

# SANDIA REPORT

SAND2007-8093  
Unlimited Release  
Printed December 2007

## The Analysis of a Sparse Grid Stochastic Collocation Method for Partial Differential Equations with High-Dimensional Random Input Data

Fabio Nobile, Raul Tempone, Clayton Webster

Prepared by  
Sandia National Laboratories  
Albuquerque, New Mexico 87185 and Livermore, California 94550

Sandia is a multiprogram laboratory operated by Sandia Corporation, a Lockheed Martin Company, for the United States Department of Energy's National Nuclear Security Administration under Contract DE-AC04-94-AL85000.

Approved for public release; further dissemination unlimited.



**Sandia National Laboratories**

Issued by Sandia National Laboratories, operated for the United States Department of Energy by Sandia Corporation.

**NOTICE:** This report was prepared as an account of work sponsored by an agency of the United States Government. Neither the United States Government, nor any agency thereof, nor any of their employees, nor any of their contractors, subcontractors, or their employees, make any warranty, express or implied, or assume any legal liability or responsibility for the accuracy, completeness, or usefulness of any information, apparatus, product, or process disclosed, or represent that its use would not infringe privately owned rights. Reference herein to any specific commercial product, process, or service by trade name, trademark, manufacturer, or otherwise, does not necessarily constitute or imply its endorsement, recommendation, or favoring by the United States Government, any agency thereof, or any of their contractors or subcontractors. The views and opinions expressed herein do not necessarily state or reflect those of the United States Government, any agency thereof, or any of their contractors.

Printed in the United States of America. This report has been reproduced directly from the best available copy.

Available to DOE and DOE contractors from  
U.S. Department of Energy  
Office of Scientific and Technical Information  
P.O. Box 62  
Oak Ridge, TN 37831

Telephone: (865) 576-8401  
Facsimile: (865) 576-5728  
E-Mail: [reports@adonis.osti.gov](mailto:reports@adonis.osti.gov)  
Online ordering: <http://www.osti.gov/bridge>

Available to the public from  
U.S. Department of Commerce  
National Technical Information Service  
5285 Port Royal Rd  
Springfield, VA 22161

Telephone: (800) 553-6847  
Facsimile: (703) 605-6900  
E-Mail: [orders@ntis.fedworld.gov](mailto:orders@ntis.fedworld.gov)  
Online ordering: <http://www.ntis.gov/help/ordermethods.asp?loc=7-4-0#online>



# The Analysis of a Sparse Grid Stochastic Collocation Method for Partial Differential Equations with High-Dimensional Random Input Data

Fabio Nobile  
MOX  
Dipartimento di Matematica  
Politecnico di Milano, Italy.  
[fabio.nobile@polimi.it](mailto:fabio.nobile@polimi.it)

Raul Tempone  
School of Computational Science  
Florida State University  
Tallahassee, FL 32606  
[rtempone@scs.fsu.edu](mailto:rtempone@scs.fsu.edu)

Clayton Webster  
Optimization and Uncertainty Quantification Department  
Sandia National Laboratories  
P.O. Box 5800  
Albuquerque, NM 87185-1318  
[cgwebst@sandia.gov](mailto:cgwebst@sandia.gov)

## Abstract

This work describes the convergence analysis of a Smolyak-type sparse grid stochastic collocation method for the approximation of statistical quantities related to the solution of partial differential equations with random coefficients and forcing terms (input data of the model). To compute solution statistics, the sparse grid stochastic collocation method uses approximate solutions, produced here by finite elements, corresponding to a deterministic set of points in the random input space. This naturally requires solving uncoupled deterministic problems and, as such, the derived strong error estimates for the fully discrete solution are used to compare the computational efficiency of the proposed method with the Monte Carlo method. Numerical examples illustrate the theoretical results and are used to compare this approach with several others, including the standard Monte Carlo.

## Acknowledgments

The first author was partially supported by M.U.R.S.T. Cofin 2005 “Numerical Modeling for Scientific Computing and Advanced Applications”. The first and second authors were partially supported by the SANDIA project # 523695. The second author wants to acknowledge the support of the Dahlquist fellowship at the Royal Institute of Technology in Stockholm, Sweden, and the support of UdelaR in Uruguay. The third author was supported by the School of Computation Science (SCS) at Florida State University and would like to thank MOX, Dipartimento di Matematica, Politecnico di Milano, Italy, for hosting a visit to complete this research. Finally, the third author would like to thank Prof. Max Gunzburger and Dr. John Burkardt for their insight, guidance and many helpful discussions.

# Contents

<b>Summary</b> .....	<b>8</b>
<b>Nomenclature</b> .....	<b>9</b>
<b>1 Introduction</b> .....	<b>11</b>
<b>2 Problem setting</b> .....	<b>14</b>
2.1 On Finite Dimensional Noise .....	15
<b>3 Collocation techniques</b> .....	<b>18</b>
3.1 Full tensor product interpolation .....	19
3.2 Smolyak approximation .....	19
3.3 Choice of interpolation abscissas .....	20
<b>4 Error analysis</b> .....	<b>23</b>
4.1 Analysis of the approximation error .....	24
4.2 Influence of truncation errors .....	35
<b>5 Application to linear elliptic PDEs with random input data</b> .....	<b>38</b>
<b>6 Numerical Examples</b> .....	<b>40</b>
<b>7 Conclusions</b> .....	<b>46</b>
<b>References</b> .....	<b>47</b>
<b>Appendix</b>	
<b>A Additional Estimates</b> .....	<b>50</b>

## Figures

1	For a two-dimensional parameter space ( $N = 2$ ) and maximum level $w = 5$ , we plot the full tensor product grid using the Clenshaw-Curtis abscissas (left) and isotropic Smolyak sparse grids $\mathcal{H}(5, 2)$ , utilizing the Clenshaw-Curtis abscissas (middle) and the Gaussian abscissas (right). . . . .	22
2	For a finite dimensional $\Gamma^N$ with $N = 5, 11$ and $21$ we plot the log of the number of distinct Clenshaw-Curtis collocation points used by the isotropic Smolyak method and the corresponding isotropic full tensor product method versus the level $w$ (or the maximum number of points $m$ employed in each direction). . . . .	22
3	The rate of convergence of the isotropic Smolyak approximation for solving problem (6.1) with correlation length $L_c = 1/64$ using both the Gaussian and Clenshaw-Curtis abscissas. For a finite dimensional probability space $\Gamma^N$ with $N = 5$ and $N = 11$ we plot the $L^2(D)$ approximation error in the expected value in the log-linear scale (left) and log-log scale (right). . . . .	42
4	The convergence of the isotropic Smolyak approximation for solving problem (6.1) with given correlation lengths $L_c = 1/2, 1/4, 1/16$ and $1/64$ using both the Gaussian and Clenshaw-Curtis abscissas. For a finite dimensional probability space $\Gamma^N$ with $N = 5$ and $N = 11$ we plot the $L^2(D)$ approximation error in the expected value versus the number of collocation points. . . . .	43
5	A 11-dimensional comparison of the isotropic Smolyak method, the anisotropic full tensor product algorithm and Monte Carlo approach for solving problem (6.1) with correlation lengths $L_c = 1/2, 1/4, 1/16$ and $1/64$ . We plot the $L^2(D)$ approximation error in the expected value versus the number of collocation points (or samples of the Monte Carlo method). . . . .	45

## Tables

1	The $N = 11$ components of the multi index $\mathbf{p}$ computed by the anisotropic full tensor product algorithm when solving problem (6.1) with a correlation length $L_c = 1/2$ .....	42
2	The $N = 11$ components of the multi index $\mathbf{p}$ computed by the anisotropic full tensor product algorithm when solving problem (6.1) with a correlation length $L_c = 1/64$ .....	43
3	For $N = 11$ , we compare the number of function evaluations required by the Anisotropic Full Tensor product method (AF) using Gaussian abscissas, Isotropic Smolyak (IS) using Clenshaw-Curtis abscissas and the Monte Carlo (MC) method using random abscissas, to reduce the original error of problem (6.1), in expectation, by a factor of $10^4$ . ....	44

## Summary

This work proposes and analyzes a Smolyak-type sparse grid stochastic collocation method for the approximation of statistical quantities related to the solution of partial differential equations with random coefficients and forcing terms (input data of the model). To compute solution statistics, the sparse grid stochastic collocation method uses approximate solutions, produced here by finite elements, corresponding to a deterministic set of points in the random input space. This naturally requires solving uncoupled deterministic problems as in the Monte Carlo method.

If the number of random variables needed to describe the input data is moderately large, full tensor product spaces are computationally expensive to use due to the *curse of dimensionality*. In this case the sparse grid approach is still expected to be competitive with the classical Monte Carlo method. Therefore, it is of major practical relevance to understand in which situations the sparse grid stochastic collocation method is more efficient than Monte Carlo. This work provides strong error estimates for the fully discrete solution using  $L^q$  norms and analyzes the computational efficiency of the proposed method. In particular, it demonstrates algebraic convergence with respect to the total number of collocation points. The derived estimates are then used to compare the method with Monte Carlo, indicating for which problems the first is more efficient than the latter.

Computational evidence complements the present theory and shows the effectiveness of the sparse grid stochastic collocation method compared to full tensor and Monte Carlo approaches.



## Nomenclature

**a.s.** almost surely

**i.i.d.** independently identically distributed

**FEM** finite element method

**MC** Monte Carlo

**meas.** measurable

**PC** polynomial chaos

**PDE** partial differential equation

**SC** stochastic collocation

**SFEM** stochastic finite element method

**SG** stochastic Galerkin

**SPDE** stochastic partial differential equation



# 1 Introduction

Mathematical modeling and computer simulations are nowadays widely used tools to predict the behavior of physical and engineering problems. Whenever a particular application is considered, the mathematical models need to be equipped with input data, such as coefficients, forcing terms, boundary conditions, geometry, etc. However, in many applications, such input data may be affected by a relatively large amount of uncertainty. This can be due to an intrinsic variability in the physical system as, for instance, in the mechanical properties of many bio-materials, polymeric fluids, or composite materials, the action of wind or seismic vibrations on civil structures, etc.

In other situations, uncertainty may come from our difficulty in characterizing accurately the physical system under investigation as in the study of groundwater flows, where the subsurface properties such as porosity and permeability in an aquifer have to be extrapolated from measurements taken only in few spatial locations.

Such uncertainties can be included in the mathematical model adopting a probabilistic setting, provided enough information is available for a complete statistical characterization of the physical system. In this framework, the input data are modeled as *random variables*, or more generally, as *random fields* with a given spatial (or temporal) correlation structure.

Therefore, the goal of the mathematical and computational analysis becomes the prediction of statistical moments of the solution (mean value, variance, covariance, etc.) or statistics of some given responses of the system (sometimes also called quantities of physical interest which are real valued functionals of the solution), given the probability distribution of the input random data. Examples of quantities of interest could be the solution values in a given region, fluxes across given boundaries, etc.

In order to parametrize the input data for a given PDE, random fields that are either coefficients or loads can often be expanded as an infinite combination of random variables by, for instance, the Karhunen-Loève [23] or Polynomial Chaos (PC) expansions [33, 37]. Although such random fields are properly described only by means of an infinite number of random variables, whenever the realizations are slowly varying in space, with a correlation length comparable to the size of the domain, only a few terms in the above mentioned expansions are typically needed to describe the random field with sufficient accuracy. Therefore, in this case, it is reasonable to limit the analysis to just a few random variables in the expansion (see e.g. [2, 16]).

In this work we focus on partial differential equations whose coefficients and forcing terms are described by a finite dimensional random vector (*finite dimensional noise assumption*, either because the problem itself can be described by a finite number of random variables or because the input coefficients are modeled as truncated random fields).

The most popular approach to solve mathematical problems in a probabilistic setting is the Monte Carlo method (see e.g. [15] and references therein). The Monte Carlo method is easy to implement and allows one to reuse available deterministic codes. Yet, the convergence rate is typically very slow, although with a mild dependence on the number on sampled random variables.

In the last few years, other approaches have been proposed, which in certain situations feature a much faster convergence rate. We mention, among others, the Spectral Galerkin method [3, 4, 16, 20, 22, 25, 27, 36], Stochastic Collocation [5, 24, 29, 35], perturbation methods or Neumann expansions [1, 17, 30, 34].

For certain classes of problems, the solution may have a very regular dependence on the input random variables. For instance, it was shown in [5] and [3] that the solution of a linear elliptic PDE with diffusivity coefficient and/or forcing term described as truncated expansions of random fields is analytic in the input random variables. In such situations, Spectral Galerkin or Stochastic Collocation methods based on orthogonal tensor product polynomials feature a very fast convergence rate.

In particular, our earlier work [5] proposed a Stochastic Collocation/Finite Element method based on standard finite element approximations in space and a collocation on a tensor grid built upon the zeros of orthogonal polynomials with respect to the joint probability density function of the input random variables. It was shown that for an elliptic PDE the error converges exponentially fast with respect to the number of points employed for each random input variable.

The Stochastic Collocation method can be easily implemented and leads naturally to the solution of uncoupled deterministic problems as in the Monte Carlo method, even in presence of input data which depend nonlinearly on the driving random variables. It can also treat efficiently the case of non independent random variables with the introduction of an auxiliary density and handle for instance cases with lognormal diffusivity coefficient, which is not bounded in  $\Omega \times D$  but has bounded realizations. When the number of input random variables is small, Stochastic Collocation is a very effective numerical tool.

On the other hand, approximations based on tensor product grids suffer from the *curse of dimensionality* since the number of collocation points in a tensor grid grows exponentially fast in the number of input random variables.

If the number of random variables is moderately large, one should rather consider sparse tensor product spaces as first proposed by Smolyak [28] and further investigated by e.g. [6, 16, 18, 35], which will be the primary focus of this paper. It is natural to expect that the use of sparse grids will reduce dramatically the number of collocation points, while preserving a high level of accuracy and thus being able to successfully compete with Monte Carlo. Our main purpose is to clarify the limitations of the previous statement and to understand in which situations the sparse grid stochastic collocation method is more efficient than Monte Carlo.

Motivated by the above, this work proposes and analyzes a Smolyak-type sparse grid stochastic collocation method for the approximation of statistical quantities related to the solution of partial differential equations whose input data are described through a finite number of random variables. The sparse tensor product grids are built upon either Clenshaw-Curtis [11] or Gaussian abscissas. After outlining the method, this work provides strong error estimates for the fully discrete solution and analyzes its computational efficiency. In particular, it proves algebraic convergence with respect to the total number of collocation points, or equivalently, the total computational work which is directly proportional to the number of collocation points. The exponent of such algebraic convergence is connected

to both the regularity of the solution and the number of input random variables,  $N$ , and essentially deteriorates with  $N$  by a  $1/\log(N)$  factor. Then, these error estimates are used to compare the method with the standard Monte Carlo, indicating for which problems the first is more efficient than the latter.

Moreover, this work addresses the case where the input random variables come from suitably truncated expansions of random fields. There it discusses how to relate the number of points in the sparse grid to the number of random variables retained in the truncated expansion in order to balance discretization error with truncation error in the input random fields. Computational evidence complements the present theory and shows the effectiveness of the sparse grid stochastic collocation method. It also includes a comparison with full tensor and Monte Carlo methods.

The outline of the work is the following: Section 2 introduces the mathematical problem, basic notations and states a regularity assumption to be used later in the error analysis. Section 3 summarizes various collocation techniques and describes the sparse approximation method under study. It also describes two types of abscissas, Clenshaw Curtis and Gaussian, that will be employed in the sparse approximation method.

Section 4 is the core of the work. We first develop strong error estimates for the fully discrete solution using  $L_P^\infty$  and  $L_P^2$  norms for Clenshaw-Curtis and Gaussian abscissas, respectively ( $P$  being the probability measure considered). These norms control the error in the approximation of expected values of smooth functionals of the solution. Then, in Section 4.2 these error estimates are used to compare the method with the standard Monte Carlo, explaining cases where the first is more efficient than the latter.

Sections 5 and 6 focus on applications to linear elliptic PDEs with random input data. In Section 5 we verify that the assumptions under which our general theory works hold in this particular case. Then we present in Section 6 some numerical results showing the effectiveness of the proposed method when compared to the full tensor and Monte Carlo methods.

## 2 Problem setting

We begin by focusing our attention on a differential operator  $\mathcal{L}$ , linear or nonlinear, on a domain  $D \subset \mathbb{R}^d$ , which depends on some coefficients  $a(\omega, x)$  with  $x \in D$ ,  $\omega \in \Omega$ , where  $(\Omega, \mathcal{F}, P)$  a complete probability space. Here  $\Omega$  is the set of outcomes,  $\mathcal{F} \subset 2^\Omega$  is the  $\sigma$ -algebra of events and  $P : \mathcal{F} \rightarrow [0, 1]$  is a probability measure. Similarly the forcing term  $f = f(\omega, x)$  can be assumed random as well.

Consider the stochastic boundary value problem: find a random function,  $u : \Omega \times \bar{D} \rightarrow \mathbb{R}$ , such that  $P$ -almost everywhere in  $\Omega$ , or in other words almost surely (a.s.), the following equation holds:

$$\mathcal{L}(a)(u) = f \quad \text{in } D \tag{2.1}$$

equipped with suitable boundary conditions. Before introducing some assumptions we denote by  $W(D)$  a Banach space of functions  $v : D \rightarrow \mathbb{R}$  and define, for  $q \in [1, \infty]$ , the stochastic Banach spaces  $L_P^q \equiv L_P^q(\Omega; W(D))$  and  $L_P^\infty \equiv L_P^\infty(\Omega; W(D))$  as

$$L_P^q = \left\{ v : \Omega \rightarrow W(D) \mid v \text{ is strongly meas. and } \int_\Omega \|v(\omega, \cdot)\|_{W(D)}^q dP(\omega) < +\infty \right\}$$

and

$$L_P^\infty = \left\{ v : \Omega \rightarrow W(D) \mid v \text{ is strongly meas. and } P - \text{ess sup}_{\omega \in \Omega} \|v(\omega, \cdot)\|_{W(D)}^2 < +\infty \right\}.$$

Of particular interest is the space  $L_P^2(\Omega; W(D))$ , consisting of Banach valued functions that have finite second moments.

