

## Sparse deterministic approximation of Bayesian inverse problems

This article has been downloaded from IOPscience. Please scroll down to see the full text article.

2012 Inverse Problems 28 045003

(<http://iopscience.iop.org/0266-5611/28/4/045003>)

View [the table of contents for this issue](#), or go to the [journal homepage](#) for more

Download details:

IP Address: 86.160.91.141

The article was downloaded on 14/03/2012 at 08:09

Please note that [terms and conditions apply](#).

# Sparse deterministic approximation of Bayesian inverse problems

C Schwab<sup>1</sup> and A M Stuart<sup>2</sup>

<sup>1</sup> Seminar for Applied Mathematics, ETH, 8092 Zurich, Switzerland

<sup>2</sup> Mathematics Institute, University of Warwick, Coventry CV4 7AL, UK

E-mail: [christoph.schwab@sam.math.ethz.ch](mailto:christoph.schwab@sam.math.ethz.ch) and [a.m.Stuart@warwick.ac.uk](mailto:a.m.Stuart@warwick.ac.uk)

Received 23 March 2011, in final form 8 February 2012

Published 13 March 2012

Online at [stacks.iop.org/IP/28/045003](http://stacks.iop.org/IP/28/045003)

## Abstract

We present a parametric deterministic formulation of Bayesian inverse problems with an input parameter from infinite-dimensional, separable Banach spaces. In this formulation, the forward problems are parametric, deterministic elliptic partial differential equations, and the inverse problem is to determine the unknown, parametric deterministic coefficients from noisy observations comprising linear functionals of the solution. We prove a generalized polynomial chaos representation of the posterior density with respect to the prior measure, given noisy observational data. We analyze the sparsity of the posterior density in terms of the summability of the input data's coefficient sequence. The first step in this process is to estimate the fluctuations in the prior. We exhibit sufficient conditions on the prior model in order for approximations of the posterior density to converge at a given algebraic rate, in terms of the number  $N$  of unknowns appearing in the parametric representation of the prior measure. Similar sparsity and approximation results are also exhibited for the solution and covariance of the elliptic partial differential equation under the posterior. These results then form the basis for efficient uncertainty quantification, in the presence of data with noise.

## 1. Introduction

Quantification of the uncertainty in predictions made by physical models, resulting from uncertainty in the input parameters to those models, is of increasing importance in many areas of science and engineering. Considerable effort has been devoted to developing numerical methods for this task. The most straightforward approach is to sample the uncertain system responses by Monte Carlo (MC) simulations. These have the advantage of being conceptually straightforward, but are constrained in terms of efficiency by their  $N^{-\frac{1}{2}}$  rate of convergence ( $N$  is the number of samples). In the 1980s, the engineering community started to develop new approaches to the problem via parametric representation of the probability space for

the input parameters [23, 24] based on the pioneering ideas of Wiener [27]. The use of sparse spectral approximation techniques [26, 22] opens the avenue toward algorithms for computational quantification of uncertainty, which beat the asymptotic complexity of MC methods, as measured by the computational cost per unit error in predicted uncertainty.

Much of the work in this area has been confined to the use of probability models on the input parameters, which are very simple, albeit leading to high-dimensional parametric representations. Typically, the randomness is described by a (possibly countably infinite) set of independent random variables representing uncertain coefficients in parametric expansions of input data, typically with the known closed-form Lebesgue densities. In many applications, such uncertainty in parameters is compensated for by (possibly noisy) observations, leading to an inverse problem. One approach to such inverse problems is via the techniques of optimal control [2]; however, this does not lead naturally to quantification of uncertainty. A Bayesian approach to the inverse problem [14, 25] allows the observations to map a possibly simple prior probability distribution on the input parameters into a posterior distribution. This posterior distribution is typically much more complicated than the prior, involving many correlations and without a useable closed form. The posterior distribution completely quantifies the uncertainty in the system's response, under given prior and structural assumptions on the system and given observational data. It allows, in particular, the Bayesian statistical estimation of unknown system parameters and responses by integration with respect to the posterior measure, which is of interest in many applications.

MC Markov chain (MCMC) methods can be used to probe this posterior probability distribution. This allows for computation of estimates of uncertain system responses conditioned on given observation data by means of approximate integration. However, these methods suffer from the same limits on computational complexity as straightforward MC methods. It is hence of interest to investigate whether sparse approximation techniques can be used to approximate the posterior density and conditional expectations given the data. In this paper, we study this question in the context of a model elliptic inverse problem. Elliptic problems with random coefficients have provided an important class of model problems for the uncertainty quantification community, see, e.g., [4, 22] and the references therein. In the context of inverse problems and noisy observational data, the corresponding elliptic problem arises naturally in the study of groundwater flow (see [19]) where hydrologists wish to determine the transmissivity (diffusion coefficient) from the head (solution of the elliptic PDE). The elliptic inverse problem hence provides natural model problem within which to study sparse representations of the posterior distribution.

In section 2, we recall the Bayesian setting for inverse problems from [25], stating and proving an infinite-dimensional Bayes' rule adapted to our inverse problem setting in theorem 2.1. Section 3 formulates the forward and inverse elliptic problem of interest, culminating in an application of Bayes' rule in theorem 3.4. The prior model is built on the work in [3, 6] in which the diffusion coefficient is represented parametrically via an infinite sum of functions, each with an independent uniformly distributed and compactly supported random variable as coefficient. Once we have shown that the posterior measure is well defined and absolutely continuous with respect to the prior, we proceed to study the analytic dependence of the posterior density in section 4, culminating in theorems 4.2 and 4.8. In section 5, we show how this parametric representation, and analyticity, may be employed to develop sparse polynomial chaos representations of the posterior density, and the key theorem 5.9 summarizes the achievable rates of convergence. In section 6, we study a variety of practical issues that arise in attempting to exploit the sparse polynomial representations as realizable algorithms for the evaluation of (posterior) expectations. Section 7 contains our concluding remarks,

and in particular, a discussion of the computational complexity of the new methodology, in comparison with that for MC-based methods.

Throughout, we concentrate on the posterior density itself. However, we also provide analysis related to the analyticity (and hence sparse polynomial representation) of various functions of the unknown input, in particular, the solution to the forward elliptic problem, and tensor products of this function. For the above class of elliptic model problems, we prove that for given data, there exist sparse,  $N$ -term ‘generalized polynomial chaos’ (gpc) approximations of this expectation with respect to the posterior (which is written as a density reweighted expectation with respect to the prior), which converge at the same rates afforded by the best  $N$ -term gpc approximations of the system response to uncertain, parametric inputs. Moreover, our analysis implies that the set  $\Lambda_N$  of the  $N$  ‘active’ gpc coefficients is identical to the set  $\Lambda_N$  of indices of a best  $N$ -term approximation of the system’s response. It was shown in [6, 7] that these rates are, in turn, completely determined by the decay rates of the input’s fluctuation expansions. We thus show that the machinery developed to describe gpc approximations of uncertain system response may be employed to study the more involved Bayesian inverse problem where the uncertainty is conditioned on observational data. Numerical algorithms, which achieve the optimal complexity implied by the sparse approximations, and numerical results demonstrating this will be given in our forthcoming work [1].

## 2. Bayesian inverse problems

Let  $G : X \rightarrow R$  denote a ‘forward’ map from some separable Banach space  $X$  of unknown parameters into another separable Banach space  $R$  of responses. We equip  $X$  and  $R$  with norms  $\|\cdot\|_X$  and  $\|\cdot\|_R$ , respectively. In addition, we are given  $\mathcal{O}(\cdot) : R \rightarrow \mathbb{R}^K$  denoting a bounded linear observation operator on the space  $R$  of system responses, which belong to the dual space  $R^*$  of the space  $R$  of system responses. We assume that the data are finite so that  $K < \infty$ , and equip  $\mathbb{R}^K$  with the Euclidean norm, denoted by  $|\cdot|$ .

We wish to determine the unknown data  $u \in X$  from the noisy observations

$$\delta = \mathcal{O}(G(u)) + \eta, \quad (1)$$

where  $\eta \in \mathbb{R}^K$  represents the noise. We assume that realization of the noise process is not known to us, but that it is a draw from the Gaussian measure  $\mathcal{N}(0, \Gamma)$ , for some positive (known) covariance operator  $\Gamma$  on  $\mathbb{R}^K$ . If we define  $\mathcal{G} : X \rightarrow \mathbb{R}^K$  by  $\mathcal{G} = \mathcal{O} \circ G$ , then we may write the equation for the observations as

$$\delta = \mathcal{G}(u) + \eta. \quad (2)$$

We define the least-squares functional (also referred to as ‘potential’ in what follows)  $\Phi : X \times \mathbb{R}^K \rightarrow \mathbb{R}$  by

$$\Phi(u; \delta) = \frac{1}{2} |\delta - \mathcal{G}(u)|_{\Gamma}^2, \quad (3)$$

where  $|\cdot|_{\Gamma} = |\Gamma^{-\frac{1}{2}} \cdot|$  so that

$$\Phi(u; \delta) = \frac{1}{2} ((\delta - \mathcal{G}(u))^{\top} \Gamma^{-1} (\delta - \mathcal{G}(u))).$$

In [25], it is shown that, under appropriate conditions on the forward and observation model  $\mathcal{G}$  and the prior measure on  $u$ , the posterior distribution on  $u$  is absolutely continuous with respect to the prior with the Radon–Nikodym derivative given by an infinite-dimensional version of Bayes’ rule. Posterior uncertainty is then determined by integration of suitably chosen functions against this posterior. At the heart of the deterministic approach proposed and analyzed here lies the *reformulation of the forward problem with stochastic input data as an infinite-dimensional, parametric deterministic problem*. We are thus interested in expressing

the posterior distribution in terms of a parametric representation of the unknown coefficient function  $u$ . To this end, we assume that, under the prior distribution, this function admits a *parametric representation* of the form

$$u = \bar{a} + \sum_{j \in \mathbb{J}} y_j \psi_j, \quad (4)$$

where  $y = \{y_j\}_{j \in \mathbb{J}}$  is an i.i.d sequence of real-valued random variables  $y_j \sim \mathcal{U}(-1, 1)$  and  $\bar{a}$ , and  $\psi_j$  are the elements of  $X$ . Here and throughout,  $\mathbb{J}$  denotes a finite or countably infinite index set, i.e. either  $\mathbb{J} = \{1, 2, \dots, J\}$  or  $\mathbb{J} = \mathbb{N}$ . All assertions proved in this paper hold in either case, and all bounds are in particular independent of the number  $J$  of parameters.

To derive the parametric expression of the prior measure  $\mu_0$  on  $y$ , we denote by

$$U = (-1, 1)^{\mathbb{J}}$$

the space of all sequences  $(y_j)_{j \in \mathbb{J}}$  of real numbers  $y_j \in (-1, 1)$ . Denoting the sub  $\sigma$ -algebra of Borel subsets on  $\mathbb{R}$ , which are also subsets of  $(-1, 1)$  by  $\mathcal{B}^1(-1, 1)$ , the pair

$$(U, \mathcal{B}) = \left( (-1, 1)^{\mathbb{J}}, \bigotimes_{j \in \mathbb{J}} \mathcal{B}^1(-1, 1) \right) \quad (5)$$

is a measurable space. We equip  $(U, \mathcal{B})$  with the uniform probability measure

$$\mu_0(\mathrm{d}y) := \bigotimes_{j \in \mathbb{J}} \frac{\mathrm{d}y_j}{2}, \quad (6)$$

which corresponds to bounded intervals for the possibly countably many uncertain parameters. Since the countable product of probability measures is again a probability measure,  $(U, \mathcal{B}, \mu_0)$  is a probability space. We assume in what follows that the prior measure on the uncertain input data, parametrized in the form (4), is  $\mu_0(\mathrm{d}y)$ . We add in passing that unbounded parameter ranges as arise, e.g., in lognormal random diffusion coefficients in models for subsurface flow [19], can be treated by the techniques developed here, at the expense of additional technicalities. We refer to [1] for details as well as for numerical experiments.

Define  $\Xi : U \rightarrow \mathbb{R}^K$  by

$$\Xi(y) = \mathcal{G}(u) \Big|_{u=\bar{a}+\sum_{j \in \mathbb{J}} y_j \psi_j}. \quad (7)$$

In the following, we view  $U$  as a bounded subset in  $\ell^\infty(\mathbb{J})$ , the Banach space of bounded sequences, and thereby introduce a notion of continuity in  $U$ .

**Theorem 2.1.** *Assume that  $\Xi : \bar{U} \rightarrow \mathbb{R}^K$  is bounded and continuous. Then,  $\mu^\delta(\mathrm{d}y)$ , the distribution of  $y$  given  $\delta$ , is absolutely continuous with respect to  $\mu_0(\mathrm{d}y)$ . Furthermore, if*

$$\Theta(y) = \exp(-\Phi(u; \delta)) \Big|_{u=\bar{a}+\sum_{j \in \mathbb{J}} y_j \psi_j}, \quad (8)$$

then

$$\frac{\mathrm{d}\mu^\delta}{\mathrm{d}\mu_0}(y) = \frac{1}{Z} \Theta(y), \quad (9)$$

where

$$Z = \int_U \Theta(y) \mu_0(\mathrm{d}y). \quad (10)$$

**Proof.** Let  $\nu_0$  denote the probability measure on  $U \times \mathbb{R}^K$  defined by  $\mu_0(dy) \otimes \pi(d\delta)$ , where  $\pi$  is the Gaussian measure  $\mathcal{N}(0, \Gamma)$ . Now, define a second probability measure  $\nu$  on  $U \times \mathbb{R}^K$  as follows. First, we specify the distribution of  $\delta$  given  $y$  to be  $\mathcal{N}(\Xi(y), \Gamma)$ . Since  $\Xi(y) : \bar{U} \rightarrow \mathbb{R}^K$  is continuous and  $\mu_0(U) = 1$ , we deduce that  $\Xi$  is  $\mu_0$  measurable. Hence, we may complete the definition of  $\nu$  by specifying that  $y$  is distributed according to  $\mu_0$ . By construction, and ignoring the constant of proportionality, which depends only on  $\delta$ ,<sup>3</sup>

$$\frac{d\nu}{d\nu_0}(y, \delta) \propto \Theta(y).$$

From the boundedness of  $\Xi$  on  $\bar{U}$ , we deduce that  $\Theta$  is bounded from below on  $\bar{U}$  by  $\theta_0 > 0$  and hence that

$$Z \geq \int_U \theta_0 \mu_0(dy) = \theta_0 > 0$$

since  $\mu_0(U) = 1$ . Noting that, under  $\nu_0$ ,  $y$  and  $\delta$  are independent, lemma 5.3 in [12] gives the desired result.  $\square$

We assume that we wish to compute the expectation of a function  $\phi : X \rightarrow S$ , for some Banach space  $S$ . With  $\phi$ , we associate the parametric mapping

$$\Psi(y) = \exp(-\Phi(u; \delta))\phi(u) \Big|_{u=\bar{a}+\sum_{j \in \mathbb{J}} y_j \psi_j} : U \rightarrow S. \quad (11)$$

From  $\Psi$ , we define

$$Z' = \int_U \Psi(y) \mu_0(dy) \in S \quad (12)$$

so that the expectation of interest is given by  $Z'/Z \in S$ . Thus, our aim is to approximate  $Z'$  and  $Z$ . Typical choices for  $\phi$  in applications might be  $\phi(u) = G(u)$ , the response of the system, or  $\phi(u) := (G(u))^{(m)} := \underbrace{G(u) \otimes \cdots \otimes G(u)}_{m \text{ times}} \in S = R^{(m)} := \underbrace{R \otimes \cdots \otimes R}_{m \text{ times}}$ . (13)

In particular, the choices  $\phi(u) = G(u)$  and  $\phi(u) = G(u) \otimes G(u)$  together facilitate computation of the mean and covariance of the response.

In the next sections, we will study the elliptic problem and deduce, from the known results concerning the parametric forward problem, the joint analyticity of the posterior density  $\Theta(y)$ , and also  $\Psi(y)$ , as a function of the parameter vector  $y \in U$ . From these results, we deduce sharp estimates on size of domain of analyticity of  $\Theta(y)$  (and  $\Psi(y)$ ) as a function of each coordinate  $y_j$ ,  $j \in \mathbb{N}$ . We concentrate on the concrete choice of  $\Psi$  defined by (13), and often the case  $p = 1$ . The analysis can be extended to other choices of  $\Psi$ .

### 3. Model parametric elliptic problem

#### 3.1. Function spaces

Our aim is to study the inverse problem of determining the diffusion coefficient  $u$  of an elliptic PDE from observation of a finite set of noisy linear functionals of the solution  $p$ , given  $u$ .

