(\vec{\xi}) \psi_{s(j)}(\vec{\xi}) \rho(\vec{\xi}) d\vec{\xi} = \delta_{ij}$$
(3.3)

where δ_{ij} is the Kronecker delta, and $\rho: I_{\Xi} \to \mathbb{R}^+$ is the probability distribution of Ξ . The exact coefficients in Eq. (3.2) can be calculated by projecting $u_0(\Xi)$ onto the basis functions $\psi_{s^{(j)}}$ via

$$c_{\boldsymbol{s}^{(j)}} = \mathrm{E}_{\rho}[u_0(\boldsymbol{\Xi})\psi_{\boldsymbol{s}^{(j)}}(\boldsymbol{\Xi})] = \int_{I_{\boldsymbol{\Xi}}} u(\vec{\xi})\psi_{\boldsymbol{s}^{(j)}}(\vec{\xi})\rho(\vec{\xi})d\vec{\xi}$$
(3.4)

where E_{ρ} is the expectation operator with respect to $\rho(\Xi)$. The QoI in Eq. (3.2) is often approximated as a truncated PCE

$$u_0(\Xi) \approx u(\Xi) = \sum_{\mathbf{s}^{(j)} \in \mathbb{N}_K^d} c_{\mathbf{s}^{(j)}} \psi_{\mathbf{s}^{(j)}}(\Xi)$$
(3.5)

where $\mathbb{N}_K^d := \{s^{(j)}: \|s^{(j)}\|_1 \leq k_d\}$, k_d is the desired degree of the PCE, and the cardinality K of the basis $\{\psi_{s^{(j)}}: s^{(j)} \in \mathbb{N}_K^d\}$ is given by $K = \frac{(k_d+d)!}{k_d!d!}$. The multidimensional integral in Eq. (3.4) can be difficult to compute for high values of d: the so-called curse of dimensionality. In such cases, the coefficients $c_{s^{(j)}}$ in Eq. (3.5) may be approximated via least-squares regression (Alemazkoor et al., 2022):

$$\vec{c} = \arg\min_{\vec{c}} \left(\left\| \vec{u}_0 - \Psi \vec{c} \right\|_2^2 \right) \tag{3.6}$$

with $\vec{c} := (c_{s^{(0)}}, c_{s^{(1)}}, ..., c_{s^{(K-1)}})^{\mathrm{T}}$ and $\vec{u}_0 := \left(u_0(\vec{\xi}^{\,(1)}), u_0(\vec{\xi}^{\,(2)}), ..., u_0(\vec{\xi}^{\,(N_p)})\right)^{\mathrm{T}}$ for a pool of N_p known input vectors $\{\vec{\xi}^{\,(i)}\}_{i=1}^{N_p}$ and the measurement matrix Ψ , given as

$$\Psi_{ij} := \psi_{\mathbf{s}^{(j)}}(\vec{\xi}^{(i)}) \text{ where } i = 1, 2, \dots, N_p \text{ and } j = 0, 1, \dots, K - 1$$
 (3.7)

In the absence of multi-fidelity strategies (see below) or other sampling techniques, achieving adequate accuracy for the above least-squares problem requires the number of samples used to scale linearly with the cardinality of the basis up to an additional logarithmic factor (Cohen and Migliorati, 2017). For a computationally demanding solver, this may become intractable.

3.1.2 Multi-Fidelity Surrogate Modeling via Kaczmarz Updating

From the variety of the multi-fidelity approaches proposed in recent years, we use the approach proposed by Alemazkoor et al. (2022), which relies on a greedy Kaczmarz algorithm (GKA), to analyze the relative contribution to the overall error from a given sample. A PCE is constructed using a large pool of randomly selected samples and their corresponding outputs from the low-fidelity model. The GKA then iteratively goes through and calculates how much the L^2 -error would change if the coefficients were updated with a given sample from the pool. This updating is achieved via an adaptation of the standard Kaczmarz algorithm:

$$\vec{c}_{\text{test}}^{(i)} = \vec{c} + \frac{u_0^{(i)} - \vec{\psi}_i \cdot \vec{c}}{\|\vec{\psi}_i\|_2^2} \vec{\psi}_i^{\text{T}}$$
(3.8)

where $\vec{\psi}_i$ is the *i*th row of the measurement matrix Ψ , and $u_0^{(i)}$ is the *i*th component of the QoI vector \vec{u}_0 . The coefficients from the Kaczmarz updating are used to compute the error value

$$\varepsilon_i = \left\| \vec{u}_0 - \mathbf{\Psi} \vec{c}_{\text{test}}^{(i)} \right\|_2 \tag{3.9}$$

and the sample that provides the largest error will be removed from the pool (equivalent to removing the corresponding ith row from Ψ and the ith entry from \vec{u}_0). The pseudocode presented in Algorithm 1 is similar to one presented in Alemazkoor et al. (2022), where a subset of the active sample pool needs to be searched to find a sufficiently poor sample for removal. It differs slightly though in that it does not compute the updated measurement matrix and is included here for the sake of completeness. The process is repeated until the pool of N_p samples is reduced to a size N_{rd} , providing the ranked set \mathcal{S}_{GKA} of $N_p - N_{rd}$ samples—in order of most error contributing to least—to be excluded from the construction of the high-fidelity surrogate; that is, $\mathcal{S}_{HiFi} = \mathcal{S}_{pool} \setminus \mathcal{S}_{GKA}$, where \mathcal{S}_{pool} is the pool of initial samples, \mathcal{S}_{HiFi} is the reduced pool of samples used to construct the high-fidelity surrogate, and \setminus is the set minus operator.

3.2 Quantities of Interest, Dimensionality, and Performance Function

The QoI function $u_0: \mathcal{D} \subseteq \mathbb{R}^d \to B \subseteq \mathbb{R}$ is a hyper surface that, for any given value of $b \in B$, defines a level set. Since level sets contain a continuum of values, u_0 is not injective, disallowing the existence of a clear inverse function, i.e., $\nexists u_0^{-1} \in L^2 \mid u_0^{-1} : B \to \mathcal{D}$. In order for there to be potential injectivity, and thus a potential inverse, the QoI must map to a space of equal or higher dimensionality than its domain. Therefore, we define $\mathbf{u}_0: \mathcal{D} \to B^{N_d} \subseteq \mathbb{R}^{N_d}$, where $\mathbf{u}_0 = (u_0^{(1)}, u_0^{(2)}, ..., u_0^{(N_d)})$ with $N_d \geq d$. For a unique candidate control vector $\vec{x}^{(k)} \in$

Algorithm 1: Greedy Kaczmarz algorithm with subset search

```
1: Inputs: \Psi, \vec{u}_0, N_{rd}, and N_{\text{sub}}
  2: Initialize N_p = \text{length}(\vec{u}_0)
  3: Initialize \Psi^{\text{opt}} = \Psi and \vec{u}_{\text{opt}} = \vec{u}_0
  4: Initialize i^{
m rows}=[1,2,...,N_p]
5: Preallocate i^{
m flagged} as an (N_p-N_{rd})	imes 1 array
  6: for i in [N_p, N_p - 1, ..., N_{rd} + 1, N_{rd}] do
                  N_{\text{sub}} \leftarrow \min(N_{\text{sub}}, i)
                 \vec{c} \leftarrow \operatorname*{arg\;min}_{\vec{c}} \left\| \vec{u}_{\; \mathrm{opt}} - \mathbf{\Psi}^{\mathrm{opt}} \vec{c} \, \right\|_{2}^{2} Preallocate \epsilon as an N_{\mathrm{sub}} 	imes 1 array
                  S \leftarrow \text{randperm}([1, 2, ..., i])
10:
                  S \leftarrow [S_1, S_2, ..., S_{N_{\text{sub}}}]
                                                                                                 \triangleright Redefine S to be the first N_{\text{sub}} values of the random
11:
         permutation
12:
                  for k = [1, 2, ..., N_{\text{sub}}] do
13:
                         \vec{c}_{\text{test}} \leftarrow \vec{c} + \frac{u_{\text{opt}}^{(j)} - \vec{\psi}_{j}^{\text{opt}} \cdot \vec{c}}{\|\vec{\psi}_{i}^{\text{opt}}\|_{2}^{2}} (\vec{\psi}_{j}^{\text{opt}})^{\text{T}} \qquad \qquad \triangleright \vec{\psi}_{j}^{\text{opt}} \text{ is the } j \text{th row of the matrix } \mathbf{\Psi}^{\text{opt}}
14:
                 arepsilon_k \leftarrow \left\| ec{u}_0 - oldsymbol{\Psi} ec{ec{c}}_{	ext{test}}^{i \ j} 
ight\|_2^1 end for
15:
16:
                  r \leftarrow \operatorname{Index}(\max(\varepsilon))
                                                                                                           ▶ Find the index of the entry with the largest error
17:
                   \begin{array}{l} r \leftarrow S_r \\ i_{N_p+1-i}^{\mathrm{flagged}} \leftarrow i_r^{\mathrm{rows}} \\ i_r^{\mathrm{rows}} \leftarrow [\ ] \\ \vec{\psi}_r^{\mathrm{opt}} \leftarrow [\ ] \end{array} 
                                                                                                                                  \triangleright Redefine r as its index regarding \Psi^{\text{opt}}
18:
19:
                                                                                                                                                 \triangleright Remove the rth entry from i^{\text{rows}}
20:
21:
                                                                                                                                                   \triangleright Remove the rth row from \Psi^{\text{opt}}
                  u_{\text{opt}}^{(r)} \leftarrow [\ ]
22:
                                                                                                                                                 \triangleright Remove the rth entry from \vec{u}_{\text{opt}}
23: end for
```