We will now make the following assumptions:

- $A_1$ ) the solution to (2.1) has realizations in the Banach space  $W(D)$ , i.e.  $u(\cdot, \omega) \in W(D)$  almost surely and  $\forall \omega \in \Omega$

$$\|u(\cdot, \omega)\|_{W(D)} \leq C \|f(\cdot, \omega)\|_{W^*(D)}$$

where we denote  $W^*(D)$  to be the dual space of  $W(D)$ , and  $C$  is a constant independent of the realization  $\omega \in \Omega$ .

- $A_2$ ) the forcing term  $f \in L_P^2(\Omega; W^*(D))$  is such that the solution  $u$  is unique and bounded in  $L_P^2(\Omega; W(D))$ .

Here we give two example problems that are posed in this setting:

**Example 2.1** *The linear problem*

$$\begin{cases} -\nabla \cdot (a(\omega, \cdot) \nabla u(\omega, \cdot)) = f(\omega, \cdot) \text{ in } \Omega \times D, \\ u(\omega, \cdot) = 0 \text{ on } \Omega \times \partial D, \end{cases} \tag{2.2}$$

with  $a(\omega, \cdot)$  uniformly bounded and coercive, i.e.

$$\exists a_{min}, a_{max} \in (0, +\infty) \text{ such that } P(\omega \in \Omega : a(\omega, x) \in [a_{min}, a_{max}] \forall x \in \bar{D}) = 1$$

and  $f(\omega, \cdot)$  square integrable with respect to  $P$ , satisfies assumptions  $A_1$  and  $A_2$  with  $W(D) = H_0^1(D)$  (see [5]).

**Example 2.2** Similarly, for  $k \in \mathbb{N}^+$ , the nonlinear problem

$$\begin{cases} -\nabla \cdot (a(\omega, \cdot) \nabla u(\omega, \cdot)) + u(\omega, \cdot)^{2k+1} = f(\omega, \cdot) & \text{in } \Omega \times D, \\ u(\omega, \cdot) = 0 & \text{on } \Omega \times \partial D, \end{cases} \quad (2.3)$$

with  $a(\omega, \cdot)$  uniformly bounded and coercive and  $f(\omega, \cdot)$  square integrable with respect to  $P$ , satisfies assumptions  $A_1$  and  $A_2$  with  $W(D) = H_0^1(D) \cap L^{2k+2}(D)$ .

**Remark 2.3 (Goals of the computation)** As said in the introduction, the goal of the mathematical and computational analysis is the prediction of statistical moments of the solution  $u$  to (2.1) (mean value, variance, covariance, etc.) or statistics of some given quantities of physical interest  $\psi(u)$ . Examples of quantities of interest could be the average value of the solution in a given region  $D_c \subset D$ ,

$$\psi(u) = \frac{1}{|D_c|} \int_{D_c} u dx,$$

and similarly average fluxes on a given direction  $n \in \mathbb{R}^d$ . In the case of Examples 2.1 and 2.2 these fluxes can be written as

$$\psi(u) = \frac{1}{|D_c|} \int_{D_c} a \frac{\partial u}{\partial n} dx.$$

## 2.1 On Finite Dimensional Noise

In some applications, the coefficient  $a$  and the forcing term  $f$  appearing in (2.1) can be described by a random vector  $[Y_1, \dots, Y_N] : \Omega \rightarrow \mathbb{R}^N$ , as in the following examples. In such cases, we will emphasize such dependence by writing  $a_N$  and  $f_N$ .

**Example 2.4 (Piecewise constant random fields)** Let us consider again problem (2.2) where the physical domain  $D$  is the union of non-overlapping subdomains  $D_i$ ,  $i = 1, \dots, N$ . We consider a diffusion coefficient that is piecewise constant and random on each subdomain, i.e.

$$a_N(\omega, x) = a_{min} + \sum_{i=1}^N \sigma_i Y_i(\omega) 1_{D_i}(x).$$

Here  $1_{D_i}$  is the indicator function of the set  $D_i$ ,  $\sigma_i, a_{min}$  are positive constants, and the random variables  $Y_i$  are nonnegative with unit variance.

In other applications the coefficients and forcing terms in (2.1) may have other type of spatial variation that is amenable to describe by an expansion. Depending on the decay of such expansion and the desired accuracy in our computations we may retain just the first  $N$  terms.

**Example 2.5 (Karhunen-Loève expansion)** We recall that any second order random field  $g(\omega, x)$ , with continuous covariance function  $cov[g] : \overline{D} \times \overline{D} \rightarrow \mathbb{R}$ , can be represented as an infinite sum of random variables, by means, for instance, of a Karhunen-Loève expansion

[23]. To this end, introduce the compact and self-adjoint operator  $T_g : L^2(D) \rightarrow L^2(D)$ , which is defined by

$$T_g v(\cdot) := \int_D \text{cov}[g](x, \cdot) v(x) dx \quad \forall v \in L^2(D).$$

Then, consider the sequence of non-negative decreasing eigenvalues of  $T_g$ ,  $\{\lambda_i\}_{i=1}^\infty$ , and the corresponding sequence of orthonormal eigenfunctions,  $\{b_i\}_{i=1}^\infty$ , satisfying

$$T_g b_i = \lambda_i b_i, \quad (b_i, b_j)_{L^2(D)} = \delta_{ij} \text{ for } i, j \in \mathbb{N}_+.$$

In addition, define mutually uncorrelated real random variables

$$Y_i(\omega) := \frac{1}{\sqrt{\lambda_i}} \int_D (g(\omega, x) - E[g](x)) b_i(x) dx, \quad i = 1, \dots$$

with zero mean and unit variance, i.e.  $E[Y_i] = 0$  and  $E[Y_i Y_j] = \delta_{ij}$  for  $i, j \in \mathbb{N}_+$ . The truncated Karhunen-Loève expansion  $g_N$ , of the stochastic function  $g$ , is defined by

$$g_N(\omega, x) := E[g](x) + \sum_{i=1}^N \sqrt{\lambda_i} b_i(x) Y_i(\omega) \quad \forall N \in \mathbb{N}_+.$$

Then by Mercer's theorem (cf [26, p. 245]), it follows that

$$\lim_{N \rightarrow \infty} \left\{ \sup_D E[(g - g_N)^2] \right\} = \lim_{N \rightarrow \infty} \left\{ \sup_D \left( \sum_{i=N+1}^{\infty} \lambda_i b_i^2 \right) \right\} = 0.$$

Observe that the  $N$  random variables in (2.5), describing the random data, are then weighted differently due to the decay of the eigen-pairs of the Karhunen-Loève expansion. The decay of eigenvalues and eigenvectors has been investigated e.g. in the works [16] and [30].

The above examples motivate us to consider problems whose coefficients are described by finitely many random variables. Thus, we will seek a random field  $u_N : \Omega \times \bar{D} \rightarrow \mathbb{R}$ , such that a.s., the following equation holds:

$$\mathcal{L}(a_N)(u_N) = f_N \quad \text{in } D, \tag{2.4}$$

We assume that equation (2.4) admits a unique solution  $u_N \in L^2_P(\Omega; W(D))$ . Therefore, following the same argument as in [5, p.1010], yields that the solution  $u_N$  of the stochastic boundary value problem (2.4) can be described by the  $[Y_1, \dots, Y_N]$  random variables, i.e.  $u_N = u_N(\omega, x) = u_N(Y_1(\omega), \dots, Y_N(\omega), x)$ .

We underline that the coefficients  $a_N$  and  $f_N$  in (2.4) may be an exact representation of the input data as in Example 2.4 or a suitable truncation of the input data as in Example 2.5. In the latter case, the solution  $u_N$  will also be an approximation of the exact solution  $u$  in (2.1) and the truncation error  $u - u_N$  has to be properly estimated (see Section 4.2).

**Remark 2.6 (Nonlinear coefficients)** *In certain cases, one may need to ensure qualitative properties on the coefficients  $a_N$  and  $f_N$  and may be worth while to describe them as*



nonlinear functions of  $Y$ . For instance, in Example 2.1 one is required to enforce positivity on the coefficient  $a_N(\omega, x)$ , say  $a_N(\omega, x) \geq a_{\min}$  for all  $x \in D$ , a.s. in  $\Omega$ . Then a better choice is to expand  $\log(a_N - a_{\min})$ . The following standard transformation guarantees that the diffusivity coefficient is bounded away from zero almost surely

$$\log(a_N - a_{\min})(\omega, x) = b_0(x) + \sum_{1 \leq n \leq N} \sqrt{\lambda_n} b_n(x) Y_n(\omega), \quad (2.5)$$

i.e. one performs a Karhunen-Loève expansion for  $\log(a_N - a_{\min})$ , assuming that  $a_N > a_{\min}$  almost surely. On the other hand, the right hand side of (2.4) can be represented as a truncated Karhunen-Loève expansion  $f_N(\omega, x) = c_0(x) + \sum_{1 \leq n \leq N} \sqrt{\mu_n} c_n(x) Y_n(\omega)$ .

For this work we denote  $\Gamma_n \equiv Y_n(\Omega)$  the image of  $Y_n$ , where we assume  $Y_n(\omega)$  to be bounded. Without loss of generality we can assume  $\Gamma_n = [-1, 1]$ . We also let  $\Gamma^N = \prod_{n=1}^N \Gamma_n$  and assume that the random variables  $[Y_1, Y_2, \dots, Y_N]$  have a joint probability density function

$$\rho : \Gamma^N \rightarrow \mathbb{R}_+, \quad \text{with } \rho \in L^\infty(\Gamma^N). \quad (2.6)$$

Thus, the plan is to approximate the function  $u_N = u_N(y, x)$ , for any  $y \in \Gamma^N$  and  $x \in \bar{D}$ . (see [5], [3])

**Remark 2.7 (Unbounded Random Variables)** *By using a similar approach to the work [5] we can easily deal with unbounded random variables, such as Gaussian or exponential ones. For the sake of simplicity in the presentation we focus our study on bounded random variables only.*

The convergence properties of the collocation techniques that will be developed in the next section depend on the regularity that the solution  $u_N$  has with respect to  $y$ . Denote  $\Gamma_n^* = \prod_{j=1, j \neq n}^N \Gamma_j$ , and let  $y_n^*$  be an arbitrary element of  $\Gamma_n^*$ . Here we require the solution to problem (2.1) to satisfy

**Assumption 2.8 (Regularity)** *For each  $y_n \in \Gamma_n$ , there exists  $\tau_n > 0$  such that the function  $u_N(y_n, y_n^*, x)$  as a function of  $y_n$ ,  $u_N : \Gamma_n \rightarrow C^0(\Gamma_n^*; W(D))$  admits an analytic extension  $u(z, y_n^*, x)$ ,  $z \in \mathbb{C}$ , in the region of the complex plane*

$$\Sigma(\Gamma_n; \tau_n) \equiv \{z \in \mathbb{C}, \text{ dist}(z, \Gamma_n) \leq \tau_n\}. \quad (2.7)$$

Moreover,  $\forall z \in \Sigma(\Gamma_n; \tau_n)$ ,

$$\|u_N(z)\|_{C^0(\Gamma_n^*; W(D))} \leq \lambda \quad (2.8)$$

with  $\lambda$  a constant independent of  $n$ .

This assumption is sound in several problems; in particular, it can be verified for the linear problem that will be analyzed in Section 5. In the more general case, this assumption should be verified for each particular application, and will have implications on the allowed regularity of the input data, e.g. coefficients, loads, etc., of the stochastic PDE under study. See also Remark 4.13 for related results based on less regularity requirements.

### 3 Collocation techniques

We seek a numerical approximation to the exact solution of (2.4) in a suitable finite dimensional subspace. To describe such a subspace properly, we introduce some standard approximation subspaces, namely:

- $W_h(D) \subset W(D)$  is a standard finite element space of dimension  $N_h$ , which contains continuous piecewise polynomials defined on regular triangulations  $\mathcal{T}_h$  that have a maximum mesh-spacing parameter  $h > 0$ . We suppose that  $W_h$  has the following deterministic approximation property: for a given function  $\varphi \in W(D)$ ,

$$\min_{v \in W_h(D)} \|\varphi - v\|_{W(D)} \leq C(s; \varphi) h^s, \quad (3.1)$$

where  $s$  is a positive integer determined by the smoothness of  $\varphi$  and the degree of the approximating finite element subspace and  $C(s; \varphi)$  is independent of  $h$ .

**Example 3.1** *Let  $D$  be a convex polygonal domain and  $W(D) = H_0^1(D)$ . For piecewise linear finite element subspaces we have*

$$\min_{v \in W_h(D)} \|\varphi - v\|_{H_0^1(D)} \leq ch \|\varphi\|_{H^2(D)}.$$

*That is,  $s = 1$  and  $C(s; \varphi) = c\|\varphi\|_{H^2(D)}$ , see for example [7].*

We will also assume that there exists a finite element operator  $\pi_h : W(D) \rightarrow W_h(D)$  with the optimality condition

$$\|\varphi - \pi_h \varphi\|_{W(D)} \leq C_\pi \min_{v \in W_h(D)} \|\varphi - v\|_{W(D)}, \quad \forall \varphi \in W(D), \quad (3.2)$$

where the constant  $C_\pi$  is independent of the mesh size  $h$ . It is worth noticing that in general the operator  $\pi_h$  will depend on the specific problem, as well as on  $y$ , i.e.  $\pi_h = \pi_h(y)$ .

- $\mathcal{P}_{\mathbf{p}}(\Gamma^N) \subset L_\rho^2(\Gamma^N)$  is the span of tensor product polynomials with degree at most  $\mathbf{p} = (p_1, \dots, p_N)$  i.e.  $\mathcal{P}_{\mathbf{p}}(\Gamma^N) = \bigotimes_{n=1}^N \mathcal{P}_{p_n}(\Gamma_n)$ , with

$$\mathcal{P}_{p_n}(\Gamma_n) = \text{span}(y_n^k, k = 0, \dots, p_n), \quad n = 1, \dots, N.$$

Hence the dimension of  $\mathcal{P}_{\mathbf{p}}(\Gamma^N)$  is  $N_p = \prod_{n=1}^N (p_n + 1)$ .

Stochastic collocation entails the evaluation of approximate values  $\pi_h u_N(y_k) = u_h^N(y_k) \in W_h(D)$ , to the solution  $u_N$  of (2.4) on a suitable set of points  $y_k \in \Gamma^N$ . Then, the fully discrete solution  $u_{h,\mathbf{p}}^N \in C^0(\Gamma^N; W_h(D))$  is a global approximation (sometimes an interpolation) constructed by linear combinations of the point values. That is

$$u_{h,\mathbf{p}}^N(y, \cdot) = \sum_{k \in \mathcal{K}} u_h^N(y_k, \cdot) l_k^{\mathbf{p}}(y), \quad (3.3)$$

where, for instance, the functions  $l_k^{\mathbf{p}}$  can be taken as the Lagrange polynomials (see Section 3.1 and 3.2). This formulation can be used to compute the mean value or variance of  $u$ ,

as described in [5, Section 2], or to approximate expected values of functionals  $\psi(u)$ , cf. Remark 2.3, by

$$E[\psi(u)] \approx E[\psi(u_{h,\mathbf{p}}^N)] \approx \sum_{k \in \mathcal{K}} \psi(u_h^N(y_k)) E[l_k^{\mathbf{P}}].$$

In the next sections we consider different choices of the evaluation points  $y_k$  and corresponding weights  $E[l_k^{\mathbf{P}}]$  in the associated quadrature formula.

### 3.1 Full tensor product interpolation

In this section we briefly recall interpolation based on Lagrange polynomials. We first introduce an index  $i \in \mathbb{N}_+$ ,  $i \geq 1$ . Then, for each value of  $i$ , let  $\{y_1^i, \dots, y_{m_i}^i\} \subset [-1, 1]$  be a sequence of abscissas for Lagrange interpolation on  $[-1, 1]$ .

For  $u \in C^0(\Gamma^1; W(D))$  and  $N = 1$  we introduce a sequence of one-dimensional Lagrange interpolation operators  $\mathcal{U}^i : C^0(\Gamma^1; W(D)) \rightarrow V_{m_i}(\Gamma^1; W(D))$

$$\mathcal{U}^i(u)(y) = \sum_{j=1}^{m_i} u(y_j^i) l_j^i(y), \quad \forall u \in C^0(\Gamma^1; W(D)), \quad (3.4)$$

where  $l_j^i \in \mathcal{P}_{m_i-1}(\Gamma^1)$  are the Lagrange polynomials of degree  $m_i-1$ , i.e.  $l_j^i(y) = \prod_{\substack{k=1 \\ k \neq j}}^{m_i} \frac{(y-y_k^i)}{(y_j^i-y_k^i)}$ , and

$$V_m(\Gamma^1; W(D)) = \left\{ v \in C^0(\Gamma^1; W(D)) : v(y, x) = \sum_{k=1}^m \tilde{v}_k(x) l_k(y), \{\tilde{v}_k\}_{k=1}^m \in W(D) \right\}.$$

Formula (3.4) reproduces exactly all polynomials of degree less than  $m_i$ . Now, in the multivariate case  $N > 1$ , for each  $u \in C^0(\Gamma^N; W(D))$  and the multi-index  $\mathbf{i} = (i_1, \dots, i_N) \in \mathbb{N}_+^N$  we define the full tensor product interpolation formulas

$$\mathcal{I}_{\mathbf{i}}^N u(y) = (\mathcal{U}^{i_1} \otimes \dots \otimes \mathcal{U}^{i_N})(u)(y) = \sum_{j_1=1}^{m_{i_1}} \dots \sum_{j_N=1}^{m_{i_N}} u(y_{j_1}^{i_1}, \dots, y_{j_N}^{i_N}) (l_{j_1}^{i_1} \otimes \dots \otimes l_{j_N}^{i_N}). \quad (3.5)$$

Clearly, the above product needs  $\prod_{n=1}^N m_{i_n}$  function evaluations. These formulas will also be used as the building blocks for the Smolyak method, described next.