Let  $D$  be a bounded Lipschitz domain in  $\mathbb{R}^d$ ,  $d = 1, 2$  or  $3$ , with the Lipschitz boundary  $\partial D$ . Let further  $(H, (\cdot, \cdot), \|\cdot\|)$  denote the Hilbert space  $L^2(D)$ , which we will identify throughout with its dual space, i.e.  $H \simeq H^*$ .

We define also the space  $V$  of variational solutions of the forward problem: specifically, we let  $(V, (\nabla \cdot, \nabla \cdot), \|\cdot\|_V)$  denote the Hilbert space  $H_0^1(D)$  (everything that follows will hold

<sup>3</sup>  $\Theta(y)$  is also a function of  $\delta$  but we suppress this for economy of notation.

for rather general, elliptic problems with affine parameter dependence and ‘energy’ space  $V$ ). The dual space  $V^*$  of all continuous, linear functionals on  $V$  is isomorphic to the Banach space  $H^{-1}(D)$ , which we equip with the dual norm to  $V$ , denoted  $\|\cdot\|_{-1}$ . We shall assume for the (deterministic) data  $f \in V^*$ .

### 3.2. Forward problem

In the bounded Lipschitz domain  $D$ , we consider the following elliptic PDE:

$$-\nabla \cdot (u \nabla p) = f \quad \text{in } D, \quad p = 0 \quad \text{in } \partial D. \quad (14)$$

Given data  $u \in L^\infty(D)$ , a weak solution of (14) for any  $f \in V^*$  is a function  $p \in V$ , which satisfies

$$\int_D u(x) \nabla p(x) \cdot \nabla q(x) \, dx = \langle q, f \rangle_{V^*} \quad \text{for all } q \in V. \quad (15)$$

Here,  $\langle \cdot, \cdot \rangle_{V^*}$  denotes the dual pairing between elements of  $V$  and  $V^*$ .

For the well-posedness of the forward problem, we shall work under the following assumption.

**Assumption 3.1.** *There exist constants  $0 < a_{\text{MIN}} \leq a_{\text{MAX}} < \infty$  so that*

$$0 < a_{\text{MIN}} \leq u(x) \leq a_{\text{MAX}} < \infty, \quad x \in D. \quad (16)$$

Under assumption 3.1, the Lax–Milgram lemma ensures the existence and uniqueness of the response  $p$  of (15). Thus, in the notation of the previous section,  $R = V$  and  $G(u) = p$ . Moreover, this variational solution satisfies the *a priori* estimate

$$\|G(u)\|_V = \|p\|_V \leq \frac{\|f\|_{V^*}}{a_{\text{MIN}}}. \quad (17)$$

We assume that the observation function  $\mathcal{O} : V \rightarrow \mathbb{R}^K$  comprises  $K$  linear functionals  $o_k \in V^*$ ,  $k = 1, \dots, K$ . In the notation of the previous section, we denote by  $X = L^\infty(D)$  the Banach space in which the unknown input parameter  $u$  takes values. It follows that

$$|\mathcal{G}(u)| \leq \frac{\|f\|_{V^*}}{a_{\text{MIN}}} \left( \sum_{k=1}^K \|o_k\|_{V^*}^2 \right)^{\frac{1}{2}}. \quad (18)$$

### 3.3. Structural assumptions on diffusion coefficient

As discussed in section 2, we introduce a parametric representation of the random input parameter  $u$  via an affine representation with respect to  $y$ , which means that the parameters  $y_j$  are the coefficients of the function  $u$  in the formal series expansion

$$u(x, y) = \bar{a}(x) + \sum_{j \in \mathbb{J}} y_j \psi_j(x), \quad x \in D, \quad (19)$$

where  $\bar{a} \in L^\infty(D)$  and  $\{\psi_j\}_{j \in \mathbb{J}} \subset L^\infty(D)$ . We are interested in the effect of approximating the solutions input parameter  $u(x, y)$ , by truncation of the series expansion (19) in the case  $\mathbb{J} = \mathbb{N}$ , and on the corresponding effect on the forward (resp. observational) map  $G(u(\cdot))$  (resp.  $\mathcal{G}(u(\cdot))$ ) to the family of elliptic equations with the above input parameters. In the decomposition (19), we have the choice to either normalize the basis (e.g., assume they all have norm one in some space) or to normalize the parameters. It is more convenient for us to do the latter. This leads us to the following assumptions which shall be made

throughout.

- (i) For all  $j \in \mathbb{J} : \psi_j \in L^\infty(D)$  and  $\psi_j(x)$  is defined for all  $x \in D$ .
- (ii)

$$y = (y_1, y_2, \dots) \in U = [-1, 1]^{\mathbb{J}}. \tag{20}$$

i.e. the parameter vector  $y$  in (19) belongs to the unit ball of the sequence space  $\ell^\infty(\mathbb{J})$ .

- (iii) For each  $u(x, y)$  to be considered, (19) holds for every  $x \in D$  and every  $y \in U$ .

We will, occasionally, use (19) with  $\mathbb{J} \subset \mathbb{N}$ , as well as with  $\mathbb{J} = \mathbb{N}$  (in the latter case the additional assumption 3.2 has to be imposed). In either case, we will work throughout under the assumption that the ellipticity condition (16) holds uniformly for  $y \in U$ .

*Uniform ellipticity assumption: There exist  $0 < a_{\text{MIN}} \leq a_{\text{MAX}} < \infty$  such that for all  $x \in D$  and for all  $y \in U$*

$$0 < a_{\text{MIN}} \leq u(x, y) \leq a_{\text{MAX}} < \infty. \tag{21}$$

We refer to assumption (21) as  $\text{UEA}(a_{\text{MIN}}, a_{\text{MAX}})$  in the following. In particular,  $\text{UEA}(a_{\text{MIN}}, a_{\text{MAX}})$  implies  $a_{\text{MIN}} \leq \bar{a}(x) \leq a_{\text{MAX}}$  for all  $x \in D$ , since we can choose  $y_j = 0$  for all  $j \in \mathbb{N}$ . Also, observe that the validity of the lower and upper inequalities in (21) for all  $y \in U$  are respectively equivalent to the conditions that

$$\sum_{j \in \mathbb{J}} |\psi_j(x)| \leq \bar{a}(x) - a_{\text{MIN}}, \quad x \in D, \tag{22}$$

and

$$\sum_{j \in \mathbb{J}} |\psi_j(x)| \leq a_{\text{MAX}} - \bar{a}(x), \quad x \in D. \tag{23}$$

We shall require in what follows a quantitative control of the relative size of the fluctuations in representation (19). To this end, we shall impose the following assumption.

**Assumption 3.2.** *The functions  $\bar{a}$  and  $\psi_j$  in (19) satisfy*

$$\sum_{j \in \mathbb{J}} \|\psi_j\|_{L^\infty(D)} \leq \frac{\kappa}{1 + \kappa} \bar{a}_{\text{MIN}},$$

with  $\bar{a}_{\text{MIN}} = \min_{x \in D} \bar{a}(x) > 0$  and  $\kappa > 0$ .

Assumption 3.1 is then satisfied by choosing

$$a_{\text{MIN}} := \bar{a}_{\text{MIN}} - \frac{\kappa}{1 + \kappa} \bar{a}_{\text{MIN}} = \frac{1}{1 + \kappa} \bar{a}_{\text{MIN}}. \tag{24}$$

### 3.4. Inverse problem

We start by proving that the forward maps  $G : X \rightarrow V$  and  $\mathcal{G} : X \rightarrow \mathbb{R}^K$  are Lipschitz.

**Lemma 3.3.** *If  $p$  and  $\tilde{p}$  are the solutions of (15) with the same right-hand side  $f$  and with coefficients  $u$  and  $\tilde{u}$ , respectively, and if these coefficients both satisfy assumption 3.1, then the forward solution map  $u \rightarrow p = G(u)$  is Lipschitz as a mapping from  $X$  into  $V$  with the Lipschitz constant defined by*

$$\|p - \tilde{p}\|_V \leq \frac{\|f\|_{V^*}}{a_{\text{MIN}}^2} \|u - \tilde{u}\|_{L^\infty(D)}. \tag{25}$$



Moreover, the forward solution map can be composed with the observation operator to prove that the map  $u \rightarrow \mathcal{G}(u)$  is Lipschitz as a mapping from  $X$  into  $\mathbb{R}^K$  with the Lipschitz constant defined by

$$|\mathcal{G}(u) - \mathcal{G}(\tilde{u})| \leq \frac{\|f\|_{V^*}}{a_{\text{MIN}}^2} \left( \sum_{k=1}^K \|o_k\|_{V^*}^2 \right)^{\frac{1}{2}} \|u - \tilde{u}\|_{L^\infty(D)}. \tag{26}$$

**Proof.** Subtracting the variational formulations for  $p$  and  $\tilde{p}$ , we find that for all  $q \in V$ ,

$$0 = \int_D u \nabla p \cdot \nabla q \, dx - \int_D \tilde{u} \nabla \tilde{p} \cdot \nabla q \, dx = \int_D u (\nabla p - \nabla \tilde{p}) \cdot \nabla q \, dx + \int_D (u - \tilde{u}) \nabla \tilde{p} \cdot \nabla q \, dx.$$

Therefore,  $w = p - \tilde{p}$  is the solution of  $\int_D u \nabla w \cdot \nabla q = L(q)$ , where  $L(v) := \int_D (\tilde{u} - u) \nabla \tilde{p} \cdot \nabla v$ . Hence,

$$\|w\|_V \leq \frac{\|L\|_{V^*}}{a_{\text{MIN}}},$$

and we obtain (25) since it follows from (17) that

$$\|L\|_{V^*} = \max_{\|v\|_V=1} |L(v)| \leq \|u - \tilde{u}\|_{L^\infty(D)} \|\tilde{p}\|_V \leq \|u - \tilde{u}\|_{L^\infty(D)} \frac{\|f\|_{V^*}}{a_{\text{MIN}}}.$$

The Lipschitz continuity of  $\mathcal{G} = \mathcal{O} \circ G : X \rightarrow \mathbb{R}^K$  is immediate since  $\mathcal{O}$  comprises the  $K$  linear functionals  $o_k$ . Thus, (25) implies (26).  $\square$

The next result may be deduced in a straightforward fashion from the preceding analysis.

**Theorem 3.4.** Under the UEA( $a_{\text{MIN}}, a_{\text{MAX}}$ ) and assumption 3.2, it follows that the posterior measure  $\mu^\delta(\text{d}y)$  on  $y$  given  $\delta$  is absolutely continuous with respect to the prior measure  $\mu_0(\text{d}y)$  with the Radon–Nikodym derivative given by (8) and (9).

**Proof.** This is a straightforward consequence of theorem 2.1 provided that we show boundedness and continuity of  $\Xi : \tilde{U} \rightarrow \mathbb{R}^K$  given by (7). Boundedness follows from (18), together with the boundedness of  $\|o_k\|_{V^*}$ , under UEA( $a_{\text{MIN}}, a_{\text{MAX}}$ ). Let  $u$  and  $\tilde{u}$  denote two diffusion coefficients generated by two parametric sequences  $y$  and  $\tilde{y}$  in  $U$ . Then, by (26) and assumption 3.2,

$$\begin{aligned} |\Xi(y) - \Xi(\tilde{y})| &\leq \frac{\|f\|_{V^*}}{a_{\text{MIN}}^2} \left( \sum_{k=1}^K \|o_k\|_{V^*}^2 \right)^{\frac{1}{2}} \|u - \tilde{u}\|_{L^\infty(D)} \\ &\leq \frac{\|f\|_{V^*}}{a_{\text{MIN}}^2} \left( \sum_{k=1}^K \|o_k\|_{V^*}^2 \right)^{\frac{1}{2}} \frac{\kappa}{1 + \kappa} \bar{a}_{\text{MIN}} \|y - \tilde{y}\|_{\ell^\infty(\mathbb{J})}. \end{aligned}$$

The result follows.  $\square$

#### 4. Complex extension of the elliptic problem

As indicated above, one main technical objective will consist in proving analyticity of the posterior density  $\Theta(y)$  with respect to the (possibly countably many) parameters  $y \in U$  in (19) defining the prior and to obtain bounds on the supremum of  $\Theta$  over the maximal domains in  $\mathbb{C}$  into which  $\Theta(y)$  can be continued analytically. Our key ingredients for getting such estimates rely on complex analysis.

It is well known that the existence theory for the forward problem (14) extends to the case where the coefficient function  $u(x)$  takes values in  $\mathbb{C}$ . In this case, the ellipticity assumption 3.1 should be replaced by the assumption that

$$0 < a_{\text{MIN}} \leq \Re(u(x)) \leq |u(x)| \leq a_{\text{MAX}} < \infty, \quad x \in D, \quad (27)$$

and all the above results remain valid with Sobolev spaces understood as spaces of complex-valued functions. Throughout what follows, we shall frequently pass to spaces of complex-valued functions, without distinguishing these notationally. It will always be clear from the context which coefficient field is implied.

#### 4.1. Notation and assumptions

We extend the definition of  $u(x, y)$  to  $u(x, z)$  for the complex variable  $z = (z_j)_{j \in \mathbb{J}}$  (by using  $z_j$  instead of  $y_j$  in the definition of  $u$  by (19)), where each  $z_j$  has modulus less than or equal to 1. Therefore,  $z$  belongs to the polydisc

$$\mathcal{U} := \bigotimes_{j \in \mathbb{J}} \{z_j \in \mathbb{C} : |z_j| \leq 1\} \subset \mathbb{C}^{\mathbb{J}}. \quad (28)$$

Note that  $\bar{\mathcal{U}} \subset \mathcal{U}$ . Using (22) and (23), when the functions  $\bar{a}$  and  $\psi_j$  are real valued, the condition  $\text{UEA}(a_{\text{MIN}}, a_{\text{MAX}})$  implies that for all  $x \in D$  and  $z \in \mathcal{U}$ ,

$$0 < a_{\text{MIN}} \leq \Re(u(x, z)) \leq |u(x, z)| \leq 2a_{\text{MAX}}, \quad (29)$$

and therefore, the corresponding solution  $p(z)$  is well defined in  $V$  for all  $z \in \mathcal{U}$  by the Lax–Milgram theorem for sesquilinear forms. More generally, we may consider an expansion of the form

$$u(x, z) = \bar{a} + \sum_{j \in \mathbb{J}} z_j \psi_j,$$

where  $\bar{a}$  and  $\psi_j$  are complex-valued functions and replace  $\text{UEA}(a_{\text{MIN}}, a_{\text{MAX}})$  by the following, complex-valued counterpart.

*Uniform ellipticity assumption in  $\mathbb{C}$  :* There exist  $0 < a_{\text{MIN}} \leq a_{\text{MAX}} < \infty$ , such that for all  $x \in D$  and all  $z \in \mathcal{U}$

$$0 < a_{\text{MIN}} \leq \Re(u(x, z)) \leq |u(x, z)| \leq a_{\text{MAX}} < \infty. \quad (30)$$

We refer to (30) as  $\text{UEAC}(a_{\text{MIN}}, a_{\text{MAX}})$ .

#### 4.2. Domains of holomorphy

The condition  $\text{UEAC}(a_{\text{MIN}}, a_{\text{MAX}})$  implies that the forward solution map  $z \mapsto p(z)$  is strongly holomorphic as a  $V$ -valued function, which is uniformly bounded in certain domains larger than  $\mathcal{U}$ . For  $0 < r \leq 2a_{\text{MAX}} < \infty$ , we define the open set

$$\mathcal{A}_r = \{z \in \mathbb{C}^{\mathbb{J}} : r < \Re(u(x, z)) \leq |u(x, z)| < 2a_{\text{MAX}} \text{ for every } x \in D\} \subset \mathbb{C}^{\mathbb{J}}. \quad (31)$$

Under  $\text{UEAC}(a_{\text{MIN}}, a_{\text{MAX}})$ , for every  $0 < r < a_{\text{MIN}}$ ,  $\mathcal{U} \subset \mathcal{A}_r$  holds.

