

FUNCTIONAL REGRESSION FOR STATE PREDICTION USING LINEAR PDE MODELS AND OBSERVATIONS*

N. C. NGUYEN[†], H. MEN[†], R. M. FREUND[‡], AND J. PERAIRE[†]

Abstract. Partial differential equations (PDEs) are commonly used to model a wide variety of physical phenomena. A PDE model of a physical problem is typically described by conservation laws, constitutive laws, material properties, boundary conditions, boundary data, and geometry. In most practical applications, however, the PDE model is only an approximation to the real physical problem due to both (i) the deliberate mathematical simplification of the model to keep it tractable and (ii) the inherent uncertainty of the physical parameters. In such cases, the PDE model may not produce a good prediction of the true state of the underlying physical problem. In this paper, we introduce a functional regression method that incorporates observations into a deterministic linear PDE model to improve its prediction of the true state. Our method is devised as follows. First, we augment the PDE model with a random Gaussian functional which serves to represent various sources of uncertainty in the model. We next derive a linear regression model for the Gaussian functional by utilizing observations and adjoint states. This allows us to determine the posterior distribution of the Gaussian functional and the posterior distribution for our estimate of the true state. Furthermore, we consider the problem of experimental design in this setting, wherein we develop an algorithm for designing experiments to efficiently reduce the variance of our state estimate. We provide several examples from the heat conduction, the convection-diffusion equation, and the reduced wave equation, all of which demonstrate the performance of the proposed methodology.

Key words. Gaussian functional, state prediction, linear PDEs, experimental design, data assimilation, nonparametric regression

AMS subject classifications. 60G15, 62G08, 65N75

DOI. 10.1137/14100275X

1. Introduction. Partial differential equations (PDEs) are widely used to model a wide variety of physical phenomena in the real world and predict associated physical states such as temperature, displacement, velocity, pressure, or density. The predictability of the PDE model depends on how well it captures the real physics. The model captures the physics and yields a good prediction if all aspects of the model, such as constitutive laws, material properties, boundary conditions, boundary data, and geometry can be accurately prescribed. In practice, the PDE model is imperfect due to (i) the deliberate mathematical simplification of the model to keep it tractable (by ignoring certain physics that pose computational difficulties) and (ii) the uncertainty of the model data (by using geometry, material properties, and boundary data that are different from those of the physical problem.). In such cases, the PDE model should be validated and the uncertainty should be quantified.

There are a number of different approaches for dealing with model discrepancy. One approach is to represent physical inputs as random parameters or random fields, thereby resulting in *stochastic PDEs*. There exist several numerical methods for solving stochastic PDEs such as Monte Carlo (MC) methods [12, 16], intrusive poly-

*Submitted to the journal's Computational Methods in Science and Engineering section January 5, 2015; accepted for publication (in revised form) November 30, 2015; published electronically March 23, 2016. This work was supported by the Air Force Office of Scientific Research under the AFOSR grant FA9550-12-1-0357 and AFOSR grant FA9550-15-1-0276.

<http://www.siam.org/journals/sisc/38-2/100275.html>

[†]MIT Department of Aeronautics and Astronautics, Cambridge, MA 02139 (cuongng@mit.edu, abbymen@mit.edu, peraire@mit.edu).

[‡]MIT Sloan School of Management, Cambridge, MA 02139 (rfreund@mit.edu).

nomial chaos [13, 38], nonintrusive polynomial chaos [37, 10, 18], stochastic collocation [37, 3, 28], response surface [5, 14], and Kriging [6, 23, 20, 39, 9].

Data assimilation is another approach that combines the PDE model with observations to estimate the physical state. There are several different data assimilation methods. In *stochastic data assimilation* [19, 24, 34], the state estimate is represented as a stochastic process and is determined by minimizing its variance. In *variational data assimilation* [21, 22, 40], the state estimate is defined as an optimal solution of a least-squares minimization principle. In *parameter estimation* [7, 1, 21, 36], a number of uncertain parameters in the PDE model are determined by matching the model outputs to the observations. In *reduced order modeling* [4, 11, 25, 26, 29, 35], the state is reconstructed by fitting the observations to the snapshots which are computed by solving a parametrized or time-varying PDE model.

In this paper, we focus on a data assimilation method recently introduced in [27]. The method is devised as follows. First, we augment the PDE model with a random Gaussian functional which serves to represent various sources of uncertainty in the model. This gives rise to a stochastic PDE model whose solution is characterized by the Gaussian functional. We next derive a linear regression model for the Gaussian functional by utilizing observations and adjoint states. This functional regression model allows us to compute the posterior distribution for our estimate of the physical state, thereby quantifying the prediction error. A crucial ingredient in our method is the *covariance operator* representing the *prior* distribution of the Gaussian functional. The bilinear covariance operators considered incorporate a number of *hyperparameters* that are determined upon the observations by maximizing the likelihood of the Gaussian functional with respect to the hyperparameters. Furthermore, we consider the problem of experimental design in this setting, wherein we develop an algorithm for designing experiments to efficiently reduce the variance of our state estimate.

We can relate our approach to existing data assimilation methods. Our approach seeks to characterize model uncertainties by introducing a Gaussian functional into the existing numerical model, while stochastic data assimilation directly represents the state estimate as a stochastic process and infers its posterior distribution from the observations, the model outputs, and the priors about the background state. We show in Appendix B that our method and the Kalman method [19] yield exactly the same posterior distribution for a judicious choice of the priors. Our approach can also be shown to yield a posterior mean that satisfies the least-square minimization principle of three-dimensional (3D) variational data assimilation [8, 22].

The proposed method shares the following features with stochastic data assimilation methods. First, it incorporates both the PDE model and the observations into the regression procedure. Second, it can handle the observations given in the form of linear functionals of the field variable. Third, it yields not only the mean estimate but also the posterior variance that quantifies the prediction error. And fourth, it provides a natural mechanism for designing experiments to reduce the posterior variance. Furthermore, the method has a distinctive feature in that it provides a systematic and rational way to construct the model structure (including the prior covariance and experimental noise) based on the observations.

There are a number of new contributions relative to our previous work [27]. Herein we propose an efficient greedy algorithm to rationally select good observations from a large number of possible experiments. We show how our method can be related to Bayesian inference, 3D variational data assimilation, and the Kalman method. Finally, we apply our approach to a variety of PDEs including convection-diffusion equations and Helmholtz equations to demonstrate its performance.

The paper is organized as follows. In section 2, we formulate the problem of interest. In section 3, we present the functional regression method for state estimation using linear PDE models and observations. In section 4, we describe a greedy algorithm for selecting the observations. In section 5, we demonstrate our approach on several examples from heat conduction, convection-diffusion equation, and reduced wave equation. In section 6, we conclude the paper with some remarks on future research. Finally, in the appendix, we show the connection our method to 3D variational data assimilation and the Kalman method.

2. Problem formulation. Let Ω denote a bounded open domain with Lipschitz boundary. Let V be an appropriate finite element approximation space, which is defined on a triangulation of the domain Ω . The Galerkin finite element formulation of a general linear PDE model can be stated as follows. We seek a solution $u^o \in V$ and an output vector $\mathbf{s}^o \in \mathbb{C}^M$ such that

$$(1a) \quad a(u^o, v) = \ell(v) \quad \forall v \in V,$$

$$(1b) \quad s_i^o = c_i(u^o), \quad i = 1, \dots, M.$$

Here $a : V \times V \rightarrow \mathbb{C}$ is a continuous bilinear form, $\ell : V \rightarrow \mathbb{C}$ is a continuous and bounded linear functional, and $c_i : V \rightarrow \mathbb{C}, i = 1, \dots, M$ are continuous and bounded linear functionals. The finite element approximation space V is of N dimensions. We shall take N to be so large that the numerical solution u is indistinguishable from the exact solution of the PDE model at any accuracy level of interest.

We assume that the underlying finite element (FE) model (1) is used to predict the true state of a deterministic and time-independent physical problem. We denote the true state by u^{true} and the associated true output vector by $\mathbf{s}^{\text{true}} \in \mathbb{C}^M$. Both u^{true} and \mathbf{s}^{true} are not known. However, we assume that we are given a vector of M observations $\mathbf{d} \in \mathbb{C}^M$, which are the measurements of the true output vector \mathbf{s}^{true} . We further assume that the observation vector differs from the true output vector \mathbf{s}^{true} by additive Gaussian noise, namely,

$$(2) \quad \mathbf{d} = \mathbf{s}^{\text{true}} + \boldsymbol{\varepsilon},$$

where $\boldsymbol{\varepsilon} = (\varepsilon_1, \dots, \varepsilon_M) \in \mathbb{C}^M$ are independent, identically distributed Gaussian random variables with zero mean and variance σ^2 .

As mentioned in the introduction, the solution u^o of the model (1) may not be a good approximation to the true state u^{true} because of various sources of uncertainty arising from the imprecise knowledge of boundary conditions, geometry, physical parameters, and governing equations. The uncertainty of our PDE model in predicting the true physical state prompts basic research questions such as the following:

- How do we combine the model (1) and the observations (2) to yield a more accurate prediction of the true state?
- How do we quantify the error in the prediction?
- How do we design experiments to improve the accuracy of the prediction?

In the remainder of this paper, we concentrate on addressing these questions. We address the first two questions by using the *functional regression* method described in the next section. We address the third question by devising an experimental design procedure in section 4.

3. Functional regression.

3.1. Gaussian functional. We suppose that we are given a bounded linear functional $g : V \rightarrow \mathbb{C}$, and we consider the following model: determine $u \in V$ and $\mathbf{s} \in \mathbb{C}^M$ such that

$$(3a) \quad a(u, v) + g(v) = \ell(v) \quad \forall v \in V,$$

$$(3b) \quad s_i = c_i(u), \quad i = 1, \dots, M.$$

Notice that the new model (3) differs from the original model (1) by the addition of the functional g . Here we can exactly determine u and \mathbf{s} only if g is known. Indeed, together with $a(u, v)$ and $\ell(v)$, the functional g determines the solution u and the output vector \mathbf{s} of the model (3). Note that the introduction of the functional g is able to capture the uncertainties in the functional ℓ and the bilinear form a , as demonstrated by numerical examples in section 5. Indeed, if we choose $g(v) = \ell(v) - a(u^{\text{true}}, v)$, then $u = u^{\text{true}}$ is the solution of (3). Unfortunately, this particular choice of g requires foreknowledge of the true state u^{true} , which we presumably do not know and which indeed we seek to accurately predict.

In order to capture various sources of uncertainty in the original model (1), we hypothesize the linear functional g as a *Gaussian process* [31]. Specifically, the Gaussian functional g is assumed to have zero mean and covariance operator k , namely,

$$(4) \quad g(v) \sim \mathcal{GP}(0, k(w, v)) \quad \forall w, v \in V.$$

We require that the covariance operator $k : V \times V \rightarrow \mathbb{C}$ is symmetric positive-definite, namely,

$$(5) \quad k(w, v) = k(v, w), \quad k(v, v) > 0 \quad \forall v \neq 0, \quad \text{and} \quad k(0, 0) = 0.$$

As the covariance operator k characterizes the space of all possible functionals prior to taking into account the observations, it plays an important role in our method. The selection of a covariance operator will be discussed later.

3.2. Functional regression model. Here we develop a linear regression model for the Gaussian functional. To this end, we first obtain the adjoint solutions $\phi_i \in V, i = 1, \dots, M$ by solving the adjoint problems:

$$(6) \quad a(v, \phi_i) = -c_i(v), \quad \forall v \in V.$$

It follows from (1), (3), and (6) that

$$(7) \quad g(\phi_i) = \ell(\phi_i) - a(u, \phi_i) = a(u^{\circ}, \phi_i) + c_i(u) = c_i(u) - c_i(u^{\circ}) = s_i - s_i^{\circ}$$

for $i = 1, \dots, M$. Moreover, we assume that the observation d_i differs from the output estimate s_i by Gaussian noise $\varepsilon_i \sim \mathcal{N}(0, \sigma^2)$, namely,

$$(8) \quad s_i = d_i - \varepsilon_i, \quad i = 1, \dots, M.$$

This equation is analogous to (2), which relates the observed data \mathbf{d} to the true outputs \mathbf{s}^{true} . We substitute (8) into (7) to obtain

$$(9) \quad d_i - s_i^{\circ} = g(\phi_i) + \varepsilon_i, \quad i = 1, \dots, M.$$

This equation can be viewed as a standard regression model, namely, a linear model with additive Gaussian noise, and we shall now show that the model (9) allows us

to determine the posterior distribution of $g(v)$ for any given test function $v \in V$ as follows.