### 3.2 Smolyak approximation

Here we follow closely the work [6] and describe the Smolyak *isotropic* formulas  $\mathcal{A}(w, N)$ . The Smolyak formulas are just linear combinations of product formulas (3.5) with the following key properties: only products with a relatively small number of points are used. With  $\mathcal{U}^0 = 0$  and for  $i \in \mathbb{N}_+$  define

$$\Delta^i := \mathcal{U}^i - \mathcal{U}^{i-1}. \quad (3.6)$$

Moreover, given an integer  $w \in \mathbb{N}_+$ , hereafter called the *level*, we define the sets

$$X(w, N) := \left\{ \mathbf{i} \in \mathbb{N}_+^N, \mathbf{i} \geq \mathbf{1} : \sum_{n=1}^N (i_n - 1) \leq w \right\}, \quad (3.7a)$$

$$\tilde{X}(w, N) := \left\{ \mathbf{i} \in \mathbb{N}_+^N, \mathbf{i} \geq \mathbf{1} : \sum_{n=1}^N (i_n - 1) = w \right\}, \quad (3.7b)$$

$$Y(w, N) := \left\{ \mathbf{i} \in \mathbb{N}_+^N, \mathbf{i} \geq \mathbf{1} : w - N + 1 \leq \sum_{n=1}^N (i_n - 1) \leq w \right\}, \quad (3.7c)$$

and for  $\mathbf{i} \in \mathbb{N}_+^N$  we set  $|\mathbf{i}| = i_1 + \dots + i_N$ . Then the isotropic Smolyak formula is given by

$$\mathcal{A}(w, N) = \sum_{\mathbf{i} \in X(w, N)} (\Delta^{i_1} \otimes \dots \otimes \Delta^{i_N}). \quad (3.8)$$

Equivalently, formula (3.8) can be written as (see [32])

$$\mathcal{A}(w, N) = \sum_{\mathbf{i} \in Y(w, N)} (-1)^{w+N-|\mathbf{i}|} \binom{N-1}{w+N-|\mathbf{i}|} \cdot (\mathcal{U}^{i_1} \otimes \dots \otimes \mathcal{U}^{i_N}). \quad (3.9)$$

To compute  $\mathcal{A}(w, N)(u)$ , one only needs to know function values on the “sparse grid”

$$\mathcal{H}(w, N) = \bigcup_{\mathbf{i} \in Y(w, N)} (\vartheta^{i_1} \times \dots \times \vartheta^{i_N}) \subset [-1, 1]^N, \quad (3.10)$$

where  $\vartheta^i = \{y_1^i, \dots, y_{m_i}^i\} \subset [-1, 1]$  denotes the set of abscissas used by  $\mathcal{U}^i$ . If the sets are nested, i.e.  $\vartheta^i \subset \vartheta^{i+1}$ , then  $\mathcal{H}(w, N) \subset \mathcal{H}(w+1, N)$  and

$$\mathcal{H}(w, N) = \bigcup_{\mathbf{i} \in \tilde{X}(w, N)} (\vartheta^{i_1} \times \dots \times \vartheta^{i_N}). \quad (3.11)$$

The Smolyak formula is actually interpolatory whenever nested points are used. This result has been proved in [6, Proposition 6 on page 277].

By comparing (3.11) and (3.10), we observe that the Smolyak approximation that employs nested points requires less function evaluations than the corresponding formula with non nested points. In the next section we introduce two particular sets of abscissas, nested and non nested, respectively. Also, Figure 1 shows, as an example, the sparse grid  $\mathcal{H}(5, 2)$  obtained in those two cases. Note that the Smolyak approximation formula, as presented in this Section, is isotropic, since all directions are treated equally. This can be seen from (3.8) observing that if a multi-index  $\mathbf{i} = (i_1, i_2, \dots, i_N)$  belongs to the set  $X(w, N)$ , then any permutation of  $\mathbf{i}$  also belongs to  $X(w, N)$  and contributes to the construction of the Smolyak approximation  $\mathcal{A}(w, N)$ .

### 3.3 Choice of interpolation abscissas

**Clenshaw-Curtis abscissas.** We first suggest to use Clenshaw-Curtis abscissas (see [11]) in the construction of the Smolyak formula. These abscissas are the extrema of Chebyshev

polynomials and, for any choice of  $m_i > 1$ , are given by

$$y_j^i = -\cos\left(\frac{\pi(j-1)}{m_i-1}\right), \quad j = 1, \dots, m_i. \quad (3.12)$$

In addition, one sets  $y_1^i = 0$  if  $m_i = 1$  and lets the number of abscissas  $m_i$  in each level to grow according to the following formula

$$m_1 = 1 \quad \text{and} \quad m_i = 2^{i-1} + 1, \quad \text{for } i > 1. \quad (3.13)$$

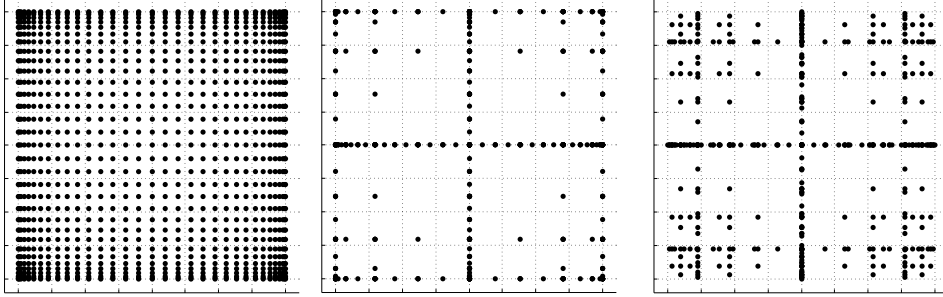
With this particular choice, one obtains nested sets of abscissas, i.e.,  $\mathcal{V}^i \subset \mathcal{V}^{i+1}$  and thereby  $\mathcal{H}(w, N) \subset \mathcal{H}(w+1, N)$ . It is important to choose  $m_1 = 1$  if we are interested in optimal approximation in relatively large  $N$ , because in all other cases the number of points used by  $\mathcal{A}(w, N)$  increases too fast with  $N$ .

**Gaussian abscissas.** We also propose to use Gaussian abscissas, i.e. the zeros of the orthogonal polynomials with respect to some positive weight. However, these Gaussian abscissas are in general not nested. Regardless, as in the Clenshaw-Curtis case, we choose the number  $m_i$  of abscissas that are used by  $\mathcal{Q}^i$  as in (3.13). See the work [31] for an insightful comparison of quadrature formulas based on Clenshaw-Curtis and Gaussian abscissas. The natural choice of the weight should be the probability density function  $\rho$  of the random variables  $Y_i(\omega)$  for all  $i$ . Yet, in the general multivariate case, if the random variables  $Y_i$  are not independent, the density  $\rho$  does not factorize, i.e.  $\rho(y_1, \dots, y_n) \neq \prod_{n=1}^N \rho_n(y_n)$ . To this end, we first introduce an auxiliary probability density function  $\hat{\rho} : \Gamma^N \rightarrow \mathbb{R}^+$  that can be seen as the joint probability of  $N$  independent random variables, i.e. it factorizes as

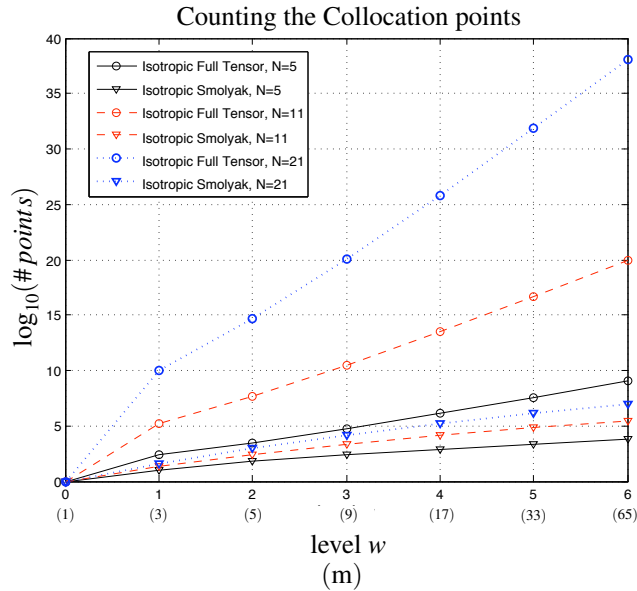
$$\hat{\rho}(y_1, \dots, y_n) = \prod_{n=1}^N \hat{\rho}_n(y_n), \quad \forall y \in \Gamma^N, \quad \text{and is such that} \quad \left\| \frac{\rho}{\hat{\rho}} \right\|_{L^\infty(\Gamma^N)} < \infty. \quad (3.14)$$

For each dimension  $n = 1, \dots, N$ , let the  $m_n$  Gaussian abscissas be the roots of the  $m_n$  degree polynomial that is  $\hat{\rho}_n$ -orthogonal to all polynomials of degree  $m_n - 1$  on the interval  $[-1, 1]$ . The auxiliary density  $\hat{\rho}$  should be chosen as close as possible to the true density  $\rho$ , so as to have the quotient  $\rho/\hat{\rho}$  not too large. Indeed, such quotient will appear in the final error estimate (see Section 4.1.2).

Examples of isotropic sparse grids, constructed from the nested Clenshaw-Curtis abscissas and the non-nested Gaussian abscissas are shown in Figure 1. There, we consider a two-dimensional parameter space and a maximum level  $w = 5$  (sparse grid  $\mathcal{H}(5, 2)$ ). To see the reduction in function evaluations with respect to full tensor product grids, we also include a plot of the corresponding Clenshaw-Curtis isotropic full tensor grid having the same maximum number of points in each direction, namely  $2^w + 1 = 33$ . Observe that if we take  $m$  points in each direction, the isotropic full tensor grid will contain  $m^N$  points while the analogous isotropic Smolyak grid  $\mathcal{H}(w, N)$  will contain much less points. Figure 2 shows the total number of points contained in the full tensor grid and in the Smolyak sparse grid as a function of the level  $w$  (or the corresponding maximum number  $m$  of points in each direction), for dimensions  $N = 5, 11, 21$ .



**Figure 1.** For a two-dimensional parameter space ( $N = 2$ ) and maximum level  $w = 5$ , we plot the full tensor product grid using the Clenshaw-Curtis abscissas (left) and isotropic Smolyak sparse grids  $\mathcal{H}(5, 2)$ , utilizing the Clenshaw-Curtis abscissas (middle) and the Gaussian abscissas (right).



**Figure 2.** For a finite dimensional  $\Gamma^N$  with  $N = 5, 11$  and  $21$  we plot the log of the number of distinct Clenshaw-Curtis collocation points used by the isotropic Smolyak method and the corresponding isotropic full tensor product method versus the level  $w$  (or the maximum number of points  $m$  employed in each direction).

## 4 Error analysis

In this section we develop error estimates that will help us compare the efficiency of the Isotropic Smolyak approximation with other alternatives, for instance the Monte Carlo method as explained in Section 4.2. Much about this has been claimed in the existing literature based on particular numerical examples. Our main goal is therefore to understand in which situations the sparse grid stochastic collocation method is more efficient than Monte Carlo.

As explained in Section 3 collocation methods can be used to approximate the solution  $u_N \in C^0(\Gamma^N; W(D))$  using finitely many function values, each of them computed by finite elements. Recall that by Assumption 2.8,  $u_N$  admits an analytic extension. Let the fully discrete numerical approximation be  $\mathcal{A}(w, N)\pi_h u_N$ . Our aim is to give a priori estimates for the total error

$$e = u - \mathcal{A}(w, N)\pi_h u_N$$

where the operator  $\mathcal{A}(w, N)$  is described by (3.8) and  $\pi_h$  is the finite element projection operator described by (3.2). We will investigate the error

$$\|u - \mathcal{A}(w, N)\pi_h u_N\| \leq \underbrace{\|u - u_N\|}_{(I)} + \underbrace{\|u_N - \pi_h u_N\|}_{(II)} + \underbrace{\|\pi_h u_N - \mathcal{A}(w, N)\pi_h u_N\|}_{(III)} \quad (4.1)$$

evaluated in the norm  $L_P^q(\Omega; W(D))$  with either  $q = 2$  or  $q = \infty$ . This yields also control of the error in the expected value of  $u$ ,  $\|E[e]\|_{W(D)} \leq E[\|e\|_{W(D)}] \leq \|e\|_{L_P^q(\Omega; W(D))}$ , and the error in the approximation of  $E[\psi(u)]$ , with  $\psi$  being a smooth functional of  $u$ . In such a case we have

$$|E[\psi(u) - \psi(\mathcal{A}(w, N)\pi_h u_N)]| \leq \left( \int_0^1 \|\delta_e \psi(u + \theta e)\|_{L_P^{q^*}(\Omega; W^*(D))} d\theta \right) \|e\|_{L_P^q(\Omega; W(D))}$$

with  $1/q + 1/q^* = 1$  and  $\delta_e \psi(u + \theta e)$  denoting the Fréchet derivative of  $\psi$  at  $u + \theta e$ .

The quantity (I) controls the truncation error for the case where the input data  $a_N$  and  $f_N$  are suitable truncations of random fields. This contribution to the total error will be considered in Section 4.2. The quantity (I) is otherwise zero if the representation of  $a_N$  and  $f_N$  is exact, as in Example 2.4. The second term (II) controls the convergence with respect to  $h$ , i.e. the finite element error, which will be dictated by standard approximability properties of the finite element space  $W_h(D)$ , given by (3.1), and the regularity in space of the solution  $u$  (see e.g. [7, 10]). For example, if we let  $q = 2$  we have

$$\|u_N - \pi_h u_N\|_{L_P^2(\Gamma^N; W(D))} \leq h^s \left( \int_{\Gamma^N} (C_\pi(y)C(s; u(y)))^2 \rho(y) dy \right)^{1/2}.$$

The full tensor product convergence results are given by [5, Theorem 1] and therefore, we will only concern ourselves with the convergence results when implementing the Smolyak approximation formula described in Section 3.2. Namely, our primary concern will be to analyze the Smolyak approximation error

$$(III) = \|\pi_h u_N - \mathcal{A}(w, N)\pi_h u_N\|_{L_P^q(\Gamma^N; W(D))}$$

for both the Clenshaw-Curtis and Gaussian versions of the Smolyak formula.

Under the very reasonable assumption that the semi-discrete finite element solution  $\pi_h u_N$  admits an analytic extension as described in Assumption 2.8 with the same analyticity region as for  $u_N$ , the behavior of the error (III) will be analogous to  $\|u_N - \mathcal{A}(w, N)u_N\|_{L^q(\Gamma^N; W(D))}$ . For this reason in the next sections we will analyze the latter.

## 4.1 Analysis of the approximation error

In this work the technique to develop error bounds for multidimensional Smolyak approximation is based on one dimensional results. Therefore, we first address the case  $N = 1$ . Let us recall the best approximation error for a function  $v : \Gamma^1 \rightarrow W(D)$  which admits an analytic extension in the region  $\Sigma(\Gamma^1; \tau) = \{z \in \mathbb{C}, \text{dist}(z, \Gamma^1) < \tau\}$  of the complex plane, for some  $\tau > 0$ . We will still denote the extension by  $v$ ; in this case,  $\tau$  is smaller than the distance between  $\Gamma^1 \subset \mathbb{R}$  and the nearest singularity of  $v(z)$  in the complex plane. Since we are considering only the case of bounded random variables, we recall the following result, whose proof can be found in [5, Lemma 7] and which is an immediate extension of the result given in [12, Chapter 7, Section 8]:

**Lemma 4.1** *Given a function  $v \in C^0(\Gamma^1; W(D))$  which admits an analytic extension in the region of the complex plane  $\Sigma(\Gamma^1; \tau) = \{z \in \mathbb{C}, \text{dist}(z, \Gamma^1) \leq \tau\}$  for some  $\tau > 0$ , there holds*

$$E_{m_i} \equiv \min_{w \in V_{m_i}} \|v - w\|_{C^0(\Gamma^1; W(D))} \leq \frac{2}{e^{\hat{\sigma}} - 1} e^{-\hat{\sigma} m_i} \max_{z \in \Sigma(\Gamma^1; \tau)} \|v(z)\|_{W(D)}$$

where  $0 < \hat{\sigma} = \log \left( \frac{2\tau}{|\Gamma^1|} + \sqrt{1 + \frac{4\tau^2}{|\Gamma^1|^2}} \right)$ .

**Remark 4.2 (Approximation with unbounded random variables)** *A related result with weighted norms holds for unbounded random variables whose probability density decays as the Gaussian density at infinity (see [5]).*

In the multidimensional case, the size of the analyticity region will depend, in general, on the direction  $n$  and it will be denoted by  $\tau_n$  (see e.g. problem considered in Section 5). The same holds for the decay coefficient  $\hat{\sigma}_n$ . In what follows, we set

$$\hat{\sigma} \equiv \min_n \hat{\sigma}_n. \tag{4.2}$$

### 4.1.1 Interpolation estimates for the Clenshaw-Curtis abscissas

In this section we develop  $L^\infty$  error estimates for the Smolyak interpolant based on Clenshaw-Curtis abscissas, cf. (3.12) and (3.13), applied to analytic functions  $u \in C^0(\Gamma^N; W(D))$  that satisfy Assumption 2.8. We remind the reader that even though in the global estimate (4.1)



it is enough to bound the approximation error (III) in the  $L^2_\rho(\Gamma^N; W(D))$  norm we will still work with the more stringent  $L^\infty(\Gamma^N; W(D))$  norm.