According to the Lax–Milgram theorem, for every  $z \in \mathcal{A}_r$ , there exists a unique solution  $p(z) \in V$  of the variational problem: given  $f \in V^*$ , for every  $z \in \mathcal{A}_r$ , find  $p \in V$  such that

$$\alpha(z; p, q) = (f, q) \quad \forall q \in V. \quad (32)$$

Here, the sesquilinear form  $\alpha(z; \cdot, \cdot)$  is defined as

$$\alpha(z; p, q) = \int_D u(x, z) \nabla p \cdot \overline{\nabla q} \, dx \quad \forall p, q \in V. \quad (33)$$

We next show that the analytic continuation of the parametric solution  $p(y)$  to the domain  $\mathcal{A}_r$  is the unique solution  $p(z)$  of (32), which satisfies the *a priori* estimate

$$\sup_{z \in \mathcal{A}_r} \|p(z)\|_V \leq \frac{\|f\|_{V^*}}{r}. \tag{34}$$

The first step of our analysis is to establish strong holomorphy of the forward solution map  $z \mapsto p(z)$  in (32) with respect to the countably many variables  $z_j$  at any point  $z \in \mathcal{A}_r$ . This follows from the observation that the function  $p(z)$  is the solution to the operator equation  $A(z)p(z) = f$ , where the operator  $A(z) \in \mathcal{L}(V, V^*)$  depends in an affine manner on each variable  $z_j$ . To prepare the argument for proving holomorphy of the functionals  $\Phi$  and  $\Theta$  appearing in (8) and (11), we give a direct proof.

Using lemma 3.3, we have proved, by means of a difference quotient argument given in [7], lemma 4.1. Lemma 4.1, together with Hartogs' theorem (see, e.g., [13]) and the separability of  $V$ , implies strong holomorphy of  $p(z)$  as a  $V$ -valued function on  $\mathcal{A}_r$ , stated as theorem 4.2. The proof of this theorem can also be found in [7]; the result will also be obtained as a corollary of the analyticity results for the functionals  $\Psi$  and  $\Theta$  proved below.

**Lemma 4.1.** *At any  $z \in \mathcal{A}_r$ , the function  $z \mapsto p(z)$  admits a complex derivative  $\partial_{z_j} p(z) \in V$  with respect to each variable  $z_j$ . This derivative is the weak solution of the problem: given  $z \in \mathcal{A}_r$ , find  $\partial_{z_j} p(z) \in V$  such that*

$$\alpha(z; \partial_{z_j} p(z), q) = L_0(q) := - \int_D \psi_j \nabla p(z) \cdot \overline{\nabla q} \, dx, \quad \text{for all } q \in V. \tag{35}$$

**Theorem 4.2.** *Under UEAC( $a_{\text{MIN}}$ ,  $a_{\text{MAX}}$ ) for any  $0 < r < a_{\text{MIN}}$  the solution  $p(z) = G(u(z))$  of the parametric forward problem is holomorphic as a  $V$ -valued function in  $\mathcal{A}_r$  and the *a priori* estimate (34) holds.*

We remark that  $\mathcal{A}_r$  also contains certain *polydiscs*: for any sequence  $\rho := (\rho_j)_{j \geq 1}$  of positive radii we define the polydisc

$$\mathcal{U}_\rho = \bigotimes_{j \in \mathbb{J}} \{z_j \in \mathbb{C} : |z_j| \leq \rho_j\} = \{z_j \in \mathbb{C} : z = (z_j)_{j \in \mathbb{J}}; |z_j| \leq \rho_j\} \subset \mathbb{C}^{\mathbb{J}}. \tag{36}$$

We say that a sequence  $\rho = (\rho_j)_{j \geq 1}$  of radii is *r-admissible* if and only if for every  $x \in D$

$$\sum_{j \in \mathbb{J}} \rho_j |\psi_j(x)| \leq \Re(\bar{a}(x)) - r. \tag{37}$$

If the sequence  $\rho$  is *r-admissible*, then the polydisc  $\mathcal{U}_\rho$  is contained in  $\mathcal{A}_r$  since on the one hand for all  $z \in \mathcal{U}_\rho$  and for almost every  $x \in D$

$$\Re(\bar{u}(x, z)) \geq \Re(\bar{a}(x)) - \sum_{j \in \mathbb{J}} |z_j \psi_j(x)| \geq \Re(\bar{a}(x)) - \sum_{j \in \mathbb{J}} \rho_j |\psi_j(x)| \geq r,$$

and on the other hand, if for every  $x \in D$

$$|u(x, z)| \leq |\bar{a}(x)| + \sum_{j \in \mathbb{J}} |z_j \psi_j(x)| \leq |\bar{a}(x)| + \Re(\bar{a}(x)) - r \leq 2|\bar{a}(x)| \leq 2a_{\text{MAX}}.$$

Here, we used  $|\bar{a}(x)| \leq a_{\text{MAX}}$  that follows from UEAC( $a_{\text{MIN}}$ ,  $a_{\text{MAX}}$ ).

Similar to (22), the validity of the lower inequality in (30) for all  $z \in \mathcal{U}$  is equivalent to the condition that

$$\sum_{j \geq 1} |\psi_j(x)| \leq \Re(\bar{a}(x)) - a_{\text{MIN}}, \quad x \in D. \tag{38}$$

This shows that the constant sequence  $\rho_j = 1$  is *r-admissible* for all  $0 < r \leq a_{\text{MIN}}$ .

**Remark 4.3.** For  $0 < r < a_{\text{MIN}}$ , there exist *r-admissible* sequences such that  $\rho_j > 1$  for all  $j \geq 1$ , i.e. such that the polydisc  $\mathcal{U}_\rho$  is strictly larger than  $\mathcal{U}$  in every variable. This will be exploited systematically below in the derivation of approximation bounds.

### 4.3. Holomorphy of response functionals

We next show that, for given data  $\delta$ , the functionals  $\mathcal{G}(\cdot)$ ,  $\Phi(u(\cdot); \delta)$  and  $\Theta(\cdot)$  depend holomorphically on the parameter vector  $z \in \mathbb{C}^{\mathbb{J}}$ , on polydiscs  $\mathcal{U}_\rho$  as in (36) for suitable  $r$ -admissible sequences of semi-axes  $\rho$ . Our general strategy for proving this will be analogous to the argument for establishing analyticity of the map  $z \mapsto G(u(z))$  as a  $V$ -valued function.

We now extend theorem 4.2 from the solution of the elliptic PDE to the posterior density and related quantities required to define expectations under the posterior, culminating in theorem 4.8 and corollary 4.9. We achieve this through a sequence of lemmas that we now derive.

The following lemma is simply a complexification of (18) and (26). It implies bounds on  $\mathcal{G}$  and its Lipschitz constant in the covariance-weighted norm.

**Lemma 4.4.** *Under UEAC( $a_{\text{MIN}}$ ,  $a_{\text{MAX}}$ ), for every  $f \in V^* = H^{-1}(D)$  and for every  $\mathcal{O}(\cdot) \in (V^*)^* \simeq V \rightarrow Y = \mathbb{R}^K$  the following holds:*

$$|\mathcal{G}(u)| \leq \frac{\|f\|_{V^*}}{a_{\text{MIN}}} \left( \sum_{k=1}^K \|o_k\|_{V^*}^2 \right)^{\frac{1}{2}}, \quad (39)$$

$$|\mathcal{G}(u) - \mathcal{G}(u)| \leq \frac{\|f\|_{V^*}}{a_{\text{MIN}}^2} \|u_1 - u_2\|_{L^\infty(D)} \left( \sum_{k=1}^K \|o_k\|_{V^*}^2 \right)^{\frac{1}{2}}. \quad (40)$$

To be concrete, we concentrate in the next lemma on computing the expected value of the pressure  $p = G(u) \in V$  under the posterior measure. To this end we define  $\Psi$  with  $\psi$  as in (13) with  $m = 1$ . We start by considering the case of a single parameter.

**Lemma 4.5.** *Let  $\mathbb{J} = \{1\}$  and take  $\phi = G : U \rightarrow V$ . With  $u(x, y)$  as in (4), under UEAC( $a_{\text{MIN}}$ ,  $a_{\text{MAX}}$ ), the functions  $\Psi : [-1, 1] \rightarrow V$  and  $\Theta : [-1, 1] \rightarrow \mathbb{R}$  and the potential  $\Phi(u(x, \cdot); \delta)$  defined by (11), (8) and (3), respectively, may be extended to functions that are strongly holomorphic on the strip  $\{y + iz : |y| < r/\kappa\}$  for any  $r \in (\kappa, 1)$ .*

**Proof.** We view  $H, V$  and  $X = L^\infty(D)$  as Banach spaces over  $\mathbb{C}$ . We extend the equation (19) to complex coefficients  $u(x, z) = \text{Re}(\bar{a}(x) + z\psi(x)) = \bar{a}(x) + y\psi(x)$  since  $z = y + i\zeta$ . Note that  $\bar{a} + z\psi$  is holomorphic in  $z$  since it is linear. Since  $\text{Re}(\bar{a} + z\psi) = \bar{a} + y\psi \geq a_{\text{MIN}}$ , it follows that, for all  $\zeta = \text{Im}(z)$ ,

$$\text{Re} \int_D u(x) |\nabla p(x) - \nabla \tilde{p}(x)|^2 dx \geq a_{\text{MIN}} \|p - \tilde{p}\|_V^2.$$

We prove that the mappings  $\Psi$  and  $\Theta$  are holomorphic by studying the properties of  $G(\bar{a} + z\psi)$  and  $\Phi(\bar{a} + z\psi)$  as the functions of  $z \in \mathbb{C}$ . Let  $h \in \mathbb{C}$  with  $|h| < \epsilon \ll 1$ . We show that

$$\lim_{|h| \rightarrow 0} h^{-1} (p(z+h) - p(z))$$

exists in  $V$  (strong holomorphy). Note first that  $\partial_z u = \psi$ . Now consider  $p$ . We have

$$\frac{1}{h} (p(z+h) - p(z)) = \frac{1}{h} (G(\bar{a} + (z+h)\psi) - G(\bar{a} + z\psi)) =: r.$$

By lemma 3.3, we deduce that

$$\|r\|_V \leq \frac{\|f\|_{H^{-1}(D)}}{a_{\text{MIN}}^2} \|\psi\|_{L^\infty(D)}.$$

From this it follows that there is a weakly convergent subsequence in  $V$ , as  $|h| \rightarrow 0$ . We proceed to deduce the existence of a strong limit. To this end, we introduce the sesquilinear form

$$b(p, q) = \int_D u \nabla p \overline{\nabla q} \, dx.$$

Then,

$$b(G(u), q) = (f, q) \quad \forall q \in V.$$

For a coefficient function  $u$  as in (19), the form  $b(\cdot, \cdot)$  is equal to the parametric sesquilinear form  $\alpha(z; p, q)$  defined in (33).

Note that for  $z = \bar{a} + y\psi \in \mathbb{R}$  and for real-valued arguments  $p$  and  $q$ , the parametric sesquilinear form  $\alpha(z; p, q)$  coincides with the bilinear form in (15). Accordingly, for every  $z \in \mathbb{C}^{\mathbb{J}}$ , the unique holomorphic extension of the parametric solution  $G(u(\bar{a} + y\psi))$  to complex parameters  $z = y + i\zeta$  is the unique variational solution of the parametric problem

$$\alpha(z; G(\bar{a} + z\psi), q) = (f, q) \quad \forall q \in V. \quad (41)$$

Assumption UEAC( $a_{\text{MIN}}$ ,  $a_{\text{MAX}}$ ) is readily seen to imply

$$\forall p \in V : \quad \text{Re}(\alpha(z; p, p)) \geq a_{\text{MIN}} \|p\|_V^2.$$

If we choose  $\delta \in (\kappa, 1)$  and choose  $z = y + i\eta$ , we obtain, for all  $\zeta$  and for  $|y| \leq \delta/\kappa$ ,

$$\text{Re}(\alpha(z; p, p)) \geq \bar{a}_{\text{MIN}}(1 - \delta) \|p\|_V^2. \quad (42)$$

From (41), we see that for such values of  $z = y + i\zeta$

$$\begin{aligned} 0 &= \alpha(z; G(\bar{a} + z\psi), q) - \alpha(z; G(\bar{a} + (z+h)\psi), q) + \alpha(z; G(\bar{a} + (z+h)\psi), q) \\ &\quad - \alpha(z+h; G(\bar{a} + (z+h)\psi), q) \\ &= \alpha(z; G(\bar{a} + z\psi) - G(\bar{a} + (z+h)\psi), q) - \int_D h\psi \nabla G(\bar{a} + (z+h)\psi) \overline{\nabla q} \, dx. \end{aligned}$$

Dividing by  $h$ , we obtain that  $r$  satisfies, for all  $z = y + i\zeta$  with  $|y| \leq \delta/\kappa$  and every  $\zeta \in \mathbb{R}$

$$\forall q \in V : \quad \alpha(z; r, q) + \int_D \psi \nabla G(\bar{a} + (z+h)\psi) \overline{\nabla q} \, dx = 0. \quad (43)$$

The second term we denote by  $s(h)$  and note that, by lemma 3.3,

$$|s(h_1) - s(h_2)| \leq \frac{1}{a_{\text{MIN}}^2} \|\psi\|_{\infty}^2 \|f\|_1 \|q\|_V |h_1 - h_2|.$$

If we denote the solution  $r$  to equation (43) by  $r_h(\bar{a}; z)$ , then we deduce from the Lipschitz continuity of  $s(\cdot)$  that  $r_h(\bar{a}; z) \rightarrow r_0(\bar{a}; z)$ , where

$$\alpha(z; r_0, q) = s(0) \quad \forall q \in V.$$

Hence,  $r_0 = \partial_z G(\bar{a} + z\psi) \in V$  and we deduce that  $G : [-1, 1] \rightarrow V$  can be extended to a complex-valued function that is strongly holomorphic on the strip  $\{y + i\zeta : |y| < \delta/\kappa, \zeta \in \mathbb{R}\}$ .

We next study the domain of holomorphy of the analytic continuation of the potential  $\Phi(\bar{a} + z\psi; d)$  to parameters  $z \in \mathbb{C}$ . It suffices to consider  $K = 1$  noting that then the unique analytic continuation of the potential  $\Phi$  is given by

$$\Phi(\bar{a} + z\psi; \delta) = \frac{1}{2\gamma^2} (\delta - \mathcal{G}(\bar{a} + z\psi))^{\top} (\delta - \mathcal{G}(\bar{a} + z\psi)). \quad (44)$$

The function  $z \mapsto \mathcal{G}(\bar{a} + z\psi)$  is holomorphic with the same domain of holomorphy as  $G(\bar{a} + z\psi)$ . Similarly, it follows that the function

$$z \mapsto (\delta - \mathcal{G}(\bar{a} + z\psi))^{\top} (\delta - \mathcal{G}(\bar{a} + z\psi))$$

is holomorphic, with the same domain of holomorphy; this is shown by composing the relevant power series expansion. From this we deduce that  $\Theta$  and  $\Psi$  are holomorphic, with the same domain of holomorphy.  $\square$

So far we have considered the case  $\mathbb{J} = \{1\}$ . We will now generalize the case. To this end, we pick an arbitrary  $m \in \mathbb{J}$  and write  $y = (y^*, y_m)$  and  $z = (z^*, z_m)$ .

**Assumption 4.6.** *There are constants  $0 < \bar{a}_{\text{MIN}} \leq \bar{a}_{\text{MAX}} < \infty$  and  $\kappa \in (0, 1)$  such that*

$$0 < \bar{a}_{\text{MIN}} \leq \bar{a} \leq \bar{a}_{\text{MAX}} < \infty, \quad \text{a.e. } x \in D, \quad \|\|\psi_j\|_{L^\infty(D)}\|_{\ell^1(\mathbb{J})} < \kappa \bar{a}_{\text{MIN}}. \quad (45)$$

For  $m \in \mathbb{J}$ , we write (19) in the form

$$u(x; y) = \bar{a}(x) + y_m \psi_m(x) + \sum_{j \in \mathbb{J} \setminus \{m\}} y_j \psi_j(x).$$

From assumption 4.6, we deduce that there are numbers  $\kappa_j \leq \kappa$  such that

$$\|\psi_j\|_{L^\infty} < \bar{a}_{\text{MIN}} \kappa_j.$$

Hence, we obtain, for every  $x \in D$  and every  $y \in U$ , the lower bound

$$\begin{aligned} u(x, y) &\geq \bar{a}_{\text{MIN}}(1 - (\kappa - \kappa_m) - \kappa_m) \\ &\geq \bar{a}_{\text{MIN}}(1 - (\kappa - \kappa_m)) \left(1 - \frac{\kappa_m}{1 - (\kappa - \kappa_m)}\right) \\ &\geq a'_{\text{MIN}}(1 - \kappa'_m) \end{aligned}$$

with  $a'_{\text{MIN}} = \bar{a}_{\text{MIN}}(1 - \kappa)$  and  $\kappa'_m = \kappa_m(1 - (\kappa - \kappa_m))^{-1} \in (0, 1)$ . With this observation, we obtain the following.