Let $\Phi = [\phi_1, \dots, \phi_M]$ be a collection of M adjoint states as determined by (6). Let $\Psi = [\psi_1 \in V, \dots, \psi_N \in V]$ be a collection of N test functions, where $\psi_j, 1 \leq j \leq N$, are basis functions of the space V . We would like to compute the posterior distribution of $\mathbf{g} \in \mathbb{C}^N$ with $g_i = g(\psi_i), i = 1, \dots, N$. According to the prior (4) and the regression model (9), the joint distribution of $(\mathbf{d} - \mathbf{s}^o)$ and \mathbf{g} is given by

$$(10) \quad \begin{bmatrix} \mathbf{d} - \mathbf{s}^o \\ \mathbf{g} \end{bmatrix} \sim \mathcal{N} \left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} \mathbf{K}(\Phi, \Phi) + \sigma^2 \mathbf{I} & \mathbf{K}(\Phi, \Psi) \\ \mathbf{K}(\Psi, \Phi) & \mathbf{K}(\Psi, \Psi) \end{bmatrix} \right),$$

where $\mathbf{K}(\Phi, \Phi) \in \mathbb{C}^{M \times M}$, $\mathbf{K}(\Phi, \Psi) \in \mathbb{C}^{M \times N}$, $\mathbf{K}(\Psi, \Phi) \in \mathbb{C}^{N \times M}$, and $\mathbf{K}(\Psi, \Psi) \in \mathbb{C}^{N \times N}$ have entries

$$(11) \quad \begin{aligned} K_{ij}(\Phi, \Phi) &= k(\phi_i, \phi_j), & i = 1, \dots, M, j = 1, \dots, M, \\ K_{ij}(\Phi, \Psi) &= k(\phi_i, \psi_j), & i = 1, \dots, M, j = 1, \dots, N, \\ K_{ij}(\Psi, \Phi) &= k(\psi_i, \phi_j), & i = 1, \dots, N, j = 1, \dots, M, \\ K_{ij}(\Psi, \Psi) &= k(\psi_i, \psi_j), & i = 1, \dots, N, j = 1, \dots, N, \end{aligned}$$

respectively. It thus follows from (10) and the conditional normal distribution formula (see Appendix A.2 in Rasmussen and Williams [31]) that the posterior distribution of \mathbf{g} is

$$(12) \quad p(\mathbf{g} | (\mathbf{d} - \mathbf{s}^o), \Phi, \Psi) \sim \mathcal{N}(\bar{\mathbf{g}}, \mathbf{G}),$$

where the mean vector $\bar{\mathbf{g}} \in \mathbb{C}^N$ and the covariance matrix $\mathbf{G} \in \mathbb{C}^{N \times N}$ are given by

$$(13) \quad \bar{\mathbf{g}} = \mathbf{K}(\Psi, \Phi) \mathbf{D}^{-1} (\mathbf{d} - \mathbf{s}^o), \quad \mathbf{G} = \mathbf{K}(\Psi, \Psi) - \mathbf{K}(\Psi, \Phi) \mathbf{D}^{-1} \mathbf{K}(\Phi, \Psi)$$

with $\mathbf{D} = \mathbf{K}(\Phi, \Phi) + \sigma^2 \mathbf{I}$. The posterior mean and covariance (13) require the adjoint states (6), the inner products (11), and the inverse of \mathbf{D} . Since N is typically much larger than M , the computational cost is dominated by the cost of computing the adjoint states.

We note that the posterior mean $\bar{\mathbf{g}}$ depends linearly on $(\mathbf{d} - \mathbf{s}^o)$, which is the difference between the observation vector and the output vector of the original model (1). Another way to look at the posterior mean is to view it as a linear combination of M inner products, each one associated with one of the adjoint states, by writing

$$(14) \quad \bar{g}_i = \sum_{j=1}^M \beta_j k(\psi_i, \phi_j), \quad i = 1, \dots, N,$$

where $\beta = \mathbf{D}^{-1} (\mathbf{d} - \mathbf{s}^o)$. We see that the discrepancy between the observed data and the original model (1) has an important effect on the posterior mean.

Also note that the posterior covariance \mathbf{G} in (13) does not explicitly depend on the values of the observed outputs but only on the adjoint states Φ and the covariance operator $k(\cdot, \cdot)$. This is based on the premise that the covariance operator $k(\cdot, \cdot)$ is somehow known and given. In reality, of course, we create/invent $k(\cdot, \cdot)$ based on notions both of general principles as well as computational expedients. As will be seen in section 3.3, we will construct/model a suitable covariance operator $k(\cdot, \cdot)$ based on data observations and maximum likelihood estimation. In this way the

posterior covariance \mathbf{G} in (13) will implicitly depend on the observed outputs and our modeling assumptions. Furthermore, notice that the covariance matrix \mathbf{G} in (13) is the difference of two terms: the first term $\mathbf{K}(\Psi, \Psi)$ is simply the prior covariance, while the second term represents a reduction in the covariance due to the observations.

Finally, we express the state estimate as $u(\mathbf{x}) = \sum_{n=1}^N u_n \psi_n(\mathbf{x})$, where $\mathbf{u} = (u_1, \dots, u_N) \in \mathbb{C}^N$ denotes the state vector. Substituting this expression into the model (3) and choosing $v = \psi_i$ for $i = 1, \dots, N$, we obtain the following linear system:

$$(15) \quad \mathbf{A}\mathbf{u} = \mathbf{l} - \mathbf{g},$$

where \mathbf{A} , \mathbf{l} , and \mathbf{g} have entries $A_{ij} = a(\psi_i, \psi_j)$, $l_i = \ell(\psi_i)$, and $g_i = g(\psi_i)$ for $i, j = 1, \dots, N$, respectively. It follows from (12) and (15) that the state vector \mathbf{u} obeys a normal distribution

$$(16) \quad \mathbf{u} | (\mathbf{d} - \mathbf{s}^o), \Phi, \Psi \sim \mathcal{N}(\bar{\mathbf{u}}, \mathbf{U}),$$

where the posterior mean vector $\bar{\mathbf{u}} \in \mathbb{C}^N$ and covariance matrix $\mathbf{U} \in \mathbb{C}^{N \times N}$ are given by

$$(17) \quad \bar{\mathbf{u}} = \mathbf{A}^{-1}(\mathbf{l} - \bar{\mathbf{g}}), \quad \mathbf{U} = \mathbf{A}^{-1} \mathbf{G} \mathbf{A}^{-\text{H}}.$$

Here $\bar{\mathbf{g}}$ and \mathbf{G} are given by (13). The superscript H denotes the complex conjugate transpose.

Let us use the notation $\boldsymbol{\psi}(\mathbf{x}) = (\psi_1(\mathbf{x}), \dots, \psi_N(\mathbf{x})) \in \mathbb{C}^N$ to denote a vector of values of the N basis functions at \mathbf{x} . It follows from (16) that

$$(18) \quad u(\mathbf{x}) \sim \mathcal{N}(\bar{u}(\mathbf{x}), \eta(\mathbf{x})),$$

where the mean value function $\bar{u}(\mathbf{x})$ and the variance function $\eta(\mathbf{x})$ are given by

$$(19) \quad \bar{u}(\mathbf{x}) = (\boldsymbol{\psi}(\mathbf{x}))^{\text{H}} \bar{\mathbf{u}}, \quad \eta(\mathbf{x}) = (\boldsymbol{\psi}(\mathbf{x}))^{\text{H}} \mathbf{U} \boldsymbol{\psi}(\mathbf{x}).$$

The mean value function represents the prediction of the true state $u^{\text{true}}(\mathbf{x})$, while the variance function quantifies the prediction error.

3.3. Covariance operator. In order to determine the posterior distribution of the state estimate, the covariance operator $k(\cdot, \cdot)$ needs to be specified. Indeed, this is a crucial aspect of our approach because it affects the estimate in (18). We propose a class of bilinear covariance operators of the form

$$(20) \quad k(w, v; \boldsymbol{\theta}) = \sum_{j=1}^J \theta_j k_j(w, v),$$

where $\boldsymbol{\theta} = (\theta_1, \theta_2, \dots, \theta_J)$ is a vector of *hyperparameters* and $k_j(\cdot, \cdot)$, $1 \leq j \leq J$ are some given symmetric bilinear forms. The particular definition of the bilinear forms k_j will depend on the type of the PDE model as well as the numerical method used to solve the PDE. We will elaborate on this point later when we discuss our numerical results in section 5.

We now note from (10) that the posterior distribution of $(\mathbf{d} - \mathbf{s}^o)$, conditional on $\boldsymbol{\theta}$, obeys the normal distribution:

$$(21) \quad (\mathbf{d} - \mathbf{s}^o) | (\boldsymbol{\theta}, \sigma) \sim \mathcal{N}(0, \mathbf{D}(\boldsymbol{\theta}, \sigma)),$$

where the matrix $\mathbf{D}(\boldsymbol{\theta}, \sigma)$ has entries

$$(22) \quad D_{mq}(\boldsymbol{\theta}, \sigma) = \sum_{j=1}^J \theta_j k_j(\phi_m, \phi_q) + \sigma^2 \Delta_{mq}, \quad m, q = 1, \dots, M,$$

where $\Delta_{mq} = 1$ if $m = q$ and $\Delta_{mq} = 0$ if $m \neq q$. Let $p((\mathbf{d} - \mathbf{s}^\circ)|(\boldsymbol{\theta}, \sigma))$ denote the probability density function of (21), and let us choose the hyperparameter vector $\boldsymbol{\theta}$ guided by the maximum likelihood principle, which by monotonicity is equivalent to log-maximum likelihood. The log likelihood function is given by

$$(23) \quad \begin{aligned} \log p((\mathbf{d} - \mathbf{s}^\circ)|(\boldsymbol{\theta}, \sigma)) &= -\frac{1}{2}(\mathbf{d} - \mathbf{s}^\circ)^H \mathbf{D}^{-1}(\boldsymbol{\theta}, \sigma)(\mathbf{d} - \mathbf{s}^\circ) \\ &\quad - \frac{1}{2} \log(\det(\mathbf{D}(\boldsymbol{\theta}, \sigma))) - \frac{M}{2} \log(2\pi). \end{aligned}$$

We therefore propose to select $\boldsymbol{\theta}$ by seeking an approximate optimal solution to the log-likelihood optimization problem, yielding

$$(24) \quad (\boldsymbol{\theta}, \sigma) \approx \arg \max_{(\boldsymbol{\theta}', \sigma') \in \Theta} \log p((\mathbf{d} - \mathbf{s}^\circ)|(\boldsymbol{\theta}', \sigma')).$$

Here the hyperparameter space Θ is a subset of \mathbb{R}^{J+1} such that the resulting covariance operator k in (20) is positive-definite. The optimization problem (24) does not have a closed form solution, and it is generically nonconvex. We therefore seek approximate solutions of (24) in cases of interest. For example, when the hyperparameter space Θ is low-dimensional, simple grid generation/enumeration and function evaluation will be sufficient for our purposes.

3.4. Bayesian interpretation of the functional regression model. In this section, we present an alternative derivation of the functional regression model via Bayesian analysis. The matrix form of the Galerkin FE formulation (1) seeks $(\mathbf{u}^\circ, \mathbf{s}^\circ) \in \mathbb{C}^N \times \mathbb{C}^M$ such that

$$(25) \quad \mathbf{A}\mathbf{u}^\circ = \mathbf{l}, \quad \mathbf{s}^\circ = \mathbf{C}^H \mathbf{u}^\circ,$$

where $\mathbf{u}^\circ \in \mathbb{C}^N$ is the vector of degrees of freedom of the FE solution u° and $C_{im} = c_m(\psi_i)$ for $i = 1, \dots, N$ and $m = 1, \dots, M$. Similarly, the matrix form of the stochastic weak formulation (3) seeks $(\mathbf{u}, \mathbf{s}) \in \mathbb{C}^N \times \mathbb{C}^M$ such that

$$(26) \quad \mathbf{A}\mathbf{u} + \mathbf{g} = \mathbf{l}, \quad \mathbf{s} = \mathbf{C}^H \mathbf{u},$$

where \mathbf{g} is a *Gaussian random vector* with the prior distribution $\mathcal{N}(\mathbf{0}, \mathbf{K})$. Note that the prior covariance matrix \mathbf{K} has entries $K_{ij} = k(\psi_i, \psi_j)$.

To define the likelihood function, we compute the adjoint state vectors stored in the matrix $\boldsymbol{\Phi} \in \mathbb{C}^{N \times M}$ by solving the adjoint problem:

$$(27) \quad \mathbf{A}^H \boldsymbol{\Phi} = -\mathbf{C}.$$

It follows from (25), (26), and (27) that

$$(28) \quad \boldsymbol{\Phi}^H \mathbf{g} = \boldsymbol{\Phi}^H (\mathbf{l} - \mathbf{A}\mathbf{u}) = \boldsymbol{\Phi}^H (\mathbf{A}\mathbf{u}^\circ - \mathbf{A}\mathbf{u}) = -\mathbf{C}^H (\mathbf{u}^\circ - \mathbf{u}) = \mathbf{s} - \mathbf{s}^\circ.$$

We now substitute $\mathbf{s} = \mathbf{d} - \boldsymbol{\varepsilon}$ into (28) to arrive at the standard linear regression model:

$$(29) \quad \mathbf{d} - \mathbf{s}^\circ = \boldsymbol{\Phi}^H \mathbf{g} + \boldsymbol{\varepsilon}.$$

Since $\varepsilon \sim \mathcal{N}(0, \sigma^2 \mathbf{I})$, $(\mathbf{d} - \mathbf{s}^\circ) | \mathbf{g} \sim \mathcal{N}(\Phi^H \mathbf{g}, \sigma^2 \mathbf{I})$. Therefore, the likelihood function, i.e., the probability density function of $\mathbf{d} - \mathbf{s}^\circ$ given the vector \mathbf{g} , is

$$(30) \quad p((\mathbf{d} - \mathbf{s}^\circ) | \mathbf{g}, \Phi) = \frac{1}{(2\pi\sigma^2)^{M/2}} \exp\left(-\frac{1}{2\sigma^2} \|\mathbf{d} - \mathbf{s}^\circ - \Phi^H \mathbf{g}\|^2\right),$$

where $\|\cdot\|$ denotes the Euclidean norm.