In our notation the norm  $\|\cdot\|_{\infty, N}$  is shorthand for  $\|\cdot\|_{L^\infty(\Gamma^N; W(D))}$  and will be used henceforth. We also define  $I_N : \Gamma^N \rightarrow \Gamma^N$  as the identity operator on an  $N$ -dimensional space. We begin by letting  $E_m$  be the error of the best approximation to functions  $u \in C^0(\Gamma^1; W(D))$  by functions  $w \in V_m$ . Similarly to [6], since the interpolation  $\mathcal{W}^i$  is exact on the subspace  $V_{m_i-1}$  we can apply the general formula

$$\|u - \mathcal{W}^i(u)\|_{\infty, 1} \leq E_{m_i-1}(u) \cdot (1 + \Lambda_{m_i}) \quad (4.3)$$

where  $\Lambda_m$  is the Lebesgue constant for our choice (3.12). It is known that

$$\Lambda_m \leq \frac{2}{\pi} \log(m-1) + 1 \quad (4.4)$$

for  $m \geq 2$ , see [13]. Using Lemma 4.1, the best approximation to functions  $u \in C^0(\Gamma^1; W(D))$  that admit an analytic extension as described by Assumption 2.8 is bounded by:

$$E_{m_i}(u) \leq \hat{C} e^{-\hat{\sigma} m_i} \quad (4.5)$$

where  $\hat{C}$  and  $\hat{\sigma} > 0$  are constants dependent on the value of  $\tau$  defined in Lemma 4.1. Hence (4.3)-(4.5) implies

$$\begin{aligned} \|(I_1 - \mathcal{W}^i)(u)\|_{\infty, 1} &\leq C \log(m_i) e^{-\sigma m_i} \leq C i e^{-\sigma 2^i}, \\ \|(\Delta^i)(u)\|_{\infty, 1} &= \|(\mathcal{W}^i - \mathcal{W}^{i-1})(u)\|_{\infty, 1} \leq \|(I_1 - \mathcal{W}^i)(u)\|_{\infty, 1} + \|(I_1 - \mathcal{W}^{i-1})(u)\|_{\infty, 1} \\ &\leq 2C i e^{-\sigma 2^{i-1}} \end{aligned}$$

for all  $i \in \mathbb{N}_+$  with positive constants  $C$  and  $\sigma = \hat{\sigma}/2$  depending on  $u$  but not on  $i$ .

The convergence proof will be split in several steps, the main results being given in Theorems 4.6 and 4.9, which state the convergence rates in terms of the level  $w$  and the total number of collocation points, respectively. We denote by  $I_d$  the identity operator applicable to functions which depend on the first  $d$  variables  $y_1, \dots, y_d$ . Then the following result holds:

**Lemma 4.3** *For functions  $u \in C^0(\Gamma^N; W(D))$  satisfying the assumption of Lemma 4.1 the isotropic Smolyak formula (3.8) based on Clenshaw Curtis abscissas satisfies:*

$$\|(I_N - \mathcal{A}(w, N))(u)\|_{\infty, N} \leq \sum_{d=1}^N R(w, d) \quad (4.6)$$

with

$$R(w, d) := \frac{1}{2} \sum_{\mathbf{i} \in \tilde{X}(w, d)} (2C)^d \left( \prod_{n=1}^d i_n \right) e^{-\sigma h(\mathbf{i}, d)} \quad (4.7)$$

and

$$h(\mathbf{i}, d) = \sum_{n=1}^d 2^{i_n-1}. \quad (4.8)$$

**Proof.** We start providing an equivalent representation of the isotropic Smolyak formula:

$$\begin{aligned}
\mathcal{A}(w, N) &= \sum_{\mathbf{i} \in X(w, N)} \bigotimes_{n=1}^N \Delta^{i_n} \\
&= \sum_{\mathbf{i} \in X(w, N-1)} \bigotimes_{n=1}^{N-1} \Delta^{i_n} \otimes \sum_{j=1}^{1+w-\sum_{n=1}^{N-1} (i_n-1)} \Delta^j \\
&= \sum_{\mathbf{i} \in X(w, N-1)} \bigotimes_{n=1}^{N-1} \Delta^{i_n} \otimes \mathcal{U}^{1+w-\sum_{n=1}^{N-1} (i_n-1)}.
\end{aligned}$$

Introducing the one-dimensional identity operator  $I_1^{(n)} : \Gamma_n \rightarrow \Gamma_n$ , for  $n = 1, \dots, N$ , the error estimate can be computed recursively using the previous representation, namely

$$\begin{aligned}
I_N - \mathcal{A}(w, N) &= I_N - \sum_{\mathbf{i} \in X(w, N-1)} \bigotimes_{n=1}^{N-1} \Delta^{i_n} \otimes \left( \mathcal{U}^{1+w-\sum_{n=1}^{N-1} (i_n-1)} - I_1^{(N)} \right) \\
&\quad - \sum_{\mathbf{i} \in X(w, N-1)} \bigotimes_{n=1}^{N-1} \Delta^{i_n} \otimes I_1^{(N)} \\
&= \sum_{\mathbf{i} \in X(w, N-1)} \bigotimes_{n=1}^{N-1} \Delta^{i_n} \otimes \left( I_1^{(N)} - \mathcal{U}^{1+w-\sum_{n=1}^{N-1} (i_n-1)} \right) \\
&\quad + (I_{N-1} - \mathcal{A}(w, N-1)) \otimes I_1^{(N)} \\
&= \sum_{d=2}^N \left[ \tilde{R}(w, d) \bigotimes_{n=d+1}^N I_1^{(n)} \right] + \left( I_1^{(1)} - \mathcal{A}(w, 1) \right) \bigotimes_{n=2}^N I_1^{(n)}
\end{aligned} \tag{4.9}$$

where, for a general dimension  $d$ , we define

$$\tilde{R}(w, d) = \sum_{\mathbf{i} \in X(w, d-1)} \bigotimes_{n=1}^{d-1} \Delta^{i_n} \otimes \left( I_1^{(d)} - \mathcal{U}^{\hat{i}_d} \right)$$

and, for any  $(i_1, \dots, i_{d-1}) \in X(w, d-1)$ , we have set  $\hat{i}_d = 1 + w - \sum_{n=1}^{d-1} (i_n - 1)$ . Observe that with this definition, the  $d$ -dimensional vector  $\mathbf{j} = (i_1, \dots, i_{d-1}, \hat{i}_d)$  belongs to the set  $\tilde{X}(w, d)$ , defined in (3.7), and the term  $\tilde{R}(w, d)$  can now be bounded as follows:

$$\begin{aligned}
\left\| \tilde{R}(w, d)(u) \right\|_{\infty, d} &\leq \sum_{\mathbf{i} \in X(w, d-1)} \prod_{n=1}^{d-1} \left\| (\Delta^{i_n})(u) \right\|_{\infty, d} \left\| \left( I_1^{(d)} - \mathcal{U}^{\hat{i}_d} \right)(u) \right\|_{\infty, d} \\
&\leq \frac{1}{2} \sum_{\mathbf{i} \in X(w, d-1)} (2C)^d \left( \prod_{n=1}^{d-1} i_n \right) \hat{i}_d e^{-\sigma(\sum_{n=1}^{d-1} 2^{i_n-1} + 2^{\hat{i}_d})} \\
&\leq \frac{(2C)^d}{2} \sum_{\mathbf{i} \in \tilde{X}(w, d)} \left( \prod_{n=1}^d i_n \right) e^{-\sigma h(\mathbf{i}, d)} =: R(w, d).
\end{aligned}$$

Hence, the Smolyak approximation error satisfies

$$\|(I_N - \mathcal{A}(w, N))(u)\|_{\infty, N} \leq \sum_{d=2}^N R(w, d) + \|(I_1^{(1)} - \mathcal{A}(w, 1))(u)\|_{\infty, 1}.$$

Observe that the last term in the previous equation can also be bounded by  $R(w, 1)$  defined in (4.7). Indeed, the set  $\tilde{X}(w, 1)$  contains only the point  $i_1 = 1 + w$  and

$$\begin{aligned} \left\| \left( I_1^{(1)} - \mathcal{A}(w, 1) \right) (u) \right\|_{\infty, 1} &= \left\| \left( I_1^{(1)} - \mathcal{U}^{1+w} \right) (u) \right\|_{\infty, 1} \\ &\leq C (1 + w) e^{-\sigma 2^{1+w}} \\ &\leq \sum_{i_1 \in \tilde{X}(w, 1)} C i_1 e^{-\sigma 2^{i_1-1}} =: R(w, 1) \end{aligned}$$

and this concludes the proof.  $\square$

**Lemma 4.4** *Let  $\delta > 0$ . Under the assumptions of Lemma 4.3 the following bound holds for the term  $R(w, d)$ ,  $d = 1, \dots, N$ :*

$$R(w, d) \leq \frac{C_1(\sigma, \delta)^d}{2} \exp\left(-\sigma d \left(2^{w/d} - \delta \tilde{C}_2(\sigma) w\right)\right) \quad (4.10)$$

where

$$\tilde{C}_2(\sigma) := 1 + \frac{1}{\log(2)} \sqrt{\frac{\pi}{2\sigma}} \quad (4.11)$$

and

$$C_1(\sigma, \delta) := \frac{4C}{e\delta\sigma} \exp\left(\delta\sigma \left\{ \frac{1}{\sigma \log^2(2)} + \frac{1}{\log(2)\sqrt{2\sigma}} + 2 \left(1 + \frac{1}{\log(2)} \sqrt{\frac{\pi}{2\sigma}}\right) \right\}\right). \quad (4.12)$$

**Proof.** The proof is divided in several steps.

1. Expand the function  $h(\mathbf{i}, d)$  up to second order around the point  $\mathbf{i}^* = (1 + w/d, \dots, 1 + w/d)$  on the subspace  $\{\mathbf{x} \in \mathbb{R}^d : |\mathbf{x} - \mathbf{1}| = w\}$ . Observe that  $\mathbf{i}^*$  is a constrained minimizer of  $h(\mathbf{i}, d)$  and

$$h(\mathbf{i}, d) \geq d2^{w/d} + \frac{\log^2(2)}{2} \sum_{n=1}^d (i_n - (1 + w/d))^2, \quad \text{for all } \mathbf{i} \in \tilde{X}(w, d). \quad (4.13)$$

2. Combining (4.7) and (4.13) estimate

$$\begin{aligned} R(w, d) &\leq \frac{(2C)^d}{2} e^{-\sigma d 2^{w/d}} \sum_{\mathbf{i} \in \tilde{X}(w, d)} \left( \prod_{n=1}^d i_n \right) e^{-\sigma \frac{\log^2(2)}{2} \sum_{n=1}^d (i_n - (1 + w/d))^2} \\ &\leq \frac{(2C)^d}{2} e^{-\sigma d 2^{w/d}} \left( \sum_{i=1}^{w+1} i e^{-\sigma \frac{\log^2(2)}{2} (i - (1 + w/d))^2} \right)^d. \end{aligned} \quad (4.14)$$

3. Next, use (A.2) from Corollary A.4 to estimate the term

$T_1 := \sum_{i=1}^{w+1} i e^{-\sigma \frac{\log^2(2)}{2} (i-(1+w/d))^2}$ . We have

$$T_1 \leq 2 \left( \frac{1}{\sigma \log^2(2)} + \frac{1}{\log(2)\sqrt{2\sigma}} \right) + 2(\text{int}\{w/d\} + 2) \left( 1 + \frac{1}{\log(2)} \sqrt{\frac{\pi}{2\sigma}} \right). \quad (4.15)$$

Combine (4.14) and (4.15), arriving at

$$R(w, d) \leq C_1(\sigma, d, w) e^{-\sigma d 2^{w/d}} \quad (4.16)$$

with

$$C_1(\sigma, d, w) \leq \frac{1}{2} (4C)^d \left\{ \left( \frac{1}{\sigma \log^2(2)} + \frac{1}{\log(2)\sqrt{2\sigma}} \right) + (\text{int}\{w/d\} + 2) \left( 1 + \frac{1}{\log(2)} \sqrt{\frac{\pi}{2\sigma}} \right) \right\}^d.$$

Now let  $\delta > 0$  and use the inequality  $x + 1 \leq e^x$ ,  $x \geq 0$ , to bound

$$\begin{aligned} C_1(\sigma, d, w) &\leq \frac{1}{2} \left( \frac{4C}{\delta\sigma} \right)^d \exp \left( d\delta\sigma \left\{ \left( \frac{1}{\sigma \log^2(2)} + \frac{1}{\log(2)\sqrt{2\sigma}} \right) + (\text{int}\{w/d\} + 2) \left( 1 + \frac{1}{\log(2)} \sqrt{\frac{\pi}{2\sigma}} \right) \right\} - d \right) \\ &\leq \frac{C_1(\sigma, \delta)^d}{2} \exp \left( \delta\sigma \left( 1 + \frac{1}{\log(2)} \sqrt{\frac{\pi}{2\sigma}} \right) w \right). \end{aligned} \quad (4.17)$$

with  $C_1(\sigma, \delta) := \frac{4C}{e\delta\sigma} \exp \left( \delta\sigma \left\{ \frac{1}{\sigma \log^2(2)} + \frac{1}{\log(2)\sqrt{2\sigma}} + 2 \left( 1 + \frac{1}{\log(2)} \sqrt{\frac{\pi}{2\sigma}} \right) \right\} \right)$  defined as in (4.12).

Estimate (4.10) follows from (4.16) and (4.17). The proof is now complete.  $\square$

**Remark 4.5 (Alternative estimate)** Observe that an alternative upper bound for  $T_1$  in (4.15) is

$$T_1 \leq \exp \left( \sigma \log^2(2) \left( 1 + \frac{w}{d} \right) \right) \frac{(2+w)^2}{2} \quad (4.18)$$

which remains bounded as  $\sigma \rightarrow 0$ . This does not happen with the bound  $C_1(\sigma, d)$ , cf. (4.12), which blows up as  $\sigma \rightarrow 0$ . As an implication of (4.18), we have

$$R(w, d) \leq \frac{(C(2+w)^2)^d}{2} e^{-\sigma d(2^{w/d} - \frac{w}{d} \log^2(2))},$$

which is an alternative to the estimate (4.10) that has an extra polynomial growth in  $w$  but remains bounded for small values of  $\sigma$ .

**Theorem 4.6** For functions  $u \in C^0(\Gamma^N; W(D))$  satisfying the assumption of Lemma 4.1. The isotropic Smolyak formula (3.8) based on Clenshaw Curtis abscissas satisfies:

$$\begin{aligned} &\|(I_N - \mathcal{A}(w, N))(u)\|_{\infty, N} \\ &\leq \inf_{\delta \in (0, \frac{\chi}{\sqrt{\pi}})} \hat{C}(\sigma, \delta, N) \times \begin{cases} e^{-\sigma w(e \log(2) - \delta \tilde{C}_2(\sigma))}, & \text{if } 0 \leq w \leq \frac{N}{\log(2)} \\ e^{-\sigma w(\frac{N}{w} 2^{w/N} - \delta \tilde{C}_2(\sigma))}, & \text{otherwise.} \end{cases} \end{aligned} \quad (4.19)$$

Here and function  $\hat{C}(\sigma, \delta, N) = \frac{C_1(\sigma, \delta)}{2} \frac{1 - C_1(\sigma, \delta)^N}{1 - C_1(\sigma, \delta)}$ . The values of  $\tilde{C}_2(\sigma)$  and  $C_1(\sigma, \delta)$  have been defined in (4.11) and (4.12), respectively.

**Proof.** From Lemmas 4.3 and 4.4 we obtain the following bound for the approximation error

$$\|(I_N - \mathcal{A}(w, N))(u)\|_{\infty, N} \leq \frac{1}{2} \sum_{d=1}^N C_1(\sigma, \delta)^d e^{-\sigma d(2^{w/d} - \frac{w}{d} \delta \tilde{C}_2(\sigma))},$$

with  $C_1(\sigma, \delta)$  defined in (4.12). Then,

$$\begin{aligned} \|(I_N - \mathcal{A}(w, N))(u)\|_{\infty, N} &\leq \frac{1}{2} \max_{1 \leq d \leq N} e^{-\sigma w(\frac{d}{w} 2^{w/d} - \delta \tilde{C}_2(\sigma))} \sum_{d=1}^N C_1(\sigma, \delta)^d, \\ &\leq \hat{C}(\sigma, \delta, N) e^{\sigma w \delta \tilde{C}_2(\sigma)} \max_{1 \leq d \leq N} e^{-\sigma w(\frac{d}{w} 2^{w/d})}, \end{aligned}$$

with

$$\begin{aligned} \hat{C}(\sigma, \delta, N) &:= \frac{1}{2} \sum_{d=1}^N C_1(\sigma, \delta)^d \\ &= \frac{C_1(\sigma, \delta)}{2} \frac{1 - C_1(\sigma, \delta)^N}{1 - C_1(\sigma, \delta)}. \end{aligned} \tag{4.20}$$

To finish the proof we further bound

$$\begin{aligned} \|(I_N - \mathcal{A}(w, N))(u)\|_{\infty, N} &\leq \hat{C}(\sigma, \delta, N) e^{\sigma w \delta \tilde{C}_2(\sigma)} e^{-\sigma w(\min_{1 \leq d \leq N} \frac{d}{w} 2^{w/d})} \\ &\leq \hat{C}(\sigma, \delta, N) e^{\sigma w \delta \tilde{C}_2(\sigma)} e^{-\sigma w(\min_{s \in [w/N, w]} \frac{1}{s} 2^s)} \end{aligned}$$

and observe that

$$\min_{s \in [w/N, w]} \frac{1}{s} 2^s = \begin{cases} e \log(2), & \text{if } 0 \leq w \leq \frac{N}{\log(2)} \\ \frac{N}{w} 2^{w/N}, & \text{otherwise.} \end{cases}$$

□

**Remark 4.7 (Alternative estimate)** Following Remark 4.5 we have an alternative to (4.19) in the estimate

$$\begin{aligned} &\|(I_N - \mathcal{A}(w, N))(u)\|_{\infty, N} \\ &\leq \frac{C(2+w)^2 (C(2+w)^2)^N - 1}{2 (C(2+w)^2 - 1)} \times \begin{cases} e^{-\sigma w \chi}, & \text{if } 0 \leq w \leq \frac{N}{\log(2)} \\ e^{-\sigma w(\frac{N}{w} 2^{w/N} - \log^2(2))}, & \text{otherwise.} \end{cases} \end{aligned} \tag{4.21}$$

Here we used the notation  $\chi = \log(2)(e - \log(2)) \approx 1.4037$ . The previous estimate can be used to produce estimates like those in Theorems 4.9 and 4.10. The alternative estimates have constants which do not blow up as  $\sigma \rightarrow 0$  but have the drawback of exhibiting additional multiplicative powers of  $\log(\eta)$ . A completely identical discussion applies to the estimates based on Gaussian abscissas, see Section 4.1.2, and will not be repeated there.