 $\{\vec{x}^{(1)}, \vec{x}^{(2)}, ..., \vec{x}^{(N_d)}\}$, each component of u_0 is itself a distinct and independent QoI derived from Q:

$$u_0^{(k)}(\vec{\xi}) = G_k(Q(\vec{x}^{(k)}; \vec{\mu}_D, F(\vec{\xi}))$$
 (3.10)

where G_k is a monotonic function and F is the one-to-one mapping mentioned in Section 3.1.1. For $k \in \{1,2,\ldots,N_d\},\ u_0^{(k)}: \mathcal{D} \to B_k \subseteq \mathbb{R},$ and $B^d:= \mathop{\times}_{k=1}^{N_d} B_k$ where $b_k \in B_k$, the level sets comprise all the points $\ell_k = \vec{\xi} \in A$ such that $u_0^{(k)}(\vec{\xi}) = b_k$. Let $S_d \subseteq \{1,2,\ldots,N_d\}$ now correspond to a particular arrangement of points from the space of candidate control vectors and the true values for unknown parameter mapped onto \mathcal{D} be denoted as $\vec{\xi}^{true}$:

$$\vec{\xi}^{true} \in \bigcap_{k \in S_d} \ell_k = \mathcal{L} \tag{3.11}$$

To start the identification exercise, given collection of level sets characterized by $b_k \in B_k$ for $k \in S_d$, the following performance function is constructed:

$$\mathcal{M}(\vec{\xi}) = \sum_{k \in S_d} \left(u^{(k)}(\vec{\xi}) - b_k \right)^2 \tag{3.12}$$

where $u^{(k)}$ is the surrogate for $u_0^{(k)}$, and the values of b_k come from empirical observations corresponding to the model Q, again disregarding the presence of noise. All potential estimates of the unknown parameters correspond to the minima of Eq. (3.12) in the domain A and will occur only at the points in \mathcal{L} . The quality of a surrogate when used for the purpose of estimating the unknown parameters via the minimization of Eq. (3.12) is not known a priori. It is, therefore, advantageous for the number of components, N_d , of u_0 to adequately exceed d. When viewed independently from each other, each surrogate $u^{(k)}$ can be constructed using only the points in $\mathcal{S}_{\mathrm{HiFi}}^{(k)} = \mathcal{S}_{\mathrm{pool}}^{(k)} \setminus \mathcal{S}_{\mathrm{GKA}}^{(k)}$. However, since the surrogates are all built on the same parameter space, we use the union of the individual sets, that is

$$S_{\text{HiFi}} = \bigcup_{k=1}^{N_d} S_{\text{HiFi}}^{(k)}$$
(3.13)

to construct the high-fidelity surrogates, which improves the overall accuracy of each surrogate without any additional computational cost.

3.3 Optimal Experiment Design via Statistical Sensitivity Measures

The methods presented in this section seek to quantify the potential of a given arrangement of control vectors S_d for estimating arbitrary unknown parameters via the minimization of Eq. (3.12), henceforth referred to as estimation potential. For brevity, we will refer to the surrogate functions $u^{(k)}$ constructed for each $u_0^{(k)}$ as potential observation functions (POFs). The estimation potential then acts as a means for the "quantification" of the likelihood of an arrangement to consistently produce a well-conditioned L^2 -cost functional based on Eq. (3.12), where by well-conditioned we mean the ability to precisely estimate arbitrary unknown parameters. In what follows we explore statistical measures of sensitivity that allow for this quantification from distinct angles. Parameter estimation problems with two unknown parameters uniformly likely to lie anywhere in their range—Legendre polynomials are used to construct PCEs—are presented. The extension to three or more parameters is the subject of our future research, where the methods proposed herein will act as the foundation for more elaborate approaches tailored to higher-dimensional problems.

3.3.1 Sobol' Indices

The first measure of sensitivity used in this work relies on analysis of variance based on Sobol' indices to quantify the estimation potential of arrangements of control vectors.

Sudret (2008) shows how the Sobol' indices can be derived from a Legendre-based PCE of arbitrary degree. With $\vec{\xi} \in [-1, 1]^2$, the kth POF is written as

$$u^{(k)}(\vec{\xi}) = \sum_{s^{(j)} \in \mathbb{N}_{\kappa}^{d}} c_{s^{(j)}}^{(k)} \psi_{s^{(j)}}(\vec{\xi})$$
(3.14)

The basis functions are grouped into sets denoted by the binary vector subscripts $i^{(1)} = (1,0)$, $i^{(2)} = (0,1)$, and $i^{(3)} = (1,1)$, and Eq. (3.14) is decomposed as

$$u^{(k)}(\vec{\xi}) = c_{s^{(0)}}^{(k)} \psi_{s^{(0)}} + \sum_{i \in \mathcal{I}} \sum_{j \in S_i} c_{s^{(j)}}^{(k)} \psi_{s^{(j)}}(\vec{\xi})$$
(3.15)

with $i \in \mathcal{I} = \{i^{(1)}, i^{(2)}, i^{(3)}\}$, and where S_i is the set of all linear indices j, which correspond to all basis functions described by i. The components of i (i_1 and i_2) indicate whether the corresponding basis function has any dependence on the first or second variable, respectively, and the value $c_{s^{(0)}}$ is the constant coefficient, while $\psi_{s^{(0)}}$ is a normalizing factor. The Sobol' indices provide the ratios of the partial variance, in terms of i, to the total variance of $u^{(k)}$, where the total variance is given by

$$D^{(k)} = \text{Var}\Big[u^{(k)}(\vec{\xi}\,)\Big] = \sum_{s^{(j)} \in \mathbb{N}_K^d \setminus s^{(0)}} \left(c_{s^{(j)}}^{(k)}\right)^2 \mathbf{E}\Big[\psi_{s^{(j)}}^2(\vec{\xi})\Big]$$
(3.16)

Sudret (2008) defines the Sobol' indices as

$$[SU_{i}]^{(k)} = \frac{1}{D^{(k)}} \sum_{j \in S_{i}} \left(c_{s^{(j)}}^{(k)} \right)^{2} E\left[\psi_{s^{(j)}}^{2}(\vec{\xi}) \right]$$
(3.17)

The first-order Sobol' indices are $[SU_{i^{(1)}}]^{(k)}$ and $[SU_{i^{(2)}}]^{(k)}$, which correspond to the basis functions dependent on only one variable, and the total Sobol' indices are

$$[SU_1^T]^{(k)} = [SU_{\boldsymbol{i}^{(1)}}]^{(k)} + [SU_{\boldsymbol{i}^{(3)}}]^{(k)} \quad \text{and} \quad [SU_2^T]^{(k)} = [SU_{\boldsymbol{i}^{(2)}}]^{(k)} + [SU_{\boldsymbol{i}^{(3)}}]^{(k)}$$
 (3.18)

which correspond to all basis functions that have any dependence on ξ_1 or ξ_2 , respectively.