According to Bayes's theorem (see (2.3) on page 17 in [33]), the posterior distribution of \mathbf{g} is given by

$$(31) \quad \text{posterior} = \frac{\text{likelihood} \times \text{prior}}{\text{marginal likelihood}}, \quad p(\mathbf{g} | (\mathbf{d} - \mathbf{s}^\circ), \Phi) = \frac{p((\mathbf{d} - \mathbf{s}^\circ) | \Phi, \mathbf{g}) p(\mathbf{g})}{p((\mathbf{d} - \mathbf{s}^\circ) | \Phi)},$$

where the marginal likelihood $p((\mathbf{d} - \mathbf{s}^\circ) | \Phi)$ is given by

$$(32) \quad p((\mathbf{d} - \mathbf{s}^\circ) | \Phi) = \int p((\mathbf{d} - \mathbf{s}^\circ) | \Phi, \mathbf{g}) p(\mathbf{g}) d\mathbf{g}.$$

Writing only the terms from the likelihood and the prior, we obtain

$$(33) \quad \begin{aligned} p(\mathbf{g} | (\mathbf{d} - \mathbf{s}^\circ), \Phi) &\propto \exp\left(-\frac{1}{2\sigma^2} \|\mathbf{d} - \mathbf{s}^\circ - \Phi^H \mathbf{g}\|^2\right) \exp\left(-\frac{1}{2} \mathbf{g}^H \mathbf{K}^{-1} \mathbf{g}\right), \\ &\propto \exp\left(-\frac{1}{2} (\mathbf{g} - \bar{\mathbf{g}})^H \mathbf{G}^{-1} (\mathbf{g} - \bar{\mathbf{g}})\right), \end{aligned}$$

where $\bar{\mathbf{g}}$ and \mathbf{G} are given by

$$(34) \quad \bar{\mathbf{g}} = \frac{1}{\sigma^2} \mathbf{G} \Phi (\mathbf{d} - \mathbf{s}^\circ), \quad \mathbf{G} = \left(\frac{1}{\sigma^2} \Phi \Phi^H + \mathbf{K}^{-1}\right)^{-1}.$$

However, since this formula requires the inverse of a matrix of size $N \times N$, we will derive an equivalent expression which is more efficient to compute than (34).

We recall the Sherman–Morrison–Woodbury (SMW) formula for the matrix inversion

$$(35) \quad (\mathbf{Z} + \mathbf{U} \mathbf{W} \mathbf{V}^H)^{-1} = \mathbf{Z}^{-1} - \mathbf{Z}^{-1} \mathbf{U} (\mathbf{W}^{-1} + \mathbf{V}^H \mathbf{Z}^{-1} \mathbf{U})^{-1} \mathbf{V}^H \mathbf{Z}^{-1}.$$

Using the SMW formula with $\mathbf{Z} = \mathbf{K}^{-1}$, $\mathbf{W} = \sigma^{-2} \mathbf{I}$, $\mathbf{U} = \mathbf{V} = \Phi$ we obtain

$$(36) \quad \mathbf{G} = \mathbf{K} - \mathbf{K} \Phi (\sigma^2 \mathbf{I} + \Phi^H \mathbf{K} \Phi)^{-1} \Phi^H \mathbf{K}.$$

It thus follows that

$$(37) \quad \begin{aligned} \mathbf{G} \Phi &= \mathbf{K} \Phi - \mathbf{K} \Phi (\sigma^2 \mathbf{I} + \Phi^H \mathbf{K} \Phi)^{-1} \Phi^H \mathbf{K} \Phi, \\ &= \mathbf{K} \Phi (\sigma^2 \mathbf{I} + \Phi^H \mathbf{K} \Phi)^{-1} ((\sigma^2 \mathbf{I} + \Phi^H \mathbf{K} \Phi) - \Phi^H \mathbf{K} \Phi), \\ &= \sigma^2 \mathbf{K} \Phi (\sigma^2 \mathbf{I} + \Phi^H \mathbf{K} \Phi)^{-1}. \end{aligned}$$

Hence, we obtain from (34)–(37) that

$$(38) \quad \bar{\mathbf{g}} = \mathbf{K} \Phi \mathbf{D}^{-1} (\mathbf{d} - \mathbf{s}^\circ), \quad \mathbf{G} = \mathbf{K} - \mathbf{K} \Phi \mathbf{D}^{-1} \Phi^H \mathbf{K},$$

where $\mathbf{D} = \sigma^2 \mathbf{I} + \Phi^H \mathbf{K} \Phi$. In this way, we need to invert the matrix \mathbf{D} of size $M \times M$. Since M is typically much smaller than N , the formula (38) is much more efficient to compute than the original one (34).

We note from the stochastic linear system (26) and the Gaussian property of \mathbf{g} that the posterior distribution of \mathbf{u} is a multivariate normal distribution

$$(39) \quad \mathbf{u} | (\mathbf{d} - \mathbf{s}^o), \Phi \sim \mathcal{N}(\bar{\mathbf{u}}, \mathbf{U}),$$

where the posterior mean and the posterior covariance matrix are given by (17), namely,

$$(40) \quad \bar{\mathbf{u}} = \mathbf{A}^{-1}(\mathbf{l} - \bar{\mathbf{g}}), \quad \mathbf{U} = \mathbf{A}^{-1} \mathbf{G} \mathbf{A}^{-\text{H}}.$$

Substituting (38) into (40) yields

$$(41) \quad \begin{aligned} \bar{\mathbf{u}} &= \mathbf{A}^{-1} (\mathbf{l} - \mathbf{K} \Phi (\sigma^2 \mathbf{I} + \Phi^{\text{H}} \mathbf{K} \Phi)^{-1} (\mathbf{d} - \mathbf{s}^o)), \\ \mathbf{U} &= \mathbf{A}^{-1} (\mathbf{K} - \mathbf{K} \Phi (\sigma^2 \mathbf{I} + \Phi^{\text{H}} \mathbf{K} \Phi)^{-1} \Phi^{\text{H}} \mathbf{K}) \mathbf{A}^{-\text{H}}. \end{aligned}$$

We observe that this result is exactly the same as the one presented in subsection 3.2.

4. Experimental design. The selection of inputs to obtain the observations is an *experimental design problem*. This problem can be posed as choosing the inputs to minimize the posterior covariance with some appropriate design criteria [32], such as A-optimality, D-optimality, E-optimality, and G-optimality [2, 30]. In this section, we develop an experimental design algorithm based on the criterion of minimizing a monotone function of the nonzero eigenvalues.

4.1. Optimizing experimental design. We assume that we are given a large set of L potential observation functionals $b_i(\cdot), 1 \leq i \leq L$, each of which corresponds to a potential experiment to be carried out. Before actually performing any of these potential experiments, we face the following decision questions. Which experiment among the L potential experiments should be implemented first? Given the results of any previous experiments, which experiment should be chosen next? How many experiments are enough to provide an accurate prediction? These decision questions can be summed up as follows: how do we optimally choose a sequence of the actualized observation functionals $c_m(\cdot), 1 \leq m \leq M$, among the potential observation functionals $b_i(\cdot), 1 \leq i \leq L$?

It will be more convenient to state our experimental design problem in the matrix form instead of in the weak form. The expressions involving $\mathbf{K}(\Psi, \Psi), \mathbf{K}(\Psi, \Phi), \mathbf{K}(\Phi, \Psi)$, etc., can be rather unwieldy, so we introduce a compact form of the notation by setting $\mathbf{K} = \mathbf{K}(\Psi, \Psi), \mathbf{K}(\Psi, \Phi) = -\mathbf{K} \mathbf{A}^{-\text{H}} \mathbf{C}, \mathbf{K}(\Phi, \Psi) = -\mathbf{C}^{\text{H}} \mathbf{A}^{-1} \mathbf{K}$, and $\mathbf{K}(\Phi, \Phi) = \mathbf{C}^{\text{H}} \mathbf{A}^{-1} \mathbf{K} \mathbf{A}^{-\text{H}} \mathbf{C}$. It follows that the posterior covariance matrix given in (17) can be written as

$$(42) \quad \mathbf{U} = \mathbf{A}^{-1} (\mathbf{K} - \mathbf{K} \mathbf{A}^{-\text{H}} \mathbf{C} (\sigma^2 \mathbf{I} + \mathbf{C}^{\text{H}} \mathbf{A}^{-1} \mathbf{K} \mathbf{A}^{-\text{H}} \mathbf{C})^{-1} \mathbf{C}^{\text{H}} \mathbf{A}^{-1} \mathbf{K}) \mathbf{A}^{-\text{H}}.$$

We can also rewrite the posterior covariance matrix \mathbf{U} as

$$(43) \quad \mathbf{U}(\mathbf{C}) = \mathbf{U}_0 - \mathbf{U}_0 \mathbf{C} (\sigma^2 \mathbf{I} + \mathbf{C}^{\text{H}} \mathbf{U}_0 \mathbf{C})^{-1} \mathbf{C}^{\text{H}} \mathbf{U}_0,$$

where $\mathbf{U}_0 = \mathbf{A}^{-1} \mathbf{K} \mathbf{A}^{-\text{H}}$ is the prior covariance matrix.

We now let $\mathbf{B} \in \mathbb{C}^{N \times L}$ be a matrix with entries $B_{nl} = b_l(\psi_n)$ for $n = 1, \dots, N$ and $l = 1, \dots, L$. Let M denote the number of experiments that we will implement, and consider the following set of matrices induced by the matrix \mathbf{B} :

$$(44) \quad \mathcal{M}_M(\mathbf{B}) := \{\mathbf{E} \in \mathbb{C}^{N \times M} : \text{each column of } \mathbf{E} \text{ is chosen among the columns of } \mathbf{B}\}.$$

In the experimental design problem, we seek an experiment matrix $\mathbf{C} \in \mathcal{M}_M(\mathbf{B})$ that minimizes the “size” of $\mathbf{U}(\mathbf{C})$, where size is measured in various ways in various design formulations such as the product of the eigenvalues of $\mathbf{U}(\mathbf{C})$, the largest eigenvalue of $\mathbf{U}(\mathbf{C})$, the sum of the eigenvalues of $\mathbf{U}(\mathbf{C})$, etc. In the interest of generality, let $h(\cdot) : \mathbb{R}^N \rightarrow \mathbb{R}$ be a given monotone function of the ordered eigenvalues $\lambda(\mathbf{U}) := (\lambda_1(\mathbf{U}), \dots, \lambda_N(\mathbf{U}))$ of \mathbf{U} , and consider the following optimization problem:

$$(45) \quad \mathbf{C}^* = \arg \min_{\mathbf{E} \in \mathcal{M}_M(\mathbf{B})} h(\lambda(\mathbf{U}(\mathbf{E}))) ,$$

where $\mathbf{U}(\mathbf{E})$ is the posterior covariance matrix as given in (43) for $\mathbf{E} = \mathbf{C}$. In other words, the optimal experiment design \mathbf{C}^* is the matrix that minimizes a particular given monotone function of the eigenvalues of the posterior covariance matrix.

Note that the optimization problem (45) has a nonlinear (and typically nonconvex) objective function and that the set of feasible solutions is an exponentially large discrete set. Furthermore, the optimal observation matrix \mathbf{C}^* is not hierarchical in the sense that increasing M might completely change the previously selected columns of \mathbf{C}^* . Because the optimal set of experiments might be very different for different values of M , one cannot compute the exact optimum by inductively amending the previous optimal set of experiments. For all of the above reasons it therefore is likely to be too computationally expensive to compute an exact optimum of (45). Below we develop a greedy iterative procedure for choosing the experiments based on the eigenvector of the largest eigenvalue of $\mathbf{U}(\mathbf{C})$.