Now we relate the number of collocation points,  $\eta = \eta(w, N) = \#\mathcal{H}(w, N)$ , to the level  $w$  of the isotropic Smolyak formula. We state the result in the following lemma:

**Lemma 4.8** *Using the isotropic Smolyak interpolant described by (3.8) with Clenshaw-Curtis abscissas, the total number of points required at level  $w$  satisfies the following bounds:*

$$N(2^w - 1) \leq \eta \leq (2eN)^w \min\{w + 1, 2eN\}, \quad (4.22)$$

Moreover, as a direct consequence of (4.22) we get that:

$$\frac{\log(\eta)}{1 + \log(2) + \log(N)} - 1 \leq w. \quad (4.23)$$

**Proof.** By using formula (3.8) and exploiting the nested structure of the Clenshaw-Curtis abscissas the number of points  $\eta = \eta(w, N) = \#\mathcal{H}(w, N)$  can be counted in the following way:

$$\eta = \sum_{\mathbf{i} \in X(w, N)} \prod_{n=1}^N r(i_n), \text{ where } r(i) := \begin{cases} 1 & \text{if } i = 1 \\ 2 & \text{if } i = 2 \\ 2^{i-2} & \text{if } i > 2 \end{cases}. \quad (4.24)$$

Now notice that for all  $n = 1, 2, \dots, N$  the following bound holds:

$$2^{i_n-2} \leq r(i_n) \leq 2^{i_n-1}. \quad (4.25)$$

We now produce a lower bound and an upper bound for  $\eta$ .

A lower bound on the number  $\eta$  of points can be obtained considering only the contribution from certain tensor grids. Indeed, for a fixed value of  $\tilde{w} = 1, \dots, w$ , let us consider the  $N$  grids with indices  $i_n = 1$ , for  $n \neq m$  and  $i_m = \tilde{w} + 1$ ,  $m = 1, \dots, N$ . Since each of those  $N$  grids has  $2^{\tilde{w}-1}$  points, we have

$$\eta \geq \sum_{\tilde{w}=1}^w N 2^{\tilde{w}-1} = N(2^w - 1).$$

On the other hand, to produce an upper bound for  $\eta$ , we recall that  $|\mathbf{i} - \mathbf{1}| = \sum_{n=1}^N (i_n - 1) \leq w$  so the following bounds hold:

$$\begin{aligned} \eta &= \sum_{\mathbf{i} \in X(w, N)} \prod_{n=1}^N r(i_n) \leq \sum_{\mathbf{i} \in X(w, N)} 2^{|\mathbf{i}-\mathbf{1}|} \leq \sum_{j=0}^w \sum_{|\mathbf{i}-\mathbf{1}|=j} 2^j \\ &\leq \sum_{j=0}^w 2^j \binom{N-1+j}{N-1} \leq \sum_{j=0}^w 2^j \prod_{s=1}^{N-1} \left(1 + \frac{j}{s}\right) \\ &\leq \sum_{j=0}^w 2^j \exp\left(\sum_{s=1}^{N-1} \frac{j}{s}\right) \leq \sum_{j=0}^w 2^j \exp(j(1 + \log(N))) \\ &\leq \sum_{j=0}^w (2eN)^j \leq \min\{(w+1)(2eN)^w, (2eN)^{w+1}\} \end{aligned}$$

and this finishes the proof.  $\square$

The next Theorem provides an error bound in terms of the total number  $\eta$  of collocation points. The proof follows directly from the results in Theorem 4.6 (taking  $\delta = (e \log(2) - 1)/\tilde{C}_2(\sigma)$ ) and Lemma 4.8; it is therefore omitted.

**Theorem 4.9 (algebraic convergence)** *For functions  $u \in C^0(\Gamma^N; W(D))$  satisfying the assumption of Lemma 4.1 the isotropic Smolyak formula (3.8) based on Clenshaw Curtis abscissas satisfies:*

$$\|(I_N - \mathcal{A}(w, N))(u)\|_{\infty, N} \leq \frac{C_1(\sigma, \delta^*)e^\sigma}{|1 - C_1(\sigma, \delta^*)|} \max\{1, C_1(\sigma, \delta^*)\}^N \eta^{-\mu_1}, \quad (4.26)$$

$$\text{with } \mu_1 = \frac{\sigma}{1 + \log(2N)}.$$

Here  $\delta^* = (e \log(2) - 1)/\tilde{C}_2(\sigma)$  and the constants  $\tilde{C}_2(\sigma)$  and  $C_1(\sigma, \delta^*)$ , defined in (4.11) and (4.12), respectively, do not depend on  $\eta$ .

Observe that the previous result indicates at least *algebraic convergence* with respect to the number of collocation points  $\eta$ . Under the same assumptions of the previous theorem and with a completely similar derivation, for large values of  $w$  we have the following sharper estimate:

**Theorem 4.10 (Subexponential convergence)** *Under the same assumptions of theorem 4.9 and for  $w > \frac{N}{\log(2)}$  it holds*

$$\|(I_N - \mathcal{A}(w, N))(u)\|_{\infty, N} \leq \frac{C_1(\sigma, \delta^*) \max\{1, C_1(\sigma, \delta^*)\}^N}{e^{\sigma \delta^* \tilde{C}_2(\sigma)} |1 - C_1(\sigma, \delta^*)|} \eta^{\mu_3} e^{-\frac{N\sigma}{2^{1/N}} \eta^{\mu_2}}, \quad (4.27)$$

$$\text{with } \mu_2 = \frac{\log(2)}{N(1 + \log(2N))} \quad \text{and} \quad \mu_3 = \frac{\sigma \delta^* \tilde{C}_2(\sigma)}{1 + \log(2N)}.$$

with constant  $C_1(\sigma, \delta^*)$  defined in (4.12) and independent of  $\eta$ .

**Proof.** We start from the result stated in Theorem 4.6 and observe that for  $w > N/\log(2)$  the function

$$g(w) = \sigma(N2^{w/N} - w\delta\tilde{C}_2(\sigma))$$

is increasing in  $w$  for all values of  $\delta < \log(2)e/\tilde{C}_2(\sigma)$ . Hence, combining (4.19) with the lower bound (4.23) we obtain the desired result.  $\square$

The previous theorem indicates at least asymptotic *subexponential convergence* with respect to the number of collocation points  $\eta$ . It should be pointed out however that the large values of  $w > N/\log(2)$  under which the bound holds are seldom used in practical computations. Therefore, from this point of view estimate (4.27) is less useful than (4.26).

**Remark 4.11 (Deterioration with respect to the dimension  $N$ )** *Depending on the distance to the singularities of the solution, related to the parameter  $\tau$  introduced in Lemma 4.1, the constant  $C_1(\sigma, \delta^*)$  may be less than 1. In such a case the only dependence of the error bounds for  $\|(I_N - \mathcal{A}(w, N))(u)\|_{\infty, N}$  is in the exponent, whose denominator slowly grows like  $\log(2N)$ .*

**Remark 4.12 (Full tensor versus Smolyak)** *An isotropic full tensor product interpolation converges roughly like  $C(\sigma, N) \exp(-\sigma p)$ , where  $p$  is the order of the polynomial space. Since the number of collocation points relates to  $p$  in this case as  $\eta = (1 + p)^N$  then  $\log(\eta) = N \log(1 + p) \leq Np$  and with respect to  $\eta$  the convergence rate can be bounded as  $C(\sigma, N)\eta^{-\sigma/N}$ . The slowdown effect that the dimension  $N$  has on the last convergence is known as the curse of dimensionality and it is the reason for not using isotropic full tensor interpolation for large values of  $N$ . On the other hand, the isotropic Smolyak approximation seems to be better suited for this case. Indeed, from the estimate (4.26) we see that the Smolyak algebraic convergence has the faster exponent  $\mathcal{O}(\frac{\sigma}{\log(2N)})$ . This is a clear advantage of the isotropic Smolyak method with respect to the full tensor and justifies our claim that the use of Smolyak approximation greatly reduces the curse of dimensionality. In Section 6 numerical results will give computational ground to this claim.*

**Remark 4.13 (Estimates based on bounded mixed derivatives)** *We can proceed in a similar way to analyze the approximation error for functions that have a bounded mixed derivative of order  $(k, \dots, k)$ . In that case, the one dimensional best approximation error is  $\|u - \mathcal{U}^i(u)\| \leq C m_i^{-k} (1 + \Lambda_{m_i})$ , with  $C$  depending on  $u$  and  $k$  but not on  $m_i$ , and using again the recursion (4.9) yields*

$$\|(I_N - \mathcal{A}(w, N))(u)\|_{\infty, N} \leq \frac{C}{|C(1 + 2^k) - 1|} (C(1 + 2^k))^N (w + 1)^{2N} 2^{-kw}. \quad (4.28)$$

Finally, the combination of (4.28) with the counting estimates in Lemma 4.8 yields

$$\|(I_N - \mathcal{A}(w, N))(u)\|_{\infty, N} \leq \frac{(C(1 + 2^k))^N}{|1 + 2^k - 1/C|} \left(1 + \log_2 \left(1 + \frac{\eta}{N}\right)\right)^{2N} \times \min \left\{ 2^k \eta^{-\frac{k \log(2)}{1 + \log(2N)}}, \eta^{-k} \left(1 + \log_2 \left(1 + \frac{\eta}{N}\right)\right)^{Nk} \right\}.$$

This estimate improves the one derived in [6]. Analogous results can be derived for gaussian abscissas and  $L^2$  norms.

#### 4.1.2 Approximation estimates for Gaussian abscissas

Similarly to the previous section, we now develop error estimates for Smolyak approximation, using Gaussian abscissas cf. Section 3.3, of  $C^0(\Gamma^N; W(D))$  analytic functions described by Assumption 2.8. As before, we remind the reader that in the global estimate (4.1) we need to bound the approximation error (III) in the norm  $L^2_\rho(\Gamma^N; W(D))$ . Yet, the Gaussian abscissas defined in Section 3.3 are constructed for the auxiliary density  $\hat{\rho} = \prod_{n=1}^N \hat{\rho}_n$ , still yielding control of the desired norm

$$\|v\|_{L^2_\rho(\Gamma^N; W(D))} \leq \left\| \frac{\rho}{\hat{\rho}} \right\|_{L^\infty(\Gamma^N)}^{1/2} \|v\|_{L^2_{\hat{\rho}}(\Gamma^N; W(D))}, \quad \text{for all } v \in C^0(\Gamma^N; W(D)).$$

In what follows we will use the shorthand notation  $\|\cdot\|_{\hat{\rho}, N}$  for  $\|\cdot\|_{L^2_{\hat{\rho}}(\Gamma^N; W(D))}$ . We now quote a useful result from Erdős and Turán [14]:



**Lemma 4.14** *For every function  $u \in C^0(\Gamma^1; W(D))$  the interpolation error with Lagrange polynomials based on gaussian abscissas satisfies*

$$\|u - \mathcal{W}^i(u)\|_{\hat{\rho},1} \leq 2 \sqrt{\int_{\Gamma^1} \hat{\rho}(y) dy} \inf_{w \in V_{m_i}} \|u - w\|_{\infty,1}. \quad (4.29)$$

Similarly to Section 4.1.1, the combination of (4.5) with (4.29) yields

$$\begin{aligned} \|(I_1 - \mathcal{W}^i)(u)\|_{\hat{\rho},1} &\leq \tilde{C} e^{-\sigma 2^i}, \\ \|(\Delta^i)(u)\|_{\hat{\rho},1} &= \|(\mathcal{W}^i - \mathcal{W}^{i-1})(u)\|_{\hat{\rho},1} \\ &\leq \|(I_1 - \mathcal{W}^i)(u)\|_{\hat{\rho},1} + \|(I_1 - \mathcal{W}^{i-1})(u)\|_{\hat{\rho},1} \\ &\leq 2\tilde{C} e^{-\sigma 2^{i-1}} \end{aligned}$$

for all  $i \in \mathbb{N}_+$  with positive constants  $\tilde{C} = \hat{C} \sqrt{(\int_{\Gamma^1} \hat{\rho}(y) dy)}$  and  $\sigma = \hat{\sigma}/2$  depending on  $u$  but not on  $i$ . We then present the following Lemma and theorem whose proofs follow, with minor changes, those given in Lemma 4.4 and Theorem 4.9, respectively. For instance, we apply (A.1) from Corollary A.4 to bound the corresponding  $T_1$  sum in the estimate of  $R(w, d)$ .

**Lemma 4.15** *For functions  $u \in C^0(\Gamma^N; W(D))$  satisfying the assumption of Lemma 4.1. The isotropic Smolyak formula (3.8) based on Gaussian abscissas satisfies:*

$$\begin{aligned} &\|(I_N - \mathcal{A}(w, N))(u)\|_{\rho,N} \\ &\leq \sqrt{\|\rho/\hat{\rho}\|_{L^\infty(\Gamma^N)}} \frac{\tilde{C}_1(\sigma)}{2} \frac{1 - \tilde{C}_1(\sigma)^N}{1 - \tilde{C}_1(\sigma)} \times \begin{cases} e^{-\sigma e \log(2) w}, & \text{if } 0 \leq w \leq \frac{N}{\log(2)} \\ e^{-\sigma N 2^{w/N}}, & \text{otherwise.} \end{cases} \end{aligned} \quad (4.30)$$

Here we have

$$\tilde{C}_1(\sigma) := 4\tilde{C} \left( 1 + \frac{1}{\log(2)} \sqrt{\frac{\pi}{2\sigma}} \right). \quad (4.31)$$

Now we relate the number of collocation points  $\eta = \eta(w, N) = \#\mathcal{H}(w, N)$  to the level  $w$  of the Smolyak formula. We state the result in the following lemma:

**Lemma 4.16** *Using the Smolyak interpolant described by (3.9) with Gaussian abscissas, the total number of points required at level  $w$  satisfies the following bounds:*

$$N(2^w + 1) \leq \eta \leq (e 2^{1+\log_2(1.5)} N)^w \min\{(w+1), e 2^{1+\log_2(1.5)} N\} \quad (4.32)$$

which implies

$$\frac{\log(\eta)}{\zeta + \log(N)} - 1 \leq w.$$

with  $\zeta := 1 + (1 + \log_2(1.5)) \log(2) \approx 2.1$ .

**Proof.** By using formula (3.9), where we collocate using the Gaussian abscissas the number of points  $\eta = \eta(w, N) = \#\mathcal{H}(w, N)$ , can be counted in the following way:

$$\eta = \sum_{i \in Y(w, N)} \prod_{n=1}^N \tilde{r}(i_n), \text{ where } 2^{i-1} \leq \tilde{r}(i) := \begin{cases} 1 & \text{if } i = 1 \\ 2^{i-1} + 1 & \text{if } i \geq 2 \end{cases}. \quad (4.33)$$

Proceeding in a similar way as for the proof of Lemma 4.8, a lower bound on the number of points  $\eta$  can be obtained as

$$\eta \geq N \sum_{\tilde{w}=w-N+1}^w (2^{\tilde{w}} + 1) \geq N(2^w + 1).$$

On the other hand, an upper bound on  $\eta$  can be obtained following the same lines as in the proof of Lemma 4.8 and observing that  $2^{i-1} \leq \tilde{r}(i) \leq 2^{(1+\epsilon)(i-1)}$ , with  $\epsilon = \log_2(1.5) \approx 0.585$ .  $\square$

Finally, the next Theorem relates the error bound (4.30) to the number of collocation points  $\eta = \eta(w, N) = \#\mathcal{H}(w, N)$ , described by Lemma 4.16.

**Theorem 4.17 (algebraic convergence)** *For functions  $u \in C^0(\Gamma^N; W(D))$  satisfying the assumption of Lemma 4.1 the isotropic Smolyak formula (3.8) based on Gaussian abscissas satisfies:*

$$\|(I_N - \mathcal{A}(w, N))(u)\|_{\rho, N} \leq \sqrt{\|\rho/\hat{\rho}\|_{L^\infty(\Gamma^N)}} e^{\sigma e \log(2)} \tilde{C}_1(\sigma) \frac{\max\{1, \tilde{C}_1(\sigma)\}^N}{|1 - \tilde{C}_1(\sigma)|} \eta^{-\tilde{\mu}_1}, \quad (4.34)$$

$$\tilde{\mu}_1 := \frac{\sigma e \log(2)}{\zeta + \log(N)},$$

with  $\zeta := 1 + (1 + \log_2(1.5)) \log(2) \approx 2.1$ . The constant  $\tilde{C}_1(\sigma)$  was defined in (4.31).

Similarly to Section 4.1.1 and with the same assumptions of the previous theorem, for large values of  $w$  we have the following sharper estimate:

**Theorem 4.18 (subexponential convergence)** *If  $w > \frac{N}{\log(2)}$  then*

$$\|(I_N - \mathcal{A}(w, N))(u)\|_{\rho, N} \leq \sqrt{\|\rho/\hat{\rho}\|_{L^\infty(\Gamma^N)}} \tilde{C}_1(\sigma) \frac{\max\{1, \tilde{C}_1(\sigma)\}^N}{|1 - \tilde{C}_1(\sigma)|} e^{-\frac{N\sigma}{2^{1/N}} \eta^{\tilde{\mu}_2}}, \quad (4.35)$$

$$\text{with } \tilde{\mu}_2 = \frac{\log(2)}{N(\zeta + \log(N))}$$

and  $\zeta := 1 + (1 + \log_2(1.5)) \log(2) \approx 2.1$ .

## 4.2 Influence of truncation errors

In this Section we consider the case where the coefficients  $a_N$  and  $f_N$  from (2.4) are suitably truncated random fields. Therefore, the truncation error  $u - u_N$  is nonzero and contributes to the total error. Such contribution should be considered as well in the error analysis. In particular, understanding the relationship of this error with the discretization error allows us to compare the efficiency of isotropic Smolyak method with other computational alternatives, for instance the Monte Carlo method.

To this end, we make the assumption that the truncation error  $u - u_N$  decays as

$$\|u - u_N\|_{L^2_P(\Omega;W(D))} \leq \zeta(N) \quad (4.36)$$

for some monotonic decreasing function  $\zeta(N)$  such that  $\zeta(N) \rightarrow 0$  as  $N \rightarrow \infty$ . For example, if one truncates the input random fields with a Karhunen-Loève expansion (see [16]), the function  $\zeta(N)$  is typically related to the decay of the eigenpairs of their covariance operators.