Figures 1 and 2 depict two arbitrary POFs, indexed by k and l. A useful visualization of the Sobol' indices can be seen in Figs. 1(b) and 2(b), where the variance from either ξ_1 or ξ_2 is shown as the total left-or-right gradient contribution or the total up-or-down gradient contribution, respectively. The amount of variance experienced along one axial direction versus the other can be quantified by taking the ratios of the Sobol' indices. We define two different ratios in terms of the first-order and total Sobol' indices as quantitative measures of the sensitivity of a POF with respect to inputs and refer to them as first-order Sobol' ratios and total Sobol' ratios, respectively:

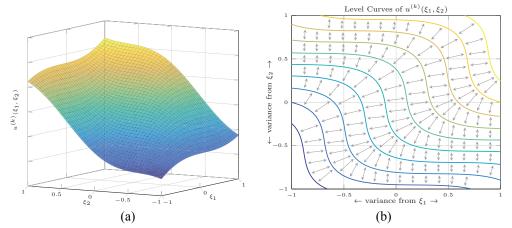


FIG. 1: The surface of the arbitrary POF $u^{(k)}$ (a) along with the two-way gradient field mapped between the POF's surface contours (b)

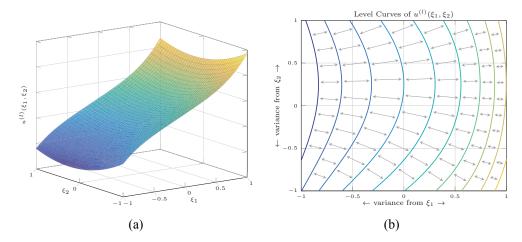


FIG. 2: The surface of the arbitrary POF $u^{(l)}$ (a) along with the two-way gradient field mapped between the POF's surface contours (b)

$$[R^1]^{(k)} := \frac{[SU_1^1]^{(k)}}{[SU_2^1]^{(k)}} \quad \text{and} \quad [R^T]^{(k)} := \frac{[SU_1^T]^{(k)}}{[SU_2^T]^{(k)}}$$
 (3.19)

Going back to Eq. (3.12), we then define measures that quantify the estimation potential of an arrangement composed of two arbitrary control vectors, $\vec{x}^{(k)}$ and $\vec{x}^{(l)}$, to be the first-order comparative Sobol' ratios and total comparative Sobol' ratios:

$$C_{kl}^{1} := \frac{\max\left(\left[R^{1}\right]^{(k)}, \left[R^{1}\right]^{(l)}\right)}{\min\left(\left[R^{1}\right]^{(k)}, \left[R^{1}\right]^{(l)}\right)} \quad \text{and} \quad C_{kl}^{T} := \frac{\max\left(\left[R^{T}\right]^{(k)}, \left[R^{T}\right]^{(l)}\right)}{\min\left(\left[R^{T}\right]^{(k)}, \left[R^{T}\right]^{(l)}\right)}$$
(3.20)

Large comparative ratios imply that a pair of POFs vary "mostly" perpendicularly to each other along the axial directions. The level curves—level sets of two-dimensional functions $u_0^{(k)}$ and $u_0^{(l)}$ —for a pair of observations corresponding to b_k and b_l are depicted in Fig. 3. These curves, by definition, contain all possible pairs of (ξ_1, ξ_2) in the domain that could produce the observed values. Moreover, the true parameter values are ideally located at an intersection of these level curves.

In the case of large comparative Sobol' ratios, the intersection of the level curves is expected, on average, to subtend angles near 90° . Figure 4 demonstrates the intersection of both POFs' level curves at an approximate angle of 90° . The resulting surface produced from the POFs and the L^2 -cost functional (3.12) is seen to contain a minimum that is pronounced more distinctly along both axial directions than surfaces produced from intersections at acute angles much less than 90° . We refer to these kinds of minima as highly isotropic (see Fig. 5) and hypothesize that maximizing the intersection angle between two level curves has the potential to lead to more pronounced minima. We further conjecture that arrangements of control vector pairs that perform well at estimating arbitrary unknown parameters is linked to maximization of comparative Sobol' ratios as a quantitative measure of estimation potential.

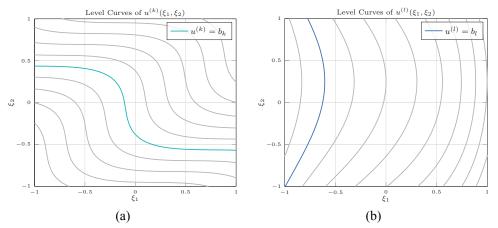


FIG. 3: Level curves for particular observations (a) $u^{(k)} = b_k$ and (b) $u^{(l)} = b_l$

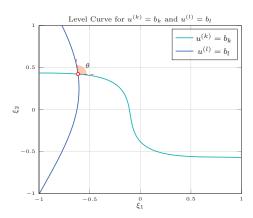


FIG. 4: The intersection of two level curves from $u^{(k)}$ and $u^{(l)}$ subtending an obtuse angle $\theta = 105^{\circ}$

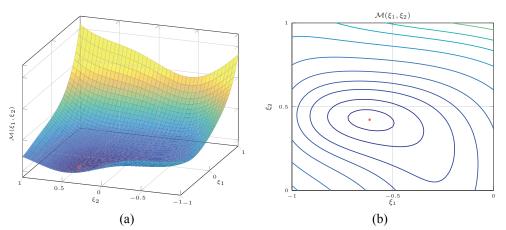


FIG. 5: The surface plot (a) and the contour plot (b) of Eq. (3.12) both show a well-pronounced minimum with high isotropy

3.3.2 Aggregated Directional Statistics

As mentioned above, for two intersecting curves, the minimum of the performance function in Eq. (3.12) is maximally pronounced when the curves subtend an angle of 90° . With multiple curves, however, the optimal angle subtended by their intersection must be updated. It turns out that for n intersecting curves, this angle is $180^{\circ}/n$. A potential drawback with the comparative Sobol' ratios, however, is that their extension to cases with more than two POFs is nontrivial. Consider a particular pair of level curves for two arbitrary POFs depicted in Fig. 6(a). The level curves intersect at three distinct coordinates. This will necessarily produce at least three distinct minima on the surface of the performance function (3.12). The occurrence of multiple minima may be circumvented by the inclusion of additional POFs. As can be seen in Fig. 6(b), this could correspond to an arrangement composed of three control vectors where the inclusion of an additional POF, indexed by m, results in a single intersection point shared by the three level curves.

A more robust quantification of estimation potential can be derived from a direct analysis of the intersection angles between multiple POF level curves. From Figs. 1(b) and 2(b), it appears that the comparative Sobol' ratios act as an indirect and imprecise aggregation of the intersection angles between the level curves. An alternative and more direct approach to compare angles of intersection is to analyze the slope fields of the level curves. The slope of a level curve $u^{(k)}(\vec{\xi}) = b_k$ at an arbitrary point $\vec{\xi} = (p_1, p_2)$ can be quantified by implicit differentiation of the equation with respect to ξ_1 or ξ_2 . Without loss of generality, let us choose ξ_1 to get

$$\frac{du^{(k)}}{d\xi_1} = \frac{\partial u^{(k)}}{\partial \xi_1} + \frac{\partial u^{(k)}}{\partial \xi_2} \frac{d\xi_2}{d\xi_1} = 0$$
(3.21)

from which the slope field (see Fig. 7) and its associated angle can be calculated as

$$m^{(k)}(\vec{\xi}) := \frac{d\xi_2}{d\xi_1} = -\frac{\partial u^{(k)}/\partial \xi_1}{\partial u^{(k)}/\partial \xi_2}$$
(3.22)

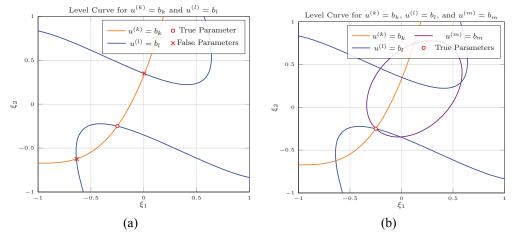


FIG. 6: The intersection using two arbitrary POFs (a) produces two false minima while the true minimum is only distinct when using three POFs (b)

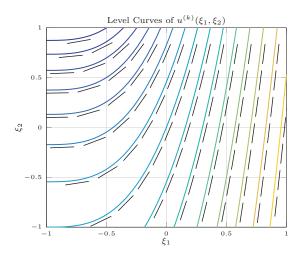


FIG. 7: Arbitrary POF $u^{(k)}$ contour lines and slope field

$$\theta^{(k)}(\vec{\xi}) := \arctan\left(m^{(k)}(\vec{\xi})\right) \tag{3.23}$$

Equation (3.23) is dependent on Eq. (3.22), which requires the first-order differentiation of the POFs written in a Legendre basis. The first derivative of the nth Legendre polynomial is given by Abramowitz and Stegun (1972)

$$\frac{d}{dx}P_n(x) = \mathcal{C}_{n-1}^{(3/2)}(x) \tag{3.24}$$

where $C_n^{(\alpha)}(x)$ is the *n*th degree Gegenbauer (ultraspherical) polynomial. The differentiation process can become tedious particularly for large cardinality PCEs, as the computation of Gegenbauer polynomials can be undesirably slow in various programming languages. Algorithm 2 is an alternative approach to computing Eq. (3.22) directly. The algorithm converts a Legendre basis, on the domain $[-1,1]^d$, into a monomial basis, on the domain $[0,1]^d$, allowing for Eq. (3.23) to be easily computed. We note that the polynomial may be mapped back to $[-1,1]^d$ after the conversion. We also note that Algorithm 2 is the multivariate extension of the one that converts univariate Legendre bases into monomial bases (Barrio and Peña, 2004).

