# An adaptive ANOVA stochastic Galerkin method for partial differential equations with high-dimensional random inputs

**Guanjie Wang · Smita Sahu · Qifeng Liao**

**Abstract** It is known that standard stochastic Galerkin methods encounter challenges when solving partial differential equations with high-dimensional random inputs, which are typically caused by the large number of stochastic basis functions required. It becomes crucial to properly choose effective basis functions, such that the dimension of the stochastic approximation space can be reduced. In this work, we focus on the stochastic Galerkin approximation associated with generalized polynomial chaos (gPC), and explore the gPC expansion based on the analysis of variance (ANOVA) decomposition. A concise form of the gPC expansion is presented for each component function of the ANOVA expansion, and an adaptive ANOVA procedure is proposed to construct the overall stochastic Galerkin system. Numerical results demonstrate the efficiency of our proposed adaptive ANOVA stochastic Galerkin method for both diffusion and Helmholtz problems.

**Keywords** Adaptive ANOVA · stochastic Galerkin methods · generalized polynomial chaos · uncertainty quantification

**Mathematics Subject Classification (2020)** 35B30 · 35R60 · 65C30 · 65D40

## 1 Introduction

Over the past few decades, there has been a significant increase in efforts to develop efficient uncertainty quantification approaches for solving partial differential equations (PDEs) with random inputs. Typically, these random inputs arise from a lack of precise measurements or a limited understanding of realistic model parameters, such as permeability coefficients in diffusion problems and refraction coefficients in acoustic problems [38,11,12].

Designing a surrogate model or calculating statistics (such as mean and variance of the solution) for partial differential equations (PDEs) with random inputs is of great interest, especially

Guanjie Wang
School of Statistics and Mathematics, Shanghai Lixin University of Accounting and Finance, Shanghai, China.
E-mail: guanjie@lixin.edu.cn
Smita Sahu
School of Mathematics and Physics, University of Portsmouth, Lion Terrace, PO1 3HF, UK.
E-mail: smita.sahu@port.ac.uk
Qifeng Liao
School of Information Science and Technology, ShanghaiTech University, Shanghai, China.
E-mail: liaoqf@shanghaitech.edu.cn

when the inputs are high-dimensional. To achieve this, extensive efforts have been made. The Monte Carlo method (MCM) and its variants are among the direct methods for computing the mean and variance [5,13]. In MCM, numerous sample points of the random inputs are generated based on their probability density functions. For each sample point, the corresponding deterministic problem can be solved using existing numerical methods. The statistics of the stochastic solution can then be estimated by aggregating the results of these deterministic solutions. While MCM is easy to implement, it converges slowly and typically requires a large number of sample points. Additionally, it does not provide a surrogate model directly, which limits its applications.

To enhance efficiency, the stochastic collocation method (SCM) and the stochastic Galerkin method (SGM) have been developed. Both the SCM and the SGM are typically more efficient than the MCM for solving partial differential equations (PDEs) with moderate dimensional random inputs [37,38,39,40,32,41]. To further accelerate the SCM and the SGM, various of methods such as the reduced basis collocation method [10], the dynamically orthogonal approximation [6,7,24], the reduced basis solver based on low-rank approximation [27] and the preconditioned low-rank projection methods [19,20] are actively studied. However, these methods still face challenges in addressing high-dimensional problems, as the number of collocation points required by the SCM and the number of unknowns in the SGM increase rapidly with an increasing number of random variables, a well-established phenomenon referred to as the curse of dimensionality.

To address the challenges posed by high-dimensional problems, novel techniques have been developed and implemented. For example, the adaptive sparse grids [1], multi-element generalized polynomial chaos [33], the compressive sensing approaches [16,17], and anchored ANOVA methods (i.e. cut-HDMR) [23,42]. In particular, the anchored ANOVA method has been extensively employed in various research studies (see for instance [30,14,29,31,21]). It is shown that the choice of the anchor point is crucial for efficient approximation [28,34]. The work [30] proposes to use the covariance decomposition to effectively evaluate the output variance of multivariate functions. An efficient approximation strategy for high-dimensional periodic functions is proposed based on the fast Fourier transform and the ANOVA decomposition in [25], and its application to PDEs with random coefficients is studied in [18]. The studies conducted in [8,35] explore the adaptive reduced basis collocation method based on ANOVA decomposition and its applications to problems involving anisotropic random inputs and stochastic Stokes-Brinkman equations.