**Lemma 4.7.** *Let assumption 4.6 hold and set  $U = [-1, 1]^{\mathbb{J}}$  and  $\phi = G : U \rightarrow V$ . Then, the functions  $\Psi : U \rightarrow V$  and  $\Theta : U \rightarrow \mathbb{R}$ , as well as the potential  $\Phi(u(x, \cdot); \delta) : U \rightarrow \mathbb{R}$ , admit unique extensions to strongly holomorphic functions on the product of strips given by*

$$\mathcal{S}_\rho := \bigotimes_{j \in \mathbb{J}} \{y_j + iz_j : |y_j| < \delta_j / \kappa'_j, \quad z_j \in \mathbb{R}\} \quad (46)$$

for any sequence  $\rho = (\rho_j)_{j \in \mathbb{J}}$  with  $\rho_j \in (\kappa'_j, 1)$ .

**Proof.** Fixing  $y^*$ , we view  $\Psi$  and  $\Theta$  as the functions of the single parameter  $y_m$ . For each fixed  $y^*$ , we extend  $y_m$  to a complex variable  $z_m$ . The estimates preceding the statement of this lemma, together with lemma 4.5, show that  $\Psi$  and  $\Theta$  are holomorphic in the strip  $\{y_m + iz_m : |y_m| < \delta_m / \kappa'_m\}$  for any  $\delta_m \in (\kappa'_m, 1)$ . Hartogs' theorem [13] and the fact that in separable Banach spaces (such as  $V$ ) weak holomorphy equals strong holomorphy extend this result onto the product of strips,  $\mathcal{S}$ .  $\square$

We note that the strip  $\mathcal{S}_\rho \subset \mathbb{C}^{\mathbb{J}}$  defined in (46) contains in particular the polydisc  $\mathcal{U}_\rho$  with  $(\rho_j)_{j \in \mathbb{J}}$ , where  $\rho_j = \delta_j / \kappa'_j$ .

#### 4.4. Holomorphy and bounds on the posterior density

So far, we have shown that the responses  $G(u)$ ,  $\mathcal{G}(u)$  and the potentials  $\Phi(u; \delta)$  depend holomorphically on the coordinates  $z \in \mathcal{A}_r \subset \mathbb{C}^{\mathbb{J}}$  in the parametric representation  $u = \bar{a} + \sum_{j \in \mathbb{J}} z_j \psi_j$ . Now, we deduce bounds on the analytic continuation of the posterior density  $\Theta(z)$  in (8) as a function of the parameters  $z$  on the domains of holomorphy. We have the following.

**Theorem 4.8.** Under UEAC( $a_{\text{MIN}}$ ,  $a_{\text{MAX}}$ ) for the analytic continuation  $\Theta(z)$  of the posterior density to the domains  $\mathcal{A}_r$  of holomorphy defined in (31), i.e. for

$$\Theta(z) = \exp(-\Phi(u; \delta)|_{u=\bar{a}+\sum_{j \in \mathbb{J}} z_j \psi_j}) \quad (47)$$

there holds for every  $0 < r < a_{\text{MIN}}$

$$\sup_{z \in \mathcal{A}_r} |\Theta(z)| = \sup_{z \in \mathcal{A}_r} |\exp(-\Phi(u(z); \delta))| \leq \exp\left(\frac{\|f\|_{V^*}^2}{r^2} \sum_{k=1}^K \|o_k\|_{V^*}^2\right). \quad (48)$$

These analyticity properties, and resulting bounds, can be extended to functions  $\phi(\cdot)$  as defined by (13), using lemma 4.7 and theorem 4.8. This gives the following result.

**Corollary 4.9.** Under UEAC( $a_{\text{MIN}}$ ,  $a_{\text{MAX}}$ ), for any  $m \in \mathbb{N}$ , the functionals  $\phi(u) = p^{(m)} \in S = V^{(m)}$  the posterior densities  $\Psi(z) = \Theta(z)\phi(u(z))$  defined in (11) admit analytic continuations as strongly holomorphic,  $V^{(m)}$ -valued functions with domains  $\mathcal{A}_r$  of holomorphy defined in (31). Moreover, for these functionals the analytic continuations of  $\Psi$  in (11) admit the bounds

$$\sup_{z \in \mathcal{A}_r} \|\Theta(z)(p(z))^{(m)}\|_{V^{(m)}} \leq \frac{\|f\|_{V^*}^m}{r^m} \exp\left(\frac{\|f\|_{V^*}^2}{r^2} \sum_{k=1}^K \|o_k\|_{V^*}^2\right). \quad (49)$$

## 5. Polynomial chaos approximations of the posterior

Building on the results of the previous section, we now proceed to approximate  $\Theta(z)$ , viewed as a holomorphic functional over  $z \in \mathbb{C}^{\mathbb{J}}$ , by the so-called *polynomial chaos representations*. Exactly the same results on analyticity and on the  $N$ -term approximation of  $\Psi(z)$  hold. We omit details for reasons of brevity of exposition and confine ourselves to establishing rates of convergence of  $N$ -term truncated representations of the posterior density  $\Theta$ . The results in this section are, in one sense, *sparsity results* on the posterior density  $\Theta$ . On the other hand, such  $N$ -term truncated gpc representations of  $\Theta$  are, as we will show in the next section, computationally accessible once sparse truncated adaptive forward solvers of the parametrized system of interest are available. Such solvers are indeed available (see, e.g., [3, 5, 22] and the references therein), so that the abstract approximation results in this section have a substantive constructive aspect. Algorithms based on Smolyak-type quadratures in  $U$ , which are designed based on the present theoretical results, will be developed and analyzed in [1]. In this section, we analyze the convergence rate of  $N$ -term truncated Legendre gpc approximations of  $\Theta$  and, with the aim of a *constructive  $N$ -term approximation* of the posterior  $\Theta(y)$  in  $U$  in section 6, we analyze also the  $N$ -term truncated monomial gpc approximations of  $\Theta(y)$ .

### 5.1. The gpc representations of $\Theta$

With the index set  $\mathbb{J}$  from the parametrization (19) of the input, we associate the countable index set

$$\mathcal{F} = \{v \in \mathbb{N}_0^{\mathbb{J}} : |v|_1 < \infty\} \quad (50)$$

of multi-indices, where  $\mathbb{N}_0 = \mathbb{N} \cup \{0\}$ . We remark that sequences  $v \in \mathcal{F}$  are finitely supported even for  $\mathbb{J} = \mathbb{N}$ . For  $v \in \mathcal{F}$ , we denote by  $\mathbb{I}_v = \{j \in \mathbb{N} : v_j \neq 0\} \subset \mathbb{N}$  the ‘support’ of  $v \in \mathcal{F}$ , i.e. the finite set of indices of entries of  $v \in \mathcal{F}$  that are nonzero and by  $\aleph(v) := \#\mathbb{I}_v < \infty$ ,  $v \in \mathcal{F}$  the ‘support size’ of  $v$ , i.e. the cardinality of  $\mathbb{I}_v$ .

For the deterministic approximation of the posterior density  $\Theta(y)$  in (8), we shall use tensorized polynomial bases similar to what is done in the so-called ‘polynomial chaos’ expansions of random fields. We shall consider two particular polynomial bases: Legendre and monomial bases.

5.1.1. *Legendre expansions of  $\Theta$ .* Since we assumed that the prior measure  $\mu_0(dy)$  is built by tensorization of the uniform probability measures on  $(-1, 1)$ , we build the bases by tensorization as follows: let  $L_k(z_j)$  denote the  $k$ th Legendre polynomial of the variable  $z_j \in \mathbb{C}$ , normalized such that

$$\int_{-1}^1 (L_k(t))^2 \frac{dt}{2} = 1, \quad k = 0, 1, 2, \dots \quad (51)$$

Note that  $L_0 \equiv 1$ . The Legendre polynomials  $L_k$  in (51) are extended to tensor-product polynomials on  $U$  via

$$L_\nu(z) = \prod_{j \in \mathbb{J}} L_{\nu_j}(z_j), \quad z \in \mathbb{C}^{\mathbb{J}}, \quad \nu \in \mathcal{F}. \quad (52)$$

The normalization (51) implies that the polynomials  $L_\nu(z)$  in (52) are well defined for any  $z \in \mathbb{C}^{\mathbb{J}}$  since the finite support of each element of  $\nu \in \mathcal{F}$  implies that  $L_\nu$  in (52) is the product of only finitely many nontrivial polynomials. It moreover implies that the set of tensorized Legendre polynomials

$$\mathbb{P}(U, \mu_0(dy)) := \{L_\nu : \nu \in \mathcal{F}\} \quad (53)$$

forms a countable orthonormal basis in  $L^2(U, \mu_0(dy))$ . This observation suggests, by virtue of lemma 5.1, the use of mean-square convergent gpc expansions to represent  $\Theta$  and  $\Psi$ . Such expansions can also serve as a basis for sampling of these quantities with draws that are equidistributed with respect to the prior  $\mu_0$ .

**Lemma 5.1.** *The density  $\Theta : U \rightarrow \mathbb{R}$  is square integrable with respect to the prior  $\mu_0(dy)$  over  $U$ , i.e.  $\Theta \in L^2(U, \mu_0(dy))$ . Moreover, if the functional  $\phi(\cdot) : U \rightarrow S$  in (11) is bounded, then*

$$\int_U \|\Psi(y)\|_S^2 \mu_0(dy) < \infty,$$

i.e.  $\Psi \in L^2(U, \mu_0(dy); S)$ .

**Proof.** Since  $\Phi$  is positive, it follows that  $\Theta(y) \in [0, 1]$  for all  $y \in U$  and the first result follows because  $\mu_0$  is a probability measure. Now, define  $K = \sup_{y \in U} |\phi(y)|$ . Then,  $\sup_{y \in U} \|\Psi(y)\|_S \leq K$  and the second result follows similarly, again using that  $\mu_0$  is a probability measure.  $\square$

**Remark 5.2.** It is a consequence of (17) that in the case where  $\phi(u) = G(u) = p \in V$ , we have  $\|\Psi(y)\|_V \leq \|f\|_{V^*}/a_{\text{MIN}}$  for all  $y \in U$ . Thus, the second assertion of lemma 5.1 holds for calculation of the expectation of the pressure under the posterior distribution on  $u$ . Indeed the assertion holds for all moments of the pressure, the concrete examples that we concentrate on here.

Since  $\mathbb{P}(U, \mu_0(dy))$  in (53) is a countable orthonormal basis of  $L^2(U, \mu_0(dy))$ , the density  $\Theta(y)$  of the posterior measure given data  $\delta \in Y$ , and the posterior reweighted pressure  $\Psi(y)$  can be represented in  $L^2(U, \mu_0(dy))$  by (parametric and deterministic) generalized Legendre polynomial chaos expansions. We start by considering the scalar-valued function  $\Theta(y)$ :

$$\Theta(y) = \sum_{\nu \in \mathcal{F}} \theta_\nu L_\nu(y) \quad \text{in } L^2(U, \mu_0(dy)), \quad (54)$$

where the gpc expansion coefficients  $\theta_\nu$  are defined by

$$\theta_\nu = \int_U \Theta(y) L_\nu(y) \mu_0(dy), \quad \nu \in \mathcal{F}. \quad (55)$$



By Parseval's equation and the normalization (51), it follows immediately from (54) and lemma 5.1 with Parseval's equality that the second moment of the posterior density with respect to the prior

$$\|\Theta\|_{L^2(U, \mu_0(dy))}^2 = \sum_{\nu \in \mathcal{F}} |\theta_\nu|^2 \quad (56)$$

is finite.

*5.1.2. Monomial expansions of  $\Theta$ .* We next consider expansions of the posterior density  $\Theta$  with respect to monomials

$$y^\nu = \prod_{j \geq 1} y_j^{\nu_j}, \quad y \in U, \quad \nu \in \mathcal{F}.$$

Once more, the infinite product is well defined since, for every  $\nu \in \mathcal{F}$ , it contains only  $\aleph(\nu)$  many nontrivial factors. By lemma 4.7 and theorem 4.8, the posterior density  $\Theta(y)$  admits an analytic continuation to the product of strips  $\mathcal{S}_\rho$  that contains, in particular, the polydisc  $\mathcal{U}_\rho$ . In  $U$ ,  $\Theta(y)$  can therefore be represented by a monomial expansion with uniquely determined coefficients  $\tau_\nu \in V$  that coincide, by uniqueness of the analytic continuation, with the Taylor coefficients of  $\Theta$  at  $0 \in U$ :

$$\forall y \in U: \quad \Theta(y) = \sum_{\nu \in \mathcal{F}} \tau_\nu y^\nu, \quad \tau_\nu := \frac{1}{\nu!} \partial_y^\nu \Theta(y) |_{y=0}. \quad (57)$$

## 5.2. Best $N$ -term approximations of $\Theta$

In our deterministic parametric approach to Bayesian estimation, the evaluation of expectations under the posterior requires evaluation of integrals (10) and (12). Our strategy is to approximate these integrals by truncating the spectral representation (54), as well as a similar expression for  $\Psi(y)$ , to a finite number  $N$  of significant terms, and to estimate the error incurred by doing so. It is instructive to compare with MC methods. Under the conditions of lemma 5.1, posterior expectation of functions  $\Psi$  have finite second moments so that MC methods exhibit the convergence rate  $N^{-1/2}$  in terms of the number  $N$  of samples, with similar extension to MCMC methods. Here, however, we will show that it is possible to derive approximations that incur error decaying more quickly than the square root of  $N$ , where  $N$  is now the number of significant terms retained in (54).

By (56), the coefficient sequence  $(\theta_\nu)_{\nu \in \mathcal{F}}$  must necessarily decay. If this decay is sufficiently strong, possibly high convergence rates of  $N$ -term approximations of integrals (10) and (12) occur. The following classical result from approximation theory [9] makes these heuristic considerations precise: denote by  $(\gamma_n)_{n \in \mathbb{N}}$  a (generally not unique) decreasing rearrangement of the sequence  $(|\theta_\nu|)_{\nu \in \mathcal{F}}$ . Then, for any summability exponents  $0 < \sigma \leq q \leq \infty$  and for any  $N \in \mathbb{N}$  holds

$$\left( \sum_{n > N} \gamma_n^q \right)^{\frac{1}{q}} \leq N^{-(\frac{1}{\sigma} - \frac{1}{q})} \left( \sum_{n \geq 1} \gamma_n^\sigma \right)^{\frac{1}{\sigma}}. \quad (58)$$

*5.2.1.  $L^2(U; \mu_0)$  approximation.* Denote by  $\Lambda_N \subset \mathcal{F}$  a set of indices  $\nu \in \mathcal{F}$  corresponding to  $N$  largest gpc coefficients  $|\theta_\nu|$  in (54), and denote by

$$\Theta_{\Lambda_N}(y) := \sum_{\nu \in \Lambda_N} \theta_\nu L_\nu(y) \quad (59)$$

the Legendre expansion (54) truncated to this set of indices. Using (58) with  $q = 2$ , Parseval's equation (56) and  $0 < \sigma \leq 1$ , we obtain for all  $N$

$$\|\Theta(z) - \Theta_{\Lambda_N}(z)\|_{L^2(U, \mu_0(\mathrm{d}y))} \leq N^{-s} \|(\theta_v)\|_{\ell^\sigma(\mathcal{F})}, \quad s := \frac{1}{\sigma} - \frac{1}{2}. \quad (60)$$

We infer from (60) that a mean-square convergence rate  $s > 1/2$  of the approximate posterior density  $\Theta_{\Lambda_N}$  can be achieved *provided that*  $(\theta_v) \in \ell^\sigma(\mathcal{F})$  for some  $0 < \sigma < 1$ .