Analyzing the behaviors of each POF's slope-angle field requires the statistics of the angle in Eq. (3.23) to be computed. However, arithmetic averages do not provide useful information for angles and other cyclic quantities. To address this challenge, we adopt a strategy based on directional statistics (Fisher, 1995; Ley and Verdebout, 2017). Let $\Theta = \{\theta_1, \theta_2, \dots, \theta_n\}$ be a collection of n angles. In directional statistics, the average is defined in terms of the angular moments. Let S_X and S_Y be

$$S_X := \{\cos(\theta_1), \cos(\theta_2), ..., \cos(\theta_n)\}$$
 (3.25)

$$S_Y := \{\sin(\theta_1), \sin(\theta_2), ..., \sin(\theta_n)\}\$$
 (3.26)

The mean direction of Θ (see Fig. 8), also referred to as the circular mean (Ley and Verdebout, 2017), denoted as $\overline{\Theta}$, is defined in terms of angular moments $\overline{S_X}$ and $\overline{S_Y}$ as

$$\overline{\Theta} := \operatorname{atan2}\left(\overline{S_Y}, \overline{S_X}\right) \tag{3.27}$$

Algorithm 2: Legendre to monomial conversion

- 1: Inputs: \vec{c}_P , \boldsymbol{I} $ho \ \vec{c}_P$ is a column vector of PCE coefficients in a Legendre basis. I is a tall matrix of \vec{c}_P 's vector indices
- 2: Initialize $n = \text{length}(\vec{c}_P) 1$
- 3: Initialize $N = \sqrt{(2I + 1)/2}$ $\triangleright N$ are the normalization factors for Legendre basis
- 4: $\ell = \text{width}(\boldsymbol{I})$
- 5: Preallocate \vec{c}_M as an $(n+1) \times 1$ array $\triangleright \vec{c}_M$ are the monomial basis coefficients
- 6: **for** i in [1, 2, ..., n + 1] **do**
- $\vec{m} \leftarrow \vec{I}_i$ $\triangleright \vec{I_i}$ is the *i*th row of \boldsymbol{I}
- $\vec{t} \leftarrow \min(\text{rows}(\boldsymbol{I})) \geq \text{rows}(\vec{m})$ $\triangleright \vec{t}$ contains all the row indices of \vec{I} that satisfy the inequality
- $m{k} \leftarrow ec{I}_{ec{t}}$
- $s \leftarrow \mathsf{height}(\boldsymbol{k})$ 10:
- $\triangleright (x)_n$ denotes the Pochhammer symbol
- $\vec{\pi}_m \leftarrow \prod_{j=1}^{\ell} (1/2)_{m_j} / (m_j!)$ $V_{mk} \leftarrow \begin{pmatrix} \mathbf{k} + \mathbf{J}_{s \times 1} \otimes \vec{m} \\ \mathbf{k} \mathbf{J}_{s \times 1} \otimes \vec{m} \end{pmatrix} \odot \vec{N}_{\vec{t}}$ \odot is the Hadamard product, \otimes is the tensor product, and J is the all-ones matrix. $\begin{pmatrix} A_{a \times b} \\ B_{a \times b} \end{pmatrix}$ is the binomial coefficients applied element-wise to

matrices A and B

- 13:
- $\vec{C}_{mk} \leftarrow \prod_{j=1}^{\ell} \operatorname{rows}(\boldsymbol{V}_{mk})$ $\vec{O}_{mk} \leftarrow (-1)^{\circ \sum_{j=1}^{\ell} (\vec{k}_j m_j)} \cdot 4^{\circ \sum_{j=1}^{\ell} m_j} \rhd \circ \text{ denotes the Hadamard power, } \vec{k}_j \text{ is } j \text{th row}$ 14:
- $\vec{c}_{M,i} \leftarrow (\vec{c}_{P,\vec{t}}) \cdot (\vec{\pi}_m \odot \vec{C}_{mk} \odot \vec{O}_{mk})$ $ightharpoonup \cdot$ denotes the scalar product of two vectors 15:
- 16: **end for**

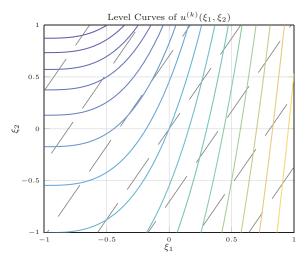


FIG. 8: Arbitrary POF $u^{(k)}$ contour lines and average (constant) slope field

where atan2 is the two-argument arctangent function. This averaging strategy, however, applies only to quantities periodic on $(-\pi, \pi]$; quantities that are cyclic but periodic on different intervals must be mapped to $(-\pi, \pi]$. Such is the case here where the angle of slope lines described by Eq. (3.23), in radians, is periodic on the interval $(-\pi/2, \pi/2]$. Therefore, to perform directional statistics, we define the boosted slope-angle function as

$$\Phi^{(k)}(\vec{\xi}) := 2\theta^{(k)}(\vec{\xi}) \tag{3.28}$$

With this new definition, any member of the sets in Eqs. (3.25) and (3.26) can be written as

$$S_X^{(k)}(\vec{\xi}) = \frac{2\left(\frac{\partial u^{(k)}}{\partial \xi_1}\right)^2}{\left(\frac{\partial u^{(k)}}{\partial \xi_1}\right)^2 + \left(\frac{\partial u^{(k)}}{\partial \xi_2}\right)^2} - 1 \tag{3.29}$$

$$S_Y^{(k)}(\vec{\xi}) = -\frac{2\frac{\partial u^{(k)}}{\partial \xi_1} \frac{\partial u^{(k)}}{\partial \xi_2}}{\left(\frac{\partial u^{(k)}}{\partial \xi_1}\right)^2 + \left(\frac{\partial u^{(k)}}{\partial \xi_2}\right)^2}$$
(3.30)

where we have used trigonometric identities to simplify things. There are several different ways to define a meaningful measure of dispersion. Here, we use the mean resultant length (MRL) (Ley and Verdebout, 2017), which is given by

$$\delta_R := \sqrt{\left(\overline{S_X}\right)^2 + \left(\overline{S_Y}\right)^2} \tag{3.31}$$

It can be shown that $\delta_R \in [0, 1]$, with $\delta_R = 1$ corresponding to zero dispersion, i.e., a collection of angles that are all equivalent, and $\delta_R = 0$ corresponding to either uniformly distributed angles or angles distributed symmetrically about any two perpendicular axes (see Fig. 9).