In this paper, we investigate the generalized polynomial chaos (gPC) expansion of component functions for the ANOVA decomposition, and present a concise form of the gPC expansion for each component function. With this formulation, we propose an adaptive ANOVA stochastic Galerkin method. The proposed method adaptively selects the effective gPC basis functions in the stochastic space, reducing the dimension of the stochastic approximation space significantly, and leveraging the orthonormality of the gPC basis to facilitate the computation of the variance of each term in the ANOVA decomposition. Note that compared with anchored ANOVA collocation methods [23,42,21,8], our proposed adaptive ANOVA stochastic Galerkin method avoids the difficulty for selecting anchor points, which are crucial for anchored ANOVA methods [28,34,14]. Additionally, the proposed method provides a straightforward approach to build a surrogate model. We conduct numerical simulations and present the results to demonstrate the effectiveness and efficiency of our proposed method.

An outline of the paper is as follows. We present our problem setting in the next section. In Section 3, we review the stochastic Galerkin method and the ANOVA decomposition for partial differential equations with random inputs. Our main theoretical results and the adaptive ANOVA stochastic Galerkin method are presented in Section 4. Numerical results are discussed in Section 5. Section 6 concludes the paper.

## 2 Problem setting

Let $D \subseteq \mathrm{R}^d$ $(d = 2, 3)$ denote a physical domain that is bounded, connected, with a polygonal boundary $\partial D$, and $\boldsymbol{x} \in \mathrm{R}^d$ denote a physical variable. Let $\boldsymbol{\mu} = (\mu_1, \dots, \mu_N)$ be a random vector of dimension of $N$, where the image of $\mu_i$ is denoted by $\Gamma_i$, and the probability density function of $\mu_i$ is denoted by $\rho_i(\mu_i)$. We further assume that the components of $\boldsymbol{\mu}$, i.e., $\mu_1, \dots, \mu_N$ are mutually independent, then the image of $\boldsymbol{\mu}$ is given by $\Gamma = \Gamma_1 \times \cdots \times \Gamma_N$, and the probability density function of $\boldsymbol{\mu}$ is given by $\rho(\boldsymbol{\mu}) = \prod_{i=1}^N \rho_i(\mu_i)$. In this work, we focus on the partial differential equations (PDEs) with random inputs, that is

$$
\begin{cases}
\mathfrak{L}(\boldsymbol{x}, \boldsymbol{\mu}, u(\boldsymbol{x}, \boldsymbol{\mu})) = f(\boldsymbol{x}) & \forall\, (\boldsymbol{x}, \boldsymbol{\mu}) \in D \times \Gamma, \\
\mathfrak{b}(\boldsymbol{x}, \boldsymbol{\mu}, u(\boldsymbol{x}, \boldsymbol{\mu})) = q(\boldsymbol{x}) & \forall\, (\boldsymbol{x}, \boldsymbol{\mu}) \in \partial D \times \Gamma,
\end{cases}
\tag{2.1}
$$

where $\mathfrak{L}$ is a linear partial differential operator with respect to physical variables, and $\mathfrak{b}$ is a boundary operator. Both operators can have random coefficients. The source function is denoted by $f(\boldsymbol{x})$, and $q(\boldsymbol{x})$ specifies the boundary conditions. Additionally, we assume that $\mathfrak{L}$ and $\mathfrak{b}$ are affinely dependent on the random inputs. Specifically, we have

$$
\mathfrak{L}(\boldsymbol{x}, \boldsymbol{\mu}, u(\boldsymbol{x}, \boldsymbol{\mu})) = \sum_{i=1}^K \Theta_{\mathfrak{L}}^{(i)}(\boldsymbol{\mu}) \mathfrak{L}_i(\boldsymbol{x}, u(\boldsymbol{x}, \boldsymbol{\mu})),
\tag{2.2}
$$