**5.2.2.  $L^1(U; \mu_0)$  and pointwise approximation of  $\Theta$ .** The analyticity of  $\Theta(y)$  in  $\mathcal{U}_\rho$  implies that  $\Theta(y)$  can be represented by the Taylor expansion (57). This expansion is unconditionally summable in  $U$  and, for any sequence  $\{\Lambda_N\}_{N \in \mathbb{N}} \subset \mathcal{F}$  that exhausts  $\mathcal{F}$ <sup>4</sup>, the corresponding sequence of  $N$ -term truncated partial Taylor sums

$$T_{\Lambda_N}(y) := \sum_{v \in \Lambda_N} \tau_v y^v \quad (61)$$

converges pointwise in  $U$  to  $\Theta$ . Since for  $y \in U$  and  $v \in \mathcal{F}$ , we have  $|y^v| \leq 1$ , for any  $\Lambda_N \subset \mathcal{F}$  of cardinality not exceeding  $N$  holds

$$\sup_{y \in U} |\Theta(y) - T_{\Lambda_N}(y)| = \sup_{y \in U} \left| \sum_{v \in \mathcal{F} \setminus \Lambda_N} \tau_v y^v \right| \leq \sum_{v \in \mathcal{F} \setminus \Lambda_N} |\tau_v|. \quad (62)$$

Similarly, we have

$$\|\Theta - T_{\Lambda_N}\|_{L^1(U, \mu_0)} = \left\| \sum_{v \in \mathcal{F} \setminus \Lambda_N} \tau_v y^v \right\|_{L^1(U, \mu_0)} \leq \sum_{v \in \mathcal{F} \setminus \Lambda_N} |\tau_v| \|y^v\|_{L^1(U, \mu_0)}.$$

For  $v \in \mathcal{F}$ , we calculate

$$\|y^v\|_{L^1(U, \mu_0)} = \int_{y \in U} |y^v| \mu_0(\mathrm{d}y) = \frac{1}{(v + \mathbf{1})!}$$

so that we find

$$\|\Theta - T_{\Lambda_N}\|_{L^1(U, \mu_0)} \leq \sum_{v \in \mathcal{F} \setminus \Lambda_N} \frac{|\tau_v|}{(v + \mathbf{1})!}. \quad (63)$$

**5.2.3. Summary.** There are, hence, two main issues to be addressed to employ the preceding approximations in practice: (i) establishing the summability of the coefficient sequences in the series (54), (57) and (ii) finding algorithms that locate sets  $\Lambda_N \subset \mathcal{F}$  of cardinality not exceeding  $N$  for which the truncated partial sums preserve the optimal convergence rates and, once these sets are localized, to determine the  $N$  'active' coefficients  $\theta_v$  or  $\tau_v$ , preferably in close to  $O(N)$  operations. In the remainder of this section, we address (i) and consider (ii) in the next section.

### 5.3. Sparsity of the posterior density $\Theta$

The analysis in the previous section shows that the convergence rate of the truncated gpc-type approximations (59) and (61) on the parameter space  $U$  is determined by the  $\sigma$ -summability of the corresponding coefficient sequences  $(|\theta_v|)_{v \in \mathcal{F}}$ ,  $(|\tau_v|)_{v \in \mathcal{F}}$ . We now show that summability (and, hence, sparsity) of Legendre and Taylor coefficient sequences in expansions (54) and

<sup>4</sup> We recall that a sequence  $\{\Lambda_N\}_{N \in \mathbb{N}} \subset \mathcal{F}$  of index sets  $\Lambda_N$  whose cardinality does not exceed  $N$  exhausts  $\mathcal{F}$  if any finite  $\Lambda \subset \mathcal{F}$  is contained in all  $\Lambda_N$  for  $N \geq N_0$  with  $N_0$  being sufficiently large.

(57) is determined by that of the sequence  $(\|\psi_j\|_{L^\infty(D)})_{j \in \mathbb{N}}$  in the input's fluctuation expansion (19). Throughout, assumptions 3.1 and 3.2 will be required to hold. We formalize the decay of  $\psi_j$  in (4) by the following.

**Assumption 5.3.** *There exists  $0 < \sigma < 1$  such that for the parametric representations (19) and (4), it holds that*

$$\sum_{j=1}^{\infty} \|\psi_j\|_{L^\infty(D)}^\sigma < \infty. \tag{64}$$

The strategy of establishing sparsity of the sequences  $(|\theta_v|)_{v \in \mathcal{F}}, (|\tau_v|)_{v \in \mathcal{F}}$  is based on estimating the sequences by Cauchy's integral formula applied to the analytic continuation of  $\Theta$ .

5.3.1. *Complex extension of the parametric problem.* To estimate  $|\theta_v|$  in (59), we shall use the holomorphy of solution to the (analytic continuation of the) parametric deterministic problem: let  $0 < K < 1$  be a constant such that

$$K \sum_{j=1}^{\infty} \|\psi_j\|_{L^\infty(D)} < \frac{a_{\min}}{8}. \tag{65}$$

Such a constant exists by assumption 5.3. For  $K$  selected in this fashion, we next choose an integer  $J_0$  such that

$$\sum_{j>J_0} \|\psi_j\|_{L^\infty(D)} < \frac{a_{\min}K}{24(1+K)}.$$

Let  $E = \{1, 2, \dots, J_0\}$  and  $F = \mathbb{N} \setminus E$ . We define

$$|v_F| = \sum_{j>J_0} |v_j|.$$

For each  $v \in \mathcal{F}$ , we define a  $v$ -dependent radius vector  $\mathbf{r} = (r_m)_{m \in \mathbb{J}}$  with  $r_m > 0$  for all  $m \in \mathbb{J}$  as follows:

$$r_m = K \text{ when } m \leq J_0 \text{ and } r_m = 1 + \frac{a_{\min}v_m}{4|v_F|\|\psi_m\|_{L^\infty(D)}} \text{ when } m > J_0, \tag{66}$$

where we make the convention that  $\frac{|v_j|}{|v_F|} = 0$  if  $|v_F| = 0$ . We consider the open discs  $\mathcal{U}_m \subset \mathbb{C}$  defined by

$$[-1, 1] \subset \mathcal{U}_m := \{z_m \in \mathbb{C} : |z_m| < 1 + r_m\} \subset \mathbb{C}. \tag{67}$$

We will extend the parametric deterministic problem (32) to parameter vectors  $z$  in the polydiscs

$$\mathcal{U}_{1+\mathbf{r}} := \bigotimes_{m \in \mathbb{J}} \mathcal{U}_m \subset \mathbb{C}^{\mathbb{J}}. \tag{68}$$

To do so, we invoke the analytic continuation of the parametric, deterministic coefficient function  $u(x, y)$  in (19) to  $z \in \mathcal{U}$  that is for such  $z$  formally given by

$$u(x, z) = \bar{a}(x) + \sum_{m \in \mathbb{J}} \psi_m(x)z_m.$$

We verify that this expression is meaningful for  $z \in \mathcal{U}_r$ : we have, for almost every  $x \in D$ ,

$$\begin{aligned} |u(x, z)| &\leq \bar{a}(x) + \sum_{m \in \mathbb{J}} |\psi_m(x)|(1 + r_m) \\ &\leq \operatorname{ess\,sup}_{x \in D} |\bar{a}(x)| + \sum_{m=1}^{J_0} \|\psi_m\|_{L^\infty(D)}(1 + K) + \sum_{m > J_0} \left(2 + \frac{a_{\min} \nu_m}{4|\nu_F| \|\psi_m\|_{L^\infty(D)}}\right) \|\psi_m\|_{L^\infty(D)} \\ &\leq \|\bar{a}\|_{L^\infty(D)} + 2 \sum_{m=1}^{\infty} \|\psi_m\|_{L^\infty(D)} + \frac{a_{\min}}{4}. \end{aligned}$$

5.3.2. Estimates of  $\theta_\nu$ .

**Proposition 5.4.** *There exists a constant  $C > 0$  such that, with the constant  $K \in (0, 1)$  in (65), for every  $\nu \in \mathcal{F}$ , the following estimate holds:*

$$|\theta_\nu| \leq C \left( \prod_{m \in \mathbb{I}(\nu)} \frac{2(1+K)}{K} \eta_m^{-\nu_m} \right), \tag{69}$$

where  $\eta_m := r_m + \sqrt{1 + r_m^2}$  with  $r_m$  as in (66).

**Proof.** For  $\nu \in \mathcal{F}$ , define  $\theta_\nu$  by (55), let  $S = \mathbb{I}(\nu)$  and define  $\bar{S} = \mathbb{J} \setminus S$ . For  $S$  denote by  $\mathcal{U}_S = \otimes_{m \in S} \mathcal{U}_m$  and  $\mathcal{U}_{\bar{S}} = \otimes_{m \in \bar{S}} \mathcal{U}_m$ , and by  $y_S = \{y_i : i \in S\}$  the extraction from  $y$ . Let  $\mathcal{E}_m$  be the ellipse in  $\mathcal{U}_m$  with foci at  $\pm 1$  and semiaxis sum  $\eta_m > 1$ . Denote also  $\mathcal{E}_S = \prod_{m \in \mathbb{I}(\nu)} \mathcal{E}_m$ . We can then write (55) as

$$\theta_\nu = \frac{1}{(2\pi i)^{|\nu|_0}} \int_U L_\nu(y) \oint_{\mathcal{E}_S} \frac{\Theta(z_S, y_{\bar{S}})}{(z_S - y_S)^{\mathbf{1}}} dz_S d\rho(y).$$

For each  $m \in \mathbb{N}$ , let  $\Gamma_m$  be a copy of  $[-1, 1]$  and  $y_m \in \Gamma_m$ . We denote by  $U_S = \prod_{m \in S} \Gamma_m$  and  $U_{\bar{S}} = \prod_{m \in \bar{S}} \Gamma_m$ . We then have

$$\theta_\nu = \frac{1}{(2\pi i)^{|\nu|_0}} \int_{U_{\bar{S}}} \oint_{\mathcal{E}_S} \Theta(z_S, y_{\bar{S}}) \int_{U_S} \frac{L_\nu(y)}{(z_S - y_S)^{\mathbf{1}}} d\rho_S(y_S) dz_S d\rho_{\bar{S}}(y_{\bar{S}}).$$

To proceed further, we recall the definitions of the Legendre functions of the second kind

$$Q_n(z) = \int_{[-1,1]} \frac{L_n(y)}{(z - y)} d\rho(y).$$

Let  $\nu_S$  be the restriction of  $\nu$  to  $S$ . We define

$$Q_{\nu_S}(z_S) = \prod_{m \in \mathbb{I}(\nu)} Q_{\nu_m}(z_m).$$

Under the Joukowski transformation  $z_m = \frac{1}{2}(w_m + w_m^{-1})$ , the Legendre polynomials of the second kind take the form

$$Q_{\nu_m} \left( \frac{1}{2}(w_m + w_m^{-1}) \right) = \sum_{k=\nu_m+1}^{\infty} \frac{q_{\nu_m k}}{w_m^k}$$

with  $|q_{\nu_m k}| \leq \pi$ . Therefore,

$$|Q_{\nu_S}(z_S)| \leq \prod_{m \in S} \sum_{k=\nu_m+1}^{\infty} \frac{\pi}{\eta_m^k} = \prod_{m \in S} \pi \frac{\eta_m^{-\nu_m-1}}{1 - \eta_m^{-1}}.$$

We then have

$$\begin{aligned}
 |\theta_v| &= \left| \frac{1}{(2\pi i)^{|v|_0}} \int_{U_{\bar{S}}} \oint_{\mathcal{E}_{\bar{S}}} \Theta(z_S, y_{\bar{S}}) \mathcal{Q}_{v_S}(z_S) dz_S d\rho_{\bar{S}}(y_S) \right| \\
 &\leq \frac{1}{(2\pi)^{|v|_0}} \int_{U_{\bar{S}}} \oint_{\mathcal{E}_{\bar{S}}} |\Theta(z_S, y_{\bar{S}})| \mathcal{Q}_{v_S}(z_S) dz_S d\rho_{\bar{S}}(y_S) \\
 &\leq \frac{1}{(2\pi)^{|v|_0}} \|\Theta(z)\|_{L^\infty(\mathcal{E}_S \times U_{\bar{S}})} \max_{\mathcal{E}_S} |\mathcal{Q}_{v_S}| \prod_{m \in S} \text{Len}(\mathcal{E}_m) \\
 &\leq \frac{1}{(2\pi)^{|v|_0}} \|\Theta(z)\|_{L^\infty(\mathcal{E}_S \times \mathcal{U}_{\bar{S}})} \prod_{m \in S} \pi \frac{\eta_m^{-v_m-1}}{1-\eta_m^{-1}} \text{Len}(\mathcal{E}_m) \\
 &\leq C \prod_{m \in S} \frac{2(1+K)}{K} \eta_m^{-v_m},
 \end{aligned}$$

as  $\text{Len}(\mathcal{E}_m) \leq 4\eta_m$ ,  $\eta_m \geq 1+K$  and as  $|\Theta(z)|$  is uniformly bounded on  $\mathcal{E}_S \times \mathcal{U}_{\bar{S}}$  by theorem 4.8.  $\square$

5.3.3. *Summability of  $\theta_v$ .* To show the  $\ell^\sigma(\mathcal{F})$  summability of  $|\theta_v|$ , we use the following result, which appears as theorem 7.2 in [6].

**Proposition 5.5.** For  $0 < \sigma < 1$  and for any sequence  $(b_v)_{v \in \mathcal{F}}$ ,

$$\left( \frac{|v|!}{v!} b^v \right)_{v \in \mathcal{F}} \in \ell^\sigma(\mathcal{F}) \iff \sum_{m \geq 1} |b_m| < 1 \quad \text{and} \quad (b_m)_{m \in \mathbb{N}} \in \ell^\sigma(\mathbb{N}).$$

This result implies the  $\sigma$ -summability of the sequence  $(\theta_v)$  of Legendre coefficients.

**Proposition 5.6.** Under assumptions 3.1 and 3.2, for  $0 < \sigma < 1$  as in assumption 5.3,  $\sum_{v \in \mathcal{F}} |\theta_v|^\sigma$  is finite.

**Proof.** We have from proposition 5.4 that

$$\begin{aligned}
 |\theta_v| &\leq C \prod_{m \in S} \frac{2(1+K)}{K} (1+r_m)^{-v_m} \\
 &\leq C \left( \prod_{m \in E, v_m \neq 0} \frac{2(1+K)}{K} \eta^{v_m} \right) \left( \prod_{m \in F, v_m \neq 0} \frac{2(1+K)}{K} \left( \frac{4|v_F| \|\psi_m\|_{L^\infty(D)}}{a_{\min} v_m} \right)^{v_m} \right),
 \end{aligned}$$

where  $\eta = 1/(1+K) < 1$ . Let  $\mathcal{F}_E = \{v \in \mathcal{F} : \mathbb{I}(v) \subset E\}$  and  $\mathcal{F}_F = \mathcal{F} \setminus E$ . From this, we have

$$\sum_{v \in \mathcal{F}} |\theta_v|^\sigma \leq C A_E A_F,$$

where

$$A_E = \sum_{v \in \mathcal{F}_E} \prod_{m \in E, v_m \neq 0} \left( \frac{2(1+K)}{K} \right)^\sigma \eta^{\sigma v_m}$$

and

$$A_F = \sum_{v \in \mathcal{F}_F} \prod_{m \in F, v_m \neq 0} \left( \frac{2(1+K)}{K} \right)^\sigma \left( \frac{4|v| \|\psi_m\|_{L^\infty(D)}}{a_{\min} v_m} \right)^{\sigma v_m}.$$

We estimate  $A_E$  and  $A_F$ : for  $A_E$ , we have

$$A_E = \left( 1 + \left( \frac{2(1+K)}{K} \right)^\sigma \sum_{m \geq 1} \eta^{pm} \right)^{J_0},$$

which is finite due to  $\eta < 1$ . For  $A_F$ , we note that for  $v_m \neq 0$ ,

$$\frac{2(1+K)}{K} \leq \left( \frac{2(1+K)}{K} \right)^{v_m}.$$

Therefore,

$$A_F \leq \sum_{v \in \mathcal{F}_F} \prod_{m \in F} \left( \frac{|v| d_m}{v_m} \right)^{\sigma v_m},$$

where

$$d_m = \frac{8(1+K) \|\psi_m\|_{L^\infty(D)}}{K a_{\min}}.$$

With the convention that  $0^0 = 1$ , we obtain from the Stirling estimate

$$\frac{n! e^n}{e \sqrt{n}} \leq n^n \leq \frac{n! e^n}{\sqrt{2\pi n}}$$

that  $|v|^{|v|} \leq |v|! e^{|v|}$ . Inserting this into the above bound for  $A_F$ , we obtain

$$\prod_{m \in F} v_m^{v_m} \geq \frac{v! e^{|v|}}{\prod_{m \in F} \max\{1, e \sqrt{v_m}\}}.$$

Hence,

$$A_F \leq \sum_{v \in \mathcal{F}_F} \left( \frac{|v|!}{v!} d^v \right)^\sigma \left( \prod_{m \in F} \max\{1, e \sqrt{v_m}\} \right)^\sigma \leq \sum_{v \in \mathcal{F}_F} \left( \frac{|v|!}{v!} \bar{d}^v \right)^\sigma,$$

where  $\bar{d}_m = e d_m$  and where we used the estimate  $e \sqrt{n} \leq e^n$ . From this, we have

$$\sum_{m \geq 1} \bar{d}_m \leq \sum_{m \in F} \frac{24(1+K) \|\psi_m\|_{L^\infty(D)}}{K a_{\min}} \leq 1.$$

Since also

$$\|\bar{d}\|_{l^\sigma(\mathbb{N})} < \infty$$

we obtain with proposition 5.5 the conclusion. □

We now show  $\sigma$ -summability of the Taylor coefficients  $\tau_v$  in (57). To this end, we proceed as in the Legendre case: first we establish sharp bounds on  $\tau_v$  by complex variable methods and then show  $\sigma$ -summability of  $(\tau_v)_{v \in \mathcal{F}}$  by a sequence factorization argument.