3.3.3 Global Approach via Mean Resultant Length

For each POF, the mean angular moments can be written as

$$\overline{S_X^{(k)}} = \frac{1}{4} \int_{-1}^1 \int_{-1}^1 S_X^{(k)}(\vec{\xi}) d\xi_1 d\xi_2$$
 (3.32)

$$\overline{S_Y^{(k)}} = \frac{1}{4} \int_{-1}^{1} \int_{-1}^{1} S_Y^{(k)}(\vec{\xi}) d\xi_1 d\xi_2$$
 (3.33)

The direct computation of the above integrals is unnecessary. Here we compute an approximation of these averages by sampling n points, $\{\vec{\xi}^{(1)}, \vec{\xi}^{(2)}, ..., \vec{\xi}^{(n)}\}$, in the domain as

$$\overline{S_X^{(k)}} \approx \frac{1}{n} \sum_{i=1}^n S_X^{(k)}(\vec{\xi}^{(i)})$$
 (3.34)

$$\overline{S_Y^{(k)}} \approx \frac{1}{n} \sum_{i=1}^n S_Y^{(k)}(\vec{\xi}^{(i)})$$
 (3.35)

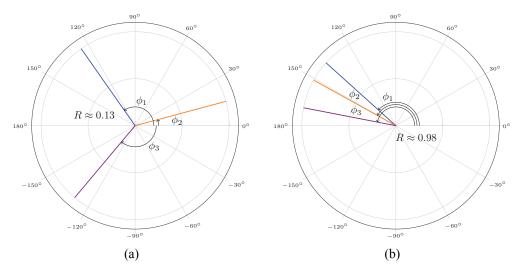


FIG. 9: More-dispersed angles (a) have a low mean resultant length $(\delta_R \to 0)$, while less-dispersed angles (b) have larger mean resultant length $(\delta_R \to 1)$

and use Eq. (3.31) as a means for the quantification of estimation potential:

$$\delta_R = \frac{1}{N_S} \sqrt{\left(\sum_{k \in S_d} \overline{S_X^{(k)}}\right)^2 + \left(\sum_{k \in S_d} \overline{S_Y^{(k)}}\right)^2}$$
(3.36)

which we refer to as global MRL, with N_S being the cardinality of S_d corresponding to the control vectors in a given arrangement. The value δ_R from Eq. (3.36) quantifies how dispersed the average boosted angles are for an arrangement of control vectors, S_d . On average, for values of $\delta_R \approx 0$, the intersection of level curves occurs at angles near $180^\circ/N_S$. We then conjecture that the surface created using Eq. (3.12) from POFs whose mean resultant length is near zero will have better-pronounced minima with greater isotropy than surfaces where $\delta_R \gg 0$. Though Eq. (3.36) is defined for an arbitrary number of POFs, only arrangements composed of $N_S = 3$ control vectors are considered in this paper when quantifying estimation potential using global MRL measure.

3.3.4 Local Approach via Mean Resultant Length

Both the comparative Sobol' ratios and global MRL rely on the evaluation of a summary statistic for each POF before quantifying the estimation potential of the arrangement. An alternative approach is to average over "point-wise" estimation potential of arrangements. To this end, we define the point-wise MRL for a combination of POFs corresponding to an arrangement as

$$\delta_R(\vec{\xi}) := \frac{1}{N_S} \sqrt{\left(\sum_{k \in S_d} S_X^{(k)}(\vec{\xi})\right)^2 + \left(\sum_{k \in S_d} S_Y^{(k)}(\vec{\xi})\right)^2}$$
(3.37)

The average of this function, $\overline{\delta}_R$, referred to as local MRL, is used for quantification of estimation potential. Values of $\overline{\delta}_R \to 0$ indicate that the point-wise average intersection of level curves

occurs at widely dispersed angles, and the values of $\overline{\delta_R} \to 1$ indicate intersections occurring at narrowly dispersed angles. A potential, yet minor, drawback of this approach is that for an arrangement of N_S control vectors chosen from a candidate set of cardinality N_d , there will be $\binom{N_d}{N_S}$ number of distinct functions (3.37) to average. If very precise averages of Eq. (3.37) are required, then this approach may demand high computational costs. However, for reasonably precise averages of the kind needed here, the computational cost of this approach would likely not exceed the cost required to create the surrogates themselves. As with the global MRL, only arrangements composed of $N_S=3$ control vectors are considered in this paper when quantifying estimation potential using local MRL.

4. NUMERICAL EXAMPLES

For all examples presented here, the PCEs constructed used an initial pool of 1500 samples fed to the low-fidelity model. Selecting samples for the high-fidelity model using GKA was done with a subset search size of 50, which reduced the pool to equal that of the cardinality of the PCEs. For each problem, 1000 random values of $\vec{\xi}^{true}$ were generated uniformly over $[-1,1]^2$ (in the domain of the Legendre polynomials), the L^2 -error was used when comparing with $\vec{\xi}^{est}$. The associated outputs for each component of the QoI vector, corresponding to Eq. (3.10), were produced, and estimates of $\vec{\xi}^{est}$ were calculated using Eq. (3.12) and a line search gradient descent algorithm with a minimum step size of $\beta=10^{-4}$ and a stopping criterion of $t_{rel}=10^{-4}$ for the gradient magnitude. To demonstrate the success of the proposed framework, scatter plots were used where we showed the correlation between the rate at which unknown parameters were estimated within the predefined L^2 -error threshold—as displayed on the horizontal axis—and the value of the indicator used to quantify the estimation potential and select the best arrangement of control vectors—displayed on the vertical axis. For indicators based on directional statistics indicators, about 2500 uniformly selected inputs were chosen in the computation of the mean resultant length.

4.1 Borehole Function

For the first example, we choose the borehole function, which is used to model groundwater flow through a borehole connecting an upper and lower aquifer. This function is given by

$$Q(\vec{x}^{(k)}; \vec{\mu}_D, \vec{\mu}) = \frac{2\pi T_u (H_u - H_l)}{\ln(r/r_w^{(k)}) \left(1 + \frac{2LT_u}{\ln(r/r_w^{(k)})(r_w^{(k)})^2 K_w^{(k)}} + \frac{T_u}{T_l}\right)}$$
(4.1)

where $\vec{x}^{(k)} = (r_w^{(k)}, K_w^{(k)})$, $\vec{\mu}_D = (T_u, H_u, r, L)$, and $\vec{\mu} = (H_l, T_l)$; see Table 1 below for the definition of parameters.

 $F_1(\cdot)$ and $F_2(\cdot)$ are linear mappings from ξ_1 and ξ_2 to H_l and T_l , respectively. The high-fidelity borehole function is written as

$$u_h^{(k)}(\vec{\xi}) = \log_{10} \left(\frac{2\pi T_u(H_u - F_1(\xi_1))}{\ln(r/r_w^{(k)}) \left(1 + \frac{2LT_u}{\ln(r/r_w^{(k)})(r_w^{(k)})^2 K_w^{(k)}} + \frac{T_u}{F_2(\xi_2)}\right)} \right)$$
(4.2)

| Var. | Description | Probability | Parameter |
|-------------|--|-------------------------------------|-----------|
| name | | distribution/Value(s) | type |
| T_u | Transmissivity of upper aquifer [m ² /yr] | 83,000 | Known |
| H_u | Potentiometric head of upper aquifer [m] | 1045 | Known |
| H_l | Potentiometric head of lower aquifer [m] | $\mathcal{U}[700,820]$ | Unknown |
| r | Radius of influence [m] | 3200 | Known |
| $r_w^{(k)}$ | Radius of borehole [m] | $\{0.05, 0.1, 0.15\}$ | Control |
| L | Length of borehole [m] | 1400 | Known |
| $K_w^{(k)}$ | Hydraulic conductivity of borehole [m/yr] | $\{500, 5\cdot 10^4, 7\cdot 10^5\}$ | Control |
| T_l | Transmissivity of lower aquifer [m²/yr] | U[63.1, 116] | Unknown |

TABLE 1: Description of each value in $\vec{x}^{(k)}$, $\vec{\mu}_D$, and $\vec{\mu}$ along with how they are evaluated

whereas the low-fidelity function is written as

$$u_l^{(k)}(\vec{\xi}) = \log_{10} \left(\frac{5T_u(H_u - F_1(\xi_1))}{\ln(r/r_w^{(k)}) \left(1.5 + \frac{2LT_u}{\ln(r/r_w^{(k)})(r_w^{(k)})^2 K_w^{(k)}} + \frac{T_u}{F_2(\xi_2)}\right)} \right)$$
(4.3)

where we have used $G_k(\cdot) = \log_{10}(\cdot)$ for all values of k as the monotonic filter function; see Eq. (3.1). Additionally, the set of all candidates for $\vec{x}^{(k)}$ is given by the following Cartesian product: $\{0.05, 0.1, 0.15\} \times \{500, 5 \cdot 10^4, 7 \cdot 10^5\}$, which yields $N_d = 9$ distinct candidates. The values of $r_w^{(k)}$ are chosen as standard values of the radius of the borehole, and the values of $K_w^{(k)}$ are chosen within the range of possible values for medium-grained sand, coarse sand, and gravel (Domenico and Schwartz, 1998). The values used in $\vec{\mu}_D$ are also selected to be consistent with the ranges reported in previous analyses using the borehole function (Alemazkoor et al., 2022; Morris et al., 1993; Zhou et al., 2011). The PCEs are constructed, after screening the initial pool of 1500 samples using GKA as laid out in Section 3.1, with sixth-degree polynomials and 28 high-fidelity model evaluations.