4.2. A greedy algorithm for choosing the experiments. Here we develop a greedy algorithm for choosing the experiments. We first consider the case where there is no noise in the observation vector \mathbf{d} , i.e., $\sigma = 0$ for the Gaussian noise. With $\sigma = 0$ it follows that \mathbf{u} must satisfy $\mathbf{C}^H \mathbf{u} = \mathbf{d}$. Let us define $\mathcal{T} := \{\mathbf{u} : \mathbf{C}^H \mathbf{u} = \mathbf{d}\}$ as the affine set containing the possible values of \mathbf{u} . Let us assume without loss of generality that the columns of \mathbf{C} are linearly independent, so $\dim(\mathcal{T}) = N - M$. Recalling $\bar{\mathbf{u}} = \mathbf{A}^{-1}(\mathbf{l} - \bar{\mathbf{g}})$ from (17), it holds that $\bar{\mathbf{u}} \in \mathcal{T}$. The posterior covariance matrix (43) can be written as

$$(46) \quad \mathbf{U} = \mathbf{U}_0 - \mathbf{U}_0 \mathbf{C} (\mathbf{C}^H \mathbf{U}_0 \mathbf{C})^{-1} \mathbf{C}^H \mathbf{U}_0 .$$

Let us also write the eigendecomposition of \mathbf{U} as $\mathbf{U} = \mathbf{Q} \mathbf{\Lambda} \mathbf{Q}^H$ for a suitable orthonormal matrix of eigenvectors \mathbf{Q} and diagonal matrix $\mathbf{\Lambda}$ of corresponding eigenvalues. Since \mathbf{U} is positive semi-definite, i.e., $\mathbf{u}^H \mathbf{U} \mathbf{u} \geq 0$ for all \mathbf{u} , all eigenvalues of \mathbf{U} are nonnegative, and we can partition the eigenvalues into those that are zero and those that are positive, and likewise partition the corresponding eigenvectors, and write

$$(47) \quad \mathbf{U} = \mathbf{Q} \mathbf{\Lambda} \mathbf{Q}^H = [\mathbf{Q}_1 \quad \mathbf{Q}_2] \begin{bmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{\Lambda}^+ \end{bmatrix} \begin{bmatrix} \mathbf{Q}_1^H \\ \mathbf{Q}_2^H \end{bmatrix} = \mathbf{Q}_2 \mathbf{\Lambda}^+ \mathbf{Q}_2^H ,$$

where the columns of \mathbf{Q}_1 are eigenvectors corresponding to eigenvalues with value 0 and the columns of \mathbf{Q}_2 are eigenvectors corresponding to positive eigenvalues, and the diagonal matrix $\mathbf{\Lambda}^+$ has positive diagonal entries. In a slight change of notation, let us denote these entries as $\Lambda_i^+ > 0$ for $i = 1, \dots, N - M$, and there is no loss of generality in assuming that $0 < \Lambda_1^+ \leq \Lambda_2^+ \leq \dots \leq \Lambda_{N-M}^+$. It is straightforward to demonstrate that

$$(48) \quad \mathbf{u}^H \mathbf{U} \mathbf{u} = 0 \iff \mathbf{U} \mathbf{u} = \mathbf{0} \iff \mathbf{u} \in \text{range}(\mathbf{C}) \iff \mathbf{u} \in \text{null}(\mathbf{Q}_2^H) .$$

For the given observation vector \mathbf{d} , \mathbf{u} is restricted to lie in the affine set \mathcal{T} , and we write $\mathbf{u}|\mathbf{d} \sim \mathcal{N}_{\mathcal{T}}(\bar{\mathbf{u}}, \mathbf{U})$ to designate that the posterior distribution of \mathbf{u} is normally distributed on \mathcal{T} with mean $\bar{\mathbf{u}}$ and covariance matrix \mathbf{U} . Even though \mathbf{U} has no inverse, the probability density function of \mathbf{u} on \mathcal{T} , which we denote as $p_{\mathcal{T}}(\mathbf{u}|\mathbf{d})$, is given by

$$(49) \quad p_{\mathcal{T}}(\mathbf{u}|\mathbf{d}) = \frac{1}{\pi^{(N-M)/2} \sqrt{\prod_{i=1}^{N-M} \Lambda_i^+}} e^{-\left[\frac{1}{2}(\mathbf{u}-\bar{\mathbf{u}})^H \mathbf{Q}_2(\Lambda^+)^{-1} \mathbf{Q}_2^H(\mathbf{u}-\bar{\mathbf{u}})\right]} .$$

Notice the precise way that the positive eigenvalues appear in (49).

Now let us suppose that we have already chosen M experiments from among the L experiments $b_i(\cdot)$, $i = 1, \dots, L$, and let \mathbf{C} denote the experiment matrix, and let us consider adding an additional experiment, which would be the $(M + 1)$ st experiment. Let this experiment be denoted as $\tilde{b}(\cdot)$, where $\tilde{b}(\cdot)$ is one of the L experiments $b_i(\cdot)$, $i = 1, \dots, L$. In the matrix form, let $\tilde{\mathbf{b}}$ be the vector with entries $\tilde{b}(\psi_n)$ for $n = 1, \dots, N$. We will assume that $\tilde{\mathbf{b}} \notin \text{range}(\mathbf{C})$, so that the new/updated experiment matrix $\tilde{\mathbf{C}} := [\mathbf{C} \ \tilde{\mathbf{b}}]$ has linearly independent columns. Then letting $\tilde{\mathbf{U}} := \mathbf{U}(\tilde{\mathbf{C}})$, we have from (46) that

$$(50) \quad \tilde{\mathbf{U}} = \mathbf{U}_0 - \mathbf{U}_0 \tilde{\mathbf{C}} (\tilde{\mathbf{C}}^H \mathbf{U}_0 \tilde{\mathbf{C}})^{-1} \tilde{\mathbf{C}}^H \mathbf{U}_0 .$$

It turns out that $\tilde{\mathbf{U}}$ simplifies to

$$(51) \quad \tilde{\mathbf{U}} = \mathbf{U} - \left(\frac{1}{\tilde{\mathbf{b}}^H \mathbf{U} \tilde{\mathbf{b}}} \right) \mathbf{U} \tilde{\mathbf{b}} \tilde{\mathbf{b}}^H \mathbf{U} .$$

(Equation (51) is derived by straightforward, but tedious, substitution using

$$\begin{aligned} & (\tilde{\mathbf{C}}^H \mathbf{U}_0 \tilde{\mathbf{C}})^{-1} \\ &= \left[\begin{array}{cc} \mathbf{C}^H \mathbf{U}_0 \mathbf{C} & \mathbf{C}^H \mathbf{U}_0 \tilde{\mathbf{b}} \\ \tilde{\mathbf{b}}^H \mathbf{U}_0 \mathbf{C} & \tilde{\mathbf{b}}^H \mathbf{U}_0 \tilde{\mathbf{b}} \end{array} \right]^{-1} \\ &= \left[\begin{array}{cc} (\mathbf{C}^H \mathbf{U}_0 \mathbf{C})^{-1} - \frac{(\mathbf{C}^H \mathbf{U}_0 \mathbf{C})^{-1} \mathbf{C}^H \mathbf{U}_0 \tilde{\mathbf{b}} \tilde{\mathbf{b}}^H \mathbf{U}_0 \mathbf{C} (\mathbf{C}^H \mathbf{U}_0 \mathbf{C})^{-1}}{\alpha} & \frac{(\mathbf{C}^H \mathbf{U}_0 \mathbf{C})^{-1} \mathbf{C}^H \mathbf{U}_0 \tilde{\mathbf{b}}}{\alpha} \\ \frac{\tilde{\mathbf{b}}^H \mathbf{U}_0 \mathbf{C} (\mathbf{C}^H \mathbf{U}_0 \mathbf{C})^{-1}}{\alpha} & \frac{1}{\alpha} \end{array} \right], \end{aligned}$$

where $\alpha := \tilde{\mathbf{b}}^H \mathbf{U} \tilde{\mathbf{b}}$, and using the block form of the inverse above in (50) and simplifying terms. Note that the assumption that $\tilde{\mathbf{b}} \notin \text{range}(\mathbf{C})$ implies via (48) that $\alpha = \tilde{\mathbf{b}}^H \mathbf{U} \tilde{\mathbf{b}} \neq 0$ so the above objects are all well-defined.)

Let us now compare the ordered eigenvalues of \mathbf{U} and $\tilde{\mathbf{U}}$, which we denote, respectively, as the arrays

$$\begin{aligned} \boldsymbol{\lambda} &:= (\lambda_1 \quad \lambda_2 \quad \cdots \quad \lambda_N) , \\ \tilde{\boldsymbol{\lambda}} &:= (\tilde{\lambda}_1 \quad \tilde{\lambda}_2 \quad \cdots \quad \tilde{\lambda}_N) . \end{aligned}$$

From the linear independence assumption regarding the columns of \mathbf{C} and $\tilde{\mathbf{C}}$ it follows using (48) that $\lambda_1 = \cdots = \lambda_M = 0$ and $\lambda_{M+1} > 0$ and also that $\tilde{\lambda}_1 = \cdots = \tilde{\lambda}_{M+1} = 0$ and $\tilde{\lambda}_{M+2} > 0$. Similar to (47), we will write

$$(52) \quad \tilde{\mathbf{U}} = [\tilde{\mathbf{Q}}_1 \quad \tilde{\mathbf{Q}}_2] \left[\begin{array}{cc} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \tilde{\Lambda}^+ \end{array} \right] \left[\begin{array}{c} \tilde{\mathbf{Q}}_1^H \\ \tilde{\mathbf{Q}}_2^H \end{array} \right] = \tilde{\mathbf{Q}}_2 \tilde{\Lambda}^+ \tilde{\mathbf{Q}}_2^H ,$$

where the columns of $\tilde{\mathbf{Q}}_1$ are eigenvectors corresponding to eigenvalues with value 0 and the columns of $\tilde{\mathbf{Q}}_2$ are eigenvectors corresponding to positive eigenvalues, and the diagonal matrix $\tilde{\Lambda}^+$ has positive diagonal entries. As with (52), we denote these entries as $\tilde{\Lambda}_i^+ > 0$ for $i = 1, \dots, N - M - 1$, and there is no loss of generality in assuming that $0 < \tilde{\Lambda}_1^+ \leq \tilde{\Lambda}_2^+ \leq \dots \leq \tilde{\Lambda}_{N-M-1}^+$. Therefore we have

$$\begin{aligned} \boldsymbol{\lambda} &= \begin{pmatrix} 0 & 0 & \cdots & 0 & \Lambda_1^+ & \Lambda_2^+ & \Lambda_3^+ & \cdots & \Lambda_{N-M}^+ \end{pmatrix}, \\ \tilde{\boldsymbol{\lambda}} &= \begin{pmatrix} 0 & 0 & \cdots & 0 & 0 & \tilde{\Lambda}_1^+ & \tilde{\Lambda}_2^+ & \cdots & \tilde{\Lambda}_{N-M-1}^+ \end{pmatrix}, \end{aligned}$$

where there are M leading zeros in $\boldsymbol{\lambda}$ and $M + 1$ leading zeros in $\tilde{\boldsymbol{\lambda}}$ above.

Observe from (51) that $\tilde{\mathbf{U}}$ is a rank-1 modification of \mathbf{U} . As a consequence, it follows from the interlacing eigenvalue theorem (see Theorem 8.1.8 of [15] and discussion therein for example) that the eigenvalues of \mathbf{U} and $\tilde{\mathbf{U}}$ are *interlaced*, namely,

$$(53) \quad \tilde{\lambda}_1 \leq \lambda_1 \leq \tilde{\lambda}_2 \leq \lambda_2 \leq \dots \leq \lambda_{N-1} \leq \tilde{\lambda}_N \leq \lambda_N .$$

It therefore follows that

$$\Lambda_1^+ \leq \tilde{\Lambda}_1^+ \leq \Lambda_2^+ \leq \tilde{\Lambda}_2^+ \leq \dots \leq \Lambda_{N-M-1}^+ \leq \tilde{\Lambda}_{N-M-1}^+ \leq \Lambda_{N-M}^+ .$$

This chain of inequalities is significant in that it provides lower bounds on the values of the updated eigenvalues $\tilde{\Lambda}_i^+$ for $i = 1, \dots, N - M - 1$:

$$(54) \quad \tilde{\Lambda}_i^+ \geq \Lambda_i^+ , \text{ for } i = 1, \dots, N - M - 1 .$$

Note that this result is true for *any* newly added experiment $\tilde{\mathbf{b}}$ satisfying the condition that $\tilde{\mathbf{C}} := [\mathbf{C} \ \tilde{\mathbf{b}}]$ has linearly independent columns.

Let us now consider if there might be an experiment $\tilde{\mathbf{b}}$ for which (54) holds at equality for all $i = 1, \dots, N - M - 1$, thus *simultaneously* providing the lowest value of the new eigenvalues $\tilde{\Lambda}_i^+$ for $i = 1, \dots, N - M - 1$. Let \mathbf{q}_N be an eigenvector of \mathbf{U} corresponding to the largest eigenvalue of \mathbf{U} , namely, $\lambda_N = \Lambda_{N-M}^+$. Suppose that we are able to set the new experiment $\tilde{\mathbf{b}}$ so that $\tilde{\mathbf{b}} = \mathbf{q}_N$. In this case it follows from (51) and (47) that

$$\begin{aligned} \tilde{\mathbf{U}} &= \mathbf{U} - \left(\frac{1}{\mathbf{q}_N^H \mathbf{U} \mathbf{q}_N} \right) \mathbf{U} \mathbf{q}_N \mathbf{q}_N^H \mathbf{U} \\ &= \mathbf{Q} \boldsymbol{\Lambda} \mathbf{Q}^H - \left(\frac{1}{\mathbf{q}_N^H \mathbf{Q} \boldsymbol{\Lambda} \mathbf{Q}^H \mathbf{q}_N} \right) \mathbf{Q} \boldsymbol{\Lambda} \mathbf{Q}^H \mathbf{q}_N \mathbf{q}_N^H \mathbf{Q} \boldsymbol{\Lambda} \mathbf{Q}^H \\ &= \mathbf{Q} \boldsymbol{\Lambda} \mathbf{Q}^H - \lambda_N \mathbf{q}_N \mathbf{q}_N^H , \end{aligned}$$

and it follows that the columns of \mathbf{Q} comprise a set of orthonormal eigenvectors of $\tilde{\mathbf{U}}$ and that the nonzero eigenvalues of $\tilde{\mathbf{U}}$ are precisely $\tilde{\Lambda}_i^+$ for $i = 1, \dots, N - M - 1$. Therefore the nonzero eigenvalues of $\tilde{\mathbf{U}}$ are in fact $\tilde{\Lambda}_i^+ = \Lambda_i^+$ for $i = 1, \dots, N - M - 1$. We therefore see that when $\tilde{\mathbf{b}} = \mathbf{q}_N$, (54) is satisfied at equality for all $i = 1, \dots, N - M - 1$. In light of the lower bounds given in (54) it holds that $\tilde{\mathbf{b}} = \mathbf{q}_N$ is the ideal next experiment as it simultaneously minimizes *any* (and hence *every*) monotone function $h(\cdot)$ of the nonzero eigenvalues of $\tilde{\mathbf{U}}$.