Now, given a desired computational accuracy to achieve,  $tol > 0$ , our aim is to choose the dimension  $N = N(tol)$  and the level  $w = w(tol)$  (or equivalently  $\eta = \eta(tol)$ , the number of collocation points) such that

$$\|u - \mathcal{A}(w, N)(u_N)\|_{L^2_P(\Omega;W(D))} \leq \zeta(N) + \|u_N - \mathcal{A}(w, N)(u_N)\|_{L^2_P(\Omega;W(D))} \approx tol.$$

More precisely, we will impose that both error contributions should be of size  $tol$ , i.e.

$$\zeta(N) \approx tol \quad (4.37)$$

and

$$\|u_N - \mathcal{A}(w, N)(u_N)\|_{L^2_P(\Omega;W(D))} \approx tol. \quad (4.38)$$

Condition (4.37) determines the dimension  $N(tol)$ , while (4.38) determines the necessary number of collocation points in the isotropic Smolyak approximation. Then, this number of collocation points is compared to the number of samples required in the standard Monte Carlo method to approximate a statistical quantity of interest with accuracy  $tol$ . The latter is  $\mathcal{O}(tol^{-2})$ .

We detail only the procedure for the choice of Clenshaw-Curtis abscissas, since the discussion for Gaussian abscissas is identical. To impose condition (4.38), we apply Theorem 4.6, with the choice  $\delta^* = (e \log(2) - 1)/\tilde{C}_2(\sigma)$  and  $\tilde{C}_2(\sigma)$  as in (4.11), yielding

$$\|u_N - \mathcal{A}(w, N)(u_N)\|_{L^2_P(\Gamma^N;W(D))} \leq \frac{C_1(\sigma, \delta^*)}{|1 - C_1(\sigma, \delta^*)|} \max\{1, C_1(\sigma, \delta^*)\}^N e^{-\sigma w}$$

where the constant  $C_1(\sigma, \delta^*)$  is defined in (4.12). Now define the constants  $C = \frac{C_1(\sigma, \delta^*)}{|1 - C_1(\sigma, \delta^*)|}$  and  $F = \max\{1, C_1(\sigma, \delta^*)\}$ . With this notation and using (4.23), we have an upper bound in terms of the number of collocation points,

$$\begin{aligned} \|u_N - \mathcal{A}(w, N)(u_N)\|_{L^2_P(\Gamma^N;W(D))} &\leq C F^N e^{-\sigma w} \\ &\leq C F^N e^\sigma \eta^{-\frac{\sigma}{1 + \log(2N)}} \approx tol. \end{aligned}$$

Then, given the value of  $N(tol)$  we can find

$$\eta(tol) \approx \left( \frac{CF^N e^\sigma}{tol} \right)^{\frac{1+\log(2N)}{\sigma}} \quad (4.39)$$

and compare with the number of samples needed to achieve accuracy  $tol$  with Monte Carlo, which is  $\eta_{MC} \approx tol^{-2}$ .

**Exponential truncation error.** Here we have  $\zeta(N) = \theta e^{-\gamma N}$ , with  $\theta$  and  $\gamma$  positive constants. Therefore the dimension depends on the required accuracy like

$$N(tol) = \frac{1}{\gamma} \log \left( \frac{\theta}{tol} \right)$$

and the number of corresponding collocation points, following (4.39), is

$$\eta(tol) \approx \left( C e^\sigma F^{\frac{1}{\gamma}} \log(\theta) \right)^{\frac{\log(2e/\gamma) + \log(\log(\theta/tol))}{\sigma}} tol^{-(1+\log(F)/\gamma) \left( \frac{\log(2e/\gamma) + \log(\log(\theta/tol))}{\sigma} \right)}.$$

From here we can see roughly that for the exponential truncation error case the isotropic Smolyak method would be more efficient than Monte Carlo only if

$$\left( 1 + \frac{\log(F)}{\gamma} \right) \frac{\log(2e/\gamma) + \log(\log(\theta/tol))}{\sigma} < 2.$$

Observe that for sufficiently stringent accuracy requirements, i.e.  $tol$  sufficiently small, the Monte Carlo method will have a better convergence rate. On the other hand, due to the very slow growth of the  $\log(\log(\theta/tol))$  term above, these values of  $tol$  may be much smaller than the ones we need in practice. Thus, the range of parameters for which the isotropic Smolyak approximation gives a better convergence rate than Monte Carlo can still be relevant in many practical problems with truncated coefficients.

Observe, moreover, that whenever the parameter  $\gamma$  is large, the behavior of the one dimensional interpolation error varies widely with respect to the different  $y$  directions. In such a case, it is likely that the isotropic Smolyak method uses too many points in the directions with fastest decay. For such a case, the isotropic Smolyak method may still be better than Monte Carlo, yet we recommend the use of an anisotropic version of the Smolyak method to obtain faster convergence. For instance, see [8, 19] where anisotropic Smolyak formulas have been proposed.

**Algebraic truncation error.** Here we have  $\zeta(N) = \theta N^{-r}$ , with  $\theta$  and  $r$  positive constants. Therefore the dimension is  $N(tol) = (tol/\theta)^{-\frac{1}{r}}$  and we have

$$\frac{tol}{F^{(\theta/tol)^{1/r}}} \approx C e^\sigma \eta^{\frac{\sigma}{\log(2e) + \frac{1}{r} \log(\theta/tol)}}.$$

After denoting  $\widehat{tol} = \frac{tol}{F^{(\theta/tol)^{1/r}}} \leq tol$ , the corresponding number of collocation points is

$$\eta(tol) \approx \left( C e^\sigma \widehat{tol} \right)^{-\frac{\log(2e) + \frac{1}{r} \log(\theta/tol)}{\sigma}}, \quad (4.40)$$

Observe that even for the case where  $F = 1$  we now have an asymptotically faster growth of  $\eta(tol)$  than in the exponential truncation case. In fact, for such a case we need to have

$$\frac{\log(2e) + \frac{1}{r} \log(\theta/tol)}{\sigma} < 2$$

for the Isotropic Smolyak method to be more efficient than Monte Carlo. If  $F > 1$  then  $\widehat{tol} < tol$  and this makes, as  $tol$  gets smaller, the comparison even more favorable to Monte Carlo, cf. (4.40).

## 5 Application to linear elliptic PDEs with random input data

In this section we apply the theory developed so far to the particular linear problem described in Example 2.1. Problem (2.2) can be written in a weak form as: find  $u \in L^2_P(\Omega; H_0^1(D))$  such that

$$\int_D E[a \nabla u \cdot \nabla v] dx = \int_D E[f v] dx \quad \forall v \in L^2_P(\Omega; H_0^1(D)). \quad (5.1)$$

A straightforward application of the Lax-Milgram theorem allows one to state the well posedness of problem (5.1) and yields

$$\|u(\omega)\|_{H_0^1(D)} \leq \frac{C_P}{a_{min}} \|f(\omega, \cdot)\|_{L^2(D)} \quad \text{a.s. and} \quad \|u\|_{L^2_P(\Omega; H_0^1(D))} \leq \frac{C_P}{a_{min}} \left( \int_D E[f^2] dx \right)^{1/2},$$

where  $C_P$  denotes the constant appearing in the Poincaré inequality:  $\|v\|_{L^2(D)} \leq C_P \|\nabla v\|_{L^2(D)}$ , for all  $v \in H_0^1(D)$ .

Once we have the input random fields described by a finite set of random variables, i.e.  $a(\omega, x) = a_N(Y_1(\omega), \dots, Y_N(\omega), x)$ , and similarly for  $f(\omega, x)$ , the “finite dimensional” version of the stochastic variational formulation (5.1) has a “deterministic” equivalent which is the following: find  $u_N \in L^2_\rho(\Gamma^N; H_0^1(D))$  such that

$$\int_{\Gamma^N} (a_N \nabla u_N, \nabla v)_{L^2(D)} \rho(y) dy = \int_{\Gamma^N} (f_N, v)_{L^2(D)} \rho(y) dy, \quad \forall v \in L^2_\rho(\Gamma^N; H_0^1(D)), \quad (5.2)$$

where  $\rho(y)$  is the joint probability density function defined by (2.6). Observe that in this work the gradient notation,  $\nabla$ , always means differentiation with respect to  $x \in D$  only, unless otherwise stated. The stochastic boundary value problem (5.1) now becomes a deterministic Dirichlet boundary value problem for an elliptic partial differential equation with an  $N$ -dimensional parameter. Then, it can be shown that problem (5.1) is equivalent to

$$\int_D a_N(y) \nabla u_N(y) \cdot \nabla \phi dx = \int_D f_N(y) \phi dx, \quad \forall \phi \in H_0^1(D), \quad \rho\text{-a.e. in } \Gamma^N. \quad (5.3)$$

For our convenience, we will suppose that the coefficient  $a_N$  and the forcing term  $f_N$  admit a smooth extension on the  $\rho$ -zero measure sets. Then, equation (5.3) can be extended a.e. in  $\Gamma^N$  with respect to the Lebesgue measure (instead of the measure  $\rho dy$ ).

It has been proved in [5] that problem (5.3) satisfies the analyticity result stated in Assumption 2.8. For instance, if we take the diffusivity coefficient as in Example 2.4 and a deterministic load the size of the analyticity region is given by

$$\tau_n = \frac{a_{min}}{4\sigma_n}. \quad (5.4)$$

On the other hand, if we take the diffusivity coefficient as a truncated expansion like in Remark 2.6, then the analyticity region  $\Sigma(\Gamma_n; \tau_n)$  is given by

$$\tau_n = \frac{1}{4\sqrt{\lambda_n} \|b_n\|_{L^\infty(D)}} \quad (5.5)$$

Observe that, in the latter case, as  $\sqrt{\lambda_n} \|b_n\|_{L^\infty(D)} \rightarrow 0$  for a regular enough covariance function (see [16]) the analyticity region increases as  $n$  increases. This fact introduces, naturally, an anisotropic behavior with respect to the “direction”  $n$ . This effect will not be exploited in the numerical methods proposed in the next sections but is the subject of ongoing research.

The finite element operator  $\pi_h$  can be introduced for this problem by projecting equation (5.3) onto the subspace  $W_h(D)$ , for each  $y \in \Gamma^N$ , i.e.  $u_h^N(y) = \pi_h u_N(y)$  satisfies

$$\int_D a_N(y) \nabla u_h^N(y) \cdot \nabla \phi_h \, dx = \int_D f_N(y) \phi_h \, dx, \quad \forall \phi_h \in W_h(D), \quad \text{for a.e. } y \in \Gamma^N. \quad (5.6)$$

Notice that the finite element functions  $u_h^N(y)$  satisfy the optimality condition (3.2), for all  $y \in \Gamma^N$ . Finally, the Smolyak formula (3.8) can be applied to  $u_h^N$  to obtain the fully discrete solution. The error estimates for the Smolyak approximation, stated in Theorems 4.9-4.10 for Clensaw-Curtis abscissas and Theorems 4.17-4.18 for Gaussian abscissas hold in this case, with parameter

$$\sigma = \frac{1}{2} \min_{n=1, \dots, N} \log \left( \frac{2\tau_n}{|\Gamma_n|} + \sqrt{1 + \frac{4\tau_n^2}{|\Gamma_n|^2}} \right). \quad (5.7)$$

## 6 Numerical Examples

This section illustrates the convergence of the sparse collocation method for the stochastic linear elliptic problem in two spatial dimensions, as described in Section 5. The computational results are in accordance with the convergence rates predicted by the theory. Actually, we observe a faster convergence than stated in Theorems 4.9 and 4.17, which hints that the current estimates may be improved.

We will also use this section to compare the convergence of the isotropic Smolyak approximation, described and analyzed in Sections 3.2 and 4.1, respectively, with other ensemble-based methods such as: the anisotropic adaptive full tensor product method described in the work [4, Section 9] and the Monte Carlo method. The problem is to solve

$$\begin{cases} -\nabla \cdot (a(\omega, \cdot) \nabla u(\omega, \cdot)) = f(\omega, \cdot) & \text{in } \Omega \times D, \\ u(\omega, \cdot) = 0 & \text{on } \Omega \times \partial D, \end{cases} \quad (6.1)$$

with  $D = [0, d]^2$  and  $d = 1$ . We consider a deterministic load  $f(\omega, x, z) = \cos(x) \sin(z)$  and construct the random diffusion coefficient  $a_N(\omega, x)$  with one-dimensional (layered) spatial dependence as

$$\log(a_N(\omega, x) - 0.5) = 1 + Y_1(\omega) \left( \frac{\sqrt{\pi}L}{2} \right)^{1/2} + \sum_{n=2}^N \zeta_n \varphi_n(x) Y_n(\omega) \quad (6.2)$$

where

$$\zeta_n := (\sqrt{\pi}L)^{1/2} \exp\left( -\frac{(\lfloor \frac{n}{2} \rfloor \pi L)^2}{8} \right), \quad \text{if } n > 1 \quad (6.3)$$

and

$$\varphi_n(x) := \begin{cases} \sin\left( \frac{\lfloor \frac{n}{2} \rfloor \pi x}{L_p} \right), & \text{if } n \text{ even,} \\ \cos\left( \frac{\lfloor \frac{n}{2} \rfloor \pi x}{L_p} \right), & \text{if } n \text{ odd.} \end{cases} \quad (6.4)$$

In this example, the random variables  $\{Y_n(\omega)\}_{n=1}^\infty$  are independent, have zero mean and unit variance, i.e.  $E[Y_n] = 0$  and  $E[Y_n Y_m] = \delta_{nm}$  for  $n, m \in \mathbb{N}_+$ , and are uniformly distributed in the interval  $[-\sqrt{3}, \sqrt{3}]$ . Consequently, the auxiliary probability density  $\hat{\rho}$  defined by (3.14) can be taken equal to the joint probability density function  $\rho$  defined by (2.6). Expression (6.2) represents the truncation of a one-dimensional random field with stationary covariance

$$\begin{aligned} & \text{cov}[\log(a_N - 0.5)](x_1, x_2) \\ &= E[(\log(a)(x_1) - E[\log(a)](x_1))((\log(a)(x_2) - E[\log(a)](x_2)))] \\ &= \exp\left( \frac{-(x_1 - x_2)^2}{L_c^2} \right). \end{aligned}$$

For  $x \in [0, d]$ , let  $L_c$  be a desired physical correlation length for the coefficient  $a$ , meaning that the random variables  $a(x)$  and  $a(y)$  become essentially uncorrelated for  $|x - y| \gg L_c$ . Then, the parameter  $L_p$  in (6.4) is  $L_p = \max\{d, 2L_c\}$  and the parameter  $L$  in (6.2) and (6.3) is  $L = L_c/L_p$ .

The rate of convergence of the isotropic Smolyak method is dictated by the decay coefficient  $\sigma$  defined by (5.7), which in this case can be bounded as

$$\sigma \geq \frac{1}{2} \log \left( 1 + \sqrt{\frac{1}{24\sqrt{\pi}L}} \right). \quad (6.5)$$



From (6.5) we notice that larger correlation lengths will have negative effects on the rate of convergence, i.e. the coefficient  $\sigma$  appearing in the estimates (4.26)-(4.27) and (4.34)-(4.35) is approaching 1 as  $L_c$  becomes large. Hence, the effect of increasing  $L_c$  is a deterioration of the rate of convergence.

Recall from Section 3.3 that the Clenshaw-Curtis abscissas are nested and therefore, in practice, we exploit this fact and construct the isotropic Smolyak interpolant using formula (3.8). Hence, the number of points  $\eta = \eta(w, N) = \#\mathcal{H}(w, N)$  can be counted as in formula (4.24). On the other hand, the Gaussian abscissas, which in this case are the roots of the Legendre polynomials, are not nested and to reduce the number of points necessary to build the isotropic Smolyak formula one utilizes the variant of (3.8), given by (3.9). Consequently, we can count the number of points  $\eta$  used by the Smolyak interpolant as in (4.33).

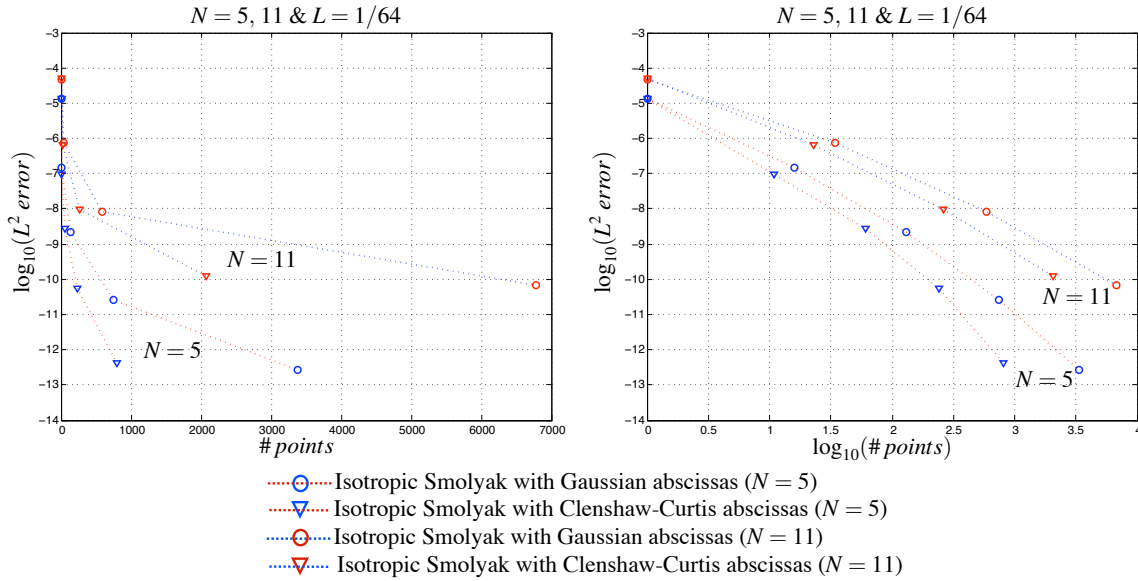
The finite element space for the spatial discretization is the span of continuous functions that are piecewise polynomials with degree two over a uniform triangulation of  $D$  with 4225 unknowns.

Observe that the collocation method only requires the solution of uncoupled deterministic problems over the set of collocation points, even in the presence of a diffusivity coefficient which depends nonlinearly on the random variables as in (6.2). This is a significant advantage that the collocation method offers compared to the classical Stochastic-Galerkin finite element method as considered, for instance, in [3, 16, 25, 36]. To study the convergence of the isotropic Smolyak approximation we consider a problem with a fixed dimension  $N$  and investigate the behavior when the level  $w$  in the Smolyak formula is increased linearly.