### 5.3.4. Bounds on the Taylor coefficients $\tau_v$ .

**Lemma 5.7.** Assume UEAC( $a_{\min}$ ,  $a_{\max}$ ) and that  $\rho = (\rho_j)_{j \geq 1}$  is an  $r$ -admissible sequence of disc radii for some  $0 < r < a_{\min}$ . Then, the Taylor coefficients  $\tau_v$  of the parametric posterior density (57) satisfy

$$\forall v \in \mathcal{F} : |\tau_v| \leq \exp \left( \frac{\|f\|_{V^*}^2}{r^2} \sum_{k=1}^K \|o_k\|_{V^*}^2 \right) \prod_{j \geq 1} \rho_j^{-v_j}. \tag{70}$$

**Proof.** For  $v = (v_j)_{j \geq 1} \in \mathcal{F}$  holds  $J = \max\{j \in \mathbb{N} : v_j \neq 0\} < \infty$ . For this  $J$ , define  $\Theta_{[J]}(z^J) := \Theta(z_1, z_2, \dots, z_J, 0, \dots)$ , i.e.  $\Theta_{[J]}(z^J)$  denotes the function of  $z^J \in \mathbb{C}^J$  obtained by setting in the posterior density  $\Theta(z)$  all coordinates  $z_j$  with  $j > J$  equal to zero. Then,

$$\partial_z^v \Theta(z)|_{z=0} = \frac{\partial^{|v|} \Theta_{[J]}}{\partial z_1^{v_1} \dots \partial z_J^{v_J}}(0, \dots, 0).$$

Since the sequence  $\rho$  is  $r$ -admissible, it follows with (48) that

$$\sup_{(z_1, \dots, z_J) \in \mathcal{U}_{\rho, J}} |\Theta_{[J]}(z_1, \dots, z_J)| \leq \exp\left(\frac{\|f\|_{V^*}^2}{r^2} \sum_{k=1}^K \|o_k\|_{V^*}^2\right), \tag{71}$$

for all  $(z_1, \dots, z_J)$  in the polydisc  $\mathcal{U}_{\rho, J} := \otimes_{1 \leq j \leq J} \{z_j \in \mathbb{C} : |z_j| \leq \rho_j\} \subset \mathbb{C}^J$ . We now prove (70) by Cauchy’s integral formula. To this end, we define  $\tilde{\rho}$  by

$$\tilde{\rho}_j := \rho_j + \epsilon \text{ if } j \leq J, \quad \tilde{\rho}_j = \rho_j \text{ if } j > J, \quad \epsilon := \frac{r}{2\|\sum_{j \leq J} |\psi_j|\|_{L^\infty(D)}}.$$

Then, the sequence  $\tilde{\rho}$  is  $r/2$ -admissible, and therefore,  $\mathcal{U}_{\tilde{\rho}} \subset \mathcal{A}_{r/2}$ . This implies that for each  $z \in \mathcal{U}_{\tilde{\rho}}$ ,  $u$  is holomorphic in each variable  $z_j$ .

It follows that  $u_j$  is holomorphic in each variable  $z_1, \dots, z_J$  on the polydisc  $\otimes_{1 \leq j \leq J} \{|z_j| < \tilde{\rho}_j\}$  that is an open neighborhood of  $\mathcal{U}_{\rho, J}$  in  $\mathbb{C}^J$ .

We may thus apply the Cauchy formula (e.g. theorem 2.1.2 of [13]) in each variable  $z_j$ :

$$u_J(z_1, \dots, z_J) = (2\pi i)^{-J} \int_{|\tilde{z}_1|=\tilde{\rho}_1} \dots \int_{|\tilde{z}_J|=\tilde{\rho}_J} \frac{u_J(\tilde{z}_1, \dots, \tilde{z}_J)}{(z_1 - \tilde{z}_1) \dots (z_J - \tilde{z}_J)} d\tilde{z}_1 \dots d\tilde{z}_J.$$

We infer

$$\frac{\partial^{|v|}}{\partial z_1^{v_1} \dots \partial z_J^{v_J}} u_J(0, \dots, 0) = v!(2\pi i)^{-J} \int_{|\tilde{z}_1|=\tilde{\rho}_1} \dots \int_{|\tilde{z}_J|=\tilde{\rho}_J} \frac{u_J(\tilde{z}_1, \dots, \tilde{z}_J)}{\tilde{z}_1^{v_1} \dots \tilde{z}_J^{v_J}} d\tilde{z}_1 \dots d\tilde{z}_J.$$

Bounding the integrand on  $\{|\tilde{z}_1| = \tilde{\rho}_1\} \times \dots \times \{|\tilde{z}_J| = \tilde{\rho}_J\} \subset \mathcal{A}_r$  with (48) implies (70).  $\square$

5.3.5.  $\sigma$ -summability of  $\tau_v$ . Proceeding in a similar fashion as in section 3 of [7], we can prove the  $\sigma$ -summability of the Taylor coefficients  $\tau_v$ .

**Proposition 5.8.** *Under assumptions 3.1, 3.2 and 5.3,  $(\|\tau_v\|_V) \in \ell^\sigma(\mathcal{F})$  for  $0 < \sigma < 1$  as in assumption 5.3.*

We remark that under the same assumptions, we also have  $\sigma$ -summability of  $(\tau_v / (v + \mathbf{1})!)_{v \in \mathcal{F}}$ , since

$$\forall v \in \mathcal{F} : |\tau_v| \geq \frac{|\tau_v|}{(v + \mathbf{1})!}.$$

5.4. Best  $N$ -term convergence rates

With (58), we infer from proposition 5.6 and from (60) convergence rates for ‘polynomial chaos’-type approximations of the posterior density  $\Theta$ .

**Theorem 5.9.** *If assumptions 3.1, 3.2 and 5.3 hold, then there is a sequence  $(\Lambda_N)_{N \in \mathbb{N}} \subset \mathcal{F}$  of index sets with cardinality not exceeding  $N$  (depending  $\sigma$  and on the data  $\delta$ ) such that the corresponding  $N$ -term truncated gpc Legendre expansions  $\Theta_{\Lambda_N}$  in (59) satisfy*

$$\|\Theta - \Theta_{\Lambda_N}\|_{L^2(U, \mu_0(dv))} \leq N^{-(\frac{1}{\sigma} - \frac{1}{2})} \|(\theta_v)\|_{\ell^\sigma(\mathcal{F}; \mathbb{R})}. \tag{72}$$

Likewise, for  $q = 1, \infty$  and for every  $N \in \mathbb{N}$ , there exist sequences  $(\Lambda_N)_{N \in \mathbb{N}} \subset \mathcal{F}$  of index sets (depending, in general, on  $\sigma, q$  and the data) whose cardinality does not exceed  $N$  such that the  $N$ -term truncated Taylor sums (61) converge with the rate  $1/\sigma - 1$ , i.e.

$$\|\Theta - T_{\Lambda_N}\|_{L^q(U, \mu_0(dy))} \leq N^{-(\frac{1}{\sigma}-1)} \|(\tau_v)\|_{\ell^\sigma(\mathcal{F}; \mathbb{R})}. \quad (73)$$

Here, for  $q = \infty$ , the norm  $\|\circ\|_{L^\infty(U; \mu_0)}$  is the supremum over all  $y \in U$ .

## 6. Approximation of expectations under the posterior

Recall that in our approach to the Bayesian estimation, the expectations under the posterior given data  $\delta$  are ratios of deterministic, infinite-dimensional parametric integrals  $Z'$  and  $Z$  with respect to the prior measure  $\mu_0$ , given by (10) and (12). For our specific elliptic inverse problem, these reduce to iterated integrals over the coordinates  $y_j \in [-1, 1]$  against a countable product of the uniform probability measures  $\frac{1}{2}dy_j$ . To render this practically feasible, the numerical evaluation of integrals of the form

$$\overline{\phi(u)}^\delta = \int_{y \in U} \phi(u(\cdot, y)) \Theta(y) \mu_0(dy) \in S \quad (74)$$

is required for functions  $\phi : U \rightarrow S$ , for a suitable state space  $S$ . Note that the choice  $\phi \equiv 1$  gives  $Z$ . For  $\phi$  not identically 1, integral (74) gives the (posterior) conditional expectation  $\mathbb{E}_{\mu^\delta}[\phi(u)]$  if normalized by  $Z$ .

For the elliptic inverse problems studied here, the choices of  $\phi(u) = u$  given by (13) with  $G(u) = p$  are of particular interest. For  $p = 1$ , this gives rise to the need to evaluate the integrals

$$\overline{p}^\delta = \int_{y \in U} p(\cdot, y) \Theta(y) \mu_0(dy) \in V \quad (75)$$

which, when normalized by  $Z$ , gives the (posterior) conditioned expectation  $\mathbb{E}_{\mu^\delta}[p]$ . We study how to approximate this integral. With the techniques developed here, and with corollary 4.9, analogous results can also be established for expectations of  $m$  point correlations of  $G(u)$  as in (13), using (74), and the normalization constant  $Z$ .

Our objective is to find constructive algorithms that achieve the high rates of convergence, in terms of number of retained terms  $N$  in a gpc expansion, implied by the theory of the previous section, and offering the potential of beating the complexity of MC-based methods. The first option to do so is to employ *sparse tensor numerical integration scheme over  $U$*  tailored to the regularity afforded by the analytic parameter dependence of the posterior density on  $y$  and of the integrands in (74). This approach is not considered here, but is considered elsewhere: we refer to [1] for details and numerical experiments. Here, we adopt an approach based on showing that integrals (74) allow *semianalytic evaluation* in log-linear<sup>5</sup> complexity with respect to  $N$ , the number of ‘active’ terms in a truncated polynomial chaos expansion of the parametric solution of the forward problem (14), (4).

To this end, we proceed as follows: based on the *assumption* that  $N$ -term gpc approximations of the parametric forward solutions  $p(x, y)$  of (14) is available, for example by the algorithms in [3, 5, 10], we show that it is possible to *construct separable  $N$ -term approximations* of the integrands in (74). The existence of such an approximate posterior density that is ‘close’ to  $\Theta$  is ensured by theorem 5.9, provided the (unknown) input data  $u$

<sup>5</sup> Meaning linear multiplied by a logarithmic factor.



satisfy certain conditions. We prove that sets  $\Lambda_N \subset \mathcal{F}$  of cardinality at most  $N$  that afford the truncation errors (72), (73) can be found in log-linear complexity with respect to  $N$  and, second, that integrals (74) with the corresponding approximate posterior density can be evaluated in such complexity and, third, we estimate the errors in the resulting conditional expectations.

### 6.1. Assumptions and notation

**Assumption 6.1.** Given a draw  $u$  of the data, an exact forward solution  $p$  of the governing equation (14) for this draw of data  $u$  is available at unit cost.

This assumption is made in order to simplify the exposition. All conclusions remain valid if this assumption is relaxed to include an additional finite-element discretization error; we refer to [1] for details. We shall use the notion of *monotone sets of multi-indices*.

**Definition 6.2.** A subset  $\Lambda_N \subset \mathcal{F}$  of finite cardinality  $N$  is called *monotone* if (M1)  $\{0\} \subset \Lambda_N$ , and if (M2)  $\forall 0 \neq v \in \Lambda_N$ , it holds that  $v - e_j \in \Lambda_N$  for all  $j \in \mathbb{I}_v$ , where  $e_j \in \{0, 1\}^{\mathbb{J}}$  denotes the index vector with 1 in the position  $j \in \mathbb{J}$  and 0 in all other positions  $i \in \mathbb{J} \setminus \{j\}$ .

Note that for monotone index sets  $\Lambda_N \subset \mathcal{F}$ , properties (M1) and (M2) in definition 6.2 imply

$$\mathbb{P}_{\Lambda_N}(U) = \text{span}\{y^v : v \in \Lambda_N\} = \text{span}\{L_v : v \in \Lambda_N\}. \quad (76)$$

Next, we will assume that a stochastic Galerkin approximation of the entire forward map of the parametric, deterministic solution with certain optimality properties is available.

**Assumption 6.3.** Given a parametric representation (19) of the unknown data  $u$ , a stochastic Galerkin approximation  $p_N \in \mathbb{P}_{\Lambda_N}(U, V)$  of the exact forward solution of the governing equation (14) is available at unit cost. Here, the set  $\Lambda_N \subset \mathcal{F}$  is a finite subset of ‘active’ gpc Legendre coefficients whose cardinality does not exceed  $N$ . In addition, we assume that the gpc approximation  $p_N \in \mathbb{P}_{\Lambda_N}(U, V)$  is quasi optimal in terms of the best  $N$ -term approximation, i.e. there exists  $C \geq 1$  independent of  $N$  such that

$$\|p - p_N\|_{L^2(U, \mu_0; V)} \leq CN^{-(1/\sigma - 1/2)} \|(\theta_v)\|_{\ell^\sigma(\mathcal{F})}. \quad (77)$$

Here,  $0 < \sigma \leq 1$  denotes the summability exponent in assumption 5.3. Note that best  $N$ -term approximations satisfy (77) with  $C = 1$ ; we may refer to (77) as a quasi-best  $N$ -term approximation property.

This best  $N$ -term convergence rate of stochastic Galerkin finite-element method (sGFEM) approximations follows from results in [6, 7], but these results do not indicate as to how sequences of sGFEM approximations that converge with this rate are actually constructed. We refer to [10] for the constructive algorithms for quasi-best  $N$ -term Legendre Galerkin approximations and to [5] for constructive algorithms for quasi-best  $N$ -term Taylor approximations and also to the references therein for details on further details for such sGFEM solvers, including space discretization. In what follows, we work under assumptions 6.1 and 6.3.

### 6.2. Best $N$ -term-based approximate conditional expectation

We first address the rates that can be achieved by the (*a priori* not accessible) best  $N$ -term approximations of the posterior density  $\Theta$  in theorem 5.9. These rates serve as benchmark rates to be achieved by any constructive procedure.

To derive these rates, we let  $\Theta_N = \Theta_{\Lambda_N}$  denote the best  $N$ -term Legendre approximations of the posterior density  $\Theta$  in theorem 5.9. With (77), we estimate

$$\begin{aligned} \|\bar{p}^\delta - \bar{p}_N^\delta\|_V &= \left\| \int_U (\Theta p - \Theta_N p_N) \mu_0(dy) \right\|_V \\ &= \left\| \int_U ((\Theta - \Theta_N)p + \Theta_N(p - p_N)) \mu_0(dy) \right\|_V \\ &\leq \int_U |\Theta - \Theta_N| \|p\|_V \mu_0(dy) + \|\Theta_N\|_{L^2(U)} \|p - p_N\|_{L^2(U, \mu_0; V)} \\ &\leq \|\Theta - \Theta_N\|_{L^2(U)} \|p\|_{L^2(U, \mu_0; V)} + \|\Theta_N\|_{L^2(U)} \|p - p_N\|_{L^2(U, \mu_0; V)} \\ &\leq CN^{-(\frac{1}{\sigma} - \frac{1}{2})}. \end{aligned}$$

With  $T_N = T_{\Lambda_N}$  denoting a best  $N$ -term Taylor approximation of  $\Theta$  in theorem 5.9, we obtain in the same fashion the bound

$$\begin{aligned} \|\bar{p}^\delta - \bar{p}_N^\delta\|_V &= \left\| \int_U (\Theta p - T_N p_N) \mu_0(dy) \right\|_V \\ &= \left\| \int_U ((\Theta - T_N)p + T_N(p - p_N)) \mu_0(dy) \right\|_V \\ &\leq \int_U |\Theta - T_N| \|p\|_V \mu_0(dy) + \|T_N\|_{L^\infty(U)} \|p - p_N\|_{L^1(U, \mu_0; V)} \\ &\leq \|\Theta - T_N\|_{L^1(U, \mu_0)} \|p\|_{L^\infty(U, \mu_0; V)} + \|T_N\|_{L^\infty(U)} \|p - p_N\|_{L^2(U, \mu_0; V)} \\ &\leq CN^{-(\frac{1}{\sigma} - 1)}. \end{aligned}$$

We now address question (ii) raised at the beginning of section 5.2, i.e. the design of practical algorithms for the construction of sequences  $(\Lambda_N)_{N \in \mathbb{N}} \subset \mathcal{F}$  such that the best- $N$  term convergence rates asserted in theorem 5.9 are attained. We develop the approximation in detail for (75); similar results for (74) may be developed for various choices of  $\phi$ .