From a practical point of view, we would not expect there to be an experiment $\tilde{\mathbf{b}}_i(\cdot)$ whose corresponding $\tilde{\mathbf{b}}_i$ vector satisfies $\tilde{\mathbf{b}}_i = \mathbf{q}_N$ (or $\tilde{\mathbf{b}}_i = -\mathbf{q}_N$). Instead, it makes good intuitive sense to seek an experiment $\tilde{\mathbf{b}}_i$ that makes the smallest angle with $\pm \mathbf{q}_N$.

This leads to the criterion to choose the next experiment by finding the index i_{next}^* such that

$$i_{\text{next}}^* = \arg \max_{i \in \{1, \dots, L\}} \frac{|\mathbf{b}_i^H \mathbf{q}_N|}{\|\mathbf{b}_i\|}.$$

This is formalized in the following algorithm.

0. (Initial values.) $M \leftarrow 0$. Given the matrix \mathbf{K} , set the initial posterior covariance matrix \mathbf{U} to be the prior covariance matrix:

$$(55) \quad \mathbf{U} = \mathbf{A}^{-1} \mathbf{K} \mathbf{A}^{-H}.$$

1. (Compute eigenvector of largest eigenvalue.) Solve the eigenvalue problem $\mathbf{U} \mathbf{v} = \lambda \mathbf{v}$ to obtain the maximum eigenvalue λ_{max} and the corresponding eigenvector \mathbf{v}_{max} .
2. (Determine the next experiment to add.) Find an index i_{next}^* such that

$$(56) \quad i_{\text{next}}^* = \arg \max_{i \in \{1, \dots, L\}} \frac{|\mathbf{b}_i^H \mathbf{v}_{\text{max}}|}{\|\mathbf{b}_i\|},$$

and append $\mathbf{b}_{i_{\text{next}}^*}$ to the current matrix \mathbf{C} to form the updated matrix $\mathbf{C} \leftarrow [\mathbf{C} \ \mathbf{b}_{i_{\text{next}}^*}]$. Perform the new experiment to obtain the experiment outcome value d_{M+1} . Update the number of experiments: $M \leftarrow M + 1$. If $M \geq M_{\text{max}}$, stop. (Optional: apply other stopping criterion to determine whether to stop.)

3. (Update posterior covariance matrix.) Given the updated matrix \mathbf{C} , calculate the updated posterior covariance matrix \mathbf{U} as

$$(57) \quad \mathbf{U} = \mathbf{A}^{-1} (\mathbf{K} - \mathbf{K} \mathbf{A}^{-H} \mathbf{C} (\sigma^2 \mathbf{I} + \mathbf{C}^H \mathbf{A}^{-1} \mathbf{K} \mathbf{A}^{-H} \mathbf{C})^{-1} \mathbf{C}^H \mathbf{A}^{-1} \mathbf{K}) \mathbf{A}^{-H}.$$

4. (Optional: update hyperparameters.) Optionally, given the updated experiments and their outcomes, reset the hyperparameters $\boldsymbol{\theta} = (\theta_1, \dots, \theta_J)$ of the covariance operator $k(\cdot, \cdot, \boldsymbol{\theta})$ in (20) by approximately solving the likelihood maximization problem (24) as described in section 3.3.
5. Go to step 1.

While the algorithm is best conceptualized for the case when $\sigma = 0$, we have found that in practice it also works well for $\sigma > 0$, using the intuitive notion that for small values of σ the mathematics developed and used in this section remains approximately valid. Indeed, as we will see in section 5, the algorithm works very well for $\sigma > 0$ in the applications of interest herein.

Notice that the computational complexity of this algorithm scales very favorably with the number of possible experiments L . Specifically, the algorithm needs to evaluate L scalar products $\mathbf{b}_i^H \mathbf{v}_{\text{max}}, i = 1, \dots, L$ only once at each of the M steps. Therefore, the algorithm can accommodate a large list of possible experiments without too much computational effort so long as M is relatively small.

Let us now discuss the optional step 4 in the above algorithm. By our way of choosing the covariance operator discussed in subsection 3.4, the matrix \mathbf{K} depends on the observations. Ultimately, the prior covariance matrix \mathbf{K} plays an important role in the posterior covariance matrix \mathbf{U} as it is given by (57). By updating the hyperparameters when a new experiment is chosen, we hope to further reduce the “size” of the posterior covariance matrix \mathbf{U} through a good choice of the prior covariance matrix \mathbf{K} . In all numerical examples presented in the next section, we carry out the optional step 4 in our greedy algorithm.

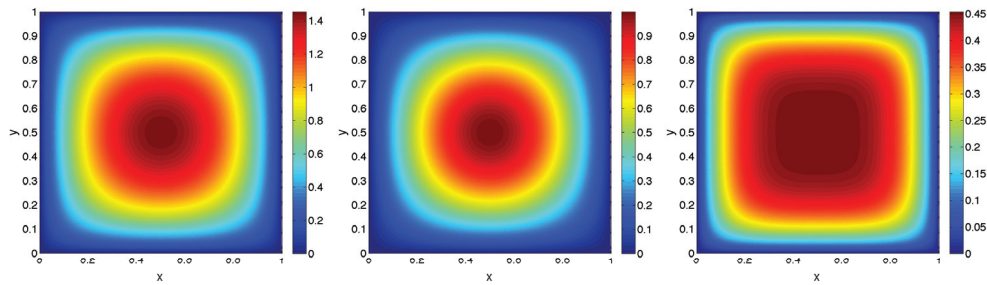


FIG. 1. Plot of u^o (left), u^{true} (middle), and their difference $u^o - u^{\text{true}}$ (right).

5. Numerical results. We will demonstrate the proposed method on numerical examples from a variety of second-order linear PDEs such as heat equations, convection-diffusion equations, and Helmholtz equations. In order to verify the performance of our method, we will specify the true state and generate the observations by adding Gaussian noise to the true outputs. Since our examples involve second-order PDEs, we will consider the covariance operator of the form

$$(58) \quad k(w, v; \theta) = \theta_1 \int_{\Omega} w v d\mathbf{x} + \theta_2 \int_{\Omega} \nabla w \cdot \nabla v d\mathbf{x} ,$$

which is related to the $H^1(\Omega)$ inner product. Note that the covariance operator is positive-semidefinite when the hyperparameters are nonnegative.

5.1. Heat conduction example. We consider the following PDE model for heat conduction in a unit square domain:

$$(59) \quad -\Delta u^o = f^o \quad \text{in } \Omega \quad \text{and} \quad u^o = 0 \quad \text{on } \partial\Omega,$$

where $\Omega \equiv (0, 1) \times (0, 1)$ and $f^o = 2\pi^2$. The Galerkin FE formulation of this model problem has the bilinear form and the linear functional

$$(60) \quad a(w, v) = \int_{\Omega} \nabla w \cdot \nabla v dx dy, \quad \ell(v) = \int_{\Omega} f^o v dx dy \quad \forall w, v \in V,$$

respectively. Here V is a FE approximation space of piecewise linear polynomials defined on a finite element mesh of 5,000 elements. Figure 1 shows a plot of the finite element solution of our PDE model (59). The possible observation functionals are specified as

$$(61) \quad b_i(v) = \int_{\Omega} \exp\left(-\frac{(x-x_i)^2 + (y-y_i)^2}{0.25^2}\right) v dx dy, \quad i = 1, \dots, L,$$

where $(x_i, y_i), i = 1, \dots, L = 961$ are the points of a 31×31 uniform grid as shown in Figure 2. These spatial points represent possible measurement locations at which we obtain the observations. We assume that the observations are noise-free.

Furthermore, we assume that the true state is $u^{\text{true}} = \sin(\pi x) \sin(\pi y)$ as depicted in Figure 1. We note that the true state u^{true} is the solution of the PDE model (59) in which the source term $f^o = 2\pi^2$ is replaced with $f^{\text{true}} = 2\pi^2 \sin(\pi x) \sin(\pi y)$. Therefore, the source of uncertainty in our PDE model (59) comes from the source term. Due to this uncertainty, as shown in Figure 1, the solution u^o of the PDE model

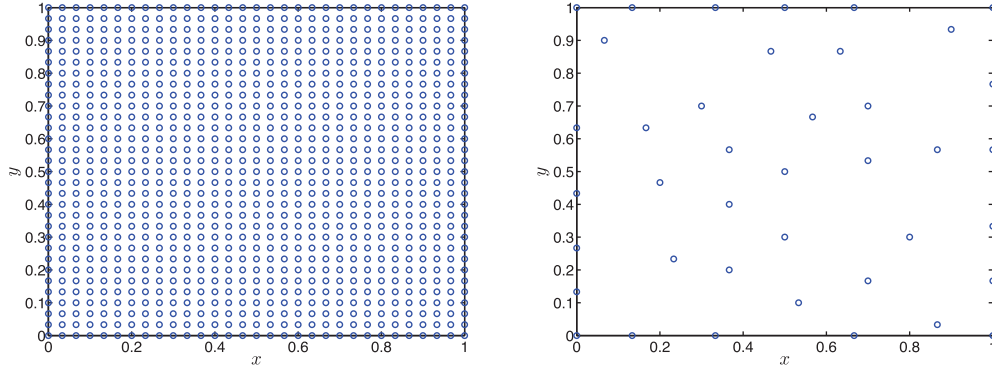


FIG. 2. Potential measurement points (left) and selected measurement locations (right).

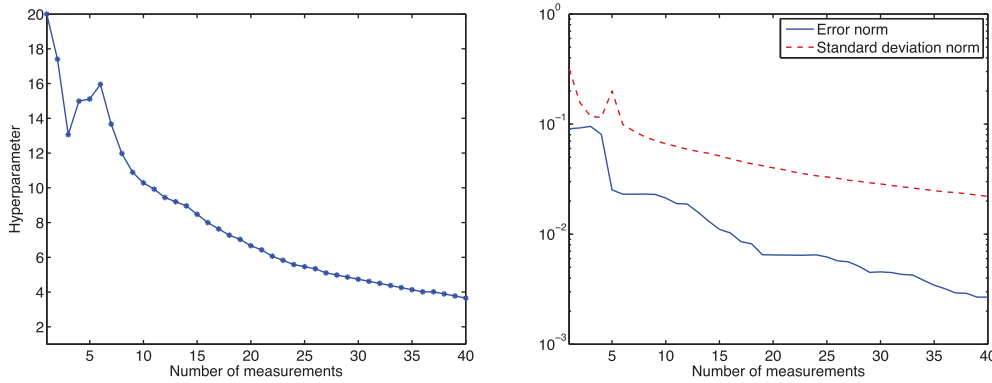


FIG. 3. The first hyperparameter θ_1 (left) and the error norm and the standard deviation norm (right) as a function of M . Note that the maximum likelihood approach yields $\theta_2 = 0$ and $\sigma = 0$.

(59) is significantly different from the true state u^{true} . In what follows, we combine the PDE model (59) with observations to improve the prediction of the true state by using the functional regression method.

The hyperparameters $(\theta_1, \theta_2, \sigma)$ are chosen to maximize the log marginal likelihood on a uniform grid $50 \times 40 \times 40$ of the hyperparameter domain $\Theta \equiv [0.1, 25] \times [0, 1] \times [0, 0.1]$. We next pursue the greedy algorithm to determine the measurement locations, which are shown in Figure 2. We observe that the selected measurement locations are distributed over the whole domain. Figure 3 shows the hyperparameter θ_1 and the $L^2(\Omega)$ norm of the standard function $\|\sqrt{\eta}\|_{\Omega}$ as well as the error norm $\|u^{\text{true}} - \bar{u}\|_{\Omega}$ as a function of the number of measurements. It is interesting to note that the maximum likelihood approach yields $\theta_2 = 0$ and $\sigma = 0$. This means that our bilinear covariance operator k is collinear with the L^2 inner product and that $\sigma = 0$ coincides with the noise-free observations. Figure 4 shows the posterior variance function and the selected measurement points at each iteration of our greedy algorithm. We observe that the variance reduces as we increase the number of measurements and that the measurement points are typically selected near the maximum peak of the posterior variance function.