The computational results for the  $L^2(D)$  approximation error to the expected value,  $E[u]$ , using the isotropic Smolyak interpolant, are shown in Figure 3. Here we consider the truncated probability space to have dimensions  $N = 5$  and  $N = 11$ . To estimate the computational error in the  $w$ -th level we approximate  $\|E[\epsilon]\| \approx \|E[\mathcal{A}(w, N)\pi_h u_N - \mathcal{A}(w+1, N)\pi_h u_N]\|$ . The results reveal, as expected, that for a small non-degenerate correlation length, i.e.  $L_c = 1/64$ , the error decreases (sub)-exponentially, as the level  $w$  increases. We also observe that the convergence rate is dimension dependent and slightly deteriorates as  $N$  increases.

To investigate the performance of the isotropic Smolyak approximation by varying the correlation length  $L_c$  we also include the cases where  $L_c = 1/16$ ,  $L_c = 1/4$  and  $L_c = 1/2$  for both  $N = 5$  and  $N = 11$ , seen in Figure 4. As predicted by (6.5), we observe that the larger correlation lengths do indeed slow down the rate of convergence. Our final interest then, is to compare our isotropic sparse tensor product method with the Monte Carlo approach and also, the anisotropic full tensor product method, proposed in [4].

The anisotropic full tensor product algorithm can be described in the following way: given a tolerance  $tol$  the method computes a multi-index  $\mathbf{p} = (p_1, p_2, \dots, p_N)$ , corresponding to the order of the approximating polynomial spaces  $\mathcal{P}_{\mathbf{p}}(\Gamma^N)$ . This adaptive algorithm increases the tensor polynomial degree with an anisotropic strategy: it increases the order of approximation in one direction as much as possible before considering the next direction. Table 1 and Table 2 show the values of components of the 11-dimensional multi-index  $\mathbf{p}$  for different values of  $tol$ , corresponding to  $L_c = 1/2$  and  $L_c = 1/64$  respectively. These tables also give insight into the anisotropic behavior of each particular problem. Notice,



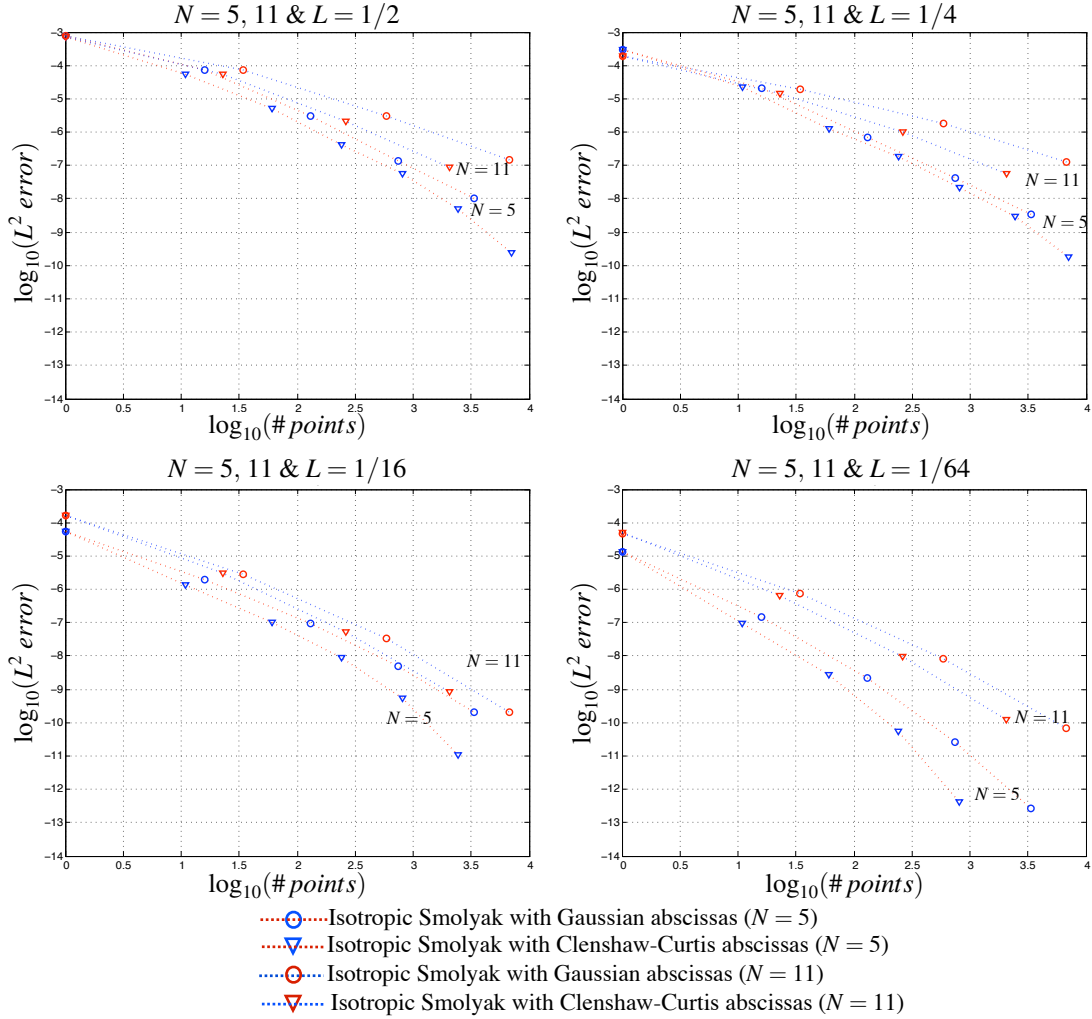
**Figure 3.** The rate of convergence of the isotropic Smolyak approximation for solving problem (6.1) with correlation length  $L_c = 1/64$  using both the Gaussian and Clenshaw-Curtis abscissas. For a finite dimensional probability space  $\Gamma^N$  with  $N = 5$  and  $N = 11$  we plot the  $L^2(D)$  approximation error in the expected value in the log-linear scale (left) and log-log scale (right).

in particular, that for the case  $L_c = 1/64$  the algorithm predicts a multi-index  $\mathbf{p}$  which is equal in all directions, i.e. an isotropic tensor product space. A convergence plot for  $L_c = 1/2$  and  $L_c = 1/64$  can be constructed by examining each row of the Table 1 and Table 2 respectively, and plotting the number of points in the tensor product grid versus the error in expectation. We estimate the error in expectation by  $\|E[\epsilon]\| \approx \|E[u_{h,\mathbf{p}}^N - u_{h,\tilde{\mathbf{p}}}^N]\|$ , with  $\tilde{\mathbf{p}} = (p_1 + 1, p_2 + 1, \dots, p_N + 1)$ . This entails an additional computational cost, which is bounded by the factor  $\exp\left(\sum_{n=1}^N 1/p_n\right)$  times the work to compute  $E[u_{h,\mathbf{p}}^N]$ .

<i>tol</i>	$N = 1$	$N = 2, 3$	$N = 4, 5$	$N = 6, 7$	$N = 8, 9$	$N = 10, 11$
1.0e-04	1	1	1	1	1	1
1.0e-05	2	1	1	1	1	1
1.0e-06	2	2	1	1	1	1
1.0e-07	3	2	2	1	1	1
1.0e-08	4	3	2	1	1	1
1.0e-09	4	4	3	1	1	1
1.0e-10	5	5	3	2	1	1
1.0e-11	5	5	4	2	1	1
1.0e-12	5	6	4	2	1	1

**Table 1.** The  $N = 11$  components of the multi index  $\mathbf{p}$  computed by the anisotropic full tensor product algorithm when solving problem (6.1) with a correlation length  $L_c = 1/2$ .

The standard Monte Carlo Finite Element Method is a popular choice for solving stochastic problems such as (6.1) (see e.g. [4, 9, 21] and the references therein). If the



**Figure 4.** The convergence of the isotropic Smolyak approximation for solving problem (6.1) with given correlation lengths  $L_c = 1/2, 1/4, 1/16$  and  $1/64$  using both the Gaussian and Clenshaw-Curtis abscissas. For a finite dimensional probability space  $\Gamma^N$  with  $N = 5$  and  $N = 11$  we plot the  $L^2(D)$  approximation error in the expected value versus the number of collocation points.

$tol$	$N = 1$	$N = 2, 3$	$N = 4, 5$	$N = 6, 7$	$N = 8, 9$	$N = 10, 11$
1.0e-03	1	1	1	1	1	1
1.0e-06	2	2	2	2	2	2
1.0e-09	3	3	3	3	3	3
1.0e-12	4	4	4	4	4	4

**Table 2.** The  $N = 11$  components of the multi index  $\mathbf{p}$  computed by the anisotropic full tensor product algorithm when solving problem (6.1) with a correlation length  $L_c = 1/64$ .

aim is to compute a functional of the solution such as the expected value, one would approximate  $E[u]$  numerically by sample averages of iid realizations of the stochastic input data. Given a number of realizations,  $M \in \mathbb{N}_+$ , we compute the sample average as follows: For each  $k = 1, \dots, M$ , sample iid realizations of  $a(\omega_k, \cdot)$  and  $f(\omega_k, \cdot)$ , solve problem

(6.1) and construct finite element approximations  $u_h^N(\omega_k, \cdot)$ . We note that once we have fixed  $\omega = \omega_k$ , the problem is completely deterministic, and may be solved by standard methods as in the collocation approach. Finally, approximate  $E[u]$  by the sample average:  $\bar{E}[u_{h,k}^N; M](\cdot) := \frac{1}{M} \sum_{k=1}^M u_h^N(\omega_k, \cdot)$ .

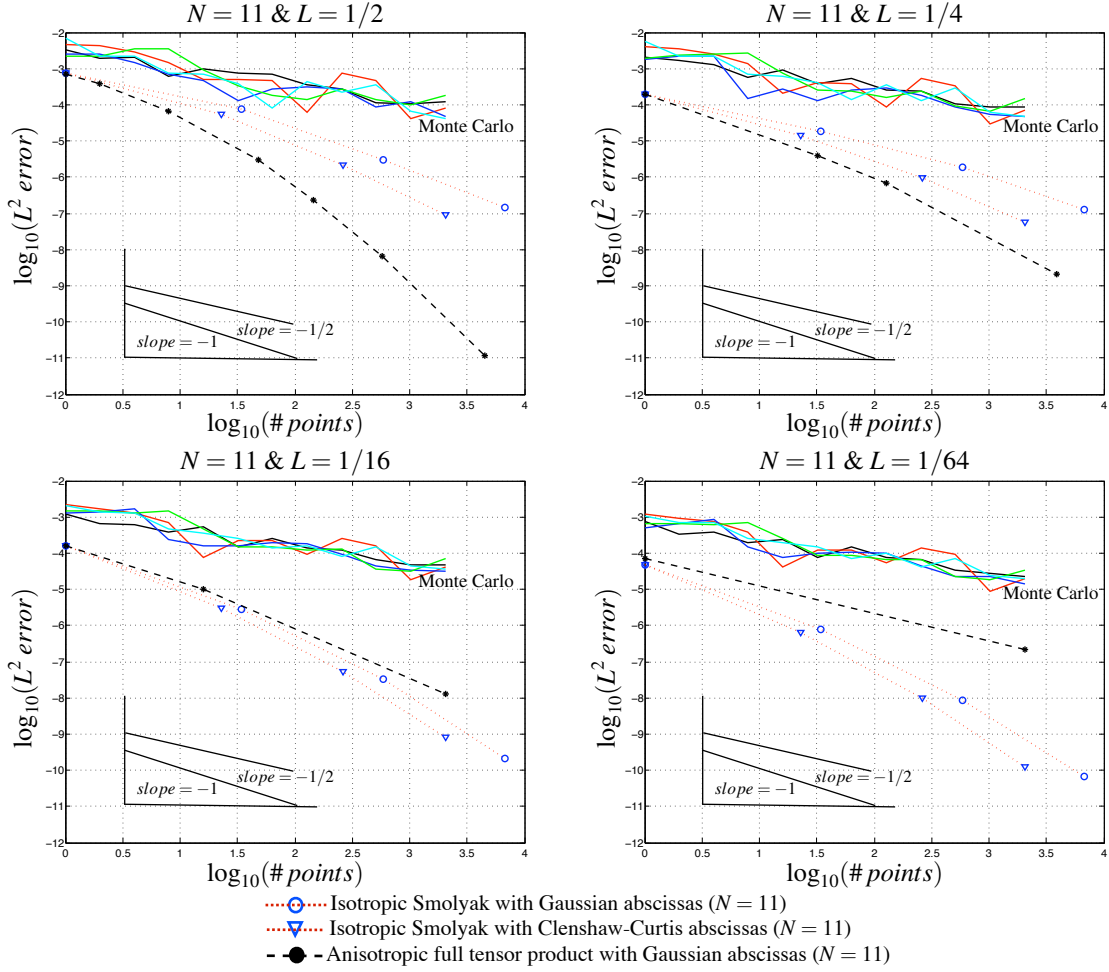
For the cases  $L_c = 1/2, 1/4, 1/16$  and  $1/64$  we take  $M = 2^i$ ,  $i = 0, 1, 2, \dots, 11$  realizations and compute the approximation to the error in expectation by  $\|E[\epsilon]\| \approx \|\bar{E}[u_{h,k}^N; M] - E[\mathcal{A}(\bar{w} + 1, N)\pi_h u_N]\|$ , where  $w = 0, 1, \dots, \bar{w}$ , with  $\bar{w} = 4$ , so that  $\mathcal{A}(5, N)$  is a highly enriched Clenshaw-Curtis isotropic sparse solution.

To study the advantages of utilizing an isotropic sparse tensor product space as opposed to an anisotropic full tensor product space we show, in Figure 4, the convergence of these methods when solving problem (6.1), using correlation lengths  $L_c = 1/2, 1/4, 1/16$  and  $L_c = 1/64$  with  $N = 11$ . We also include 5 ensembles of the Monte Carlo method described previously. Figure 4 reveals that for the isotropic case with  $L_c = 1/64$  the isotropic Smolyak method obtains a faster rate of convergence than the anisotropic full tensor product method. This is due to a slower decay of the eigenvalues expansion (6.2) and hence, an almost equal weighing of all  $N = 11$  random variables. On the contrary, opposite behavior can be observed for  $L_c = 1/2$ . Since, in this case, the rate of decay of the expansion is faster, the anisotropic full tensor method weighs heavily the important modes and, therefore, achieves a faster convergence than the isotropic Smolyak method.

In all four cases we observe that the 2 methods out-perform the Monte Carlo method. We know that the amount of work to reach the accuracy  $\epsilon$  in the Monte Carlo approach can be approximated by  $\epsilon \approx O(M^{-1/2})$  times the amount of work per sample, where  $M$  is the number of samples. This is only affected by the problem dimension through the eventual increase of the work per sample. Nevertheless, the convergence rate is quite slow and a high level of accuracy is only achieved when an large amount of function evaluations are required. This can be seen from Figure 4 where we include reference lines with slopes  $-1/2$  and  $-1$ , respectively, or in Table 3 where, for  $N = 11$ , we compare the work, proportional to the number of samples, which is the number of collocation points, required by each method to decrease the original error by a factor of  $10^4$ , for all four correlation lengths  $L_c = 1/2, 1/4, 1/16$  and  $L_c = 1/64$ .

$L_c$	AF	IS	MC
1/2	$2.5 \times 10^2$	$2.5 \times 10^3$	$5.0 \times 10^9$
1/4	$1.2 \times 10^3$	$4.0 \times 10^3$	$2.0 \times 10^9$
1/16	$2.0 \times 10^3$	$5.0 \times 10^2$	$1.6 \times 10^9$
1/64	$2.0 \times 10^5$	$3.6 \times 10^2$	$1.3 \times 10^9$

**Table 3.** For  $N = 11$ , we compare the number of function evaluations required by the Anisotropic Full Tensor product method (AF) using Gaussian abscissas, Isotropic Smolyak (IS) using Clenshaw-Curtis abscissas and the Monte Carlo (MC) method using random abscissas, to reduce the original error of problem (6.1), in expectation, by a factor of  $10^4$ .



**Figure 5.** A 11-dimensional comparison of the isotropic Smolyak method, the anisotropic full tensor product algorithm and Monte Carlo approach for solving problem (6.1) with correlation lengths  $L_c = 1/2, 1/4, 1/16$  and  $1/64$ . We plot the  $L^2(D)$  approximation error in the expected value versus the number of collocation points (or samples of the Monte Carlo method).

## 7 Conclusions

In this work we proposed and analyzed a sparse grid stochastic collocation method for solving partial differential equations whose coefficients and forcing terms depend on a finite number of random variables. The sparse grids are constructed from the Smolyak formula, utilizing either Clenshaw-Curtis or Gaussian abscissas. The method leads to the solution of uncoupled deterministic problems and, as such, it is simple to implement, allows for the use of legacy codes and is fully parallelizable like a Monte Carlo method.

This method is an improvement of the stochastic collocation method on tensor product grids proposed in [5]. The use of sparse grids considered in the present work (as opposed to full tensor grids), reduces considerably the *curse of dimensionality* and allows us to treat effectively problems that depend on a moderately large number of random variables, while keeping a high level of accuracy.

Upon assumption that the solution depends analytically on each random variable (which is a reasonable assumption for a certain class of applications, see [3, 5]), we derived strong error estimates for the fully discrete sparse grid stochastic collocation solution and analyzed its computational efficiency. In particular, the main result is the algebraic convergence with respect to the total number of collocation points, cf. Theorem 4.9 and Theorem 4.17. The exponent of such algebraic convergence depends on both the regularity of the solution and the number of input random variables,  $N$ . The exponent essentially deteriorates with  $N$  by a factor of  $1/\log(N)$ . The theory is confirmed numerically by the examples presented in Section 6. We also utilized the error estimates to compare the method with Monte Carlo in terms of computational work to achieve a given accuracy, indicating for which problems the first is more efficient than the latter. To this effect, in Section 4.2 we considered a case where the input random variables come from suitably truncated expansions of random fields and related the number of collocation points in the sparse grid to the number of random variables retained in the truncated expansion. We also developed error estimates with less regularity requirements in Remark 4.13.