### 6.3. Constructive $N$ -term approximation of the potential $\Phi$

We show that, from the quasi-best  $N$ -term optimal stochastic Galerkin approximation  $u_N \in \mathbb{P}_{\Lambda_N}(U, V)$ , and in particular, from its (monotone) index set  $\Lambda_N$ , a corresponding  $N$ -term approximation  $\Phi_N$  of the potential  $\Phi$  in (3) can be computed. We denote the observation corresponding to the stochastic Galerkin approximation of the system response  $p_N$  by  $\mathcal{G}_N$ , i.e. the mapping

$$U \ni y \mapsto \mathcal{G}_N(u)|_{u=\bar{a}+\sum_{j \in \mathbb{J}} y_j \psi_j} = (\mathcal{O} \circ G_N)(u)|_{u=\bar{a}+\sum_{j \in \mathbb{J}} y_j \psi_j}, \tag{78}$$

where  $G_N(u) = p_N \in \mathbb{P}_{\Lambda_N}(U; V)$ . By the linearity and boundedness of the observation functional  $\mathcal{O}(\cdot)$  then  $\mathcal{G}_N \in \mathbb{P}_{\Lambda_N}(U; \mathbb{R}^K)$ ; in the following, we assume for simplicity  $K = 1$  so that  $\mathcal{G}_N|_{u=\bar{a}+\sum_{j \in \mathbb{J}} y_j \psi_j} \in \mathbb{P}_{\Lambda_N}(U)$ . We then denote by  $U \ni u \mapsto \Phi$  the potential in (3) and by  $\Phi_N$  the potential of the stochastic Galerkin approximation  $\mathcal{G}_N$  of the forward observation map. For notational convenience, we suppress the explicit dependence on the data  $\delta$  in the following and assume that the Gaussian covariance  $\Gamma$  of the observational noise  $\eta$  in (1) is the identity:  $\Gamma = I$ . Then, for every  $y \in U$ , with  $u = \bar{a} + \sum_{j \in \mathbb{J}} y_j \psi_j$ , the exact potential  $\Phi$  and the potential  $\Phi_N$  based on  $N$ -term approximation  $p_N$  of the forward solution take the form

$$\Phi(y) = \frac{1}{2} (\delta - \mathcal{G}(u))^2, \quad \Phi_N(y) = \frac{1}{2} (\delta - \mathcal{G}_N(u))^2. \tag{79}$$

By lemma 4.7, these potentials admit extensions to holomorphic functions of the variables  $z \in \mathcal{S}_\rho$  in the strip  $\mathcal{S}_\rho$  defined in (46). Since  $\Lambda_N$  is monotone, we may write  $p_N \in \mathbb{P}_{\Lambda_N}(U, V)$  and  $\mathcal{G}_N \in \mathbb{P}_{\Lambda_N}(U)$  in terms of their (uniquely defined) Taylor expansions about  $y = 0$ :

$$\mathcal{G}_N(u) = \sum_{v \in \Lambda_N} g_v y^v. \tag{80}$$

This implies, for every  $y \in U$ ,  $\Phi_N(y) = \delta^2 - 2\delta\mathcal{G}_N(y) + (\mathcal{G}_N(y))^2$ , where

$$(\mathcal{G}_N(y))^2 = \sum_{v, v' \in \Lambda_N} g_v g_{v'} y^{v+v'} \in \mathbb{P}_{\Lambda_N + \Lambda_N}(U)$$

has a higher polynomial degree and possibly  $O(N^2)$  coefficients. Therefore, an exact evaluation of a gpc approximation of the potential  $\Phi_N$  might incur loss of linear complexity with respect to  $N$ . To preserve log-linear in  $N$  complexity, we perform an  $N$ -term truncation  $[\Phi_N]_{\#N}$  of  $\Phi_N$ , thereby introducing an additional error which, as we show next, is of the same order as the error of gpc approximation of the system’s response. The following lemma is stated in slightly more general form than is currently needed, since it will also be used for the error analysis of the posterior density ahead.

**Lemma 6.4.** *Consider the two sequences  $(g_v) \in \ell^\sigma(\mathcal{F})$  and  $(g'_{v'}) \in \ell^\sigma(\mathcal{F}')$ ,  $0 < \sigma \leq 1$ . Then,*

$$(g_v g'_{v'})_{(v, v') \in \mathcal{F} \times \mathcal{F}'} \in \ell^\sigma(\mathcal{F} \times \mathcal{F}')$$

and there holds

$$\|(g_v g'_{v'})\|_{\ell^\sigma(\mathcal{F} \times \mathcal{F}')}^\sigma \leq \|g_v\|_{\ell^\sigma(\mathcal{F})}^\sigma \|g'_{v'}\|_{\ell^\sigma(\mathcal{F}')}^\sigma. \tag{81}$$

Moreover, a best  $N$ -term truncation  $[\circ]_{\#}$  of products of corresponding best  $N$ -term truncated Taylor polynomials, defined by

$$\left[ \left( \sum_{v \in \Lambda_N} g_v y^v \right) \left( \sum_{v' \in \Lambda'_N} g'_{v'} y^{v'} \right) \right]_{\#N} := \sum_{(v, v') \in \Lambda_N^1} g_v g'_{v'} y^{v+v'} \in \mathbb{P}_{\Lambda_N^1}(U) \tag{82}$$

where  $\Lambda_N^1 \subset \mathcal{F} \times \mathcal{F}'$  is the set of sums of index pairs  $(v, v') \in \mathcal{F} \times \mathcal{F}'$  of at most  $N$  largest (in absolute value) products  $g_v g_{v'}$ , has a pointwise error in  $U$  bounded by

$$N^{-(\frac{1}{\sigma}-1)} \|g_v\|_{\ell^\sigma(\mathcal{F})} \|g'_{v'}\|_{\ell^\sigma(\mathcal{F}')}. \tag{83}$$

Moreover, if the index sets  $\Lambda_N \subset \mathcal{F}$  and  $\Lambda'_N \subset \mathcal{F}'$  are each monotone, the index set  $\bar{\Lambda}_N := \{v + v' : (v, v') \in \Lambda_N^1\} \subset \mathcal{F}$  can be chosen monotone with cardinality at most  $2N$ .

**Proof.** We calculate

$$\begin{aligned} \|g_v g'_{v'}\|_{\ell^\sigma(\mathcal{F} \times \mathcal{F}')}^\sigma &= \sum_{v \in \mathcal{F}} \sum_{v' \in \mathcal{F}'} |g_v g'_{v'}|^\sigma = \sum_{v \in \mathcal{F}} \left( |g_v|^\sigma \sum_{v' \in \mathcal{F}'} |g'_{v'}|^\sigma \right) \\ &= \|g_v\|_{\ell^\sigma(\mathcal{F})}^\sigma \|g'_{v'}\|_{\ell^\sigma(\mathcal{F}')}^\sigma. \end{aligned}$$

Since  $(g_v g'_{v'}) \in \ell^\sigma(\mathcal{F} \times \mathcal{F}')$ , we may apply (58) with (81) as follows:

$$\begin{aligned} &\left\| \left[ \sum_{v \in \Lambda_N} \sum_{v' \in \Lambda'_N} g_v g'_{v'} y^{v+v'} \right] - \left[ \sum_{v \in \Lambda_N} \sum_{v' \in \Lambda'_N} g_v g'_{v'} y^{v+v'} \right]_{\#N} \right\|_{L^\infty(U)} \\ &\leq \sum_{(v, v') \in \mathcal{F} \times \mathcal{F}' \setminus \Lambda_N^1} |g_v g'_{v'}| \leq N^{-(\frac{1}{\sigma}-1)} \|g_v\|_{\ell^\sigma(\mathcal{F})} \|g'_{v'}\|_{\ell^\sigma(\mathcal{F}')}. \end{aligned}$$

Evidently,  $\bar{\Lambda}_N \subseteq \Lambda_N + \Lambda'_N$  and the cardinality of the set  $\Lambda_N + \Lambda'_N$  is at most  $2N$ . If  $\Lambda_N$  and  $\Lambda'_N$  are monotone, then  $\Lambda_N + \Lambda'_N$  is monotone. To see it, let  $\mu \in \Lambda_N + \Lambda'_N$ . Then,  $\mu = \nu + \nu'$  for some  $\nu \in \Lambda_N$ ,  $\nu' \in \Lambda'_N$ , and  $\mathbb{I}_\mu = \mathbb{I}_\nu \cup \mathbb{I}_{\nu'}$ . Let  $0 \neq \mu$ ,  $j \in \mathbb{I}_\mu$  and assume w.l.o.g. that  $j \in \mathbb{I}_\nu$ . Then,  $\mu - e_j = (\nu - e_j) + \nu' \in \Lambda_N + \Lambda'_N$  by the assumed monotonicity of the set  $\Lambda_N$ . If  $j \in \mathbb{I}_{\nu'}$ , the argument is analogous. Therefore,  $\mu - e_j \in \Lambda_N + \Lambda'_N$  for every  $j \in \mathbb{I}_\mu$ . Hence,  $\Lambda_N + \Lambda'_N \subset \mathcal{F}$  is monotone.  $\square$

Lemma 6.4 is key to the analysis of consistency errors in the approximate evaluation of  $N$ -term truncated power series and, in particular, of the potential  $\exp(-\Phi(u; \delta))$ , which appears in the posterior density  $\Theta$ . It crucially involves Taylor-type polynomial chaos expansions. Expansions based on Legendre (or other) univariate polynomial bases can be covered by lemma 6.4 by conversion to monomial bases, using (76), as long as  $N$ -term truncations are restricted to monotone index sets  $\Lambda_N \subset \mathcal{F}$ .

Applying lemma 6.4 with  $\mathcal{F}' = \mathcal{F}$  and with  $(g_{\nu'})_{\nu' \in \mathcal{F}'} = (g_\nu)_{\nu \in \mathcal{F}}$ , we find

$$\begin{aligned} \sup_{y \in U} |\Phi_N(y) - [\Phi_N(y)]_{\#N}| &= \sup_{y \in U} |(\mathcal{G}_N(y))^2 - [(\mathcal{G}_N(y))^2]_{\#N}| \\ &\leq N^{-(\frac{1}{\sigma}-1)} \|g_\nu\|_{\ell^\sigma(\mathcal{F})}^2. \end{aligned} \tag{84}$$

6.4. Constructive  $N$ -term approximation of  $\Theta = \exp(-\Phi)$

With the  $N$ -term approximation  $[\Phi_N]_{\#N}$ , we now define the *constructive  $N$ -term approximation  $\Theta_N$  of the posterior density*. We continue to work under assumption 6.3, i.e. that  $N$ -term truncated gpc approximations  $p_N$  of the forward solution  $p(y) = G(u(y))$  of the parametric problem are available that satisfy (77). For an integer  $K(N) \in \mathbb{N}$  to be selected below, we define

$$\Theta_N = \sum_{k=0}^{K(N)} \frac{(-1)^k}{k!} [([\Phi_N]_{\#N})^k]_{\#N}. \tag{85}$$

We then estimate (all integrals are with respect to the prior measure  $\mu_0(dy)$ )

$$\begin{aligned} \|\Theta - \Theta_N\|_{L^1(U)} &= \left\| e^{-\Phi} - e^{-[\Phi_N]_{\#N}} + e^{-[\Phi_N]_{\#N}} - \sum_{k=0}^{K(N)} \frac{(-1)^k}{k!} [([\Phi_N]_{\#N})^k]_{\#N} \right\|_{L^1(U)} \\ &\leq \|e^{-\Phi} - e^{-[\Phi_N]_{\#N}}\|_{L^1(U)} + \left\| e^{-[\Phi_N]_{\#N}} - \sum_{k=0}^{K(N)} \frac{(-1)^k}{k!} [([\Phi_N]_{\#N})^k]_{\#N} \right\|_{L^1(U)} \\ &=: \text{I} + \text{II}. \end{aligned}$$

We estimate both terms separately.

For term I, we observe that due to  $x = [\Phi_N]_{\#N} - \Phi \geq 0$  for sufficiently large values of  $N$ , it holds  $0 \leq 1 - e^{-x} \leq x$ , so that by the triangle inequality and the bound (84):

$$\begin{aligned} \text{I} &= \|e^{-\Phi} (1 - e^{\Phi - [\Phi_N]_{\#N}})\|_{L^1(U)} \leq \|\Theta\|_{L^\infty(U)} \|1 - e^{-(\Phi_N]_{\#N} - \Phi)}\|_{L^1(U)} \\ &\leq \|\Theta\|_{L^\infty(U)} \|\Phi - [\Phi_N]_{\#N}\|_{L^1(U)} \leq C(\|\Phi - \Phi_N\|_{L^1(U)} + \|\Phi_N - [\Phi_N]_{\#N}\|_{L^1(U)}) \\ &\leq \|p - p_N\|_{L^2(U,V)} + CN^{-(\frac{1}{\sigma}-1)} \leq CN^{-(\frac{1}{\sigma}-1)}, \end{aligned}$$

where  $C$  depends on  $\delta$ , but is independent of  $N$ . In the preceding estimate, we used that  $\Phi > 0$  and  $0 \leq \Theta = \exp(-\Phi) < 1$  imply

$$\|\Phi - \Phi_N\|_{L^1(U)} \leq \|\mathcal{O}\|_{V^*} \|p - p_N\|_{L^2(U,V)} (2|\delta| + \|\mathcal{O}\|_{V^*} \|p + p_N\|_{L^2(U,V)}).$$

We turn to term II. Using the (globally convergent) series expansion of the exponential function, we may estimate with the triangle inequality

$$\Pi \leq \|R_{K(N)}\|_{L^1(U)} + \sum_{k=0}^{K(N)} \frac{1}{k!} \|([\Phi_N]_{\#N})^k - [([\Phi_N]_{\#N})^k]_{\#N}\|_{L^1(U)}, \tag{86}$$

where the remainder  $R_{K(N)}$  equals

$$R_{K(N)} = \sum_{k=K(N)+1}^{\infty} \frac{(-1)^k}{k!} ([\Phi_N]_{\#N})^k. \tag{87}$$

To estimate the second term in the bound (86), we claim that for every  $k, N \in \mathbb{N}_0$  the following holds:

$$\|([\Phi_N]_{\#N})^k - [([\Phi_N]_{\#N})^k]_{\#N}\|_{L^\infty(U)} \leq N^{-(\frac{1}{\sigma}-1)} \|(g_\nu)\|_{\ell^\sigma(\mathcal{F})}^{2k\sigma}. \tag{88}$$

We prove (88) for arbitrary, fixed  $N \in \mathbb{N}$  by induction with respect to  $k$ . For  $k = 0, 1$ , the bound is obvious. Assume now that the bound has been established for all powers up to some  $k \geq 2$ . Writing  $([\Phi_N]_{\#N})^{k+1} = ([\Phi_N]_{\#N})^k [\Phi_N]_{\#N}$  and denoting the sequence of Taylor coefficients of  $[\Phi_N]^k$  by  $g'_{\nu^k}$  with  $\nu^k \in (\mathcal{F} \times \mathcal{F})^k \simeq \mathcal{F}^{2k}$ , we note that by the  $k$ -fold application of (81) it follows that  $\|(g'_{\nu^k})\|_{\ell^\sigma(\mathcal{F}^{2k})} \leq \|(g_\nu)\|_{\ell^\sigma(\mathcal{F})}^{2k\sigma}$ . By the definition of  $[\Phi_N]_{\#N}$ , the same bound also holds for the coefficients of  $([\Phi_N]_{\#N})^k$ , for every  $k \in \mathbb{N}$ . We may therefore apply lemma 6.4 to the product  $([\Phi_N]_{\#N})^k [\Phi_N]_{\#N}$  and obtain the estimate (88) with  $k + 1$  in place of  $k$  from (83). Inserting (88) into (86), we find

$$\begin{aligned} \sum_{k=0}^{K(N)} \frac{1}{k!} \|([\Phi_N]_{\#N})^k - [([\Phi_N]_{\#N})^k]_{\#N}\|_{L^1(U)} &\leq N^{-(\frac{1}{\sigma}-1)} \sum_{k=0}^{K(N)} \frac{1}{k!} \|(g_\nu)\|_{\ell^\sigma(\mathcal{F})}^{2k\sigma} \\ &\leq N^{-(\frac{1}{\sigma}-1)} \exp(\|(g_\nu)\|_{\ell^\sigma(\mathcal{F})}^{2\sigma}). \end{aligned} \tag{89}$$

In a similar fashion, we estimate the remainder  $R_{K(N)}$  in (86): as the truncated Taylor expansion  $[\Phi_N]_{\#N}$  converges pointwise to  $\Phi_N$  and to  $\Phi > 0$ , for sufficiently large  $N$ , we have  $[\Phi_N]_{\#N} > 0$  for all  $y \in U$ , so that the series (87) is alternating and converges pointwise. Hence, its truncation error is bounded by the leading term of the tail sum:

$$\|R_{K(N)}\|_{L^\infty(U)} \leq \frac{\|[\Phi_N]_{\#N}\|_{L^\infty(U)}^{K(N)+1}}{(K(N) + 1)!} \leq \frac{\|(g_\nu)\|_{\ell^1(\mathcal{F})}^{2(K(N)+1)}}{(K(N) + 1)!}. \tag{90}$$

Now, given  $N$  sufficiently large, we choose  $K(N)$  so that the bound (90) is smaller than (89), which leads with Stirling's formula in (90) to the requirement

$$(K + 1) \ln \left( \frac{Ae}{K} \right) \leq \ln B - \left( \frac{1}{\sigma} - 1 \right) \ln N \tag{91}$$

for some constants  $A, B > 0$  independent of  $K$  and  $N$  (depending on  $p$  and on  $(g_\nu)$ ). One verifies that (91) is satisfied by selecting  $K(N) \simeq \ln N$ .