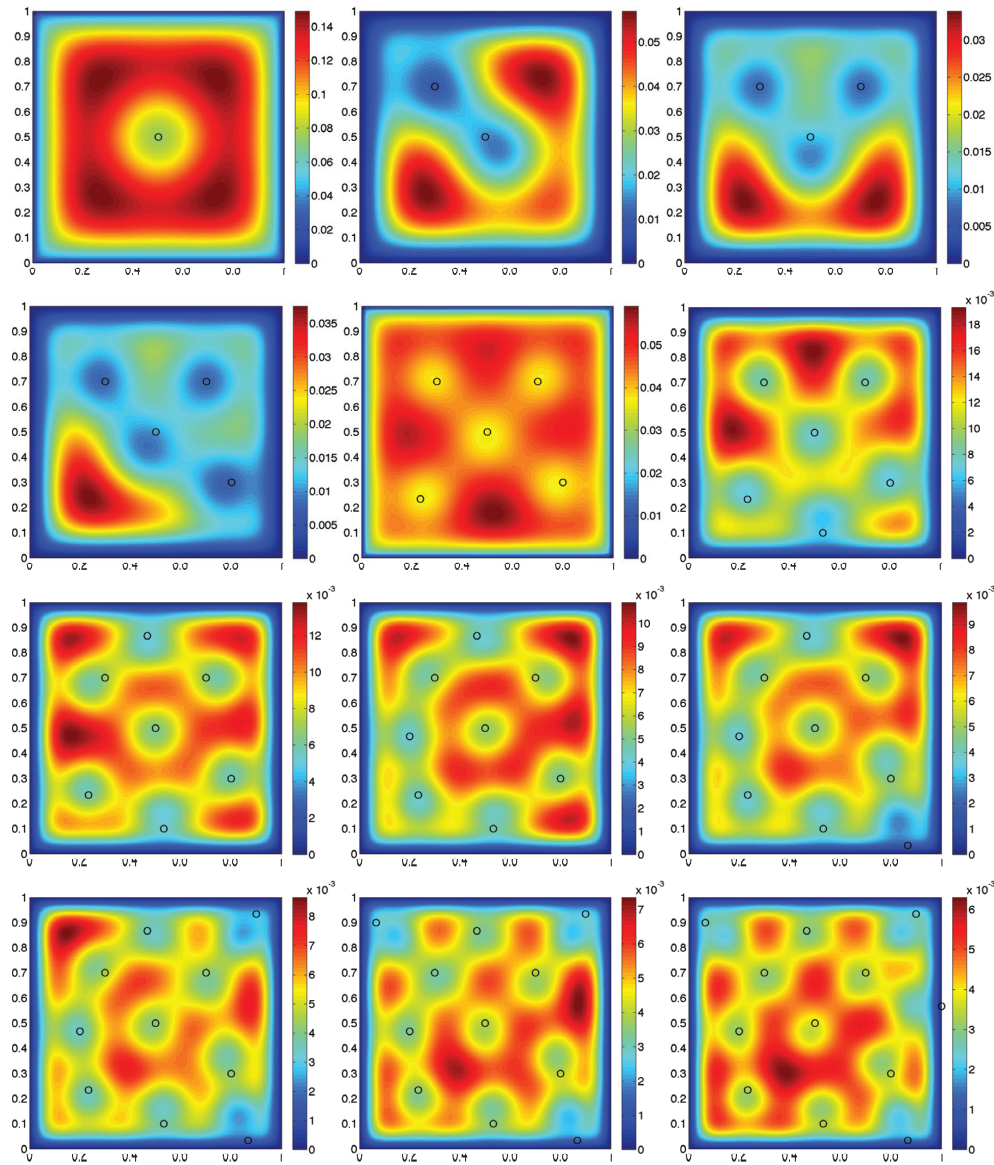


FIG. 4. The variance function $\eta(\mathbf{x})$ and the measurement points for $M = 1, \dots, 12$, where the order goes from left to right and top to bottom.

5.2. Convection-diffusion example. We next consider a convection-diffusion problem on a unit square $\Omega \equiv (0, 1) \times (0, 1)$:

$$(62) \quad -\Delta u^\circ + \mathbf{c}^\circ \cdot \nabla u^\circ = f^\circ \quad \text{in } \Omega \quad \text{and} \quad u^\circ = 0 \quad \text{on } \partial\Omega,$$

where $\mathbf{c}^\circ = (10, 10)$ and $f^\circ = 10$. The Galerkin FE formulation of the above PDE model is derived with the bilinear form and the linear functional

$$(63) \quad a(w, v) = \int_{\Omega} \nabla w \cdot \nabla v dx dy + \int_{\Omega} (\mathbf{c}^\circ \cdot \nabla w) v dx dy, \quad \ell(v) = \int_{\Omega} f^\circ v dx dy \quad \forall w, v \in V,$$

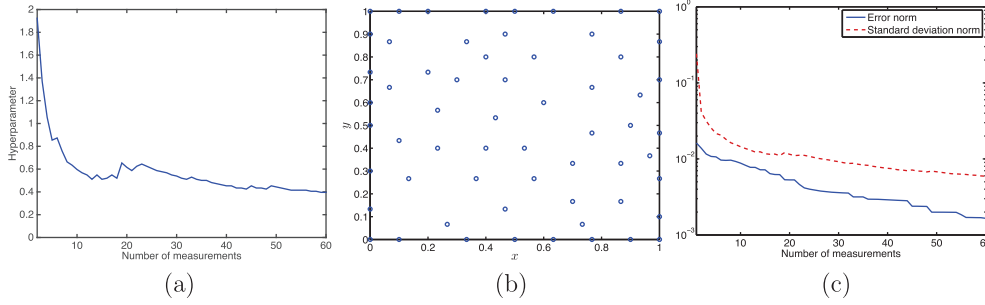


FIG. 5. The hyperparameter θ_1 versus M (a), the measurement locations (b), and the error and the standard deviation versus M (c). Note that $\theta_2 = 0$ and $\sigma = 0$.

where V is a FE approximation space of piecewise polynomials of degree $p = 3$ and is defined on a mesh of 800 elements. The possible measurement points and observation functionals are the same as in the previous example. We assume that the observations are noise-free.

To assess our method, we assume that the true state u^{true} is the solution of the PDE model (62) in which the convective velocity $\mathbf{c}^o = (10, 10)$ is replaced with $\mathbf{c}^{\text{true}} = (15, 15)$. Hence, the source of uncertainty in our PDE model (62) comes from the convective velocity. In this example, the uncertainty is present in the bilinear form a , whereas in the previous example the uncertainty is present in the linear functional ℓ . This example serves to demonstrate the effectiveness of our method on cases where uncertainty affects PDE operators.

The covariance operator has the form (58), where (θ_1, θ_2) and σ are computed by maximizing the likelihood on a uniform grid $200 \times 50 \times 50$ of the domain $\Theta \equiv [0.1, 10] \times [0, 1] \times [0, 0.1]$. Note that $\theta_2 = 0$ and $\sigma = 0$, while θ_1 is shown in Figure 5(a). Figure 5 also shows the measurement locations and the norm of the standard function $\|\sqrt{\eta}\|_{\Omega}$ as well as the error norm $\|u^{\text{true}} - \bar{u}\|_{\Omega}$ as a function of the number of measurements. As expected, both the error norm and the standard deviation norm decrease as M increases. Hence, increasing the number of measurements reduces the uncertainty. This can be seen more clearly in Figure 6, which shows the absolute error function and the standard deviation function for $M = 20, 40$, and 60 . We observe that the error near the top right corner of the domain is much larger than the one near the bottom left corner of the domain, while the standard deviation is spread over the physical domain much more evenly than the error.

5.3. A Helmholtz example. We consider the sound-hard scattering of an incident plane wave $u^{\text{inc}} = \exp(ikx)$ by a circle of radius $a = 1$, where k is the wave number. The scattered field satisfies the exterior Helmholtz problem

$$\begin{aligned}
 \Delta u^{\text{true}} + k^2 u^{\text{true}} &= 0 \quad \text{in } \mathbb{R}^2 \setminus \Omega_c, \\
 \nabla u^{\text{true}} \cdot \mathbf{n} + \nabla u^{\text{inc}} \cdot \mathbf{n} &= 0 \quad \text{on } \Gamma_c, \\
 \lim_{r \rightarrow \infty} \sqrt{r} \left(\frac{\partial u^{\text{true}}}{\partial r} - ik u^{\text{true}} \right) &= 0,
 \end{aligned}
 \tag{64}$$

where Ω_c is the domain of the unit circle and Γ_c is the boundary of the unit circle. The last condition is known as Sommerfeld radiation condition. It is known that the

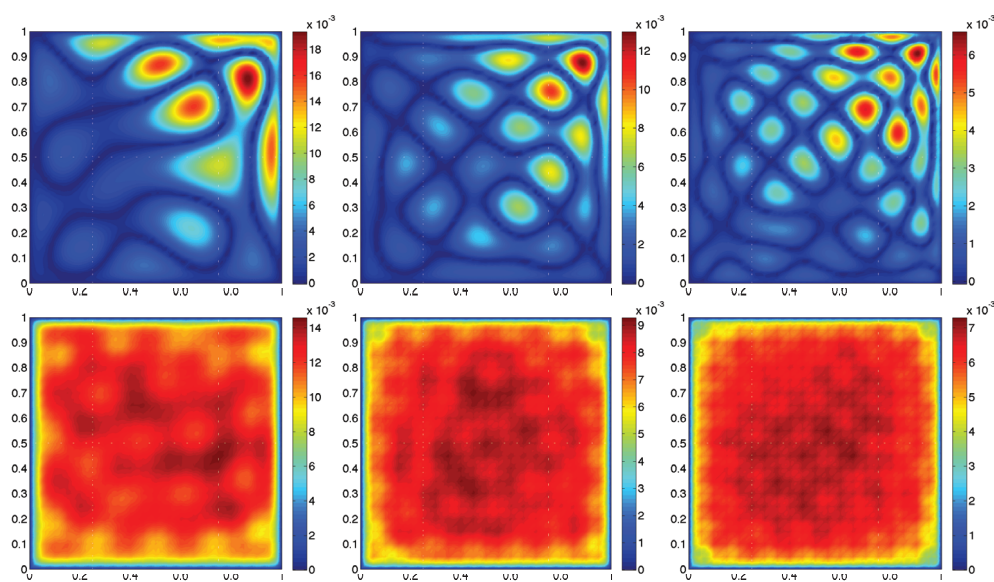


FIG. 6. The absolute error function $|\bar{u}(\mathbf{x}) - u^{\text{true}}(\mathbf{x})|$ (top row) and the standard deviation function $\sqrt{\eta(\mathbf{x})}$ (bottom row) for $M = 20$ (left), $M = 40$ (middle), and $M = 60$ (right).

solution of the above Helmholtz problem has the analytical form

$$(65) \quad u^{\text{true}}(x, y) = -\frac{J'_0(ka)}{H_0^1(ka)} H_0^1(kr) - 2 \sum_{n=1}^{\infty} i^n \frac{J'_n(ka)}{H_n^1(ka)} H_n^1(kr) \cos(n\theta),$$

where (r, θ) are polar coordinates of (x, y) , and J_n and H_n^1 are the Bessel function and the Hankel function of the first kind, respectively [17]. Note that the prime denotes the derivative of a function with respect to its argument.

Since the scatterer is simple the solution of the exterior Helmholtz problem (64) is known analytically. For more complicated scatterers, however, the problem is usually solved by using a numerical method. In the case of a finite element method, the unbounded domain needs to be truncated to a bounded domain and the Sommerfeld radiation condition should be replaced with a suitable absorbing boundary condition at the exterior boundary of the truncated domain. This gives rise to the following PDE model:

$$(66) \quad \begin{aligned} \Delta u^o + k^2 u^o &= 0 && \text{in } \Omega, \\ \nabla u^o \cdot \mathbf{n} + \nabla u^{\text{inc}} \cdot \mathbf{n} &= 0 && \text{on } \Gamma_c, \\ \nabla u^o \cdot \mathbf{n} - ik u^o &= 0 && \text{on } \Gamma_o, \end{aligned}$$

where Ω is the truncated domain as shown in Figure 7 and Γ_o is the exterior boundary of the truncated domain. Note here that we approximate the Sommerfeld radiation condition by the first-order absorbing boundary condition. The weak formulation of the PDE model (66) is given by

$$(67) \quad a(u^o, v) = \ell(v) \quad \forall v \in X,$$

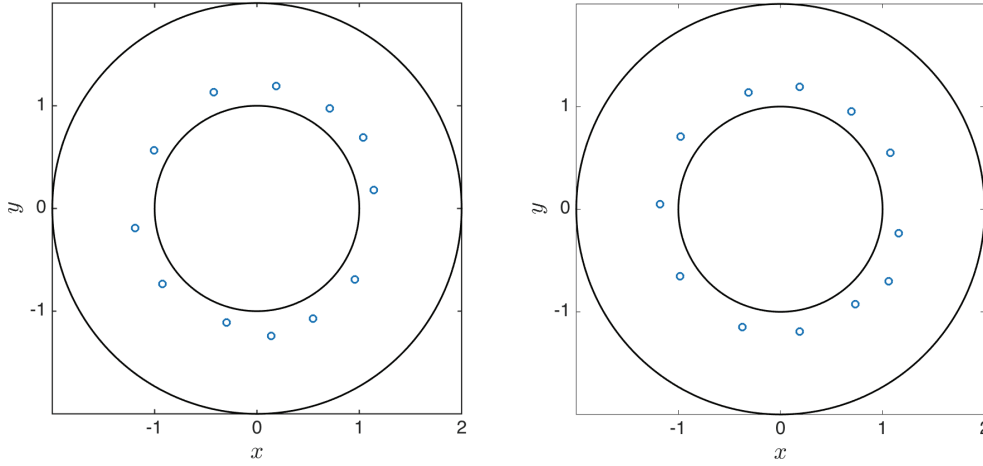


FIG. 7. The measurement locations obtained using the greedy algorithm for two different noise levels: $\sigma_{\text{noise}} = 0.005$ (left) and $\sigma_{\text{noise}} = 0.01$ (right).

where, for all $w, v \in X$,

$$(68) \quad \begin{aligned} a(w, v) &= \int_{\Omega} (\nabla w \cdot \nabla v - k^2 wv) dx dy - \int_{\Gamma_o} ikwv dx dy, \\ \ell(v) &= - \int_{\Gamma_c} (\nabla u^{\text{inc}} \cdot \mathbf{n}) v dx dy, \end{aligned}$$

and $X \equiv H^1(\Omega)$. In practice, we replace X with a finite-dimensional space V which is a FE approximation space of piecewise polynomials of degree $p = 3$ and is defined on a mesh of 800 elements. This FE discretization is fine enough such that the numerical error is negligible.