The sparse grid method is very effective for problems whose input data depend on a moderate number of random variables, which “weigh equally” in the solution. For such an isotropic situation the displayed convergence is faster than standard collocation techniques built upon full tensor product spaces.

On the other hand, the convergence rate deteriorates when we attempt to solve highly anisotropic problems, such as those appearing when the input random variables come e.g. from Karhunen-Loève truncated expansions of “smooth” random fields. In such cases, a full anisotropic tensor product approximation, as proposed in [4, 5], may still be more effective for a small or moderate number of random variables.

Future directions of this research will include the development and analysis of an anisotropic version of the Sparse Grid Stochastic Collocation method, which will combine an optimal treatment of the anisotropy of the problem while reducing the *curse of dimensionality* via the use of sparse grids.

## References

- [1] I. Babuška and P. Chatzipantelidis. On solving elliptic stochastic partial differential equations. *Comput. Methods Appl. Mech. Engrg.*, 191(37-38):4093–4122, 2002. [12](#)
- [2] I. M. Babuška, K. M. Liu, and R. Tempone. Solving stochastic partial differential equations based on the experimental data. *Math. Models Methods Appl. Sci.*, 13(3):415–444, 2003. Dedicated to Jim Douglas, Jr. on the occasion of his 75th birthday. [11](#)
- [3] I. M. Babuška, R. Tempone, and G. E. Zouraris. Galerkin finite element approximations of stochastic elliptic partial differential equations. *SIAM J. Numer. Anal.*, 42(2):800–825, 2004. [12](#), [17](#), [41](#), [46](#)
- [4] I. M. Babuška, R. Tempone, and G. E. Zouraris. Solving elliptic boundary value problems with uncertain coefficients by the finite element method: the stochastic formulation. *Comput. Methods Appl. Mech. Engrg.*, 194(12-16):1251–1294, 2005. [12](#), [40](#), [41](#), [42](#), [46](#)
- [5] I.M. Babuška, F. Nobile, and R. Tempone. A stochastic collocation method for elliptic partial differential equations with random input data. *SIAM J. Numer. Anal.*, 43(3):1005–1034, 2007. [12](#), [14](#), [16](#), [17](#), [19](#), [23](#), [24](#), [38](#), [46](#)
- [6] V. Barthelmann, E. Novak, and K. Ritter. High dimensional polynomial interpolation on sparse grids. *Adv. Comput. Math.*, 12(4):273–288, 2000. Multivariate polynomial interpolation. [12](#), [19](#), [20](#), [25](#), [32](#)
- [7] S. C. Brenner and L. R. Scott. *The Mathematical Theory of Finite Element Methods*. Springer-Verlag, New York, 1994. [18](#), [23](#)
- [8] H-J. Bungartz and M. Griebel. Sparse grids. *Acta Numer.*, 13:147–269, 2004. [36](#)
- [9] J. Burkardt, M. Gunzburger, and C. Webster. Reduced order modeling of some nonlinear stochastic partial differential equations. *International Journal of Numerical Analysis and Modeling*, 4(3-4):368–391, 2007. [42](#)
- [10] P. G. Ciarlet. *The Finite Element Method for Elliptic Problems*. North-Holland, New York, 1978. [23](#)
- [11] C. W. Clenshaw and A. R. Curtis. A method for numerical integration on an automatic computer. *Numer. Math.*, 2:197–205, 1960. [12](#), [20](#)
- [12] R. A. DeVore and G. G. Lorentz. *Constructive approximation*, volume 303 of *Grundlehren der Mathematischen Wissenschaften [Fundamental Principles of Mathematical Sciences]*. Springer-Verlag, Berlin, 1993. [24](#)
- [13] V. K. Dzjadyk and V. V. Ivanov. On asymptotics and estimates for the uniform norms of the Lagrange interpolation polynomials corresponding to the Chebyshev nodal points. *Analysis Mathematica*, 9(11):85–97, 1983. [25](#)
- [14] P. Erdős and P. Turán. On interpolation. I. Quadrature- and mean-convergence in the Lagrange-interpolation. *Ann. of Math. (2)*, 38(1):142–155, 1937. [32](#)

- [15] G.S. Fishman. *Monte Carlo: Concepts, algorithms, and applications*. Springer Series in Operations Research. Springer-Verlag, New York, 1996. [11](#)
- [16] P. Frauenfelder, C. Schwab, and R. A. Todor. Finite elements for elliptic problems with stochastic coefficients. *Comput. Methods Appl. Mech. Engrg.*, 194(2-5):205–228, 2005. [11](#), [12](#), [16](#), [35](#), [39](#), [41](#)
- [17] A. Gaudagnini and S. Neumann. Nonlocal and localized analysis of conditional mean steady state flow in bounded, randomly nonuniform domains. 1. Theory and computational approach. 2. Computational examples. *Water Resources Research*, 35(10):2999–3039, 1999. [12](#)
- [18] T. Gerstner and M. Griebel. Numerical integration using sparse grids. *Numer. Algorithms*, 18(3-4):209–232, 1998. [12](#)
- [19] T. Gerstner and M. Griebel. Dimension-adaptive tensor-product quadrature. *Computing*, 71(1):65–87, 2003. [36](#)
- [20] R. G. Ghanem and P. D. Spanos. *Stochastic finite elements: a spectral approach*. Springer-Verlag, New York, 1991. [12](#)
- [21] M. Grigoriu. *Stochastic calculus*. Birkhäuser Boston Inc., Boston, MA, 2002. Applications in science and engineering. [42](#)
- [22] O. P. Le Maître, O. M. Knio, H. N. Najm, and R. G. Ghanem. Uncertainty propagation using Wiener-Haar expansions. *J. Comput. Phys.*, 197(1):28–57, 2004. [12](#)
- [23] M. Loève. *Probability theory*. Springer-Verlag, New York, fourth edition, 1977. Graduate Texts in Mathematics, Vol. 45 and 46. [11](#), [16](#)
- [24] L. Mathelin, M. Y. Hussaini, and T. A. Zang. Stochastic approaches to uncertainty quantification in CFD simulations. *Numer. Algorithms*, 38(1-3):209–236, 2005. [12](#)
- [25] H. G. Matthies and A. Keese. Galerkin methods for linear and nonlinear elliptic stochastic partial differential equations. *Comput. Methods Appl. Mech. Engrg.*, 194(12-16):1295–1331, 2005. [12](#), [41](#)
- [26] F. Riesz and B. Sz-Nagy. *Functional Analysis*. Dover, 1990. [16](#)
- [27] L.J. Roman and M. Sarkis. Stochastic galerkin method for elliptic SPDES: a white noise approach. *DCDS-B Journal*, 6(4):941–955, 2006. [12](#)
- [28] S.A. Smolyak. Quadrature and interpolation formulas for tensor products of certain classes of functions. *Dokl. Akad. Nauk SSSR*, 4:240–243, 1963. [12](#)
- [29] M.A. Tatang. *Direct incorporation of uncertainty in chemical and environmental engineering systems*. PhD thesis, MIT, 1995. [12](#)
- [30] R. A. Todor. *Sparse Perturbation Algorithms for Elliptic PDE's with Stochastic Data*. PhD thesis, Dipl. Math. University of Bucharest, 2005. [12](#), [16](#)
- [31] L. N. Trefethen. Is Gauss quadrature better than Clenshaw-Curtis? *SIAM Review*, to appear, 2006. [21](#)



- [32] G. W. Wasilkowski and H. Woźniakowski. Explicit cost bounds of algorithms for multivariate tensor product problems. *Journal of Complexity*, 11:1–56, 1995. [20](#)
- [33] N. Wiener. The homogeneous chaos. *Amer. J. Math.*, 60:897–936, 1938. [11](#)
- [34] C. L. Winter and D. M. Tartakovsky. Groundwater flow in heterogeneous composite aquifers. *Water Resour. Res.*, 38(8):23.1 (doi:10.1029/2001WR000450), 2002. [12](#)
- [35] D. Xiu and J.S. Hesthaven. High-order collocation methods for differential equations with random inputs. *SIAM J. Sci. Comput.*, 27(3):1118–1139 (electronic), 2005. [12](#)
- [36] D. Xiu and G. E. Karniadakis. Modeling uncertainty in steady state diffusion problems via generalized polynomial chaos. *Comput. Methods Appl. Mech. Engrg.*, 191(43):4927–4948, 2002. [12](#), [41](#)
- [37] D. Xiu and G. E. Karniadakis. The Wiener-Askey polynomial chaos for stochastic differential equations. *SIAM J. Sci. Comput.*, 24(2):619–644, 2002. [11](#)

## A Additional Estimates

Here we present auxiliary results that are used in Section 4. Let us recall the definition for the integer and fractional parts of a non-negative real number  $x$ , that satisfy  $x = \text{frac}\{x\} + \text{int}\{x\}$ ,  $\forall x \in \mathbb{R}^+$  with  $\text{int}\{x\}$  being the largest natural number that is smaller or equal than  $x$ .

**Lemma A.1** *Given  $w \in \mathbb{N}_+$ , for any  $\alpha > 0$  and  $0 \leq \beta < w$ , we have*

$$\sum_{i=0}^w e^{-\alpha(i-\beta)^2} \leq 2e^{-\alpha(\min\{\text{frac}\{\beta\}, 1-\text{frac}\{\beta\}\})^2} \left(1 + \frac{1}{2}\sqrt{\frac{\pi}{\alpha}}\right)$$

**Proof.** Let us write  $\beta = \text{int}\{\beta\} + \text{frac}\{\beta\}$ . Then

$$\begin{aligned} \sum_{i=0}^w e^{-\alpha(i-\beta)^2} &= \sum_{j=-\text{int}\{\beta\}}^{w-\text{int}\{\beta\}} e^{-\alpha(j-\text{frac}\{\beta\})^2} \\ &= \sum_{j=-\text{int}\{\beta\}}^0 e^{-\alpha(j-\text{frac}\{\beta\})^2} + \sum_{j=1}^{w-\text{int}\{\beta\}} e^{-\alpha(j-\text{frac}\{\beta\})^2} \\ &= \sum_{j=0}^{\text{int}\{\beta\}} e^{-\alpha(j+\text{frac}\{\beta\})^2} + \sum_{j=0}^{w-\text{int}\{\beta\}-1} e^{-\alpha(j+(1-\text{frac}\{\beta\}))^2}. \end{aligned}$$

Finally, estimate

$$\begin{aligned} \sum_{i=0}^w e^{-\alpha(i-\beta)^2} &\leq 2e^{-\alpha(\min\{\text{frac}\{\beta\}, 1-\text{frac}\{\beta\}\})^2} \sum_{j=0}^{\infty} e^{-\alpha j^2} \\ &\leq 2e^{-\alpha(\min\{\text{frac}\{\beta\}, 1-\text{frac}\{\beta\}\})^2} \left(1 + \sum_{j=1}^{\infty} e^{-\alpha j^2}\right) \\ &\leq 2e^{-\alpha(\min\{\text{frac}\{\beta\}, 1-\text{frac}\{\beta\}\})^2} \left(1 + \int_0^{\infty} e^{-\alpha x^2} dx\right) \\ &\leq 2e^{-\alpha(\min\{\text{frac}\{\beta\}, 1-\text{frac}\{\beta\}\})^2} \left(1 + \frac{1}{2}\sqrt{\frac{\pi}{\alpha}}\right). \end{aligned}$$

□

Now we state and prove an auxiliary estimate, to be used later on in the proof of Lemma A.3.

**Lemma A.2** *If  $\alpha > 0$  we have*

$$\sum_{k=1}^{\infty} k e^{-\alpha k^2} \leq \frac{1}{\alpha} + \frac{1}{\sqrt{2\alpha}}$$

**Proof.** Observe first that  $x e^{-\alpha x^2} \leq \frac{1}{\sqrt{2e\alpha}}$  for all  $x \geq 0$  and that the bound is attained at  $x^* = \frac{1}{\sqrt{2\alpha}}$ . Then, for any integer  $k_0 \geq x^*$  we can estimate

$$\sum_{k=1}^{\infty} k e^{-\alpha k^2} \leq \frac{k_0}{\sqrt{2e\alpha}} + \int_{k_0}^{+\infty} x e^{-\alpha x^2} dx \leq \frac{k_0}{\sqrt{2e\alpha}} + \frac{e^{-\alpha k_0^2}}{2\alpha}.$$

Finally, choosing  $k_0 = \text{int}\{x^*\} + 1 = \text{int}\{1/\sqrt{2\alpha}\} + 1$  the desired result follows.  $\square$

**Lemma A.3** *Given  $w \in \mathbb{N}_+$ , for any  $\alpha > 0$  and  $0 \leq \beta < w$ , we have*

$$\begin{aligned} \sum_{i=1}^w i e^{-\alpha(i-\beta)^2} &\leq 2 e^{-\alpha(1-\text{frac}\{\beta\})^2} \left( \frac{1}{\alpha} + \frac{1}{\sqrt{2\alpha}} \right) \\ &\quad + (\text{int}\{\beta\} + 1) 2 e^{-\alpha(\min\{\text{frac}\{\beta\}, 1-\text{frac}\{\beta\}\})^2} \left( 1 + \frac{1}{2} \sqrt{\frac{\pi}{\alpha}} \right) \end{aligned}$$

**Proof.** Write

$$\begin{aligned} \sum_{i=0}^w i e^{-\alpha(i-\beta)^2} &= \sum_{j=-\text{int}\{\beta\}}^{w-\text{int}\{\beta\}} (j + \text{int}\{\beta\}) e^{-\alpha(j-\text{frac}\{\beta\})^2} \\ &= \sum_{j=-\text{int}\{\beta\}}^{w-\text{int}\{\beta\}} (j-1) e^{-\alpha(j-\text{frac}\{\beta\})^2} \\ &\quad + (\text{int}\{\beta\} + 1) \sum_{j=-\text{int}\{\beta\}}^{w-\text{int}\{\beta\}} e^{-\alpha(j-\text{frac}\{\beta\})^2} \end{aligned}$$

and bound

$$\begin{aligned} \sum_{i=0}^w i e^{-\alpha(i-\beta)^2} &\leq \sum_{j=1}^{w-\text{int}\{\beta\}} (j-1) e^{-\alpha(j-\text{frac}\{\beta\})^2} \\ &\quad + (\text{int}\{\beta\} + 1) 2 e^{-\alpha(\min\{\text{frac}\{\beta\}, 1-\text{frac}\{\beta\}\})^2} \left( 1 + \frac{1}{2} \sqrt{\frac{\pi}{\alpha}} \right) \\ &\leq e^{-\alpha(1-\text{frac}\{\beta\})^2} \sum_{j=1}^{w-\text{int}\{\beta\}} (j-1) e^{-\alpha(j-1)^2} \\ &\quad + (\text{int}\{\beta\} + 1) 2 e^{-\alpha(\min\{\text{frac}\{\beta\}, 1-\text{frac}\{\beta\}\})^2} \left( 1 + \frac{1}{2} \sqrt{\frac{\pi}{\alpha}} \right) \end{aligned}$$

Finally, use the auxiliary Lemma A.2 to estimate

$$\sum_{j=0}^{w-\text{int}\{w/d\}} j e^{-\alpha j^2} \leq \sum_{j=1}^{\infty} j e^{-\alpha j^2} \leq \left( \frac{1}{\alpha} + \frac{1}{\sqrt{2\alpha}} \right)$$

$\square$

We have, as a direct consequence of Lemmas A.1 and A.3 the following estimates:

**Corollary A.4** *There holds*

$$\sum_{i=0}^w e^{-\sigma \frac{\log^2(2)}{2}(i-w/d)^2} \leq 2 \left( 1 + \frac{1}{\log(2)} \sqrt{\frac{\pi}{2\sigma}} \right) \quad (\text{A.1})$$

*and*

$$\begin{aligned} \sum_{i=0}^w (1+i) e^{-\sigma \frac{\log^2(2)}{2}(i-w/d)^2} &\leq 2 \left( \frac{1}{\sigma \log^2(2)} + \frac{1}{\log(2)\sqrt{2\sigma}} \right) \\ &\quad + 2(\text{int}\{w/d\} + 2) \left( 1 + \frac{1}{\log(2)} \sqrt{\frac{\pi}{2\sigma}} \right). \end{aligned} \quad (\text{A.2})$$

## DISTRIBUTION:

### External Distribution

- 1 Prof. I.M. Babuška  
The University of Texas at Austin  
Department of Aerospace Engineering and Engineering Mechanics  
1 University Station, C0200  
201 E. 21st Street, ACE 6.416  
Austin, TX 78712 USA
- 1 Dr. J. Burkardt  
Interdisciplinary Center for Applied Mathematics (ICAM)  
Wright House 0531  
Virginia Tech  
Blacksburg, VA, 24061 USA
- 1 Prof. Q. Du  
Department of Mathematics  
Pennsylvania State University  
University Park, PA 16802, USA
- 1 Prof. H. Elman  
Department of Computer Science  
University of Maryland  
College Park, MD 20742 USA
- 1 Prof. M. Gunzburger  
School of Computational Science  
Florida State University  
400 Dirac Science Library  
Tallahassee, FL 32306-4120 USA
- 4 Prof. R. Tempone  
School of Computational Science  
Florida State University  
400 Dirac Science Library  
Tallahassee, FL 32306-4120 USA
- 4 Prof. F. Nobile  
MOX, Dipartimento di Matematica  
Politecnico di Milano  
via Bonardi 9, 20133 Milano, ITALY
- 1 Prof. C. Trenchea  
Department of Mathematics  
301 Thackeray Hall  
University of Pittsburgh  
Pittsburgh, Pennsylvania 15260 USA

## Internal Distribution

1	MS 1320	Pavel Bochev, 01414
1	MS 1320	Scott Collis, 01416
1	MS 1318	Mike Eldred, 01411
1	MS 0346	Richard Field, 01526
1	MS 0828	Tony Giunta, 01544
1	MS 1320	Rich Lehoucq, 01416
1	MS 0321	James Peery, 01400
1	MS 1318	Eric Phipps, 01411
1	MS 0828	John Red-Horse, 01544
1	MS 1320	Denis Ridzal, 01411
1	MS 1318	James Stewart, 01411
1	MS 0370	Tim Trucano, 01411
1	MS 1318	David Womble, 01410
2	MS 9018	Central Technical Files, 8944
2	MS 0899	Technical Library, 4536