The plots in Fig. 10 use comparative Sobol' ratios and provide a good indication to select the best control vectors. The plots in Fig. 11, which use noncorrected mean resultant length values, however, provide a much better indication, through a strong correlation between success and indicator value, basically pinning down the best candidates for control vectors.

4.2 Advection-Diffusion-Reaction Equation: River Contaminant

For the second example, we use a partial differential equation (PDE) that models the steady-state concentration of pollutants in a straight shallow river (Hamdi, 2007). The PDE is given as

$$\vec{\nabla} \cdot \left(\mathbf{D} \vec{\nabla} C + \vec{v}(x_2) C \right) + rC = S(\vec{x})$$
(4.4)

where

$$\boldsymbol{D} := \begin{pmatrix} D_1 & 0 \\ 0 & D_2 \end{pmatrix} \quad \vec{v} := \begin{pmatrix} v_1(x_2) \\ 0 \end{pmatrix} \tag{4.5}$$

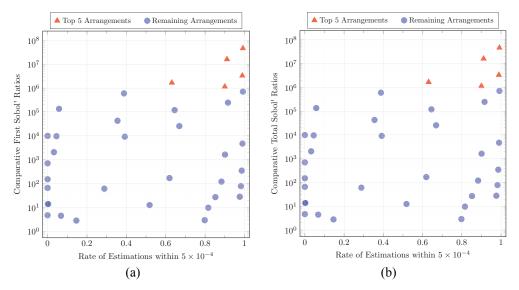


FIG. 10: Comparative first-order and total ratios as an indicator for 36 pairs of control vectors

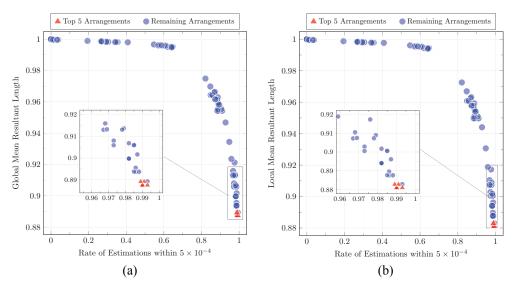


FIG. 11: Global and local mean resultant length as an indicator for 84 triples of control vectors

and on the domain $(x_1, x_2) \in [0, L] \times [-w/2, w/2]$ with the following boundary conditions:

$$C=0$$
 on the left boundary $x_1=0$
$$\frac{\partial C}{\partial x_0}=0$$
 on the right boundary $x_1=L$ and river edges $x_2=\pm w/2$

and $D_2 = (w/L)^{4/3}D_1$ based on Richardson's turbulent pair diffusion rule (Batchelor, 1952; Hamdi, 2007). A simple finite difference scheme is applied to solve this PDE, where the domain

in x_1 is divided into m-1 subintervals of length Δx_1 , and the domain in x_2 was divided into n-1 subintervals of length Δx_2 with a convergence rate of $\mathcal{O}(\Delta x_1^2, \Delta x_2^2)$.

We set the parameter estimation problem here to be that of source localization for a point source that is approximated as

$$S(x,y) = \lambda \delta(x_1 - x_{1,s}) \delta(x_2 - x_{2,s}) \approx \frac{\lambda}{2\pi\sigma^2} \exp\left(-\frac{1}{2\sigma^2} \left((x_1 - x_{1,s})^2 + (x_2 - x_{2,s})^2 \right) \right)$$
(4.6)

where $\sigma = \max(\Delta x_1^l, \Delta x_2^l)$ with Δx_1^l and Δx_2^l , corresponding to the mesh size in low-fidelity model. The mathematical model is the solution to the PDE on the $m \times n$ grid:

$$Q(\vec{x}^{(k)}; \vec{\mu}_D, \vec{\mu}) = C_{(m \times n)}(x_1^{(k)}, x_2^{(k)}; D_1, L, w, r, v_0, \lambda, x_s, y_s)$$
(4.7)

where $\vec{x}^{(k)} = (x_1^{(k)}, x_2^{(k)})$, $\vec{\mu}_D = (D_1, L, w, r, v_0, \lambda)$, and $\vec{\mu} = (x_{1,s}, x_{2,s})$; see Table 2 below for the definition of parameters.

Again with $F_1(\cdot)$ and $F_2(\cdot)$, the linear mappings from ξ_1 and ξ_2 to x_s and y_s , respectively, we have the high-fidelity solution as

$$u_h^{(k)}(\vec{\xi}\,) = \log_{10} \left(C_{(253 \times 41)}(x_1^{(k)}, x_2^{(k)}; D_1, L, w, r, v_0, \lambda, F_1(\xi_1), F_2(\xi_2)) \right) \tag{4.8}$$

whereas the low-fidelity solution is

$$u_l^{(k)}(\vec{\xi}\,) = \log_{10} \left(C_{(64\times11)}(x_1^{(k)}, x_2^{(k)}; D_1, L, w, r, v_0, \lambda, F_1(\xi_1), F_2(\xi_2)) \right) \tag{4.9}$$

with $G_k(\cdot) = \log_{10}(\cdot)$ being the monotonic filter function for all values of k. Four different detector arrays are used as hypothetical in situ measurement apparatuses. Each detector array is comprised of $N_d = 15$ detectors, where each detector corresponds to a unique control vector—a unique coordinate in (x_1, x_2) —with its associated POF. The constructed PCE is a 15-degree polynomial built eventually with 136 high-fidelity runs.

TABLE 2: Description of each value in $\vec{x}^{(k)}$, $\vec{\mu}_D$, and $\vec{\mu}$ along with how they are evaluated

| Var. | Description | Probability | Parameter |
|---------------|--|----------------------------|-----------|
| name | | distribution/Value(s) | type |
| D_1 | Downriver diffusion constant [m ² /s] | 8 | Known |
| L | Length of the river segment [m] | 1000 | Known |
| w | Width of the river segment [m] | 100 | Known |
| r | Reaction coefficient $[s^{-1}]$ | $2.2\cdot 10^{-6}$ | Known |
| v_0 | Center velocity [m/s] | 0.12 | Known |
| $x_1^{(k)}$ | Length-wise coordinate of concentration [m] | Array Dependent | Control |
| $x_{1,s}$ | Length-wise coordinate of source location [m] | $\mathcal{U}[333.3,666.7]$ | Unknown |
| $x_{2}^{(k)}$ | Width-wise coordinate of concentration [m] | Array Dependent | Control |
| $x_{2,s}$ | Width-wise coordinate of source location [m] | $\mathcal{U}[-50, 50]$ | Unknown |
| λ | Source rate [kg/s] | 2.3 | Known |

4.2.1 Upriver-Situated Detector Array

A conceivable configuration for a detector array would be to place detectors uniformly upriver within the river segment, as seen in Fig. 12.

The plots in Figs. 13 and 14, based on comparative Sobol' ratios and noncorrected mean resultant length values, respectively, provide strong indications for the selection of best control vectors all with reasonably strong correlation between success and indicator value.

4.2.2 Mid-River-Situated Detector Array

Another reasonable configuration is that the detector array is placed in the center of the river segment, as seen in Fig. 15.

Barring Fig. 16(b), which provides reasonably good insight, the plots in Fig. 17 and the plot in Fig. 16(a) provide poor indications for the choice of control vectors.

4.2.3 Downriver-Situated Detector Array

A third reasonable configuration is that the detector array is placed downriver, as seen in Fig. 18.