Therefore, under assumptions 6.1 and 6.3, we have shown how to construct an  $N$ -term approximate posterior density  $\Theta_N$  by summing  $K = O(\ln N)$  many terms in (85). The approximate posterior density has at most  $O(N)$  nontrivial terms, which can be integrated exactly against the separable prior  $\mu_0$  over  $U$  in complexity that behaves log-linearly with respect to  $N$ , under assumptions 6.1 and 6.3: the construction of  $\Theta_N$  requires  $K$ -fold performance of the  $[\cdot]_{\#N}$ -truncation operation in (82) of products of Taylor expansions, with each factor having at most  $N$  nontrivial entries, amounting altogether to solving (possibly approximately)  $O(KN \ln N) = O(N(\ln N)^2)$  forward problems.

**Remark 6.5.** Inspecting the (constructive) proof of lemma 6.4 and the definition of the  $N$ -term approximation  $\Theta_N$  of the posterior density (85), we see that the index set  $\Lambda_N^\Theta$  of active Taylor gpc coefficients of  $\Theta_N$  satisfies

$$\Lambda_N^\Theta \subset \overline{\Lambda_N^\Theta} := (\Lambda_N + \Lambda_N) + \cdots (K(N) - \text{times}) \cdots + (\Lambda_N + \Lambda_N) \subset \mathcal{F},$$

where  $\Lambda_N \subset \mathcal{F}$  is the set of  $N$  active gpc coefficients in the approximate forward solver in assumption 6.3.

If, in particular,  $\Lambda_N$  is monotone, so is the set  $\overline{\Lambda_N^\Theta}$ . This follows by induction over  $K$  with the argument in the last part of the proof of lemma 6.4. Moreover, the cardinality of  $\Lambda_N^\Theta$  is bounded by  $2NK(N) \lesssim N \log(N)$ .

## 7. Conclusions

This paper is concerned with the formulation of Bayesian inversion as a problem in infinite-dimensional parametric integration and the construction of algorithms that exploit analyticity of the forward map from state space to data space to approximate these integration problems. In this section, we make some concluding remarks about the implications of our analysis. We discuss computational complexity for such problems, and we discuss further directions for research.

### 7.1. Computational cost: idealized analysis

Throughout, we have been guided by the desire to create algorithms that outperform MC-based methods. To gain insight into this issue, we first proceed under the (idealized) setting of assumptions 6.1 and 6.3, which imply that the PDE (14), for fixed parameter  $u$ , and its parametric solution, for all  $u \in U$ , can both be approximated at unit cost. In this situation, we can study the *cost per unit error* of MC and gpc methods as follows. We neglect logarithmic corrections for clarity of exposition. The MC method will require  $\mathcal{O}(N)$  work to achieve an error of size  $N^{-\frac{1}{2}}$ , where  $N$  is a number of samples from the prior. To obtain error  $\epsilon$  thus requires work of order  $\mathcal{O}(\epsilon^{-2})$ . Recall the parameter  $\sigma$  from assumption 5.3 that measures the rate of decay of the input fluctuations and, as we have shown, governs the smoothness properties of the analytic map from unknown to data. The gpc method based on the best  $N$ -term approximation requires work which is linear in  $N$  to obtain an error of size  $N^{-(1/\sigma-1)}$ . Thus, to obtain error  $\epsilon$  requires work of order  $\mathcal{O}(\epsilon^{\sigma/(1-\sigma)})$ . For all  $\sigma < 2/3$ , the complexity of the new gpc methods, under our idealized assumptions, is superior to that of MC-based methods.

### 7.2. Computational cost: practical issues

The analysis of the previous subsection provides a clear way to understand the potential of the methods introduced in this paper and is useful for communicating the central idea. However, by working under the stated assumptions 6.1 and 6.3, some aspects of the true computational complexity of the problem are hidden. In this subsection, we briefly discuss further issues that arise. Throughout, we assume that the desired form of the unknown diffusion coefficient for the forward PDE (14) is given by (19) in the case where  $\mathbb{J} = \mathbb{N}$ :

$$u(x, y) = \bar{a}(x) + \sum_{j \in \mathbb{N}} y_j \psi_j(x), \quad x \in D. \quad (92)$$

To quantify the complexity of the problem, we assume that, for some  $b > 0$ ,

$$\|\psi_j\|_{L^\infty(D)} \asymp j^{-(1+b)}. \quad (93)$$

Then, assumption 5.3 holds for any  $\sigma > (1 + b)^{-1}$ . In practice, to implement either MC- or gpc-based methods, it is necessary to truncate the series (92) to  $J$  terms to obtain

$$u^J(x, y) = \bar{a}(x) + \sum_{1 \leq j \leq J} y_j \psi_j(x), \quad x \in D. \quad (94)$$

To quantify the computational cost of the problem, we assume that the non-parametric forward problem (14) with fixed  $u \in U$  incurs costs  $\text{pde}(J, \epsilon)$  to make an error of size  $\epsilon$  in  $V$ . Likewise, we assume that the parametric forward problem (14), for all  $u \in U$ , incurs costs  $\text{ppde}(N, J, \epsilon)$  to make an error of  $\epsilon$  in  $L^2(U, \mu_0(\text{d}y); V)$  via computation of an approximation to a quasi-optimal best  $N$ -term gpc approximation.

Both MC- and gpc-based methods will incur an error caused by truncation to  $J$  terms. Using the Lipschitz property of  $\mathcal{G}$  expressed in (26), together with the arguments developed in [8],<sup>6</sup> we deduce that the error in computing expectations caused by truncation of the input data to  $J$  terms is proportional to

$$\sum_{j=J+1}^{\infty} \|\psi_j\|_{L^\infty(D)}.$$

Under assumption (93), this is of order  $\mathcal{O}(J^{-b})$  and since  $b$  may be chosen arbitrarily close to  $1/\sigma - 1$  we obtain an error  $\mathcal{O}(J^{1-1/\sigma})$  from truncation.

The total error for MC-based methods using  $N$  samples is then of the form

$$E_{\text{mc}} = \frac{C(J)}{N^{1/2}} + \mathcal{O}(J^{1-1/\sigma}) + \epsilon.$$

In the case where  $C(J)$  is independent of  $J$ , which arises for pure MC methods based on prior sampling and for the independence MCMC sampler [15, 20], choosing  $N$  and  $J$  to balance the error gives  $N = \mathcal{O}(\epsilon^{-2})$  and  $J = \mathcal{O}(\epsilon^{-\sigma/(1-\sigma)})$  and, with these relationships imposed, the cost is  $N \times \text{pde}(J, \epsilon)$  since one forward PDE solve is made at each step of any MC method. In practice, the standard MC sampling may be ineffective, because samples from the prior are not well distributed with respect to the posterior density; this is especially true for problems with large numbers of observations and/or small observational noise. In this case, MCMC methods may be favored and it is possible that  $C(J)$  will grow with  $J$ ; see [21] for an analysis of this effect for random walk Metropolis algorithms. Balancing the error terms will then lead to a further increase in computational cost.

For gpc methods based on  $N$ -term truncation, the error is of the form

$$E_{\text{gpc}} = \mathcal{O}(N^{1-1/\sigma}) + \mathcal{O}(J^{1-1/\sigma}) + \epsilon$$

implying that  $N = J = \mathcal{O}(\epsilon^{-\sigma/(1-\sigma)})$  to balance errors. This expressions must be substituted into  $\text{ppde}(N, J, \epsilon)$  to deduce the asymptotic cost.

In practice, however, the gpc methods can also suffer when the number of observed data is high, or when the observational noise is small. To see this, note that the choice of active terms in the expansion (55) is independent of the data and is determined by the prior. For these reasons, it may be computationally expedient in practice to study methods that marry MCMC and gpc [16–18]. In a forthcoming paper [11], we will investigate the performance of the gpc-based posterior approximations, in particular in the case of values of  $\sigma$ , which are

<sup>6</sup> The key idea in [8] is that error in the forward problem transfers to error in the Bayesian inverse problem, as measured in the Hellinger metric and hence for a wide class of expectations; the analysis in [8] is devoted to Gaussian priors and situations where the Lipschitz constant of the forward model depends on the realization of the input data  $u$  and the Fernique theorem is used to control this dependence; this is more complex than required here because the Lipschitz constants in (26) here do not depend on the realization of the input data  $u$ . For these reasons, we do not feel that it is necessary to provide a proof of the error incurred by truncation.

close to  $\sigma = 1$ , i.e. in the case of little or no sparsity in the expansion of the unknown  $u$ , for parametric precomputation of an approximation of the law of the forward model, removing the necessity to compute a forward solution at each step, and by extending this idea further to multi-level LMCMC.

### 7.3. Outlook

We have proved that for a class of inverse diffusion problems with an unknown diffusion coefficient  $u$ , that in the context of a Bayesian approach to the solution of these inverse problems, given the data  $\delta$ , for a class of diffusion coefficients  $u$  that are spatially heterogeneous and uncertainty parametrized by a countable number of random coordinate variables, *sparsity in the gpc expansion of  $u$  entails the same sparsity in the density of the Bayesian posterior with respect to the prior measure.*

We have provided a constructive proof of *how to obtain an approximate posterior density by an  $O(N)$ -term truncated gpc expansion, based on a set  $\Lambda_N \subset \mathcal{F}$  of  $N$  active gpc coefficients in the parametric system's forward response.* We have indicated that several algorithms for the linear complexity computation of approximate parametrizations including prediction of the sets  $\Lambda_N$  with quasi-optimality properties (in the sense of best  $N$ -term approximations) are now available.

In [1], based on this work, we present a detailed analysis including the error incurred through finite-element discretization of the forward problem in the physical domain  $D$ , under slightly stronger hypotheses on the data  $u$  and  $f$  than studied here. Implementing these methods, and comparing them with other methods, such as those studied in [11], will provide further guidance for the development of the promising ideas introduced in this paper and variants on them.

Furthermore, we have assumed in this paper that the observation functional  $\mathcal{O}(\cdot) \in V^*$  that precludes, in space dimensions 2 and higher, point observations. Once again, results that are completely analogous to those in this paper hold also for such  $\mathcal{O}$ , albeit again under stronger hypotheses on  $u$  and on  $f$ . This will also be elaborated in [1].

As indicated in [3, 5–7, 10, 22], the gpc parametrizations (by either Taylor-type or Legendre-type polynomial chaos representations) of the laws of these quantities allow a choice of discretization of each gpc coefficient of the quantity of interest by sparse tensorization of hierarchic bases in the physical domain  $D$  and the gpc basis functions  $L_\nu(y)$  resp.  $y^\nu$  so that the additional discretization error incurred by the discretization in  $D$  can be kept of the order of the gpc truncation error with an overall computational complexity, which does not exceed that of a single, deterministic solution of the forward problem. These issues will be addressed in [1] as well.

### Acknowledgments

CS was supported by SNF and by ERC under FP7 grant AdG247277. AMS was supported by EPSRC and ERC.

### References

- [1] Andreev R, Schwab C and Stuart A M in preparation
- [2] Banks H T and Kunisch K 1989 *Estimation Techniques for Distributed Parameter Systems* (Basel: Birkhäuser)
- [3] Bieri M, Andreev R and Schwab C 2009 Sparse tensor discretization of elliptic SPDEs *SIAM J. Sci. Comput.* **31** 4281



- [4] Babuška I, Tempone R and Zouraris G E 2004 Galerkin finite element approximations of stochastic elliptic partial differential equations *SIAM J. Numer. Anal.* **42** 800–25
- [5] Chkifa A, Cohen A, DeVore R and Schwab C 2011 Sparse adaptive Taylor approximation algorithms for parametric and stochastic elliptic PDEs *Seminar for Applied Mathematics (ETH, Zürich, Switzerland)* Report 2011-44 (<http://www.sam.math.ethz.ch/reports/2011/44>)
- [6] Cohen A, DeVore R and Schwab C 2010 Convergence rates of best  $N$ -term Galerkin approximations for a class of elliptic SPDEs *J. Found. Comput. Math.* **10** 615–46
- [7] Cohen A, DeVore R and Schwab C 2011 Analytic regularity and polynomial approximation of parametric and stochastic elliptic PDEs *Anal. Appl.* at press
- [8] Cotter S L, Dashti M and Stuart A M 2010 Approximation of Bayesian inverse problems in differential equations *SIAM J. Numer. Anal.* **48** 322–45
- [9] DeVore R 1998 Nonlinear approximation *Acta Numer.* **7** 51–150
- [10] Gittelson C J 2011 Adaptive wavelet methods for elliptic partial differential equations with random operators *Seminar for Applied Mathematics (ETH, Zürich, Switzerland)* Report 2011-37 (<http://www.sam.math.ethz.ch/reports/2011/37>)
- [11] Hoang V Ha, Schwab C and Stuart A M 2012 in preparation
- [12] Hairer M, Stuart A M and Voss J 2007 Analysis of SPDEs arising in path sampling, part II: The nonlinear case *Ann. Appl. Probab.* **17** 1657–706
- [13] Hörmander L 1990 *An Introduction to Complex Analysis in Several Variables (North Holland Mathematical Library)* 3rd edn (Amsterdam: North-Holland)
- [14] Kaipio J and Somersalo E 2005 *Statistical and Computational Inverse Problems (Applied Mathematical Sciences vol 160)* (Berlin: Springer)
- [15] Liu J 2001 *Monte Carlo Strategies in Scientific Computing (Springer Texts in Statistics)* (New York: Springer)
- [16] Marzouk Y M, Najm H N and Rahn L A 2007 Stochastic spectral methods for efficient Bayesian solution of inverse problems *J. Comput. Phys.* **224** 560–86
- [17] Marzouk Y M and Xiu D 2009 A stochastic collocation approach to Bayesian inference in inverse problems *Commun. Comput. Phys.* **6** 826–47
- [18] Marzouk Y M and Najm H N 2009 Dimensionality reduction and polynomial chaos acceleration of Bayesian inference in inverse problems *J. Comput. Phys.* **228** 1862–902
- [19] McLaughlin D and Townley L R 1996 A reassessment of the groundwater inverse problem *Water Resources Res.* **32** 1131–61
- [20] Robert C P and Casella G C 1999 *Monte Carlo Statistical Methods (Springer Texts in Statistics)* (Berlin: Springer)
- [21] Roberts G O and Sherlock C 2012 Optimal scaling of random walk metropolis algorithms with discontinuous target densities *Ann. Appl. Probab.* at press
- [22] Schwab C and Gittelson C J 2011 Sparse tensor discretizations of high-dimensional parametric and stochastic PDEs *Acta Numer.* **20** 291–467
- [23] Spanos P D and Ghanem R 1989 Stochastic finite element expansion for random media *J. Eng. Mech.* **115** 1035–53
- [24] Spanos P D and Ghanem R 2003 *Stochastic Finite Elements: A Spectral Approach* (New York: Dover)
- [25] Stuart A M 2010 Inverse problems: a Bayesian approach *Acta Numer.* **19** 451–559
- [26] Todor R A and Schwab C 2007 Convergence rates for sparse chaos approximations of elliptic problems with stochastic coefficients *IMA J. Numer. Anal.* **27** 232–61
- [27] Wiener N 1938 The homogeneous chaos *Am. J. Math.* **60** 897–936