This example serves to demonstrate our approach when the uncertainty sources originate from the domain truncation and the boundary condition. We aim to improve the prediction of the scattered field u^{true} by combining the PDE model (66) with observations of the scattered field. Of particular interest is the value of the scattered field on the unit circle, as it ultimately determines the radar cross section. The possible observation functionals are specified as

$$(69) \quad b_i(v) = \int_{\Omega} \exp\left(-\frac{(x - x_i)^2 + (y - y_i)^2}{0.25^2}\right) v dx dy, \quad i = 1, \dots, L,$$

where $(x_i, y_i), i = 1, \dots, L = 3720$ coincide with the mesh points. We will consider noisy observations which are obtained by adding Gaussian noise of zero mean and standard deviation σ_{noise} to the noise-free observations.

The hyperparameters $(\theta_1, \theta_2, \sigma)$ are computed by maximizing the likelihood on a uniform grid $100 \times 50 \times 50$ of the domain $\Theta \equiv [0.01, 1] \times [0, 1] \times [0, 0.02]$. We present in Table 1 the hyperparameters for $M = 4, 8, 12$. We observe that σ is quite close to σ_{noise} for $M = 12$. Hence, the likelihood maximization approach yields a good estimate of the noise level when the number of measurements is sufficient. We present in Figure 7 the selected measurement locations for two different noise levels in the observations: $\sigma_{\text{noise}} = 0.005$ and $\sigma_{\text{noise}} = 0.01$. It is interesting to see that in all cases

TABLE 1

The hyperparameters $(\theta_1, \theta_2, \sigma)$ for the two different noise levels in the observations.

M	$\sigma_{\text{noise}} = 0.005$			$\sigma_{\text{noise}} = 0.01$		
	θ_1	θ_2	σ	θ_1	θ_2	σ
4	0.0300	0.0000	0.0000	0.0100	0.0000	0.0125
8	0.0282	0.0000	0.0081	0.0400	0.0000	0.0070
12	0.0455	0.0000	0.0049	0.0336	0.0000	0.0086

the measurement locations are quite close to the unit circle and distributed quite uniformly along the unit circle.

We next show the true state, the mean prediction, and the 95% confidence region (shaded area) for $\sigma_{\text{noise}} = 0.005$ in Figure 8, and for $\sigma_{\text{noise}} = 0.01$ in Figure 9. Here the 95% confidence region is an area bounded by the mean prediction plus and minus two times the standard deviation function. We see that increasing M can improve the accuracy of the mean prediction for $\sigma_{\text{noise}} = 0.005$. However, for $\sigma_{\text{noise}} = 0.01$, the mean prediction does not get better as we increase M . This means that when σ_{noise} is equal or greater than 0.01, the observations do not add a positive contribution to the prediction of the true state. For $\sigma_{\text{noise}} = 0.005$ the 95% confidence region shrinks quite rapidly as M increases, whereas for $\sigma_{\text{noise}} = 0.01$ the 95% confidence region does not shrink as M increases. Therefore, when the observations become too noisy, adding more observations does not help improve the prediction. Clearly, the mean prediction and error estimate get worse as the noise level increases.

6. Conclusions. We have presented a statistical method that combines a linear PDE model with observations to predict the state of a physical problem. First, a random functional is introduced into the PDE model to account for various sources of uncertainty in the model. This random functional is posed as a Gaussian process with zero mean and prior covariance operator. Next, a linear regression model for the Gaussian functional is derived by utilizing the adjoint states and the observations. This regression model allows us to compute the posterior distribution of both the Gaussian functional and the state estimate. A key ingredient of our method is the prior covariance operator of the Gaussian functional. We propose a class of affine bilinear forms for the prior covariance operator and determine the associated hyperparameters based on the observations. Furthermore, we devise a greedy algorithm to select observations among a large number of possible measurements.

We would like to extend our approach to nonlinear PDEs. Two new challenges arise in the nonlinear case. First, although the functional g is Gaussian, the state u is no longer a Gaussian random field because of the nonlinearity. Second, the nonlinearity prevents a direct link between the adjoint solutions and the observations, thereby complicating the regression model. Our research will focus on addressing these issues.

Appendix A. Relationship with variational data assimilation. Let $\beta = D^{-1}(d - s^o)$ and $\bar{q} = \Phi\beta$. We can easily show that $(\bar{u}, \bar{q}, \beta)$ is the optimal solution of the following least squares minimization problem

$$(70) \quad \min_{(z, w, \gamma) \in \mathbb{C}^N \times \mathbb{C}^N \times \mathbb{C}^M} \frac{1}{2} w^H K w + \frac{1}{2} \sigma^2 \gamma^H \gamma$$

$$\text{s.t.} \quad A z + K w = l,$$

$$C^H z + \sigma^2 \gamma = d.$$

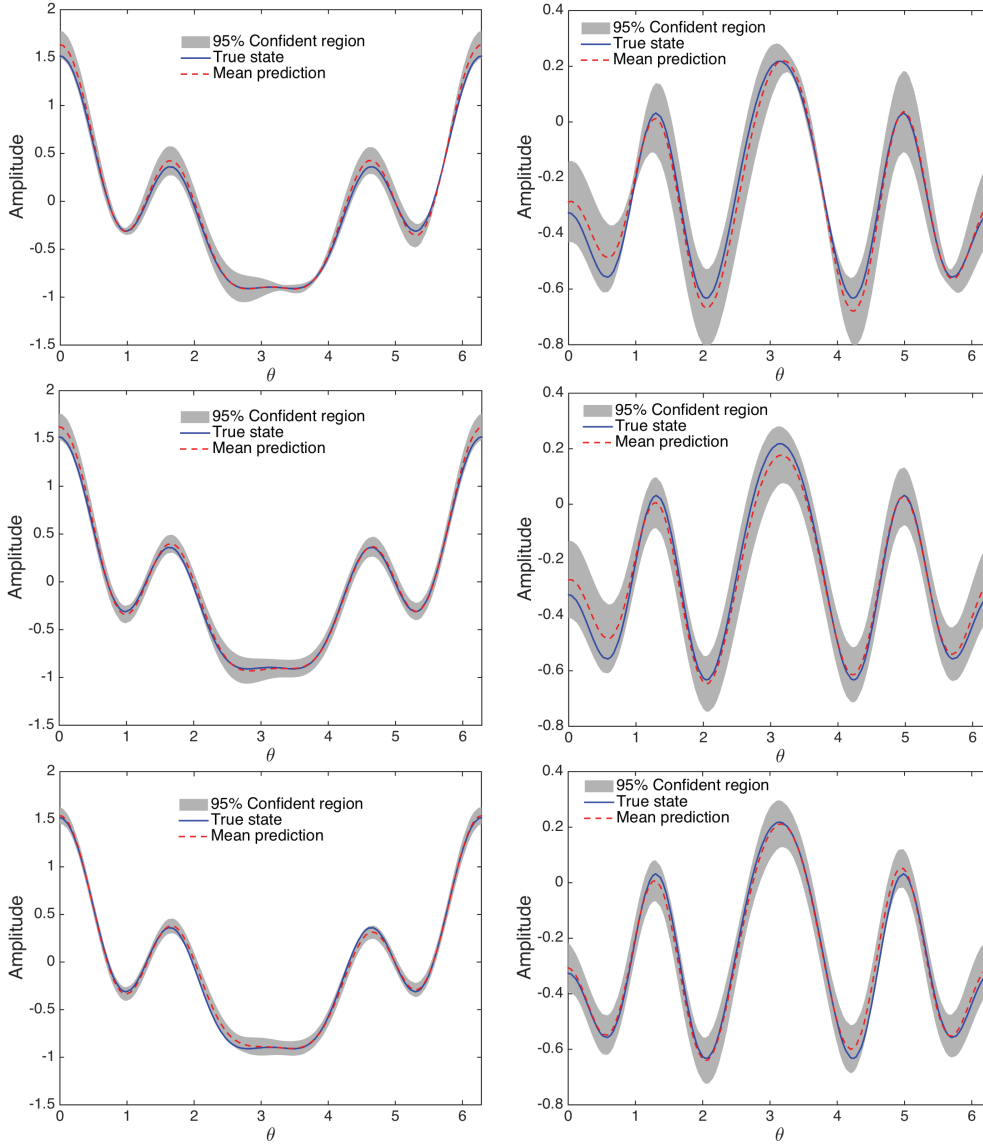


FIG. 8. Panels show the true state and the mean prediction as a function of θ : the real part (left) and the imaginary part (right). These results are obtained using $M = 4$ (top row), $M = 8$ (middle row), and $M = 12$ (bottom row) observations with $\sigma_{\text{noise}} = 0.005$. In these plots the shaded area represents the mean prediction plus and minus two times the standard deviation function (corresponding to the 95% confidence region).

By eliminating the constraints in the least-squares problem (70) we obtain the following result:

$$(71) \quad \bar{\mathbf{u}} := \arg \min_{\mathbf{z} \in \mathbb{C}^N} \frac{1}{2} (\mathbf{A}\mathbf{z} - \mathbf{l})^H \mathbf{K}^{-1} (\mathbf{A}\mathbf{z} - \mathbf{l}) + \frac{1}{2} \sigma^{-2} (\mathbf{C}^H \mathbf{z} - \mathbf{d})^H (\mathbf{C}^H \mathbf{z} - \mathbf{d}).$$

We see that the posterior mean is optimal in the sense that it minimizes an error objective function, which is defined as the sum of the model error and the output

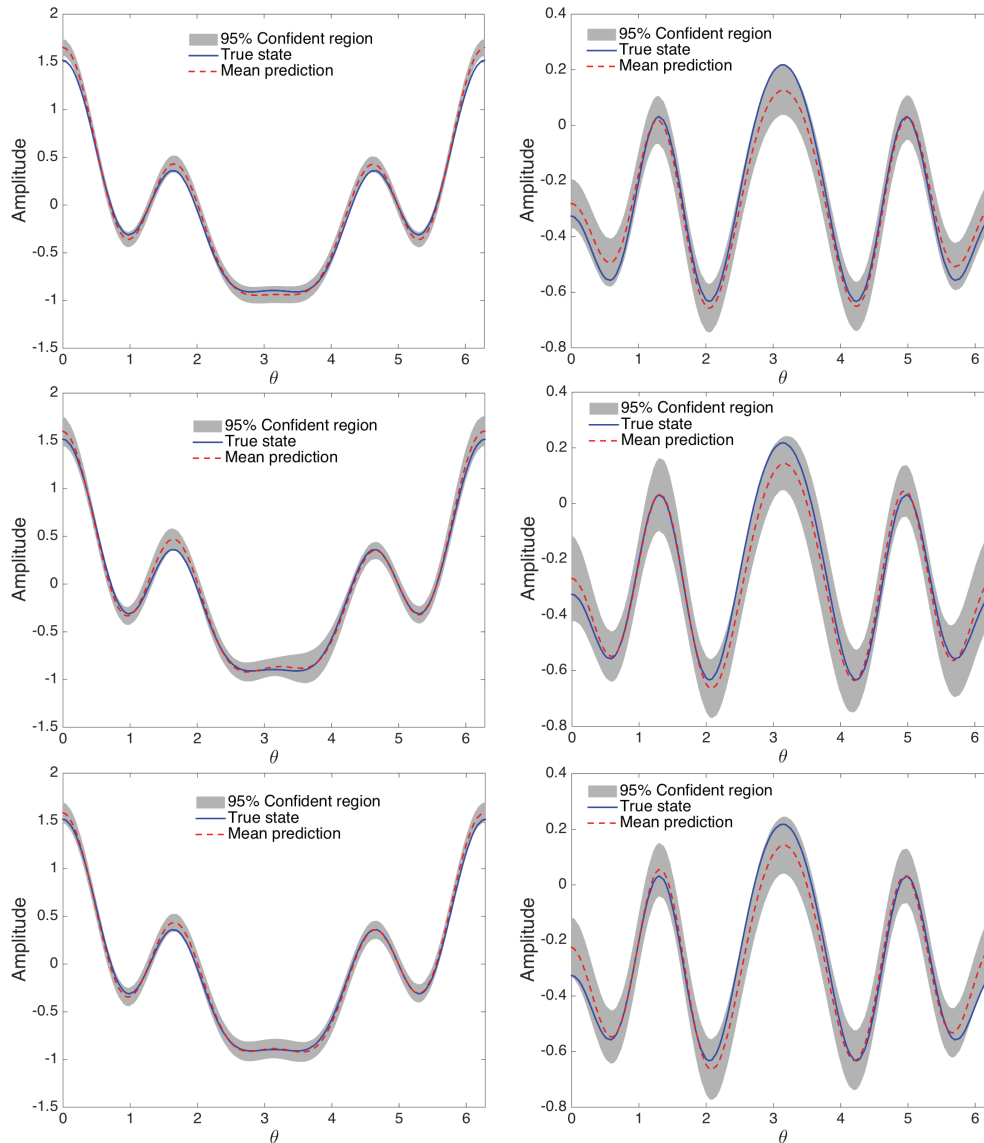


FIG. 9. Panels show the true state and the mean prediction as a function of θ : the real part (left) and the imaginary part (right). These results are obtained using $M = 4$ (top row), $M = 8$ (middle row), and $M = 12$ (bottom row) observations with $\sigma_{\text{noise}} = 0.01$. In these plots the shaded area represents the mean prediction plus and minus two times the standard deviation function (corresponding to the 95% confidence region).

error. Furthermore, we can write (71) as

$$(72) \quad \bar{\mathbf{u}} := \arg \min_{\mathbf{z} \in \mathbb{C}^N} \frac{1}{2} (\mathbf{z} - \mathbf{u}^o)^H \mathbf{U}_0^{-1} (\mathbf{z} - \mathbf{u}^o) + \frac{1}{2} \sigma^{-2} (\mathbf{C}^H \mathbf{z} - \mathbf{d})^H (\mathbf{C}^H \mathbf{z} - \mathbf{d}),$$

where $\mathbf{U}_0 = \mathbf{A}^{-1} \mathbf{K} \mathbf{A}^{-H}$. The optimization formulation (72) is known as 3D variational data assimilation [8, 22]. It shows the relation between our method and 3D variational data assimilation.