FIG. 12: The upriver detector array (L): the marked locations show where hypothetical *in situ* detectors would be in the river with respect to the indices (i, j) of the low-fidelity mesh coordinates $(x_{1,i}, x_{2,j})$

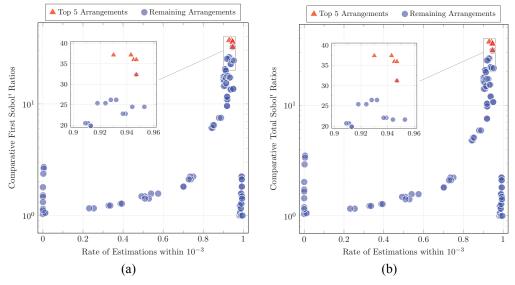


FIG. 13: Comparative first-order and total ratios as an indicator for 105 pairs of control vectors

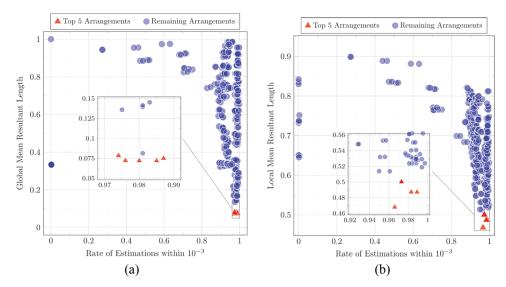


FIG. 14: Global and local mean resultant length as an indicator for 455 triples of control vectors

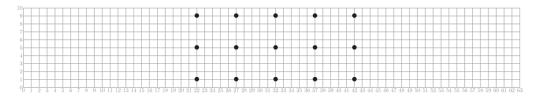


FIG. 15: The mid-river detector array (M): the marked locations show where hypothetical *in situ* detectors would be in the river with respect to the indices (i, j) of the low-fidelity mesh coordinates $(x_{1,i}, x_{2,j})$

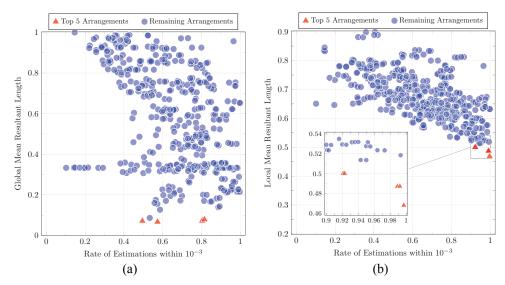


FIG. 16: Global and local mean resultant length as an indicator for 455 triples of control vectors

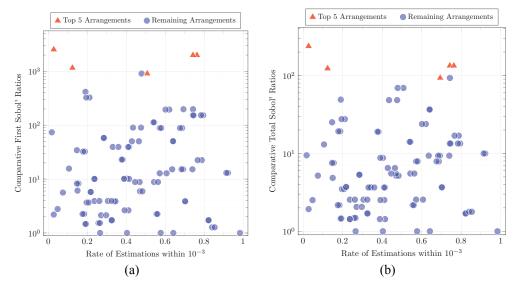


FIG. 17: Comparative first-order and total ratios as an indicator for 105 pairs of control vectors

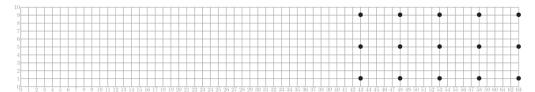


FIG. 18: The downriver detector array (R): the marked locations show where hypothetical *in situ* detectors would be in the river with respect to the indices (i, j) of the low-fidelity mesh coordinates $(x_{1,i}, x_{2,j})$

Just like the mid-river detector array, the downriver array provides reasonably good insight through Fig. 19(b), yet the plots in Fig. 20 and the plot in Fig. 19(a) provide poor indications for the choice of control vectors.

4.2.4 Whole-River-Situated Detector Array

A fourth reasonable configuration is that the detector array is uniformly spread across the whole river segment as seen in Fig. 21.

Similar to the previous two arrays, the whole-river array provides reasonably good insight in Fig. 22(b). The plots in Fig. 23 provide a poor indication for the choice of control vector, while the plot in Fig. 22(a) provides mediocre insight.

5. CONCLUDING REMARKS

We presented a novel framework for the identification of optimal experimental arrangement, from a finite set of candidates, for precise parameter estimation. The framework is particularly useful when minimal data is available through experiments due to limited experimental budget or practical constraints and rests on the following pillars: (i) efficient emulators built through multifidelity resolution of a parent mathematical model of a physical phenomenon and (ii) statistical

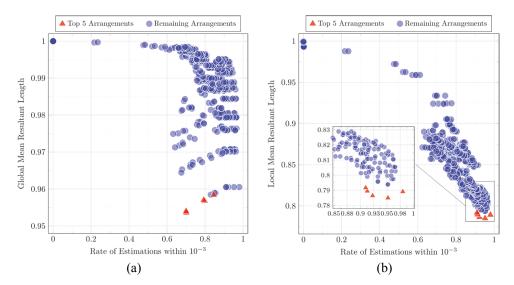


FIG. 19: Global and local mean resultant length as an indicator for 455 triples of control vectors

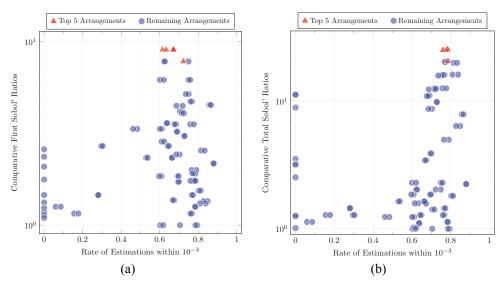


FIG. 20: Comparative first-order and total ratios as an indicator for 105 pairs of control vectors

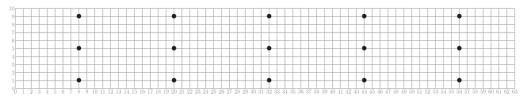


FIG. 21: The whole-river detector array (W): the marked locations show where hypothetical *in situ* detectors would be in the river with respect to the indices (i, j) of the low-fidelity mesh coordinates $(x_{1,i}, x_{2,j})$

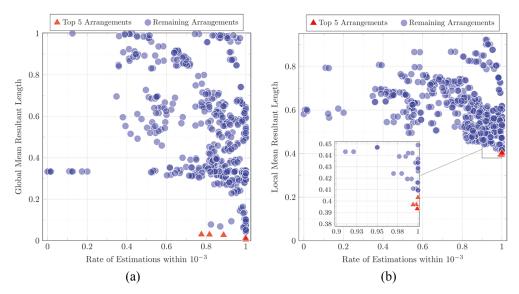


FIG. 22: Global and local mean resultant length as an indicator for 455 triples of control vectors

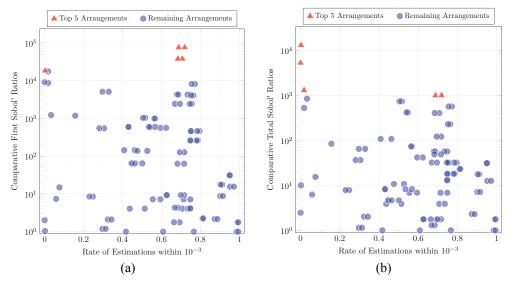


FIG. 23: Comparative first-order and total ratios as an indicator for 105 pairs of control vectors

measures of sensitivity that systematically inform the choice of experimental arrangement. The surrogate models provide potential observation functions (POFs) for the parent mathematical model as a function of the model's unknown parameters for a particular choice of control parameters, referred to as a control vector. An arrangement is a choice of several control vectors from a set of candidate control vectors and is optimal when it provides the best route to parameter estimation among all possible choices for a given predefined number of control vectors. The outputs from the POFs are used to define an L^2 -cost functional, which is then used to quantify

the estimation potential of an arrangement based on how well-conditioned the minima is for arbitrary observations. This is done via strategies derived from variance-based sensitivity analysis and directional statistics including comparative Sobol' ratios and measures of dispersion based on both global and local mean resultant lengths that allow for the examination of the average angle of intersection between the POFs' level curves. Two numerical examples were presented involving the estimation of deep underground aquifer properties and the source localization of a contaminant in a river, collectively indicating the potential of the proposed framework to identify the optimal experimental design.

ACKNOWLEDGMENTS

MT and YC acknowledge financial support from the Office of Naval Research, United States under Grants N00014-20-1-2849 and N00014-22-1-2012 (through MUST program at UMassD). Moreover, MT acknowledges support from the Office of Naval Research, United States under Grant N00014-21-1-2570. The computational resources were provided by the Center for Scientific Computing and Data Science Research (CSCDR) at the University of Massachusetts Dartmouth.