Appendix B. Relationship with the Kalman method. Let us review the Kalman method [19] for state estimation. We begin by assuming that the state vector \mathbf{u} is Gaussian with the prior distribution

$$(73) \quad p(\mathbf{u}) \propto \exp\left(-\frac{1}{2}(\mathbf{u} - \mathbf{u}^b)^H \mathbf{B}^{-1}(\mathbf{u} - \mathbf{u}^b)\right),$$

where \mathbf{u}^b is the background state vector and \mathbf{B} is the background covariance matrix. Furthermore, the data \mathbf{d} is assumed to have a Gaussian probability density function with covariance \mathbf{R} and mean $\mathbf{C}^H \mathbf{u}$. The likelihood function is thus given by

$$(74) \quad p(\mathbf{d}|\mathbf{u}) \propto \exp\left(-\frac{1}{2}(\mathbf{C}^H \mathbf{u} - \mathbf{d})^H \mathbf{R}^{-1}(\mathbf{C}^H \mathbf{u} - \mathbf{d})\right).$$

The posterior distribution of \mathbf{u} is given by Bayes' theorem (see equation (2.3) on page 17 in [33]):

$$(75) \quad p(\mathbf{u}|\mathbf{d}) \propto p(\mathbf{d}|\mathbf{u})p(\mathbf{u}).$$

It can be shown by algebraic manipulations that the posterior distribution is also Gaussian

$$(76) \quad p(\mathbf{u}|\mathbf{d}) \propto \exp\left(-\frac{1}{2}(\mathbf{u} - \hat{\mathbf{u}})^H \hat{\mathbf{U}}^{-1}(\mathbf{u} - \hat{\mathbf{u}})\right).$$

Here the posterior mean $\hat{\mathbf{u}}$ and covariance $\hat{\mathbf{U}}$ are given by the Kalman formulas

$$(77) \quad \hat{\mathbf{u}} = \mathbf{u}^b + \mathbf{H}(\mathbf{d} - \mathbf{C}^H \mathbf{u}^b), \quad \hat{\mathbf{U}} = \mathbf{B} - \mathbf{H} \mathbf{C}^H \mathbf{B},$$

where

$$(78) \quad \mathbf{H} = \mathbf{B} \mathbf{C} (\mathbf{C}^H \mathbf{B} \mathbf{C} + \mathbf{R})^{-1}$$

is the so-called Kalman gain matrix.

We now show that our method is related to the Kalman method for a particular choice of the background state vector and the background covariance matrix. In particular, let us choose the prior information as

$$(79) \quad \mathbf{u}^b = \mathbf{u}^o, \quad \mathbf{B} = \mathbf{A}^{-1} \mathbf{K} \mathbf{A}^{-H}, \quad \mathbf{R} = \sigma^2 \mathbf{I}.$$

It then follows from (77)–(79) that

$$(80) \quad \begin{aligned} \hat{\mathbf{u}} &= \mathbf{u}^o + \mathbf{A}^{-1} \mathbf{K} \mathbf{A}^{-H} \mathbf{C} (\mathbf{C}^H \mathbf{A}^{-1} \mathbf{K} \mathbf{A}^{-H} \mathbf{C} + \sigma^2 \mathbf{I})^{-1} (\mathbf{d} - \mathbf{s}^o), \\ \hat{\mathbf{U}} &= \mathbf{A}^{-1} \mathbf{K} \mathbf{A}^{-H} - \mathbf{A}^{-1} \mathbf{K} \mathbf{A}^{-H} \mathbf{C} (\mathbf{C}^H \mathbf{B} \mathbf{C} + \sigma^2 \mathbf{I})^{-1} \mathbf{C}^H \mathbf{A}^{-1} \mathbf{K} \mathbf{A}^{-H}. \end{aligned}$$

We see that the Kalman method yields exactly the same posterior distribution as our method when the priors are chosen by (79). Our method provides a systematic way to choose appropriate priors which are very important in the context of the Kalman method.

REFERENCES

- [1] A. P. DEMSTER, N. M. LAIRD, AND D. B. RUBIN, *Maximum Likelihood from Incomplete Data via the EM Algorithm*. J. Roy. Statist. Soc. Ser. B, 39 (1977), pp. 1–38.
- [2] A. ATKINSON, A. DONEV, AND R. TOBIAS, *Optimum Experimental Designs, with SAS*, Oxford University Press, Oxford, 2007.
- [3] I. BABUŠKA, F. NOBILE, AND R. TEMPONE, *A stochastic collocation method for elliptic partial differential equations with random input data*, SIAM J. Numer. Anal., 45 (2007), pp. 1005–1034.
- [4] M. BARRAULT, Y. MADAY, N. C. NGUYEN, AND A. T. PATERA, *An “empirical interpolation” method: Application to efficient reduced-basis discretization of partial differential equations*, C. R. Math., 339 (2004), pp. 667–672.
- [5] H.-S. CHUNG AND J. J. ALONSO, *Using Gradients to Construct Response Surface Models for High-Dimensional Design Optimization Problems*, AIAA paper 2001-0922, 2001.
- [6] H.-S. CHUNG AND J. J. ALONSO, *Using Gradients to Construct Cokriging Approximation Models for High-Dimensional Design Optimization Problems*, AIAA paper 2002-0317, 2002.
- [7] A. COHEN, M. A. DAVENPORT, AND D. LEVIATAN, *On the stability and accuracy of least squares approximations*, Found. Comput. Math., 13 (2013), pp. 819–834.
- [8] P. COURTIER, E. ANDERSSON, W. HECKLEY, D. VASILJEVIC, M. HAMRUD, A. HOLLINGSWORTH, F. RABIER, M. FISHER, AND J. PAILLEUX, *The ECMWF implementation of three-dimensional variational assimilation (3D-Var). I: Formulation*, Quart. J. Roy. Meteor. Soc., 124 (1998), pp. 1783–1807.
- [9] R. P. DWIGHT AND Z.-H. HAN, *Efficient Uncertainty Quantification Using Gradient-Enhanced Kriging*, AIAA paper 2009-2276, 2009.
- [10] M. S. ELDERED, C. G. WEBSTER, AND P. G. CONSTANTINE, *Design Under Uncertainty Employing Stochastic Expansion Methods*, AIAA paper 2008-6001, 2008.
- [11] R. EVERSON AND L. SIROVICH, *Karhunen–Loeve procedure for gappy data*, Opt. Soc. Am. A, 12 (1995), pp. 1657–1664.
- [12] G. FISHMAN, *Monte Carlo: Concepts, Algorithms, and Applications*, Springer, New York, 1996.
- [13] R. GHANEM AND J. RED-HORSE, *Propagation of probabilistic uncertainty in complex physical systems using a stochastic finite element approach*, Phys. D, 133 (1999), pp. 137–144.
- [14] A. A. GIUNTA, M. S. ELDERED, AND J. P. CASTRO, *Uncertainty quantification using response surface approximations*, in Proceedings of the 9th ASCE Joint Specialty Conference on Probabilistic Mechanics and Structural Reliability, Albuquerque, NM., 2004.
- [15] G. GOLUB AND C. F. VAN LOAN, *Matrix Computations*, Johns Hopkins University Press, Baltimore, MD, 996.
- [16] J. M. HAMMERSLEY, D. C. HANDSCOMB, AND G. WEISS, *Monte Carlo methods*, Phys. Today, 18 (1965), p. 55.
- [17] I. HARARI AND T. J. R. HUGHES, *Galerkin/least-squares finite element methods for the reduced wave equation with non-reflecting boundary conditions in unbounded domains*, Comput. Methods Appl. Mech. Engrg., 98 (1992), pp. 411–454.
- [18] J. JAKEMAN, M. ELDERED, AND D. XIU, *Numerical approach for quantification of epistemic uncertainty*, J. Comput. Phys., doi:10.1016/j.jcp.2010.03.003.
- [19] R. E. KALMAN, *A new approach to linear filtering and prediction problems*, Trans. ASME J. Basic Eng., 82 (1960), pp. 35–45.
- [20] J. LAURENCEAU AND P. SAGAUT, *Building efficient response surfaces of aerodynamic functions with Kriging and cokriging*, AIAA Journal, 46 (2008), pp. 498–507.
- [21] J. M. LEWIS, S. LAKSHMIVARAHAN, AND S. DHALL, *Dynamic Data Assimilation: A Least Squares Approach*, Cambridge University Press, Cambridge, 2006.
- [22] Z. LI AND I. M. NAVON, *Optimality of variational data assimilation and its relationship with the Kalman filter and smoother*, Quart. J. Roy. Meteor. Soc., 127 (2001), pp. 661–683.
- [23] W. LIU AND S. M. BATILL, *Gradient-Enhanced Response Surface Approximations Using Kriging Models*, AIAA paper 2002-5456, 2002.
- [24] A. C. LORENC, *A global three-dimensional multivariate statistical interpolation scheme*, Monthly Weather Rev., 109 (1981), pp. 701–721.
- [25] Y. MADAY, N. C. NGUYEN, A. T. PATERA, AND G. S. H. PAU, *A general multipurpose interpolation procedure: The magic points*, Comm. Pure Appl. Anal., 8 (2008), pp. 383–404.
- [26] N. C. NGUYEN, A. T. PATERA, AND J. PERAIRE, *A “best points” interpolation method for efficient approximation of parametrized functions*, Internat. J. Numer. Methods Engrg., 73 (2008), pp. 521–543.
- [27] N. C. NGUYEN AND J. PERAIRE, *Gaussian functional regression for linear partial differential equations*, Comput. Methods Appl. Mech. Engrg., 287 (2015), pp. 69–89.

- [28] F. NOBILE, R. TEMPONE, AND C. G. WEBSTER, *A sparse grid stochastic collocation method for partial differential equations with random input data*, SIAM J. Numer. Anal., 46 (2008), pp. 2309–2345.
- [29] A. T. PATERA AND E. M. RONQUIST, *Regression on parametric manifolds: Estimation of spatial fields, functional outputs, and parameters from noisy data*, C. R. Math., 350 (2012), pp. 543–547.
- [30] F. PUKELSHEIM, *Optimal Design of Experiments*, Wiley & Sons, New York, 1993.
- [31] C. E. RASMUSSEN AND C. K. I. WILLIAMS, *Gaussian Processes for Machine Learning*, MIT Press, Cambridge, MA, 2006.
- [32] J. SACKS, W. J. WELCH, T. J. MITCHELL, AND H. P. WYNN, *Design and analysis of computer experiments*, Statist. Sci., 4 (1989), pp. 409–423.
- [33] S. SÄRKKÄ, *Bayesian Filtering and Smoothing*, Cambridge University Press, Cambridge, MA, 2013.
- [34] T. SONDERGAARD AND P. F. J. LERMUSIAUX, *Data assimilation with Gaussian mixture models using the dynamically orthogonal field equations. Part I: Theory and scheme*, Monthly Weather Rev., 141 (2013), pp. 1737–1760.
- [35] K. WILLCOX, *Unsteady flow sensing and estimation via the gappy proper orthogonal decomposition*, Comput. & Fluids, 35 (2006), pp. 208–226.
- [36] J. WOLBERG, *Data Analysis Using the Method of Least Squares: Extracting the Most Information from Experiments*, Springer, New York, 2005.
- [37] D. XIU AND J. A. HESTHAVEN, *High-order collocation methods for differential equations with random inputs*, SIAM J. Sci. Comput., 27 (2005), pp. 1118–1139.
- [38] D. XIU AND G. E. KARNIAKAKIS, *Modeling uncertainty in flow simulations via generalized polynomial chaos*, J. Comput. Phys., 187 (2003), pp. 137–167.
- [39] W. YAMAZAKI, M. P. RUMPFKEIL, AND D. J. MAVRIPLIS, *Design Optimization Utilizing Gradient/Hessian Enhanced Surrogate Model*, AIAA paper 2010-4363, 2010.
- [40] M. YANO, J. D. PENN, AND A. T. PATERA, *A model-data weak formulation for simultaneous estimation of state and model bias*, C. R. Math., 351 (Dec. 2013), pp. 937–941.