REFERENCES

- Abramowitz, M. and Stegun, I.A., *Orthogonal Polynomials*, Applied Math Series, Washington, DC: U.S. Dept. of Commerce, National Bureau of Standards, 1972.
- Alemazkoor, N., Louhghalam, A., and Tootkaboni, M., A Multi-Fidelity Polynomial Chaos-Greedy Kaczmarz Approach for Resource-Efficient Uncertainty Quantification on Limited Budget, Comput. Methods Appl. Mech. Eng., vol. 389, p. 114290, 2022.
- Ali, A.M., Yao, K., Collier, T.C., Taylor, C.E., Blumstein, D.T., and Girod, L., An Empirical Study of Collaborative Acoustic Source Localization, *Proc. of the 6th Int. Conf. on Information Processing in Sensor Networks*, Cambridge, MA, pp. 41–50, 2007.
- Bandara, S., Schlöder, J.P., Eils, R., Bock, H.G., and Meyer, T., Optimal Experimental Design for Parameter Estimation of a Cell Signaling Model, *PLoS Comput. Biol.*, vol. 5, no. 11, p. e1000558, 2009.
- Barrio, R. and Peña, J.M., Basis Conversions Among Univariate Polynomial Representations, *Comptes Rendus Mathematique*, vol. **339**, no. 4, pp. 293–298, 2004.
- Batchelor, G.K., Diffusion in a Field of Homogeneous Turbulence: II. The Relative Motion of Particles, *Math. Proc. Cambridge Phil. Soc.*, vol. **48**, pp. 345–362, 1952.
- Cohen, A. and Migliorati, G., Optimal Weighted Least-Squares Methods, SMAI J. Comput. Math., vol. 3, pp. 181–203, 2017.
- Domenico, P.A. and Schwartz, F.W., *Hydraulic Conductivity and Permeability of Geologic Materials*, 2nd ed., Hoboken, NJ: Wiley, 1998.
- Fisher, N.I., Statistical Analysis of Circular Data, Cambridge, UK: Cambridge University Press, 1995.
- Gasca-Ortiz, T., Domínguez-Mota, F.J., and Pantoja, D.A., Determination of Optimal Diffusion Coefficients in Lake Zirahuén through a Local Inverse Problem, *Mathematics*, vol. 9, no. 14, p. 1695, 2021.
- Ghanem, R.G. and Spanos, P.D., *Stochastic Finite Elements: A Spectral Approach*, North Chelmsford, MA: Courier Corporation, 2003.
- Hamdi, A., Identification of Point Sources in Two-Dimensional Advection-Diffusion-Reaction Equation: Application to Pollution Sources in a River. Stationary Case, *Inverse Prob. Sci. Eng.*, vol. 15, no. 8, pp. 855–870, 2007.

Jumabekova, A. and Berger, J., Optimal Experiment Design for the Estimation of Building Wall Material Thermal Properties, J. Phys.: Conf. Ser., vol. 2444, p. 012007, 2023.

- Jumabekova, A., Berger, J., Foucquier, A., and Dulikravich, G.S., Searching an Optimal Experiment Observation Sequence to Estimate the Thermal Properties of a Multilayer Wall under Real Climate Conditions, *Int. J. Heat Mass Transf.*, vol. 155, p. 119810, 2020.
- Kowalsky, M., Finsterle, S., Williams, K.H., Murray, C., Commer, M., Newcomer, D., Englert, A., Steefel, C.I., and Hubbard, S., On Parameterization of the Inverse Problem for Estimating Aquifer Properties Using Tracer Data, *Water Resour. Res.*, vol. 48, no. 6, 2012.
- Le Gratiet, L. and Cannamela, C., Cokriging-Based Sequential Design Strategies Using Fast Cross-Validation Techniques for Multi-Fidelity Computer Codes, *Technometrics*, vol. 57, no. 3, pp. 418–427, 2015.
- Ley, C. and Verdebout, T., Modern Directional Statistics, Boca Raton, FL: CRC Press, p. 11, 2017.
- Lile, O.B., Morris, M., and Rønning, J.S., Estimating Groundwater Flow Velocity from Changes in Contact Resistance during a Saltwater Tracer Experiment, *J. Appl. Geophys.*, vol. **38**, no. 2, pp. 105–114, 1997.
- Morris, M.D., Mitchell, T.J., and Ylvisaker, D., Bayesian Design and Analysis of Computer Experiments: Use of Derivatives in Surface Prediction, *Technometrics*, vol. **35**, no. 3, pp. 243–255, 1993.
- Narayan, A., Gittelson, C., and Xiu, D., A Stochastic Collocation Algorithm with Multifidelity Models, *SIAM J. Sci. Comput.*, vol. **36**, no. 2, pp. A495–A521, 2014.
- Parussini, L., Venturi, D., Perdikaris, P., and Karniadakis, G.E., Multi-Fidelity Gaussian Process Regression for Prediction of Random Fields, J. Comput. Phys., vol. 336, pp. 36–50, 2017.
- Piazzola, C., Tamellini, L., Pellegrini, R., Broglia, R., Serani, A., and Diez, M., Uncertainty Quantification of Ship Resistance via Multi-Index Stochastic Collocation and Radial Basis Function Surrogates: A Comparison, AIAA Aviation 2020 Forum, Virtual, p. 3160, 2020.
- Pronzato, L., Optimal Experimental Design and Some Related Control Problems, *Automatica*, vol. 44, no. 2, pp. 303–325, 2008.
- Rainwater, K.A., Wise, W.R., and Charbeneau, R.J., Parameter Estimation through Groundwater Tracer Tests, Water Resour. Res., vol. 23, no. 10, pp. 1901–1910, 1987.
- Rao, N.S., Shankar, M., Chin, J.C., Yau, D.K., Ma, C.Y., Yang, Y., Hou, J.C., Xu, X., and Sahni, S., Localization under Random Measurements with Application to Radiation Sources, 2008 11th Int. Conf. on Information Fusion, Cologne, Germany, pp. 1–8, 2008.
- Romano, N. and Santini, A., Determining Soil Hydraulic Functions from Evaporation Experiments by a Parameter Estimation Approach: Experimental Verifications and Numerical Studies, *Water Resour. Res.*, vol. **35**, no. 11, pp. 3343–3359, 1999.
- Rózsás, Á., Slobbe, A., Martini, G., and Jansen, R., Structural and Load Parameter Estimation of a Real-World Reinforced Concrete Slab Bridge Using Measurements and Bayesian Statistics, *Struct. Concrete*, vol. 23, no. 6, pp. 3569–3600, 2022.
- Sheng, X. and Hu, Y.H., Maximum Likelihood Multiple-Source Localization Using Acoustic Energy Measurements with Wireless Sensor Networks, *IEEE Trans. Signal Process.*, vol. 53, no. 1, pp. 44–53, 2004.
- Shin, Y. and Xiu, D., Nonadaptive Quasi-Optimal Points Selection for Least Squares Linear Regression, SIAM J. Sci. Comput., vol. 38, no. 1, pp. A385–A411, 2016.
- Song, X., Lv, L., Sun, W., and Zhang, J., A Radial Basis Function-Based Multi-Fidelity Surrogate Model: Exploring Correlation between High-Fidelity and Low-Fidelity Models, *Struct. Multidisc. Opt.*, vol. 60, pp. 965–981, 2019.
- Sudret, B., Global Sensitivity Analysis Using Polynomial Chaos Expansions, *Reliab. Eng. Syst. Safety*, vol. **93**, no. 7, pp. 964–979, 2008.

- Wu, C.Q., Berry, M.L., Grieme, K.M., Sen, S., Rao, N.S., Brooks, R.R., and Cordone, G., Network Detection of Radiation Sources Using Localization-Based Approaches, *IEEE Trans. Indust. Inf.*, vol. 15, no. 4, pp. 2308–2320, 2019.
- Wu, X., Liu, M., and Wu, Y., *In-Situ* Soil Moisture Sensing: Optimal Sensor Placement and Field Estimation, *ACM Trans. Sensor Netw.*, vol. **8**, no. 4, pp. 1–30, 2012.
- Xiu, D., Numerical Methods for Stochastic Computations: A Spectral Method Approach, Princeton, NJ: Princeton University Press, 2010.
- Xiu, D. and Karniadakis, G.E., The Wiener–Askey Polynomial Chaos for Stochastic Differential Equations, *SIAM J. Sci. Comput.*, vol. **24**, no. 2, pp. 619–644, 2002.
- Yao, K., Chen, J.C., and Hudson, R.E., Maximum-Likelihood Acoustic Source Localization: Experimental Results, 2002 IEEE Int. Conf. on Acoustics, Speech, and Signal Processing, Orlando, FL, pp. III–2949, 2002.
- Zhang, Q., Volker, R.E., and Lockington, D.A., Experimental Investigation of Contaminant Transport in Coastal Groundwater, *Adv. Environ. Res.*, vol. 6, no. 3, pp. 229–237, 2002.
- Zhou, Q., Qian, P.Z., and Zhou, S., A Simple Approach to Emulation for Computer Models with Qualitative and Quantitative Factors, *Technometrics*, vol. **53**, no. 3, pp. 266–273, 2011.
- Zhu, X., Narayan, A., and Xiu, D., Computational Aspects of Stochastic Collocation with Multifidelity Models, SIAM/ASA J. Uncertainty Quant., vol. 2, no. 1, pp. 444–463, 2014.