$$
\mathfrak{b}(\boldsymbol{x}, \boldsymbol{\mu}, u(\boldsymbol{x}, \boldsymbol{\mu})) = \sum_{i=1}^K \Theta_{\mathfrak{b}}^{(i)}(\boldsymbol{\mu}) \mathfrak{b}_i(\boldsymbol{x}, u(\boldsymbol{x}, \boldsymbol{\mu})),
\tag{2.3}
$$

where $\{\mathfrak{L}_i\}_{i=1}^K$ are parameter-independent linear differential operators, and $\{\mathfrak{b}_i\}_{i=1}^K$ are parameter-independent boundary operators. Both $\Theta_{\mathfrak{L}}^{(i)}(\boldsymbol{\mu})$ and $\Theta_{\mathfrak{b}}^{(i)}(\boldsymbol{\mu})$ take values in R for $i = 1, \dots, K$.

It is of interest to design a surrogate model for the problem (2.1) or calculate statistics of the stochastic solution $u(\boldsymbol{x}, \boldsymbol{\mu})$, such as the mean and the variance.

## 3 Stochastic Galerkin method and ANOVA decomposition

In this section, we introduce the stochastic Galerkin methods for solving problem (2.1), and we review the ANOVA decomposition for multi-variable functions. For the sake of presentation simplicity, we consider problems that satisfy homogeneous Dirichlet boundary conditions. However, it is noteworthy that the approach we present can be readily extended to other arbitrary (well-posed) boundary conditions.

3.1 Variational formulation

To introduce the variational form of (2.1), some notations are required. We first define the Hilbert spaces $L^2(D)$ and $L^2_\rho(\Gamma)$ via

$$
L^2(D) := \left\{ v(\boldsymbol{x}) : D \to \mathrm{R} \,\middle|\, \int_D v^2(\boldsymbol{x})\, \mathrm{d}\boldsymbol{x} < \infty \right\},
$$

$$
L^2_\rho(\Gamma) := \left\{ g(\boldsymbol{\mu}) : \Gamma \to \mathrm{R} \,\middle|\, \int_\Gamma \rho(\boldsymbol{\mu}) g^2(\boldsymbol{\mu})\, \mathrm{d}\boldsymbol{\mu} < \infty \right\},
$$

which are equipped with the inner products

$$\langle v(\boldsymbol{x}), \hat{v}(\boldsymbol{x})\rangle_{L^2} := \int_D v(\boldsymbol{x})\hat{v}(\boldsymbol{x})\,\mathrm{d}\boldsymbol{x},$$

$$\langle g(\boldsymbol{\mu}), \hat{g}(\boldsymbol{\mu})\rangle_{L_\rho^2} := \int_\Gamma \rho(\boldsymbol{\mu})g(\boldsymbol{\mu})\hat{g}(\boldsymbol{\mu})\,\mathrm{d}\boldsymbol{\mu}.$$

Following presentation from Babuška et al.[3], we define the tensor space of $L^2(D)$ and $L_\rho^2(\Gamma)$ as

$$L^2(D) \otimes L_\rho^2(\Gamma) := \left\{ w(\boldsymbol{x}, \boldsymbol{\mu}) \middle| w(\boldsymbol{x}, \boldsymbol{\mu}) = \sum_{i=1}^n v_i(\boldsymbol{x})g_i(\boldsymbol{\mu}), v_i(\boldsymbol{x}) \in L^2(D), g_i(\boldsymbol{x}) \in L_\rho^2(\Gamma), n \in \mathrm{N} \right\},$$

which is equipped with the inner product

$$\langle w(\boldsymbol{x}, \boldsymbol{\mu}), \hat{w}(\boldsymbol{x}, \boldsymbol{\mu})\rangle_{L^2 \otimes L_\rho^2} = \int_\Gamma \int_D w(\boldsymbol{x}, \boldsymbol{\mu})\hat{w}(\boldsymbol{x}, \boldsymbol{\mu})\rho(\boldsymbol{\mu})\mathrm{d}\boldsymbol{x}\mathrm{d}\boldsymbol{\mu}.$$

We next define the space

$$H_0^1(D) := \left\{ v \in H^1(D) \,|\, v = 0 \text{ on } \partial D \right\},$$

where $H^1(D)$ is the Sobolev space

$$H^1(D) := \left\{ v \in L^2(D)\,,\, \partial v / \partial x_i \in L^2(D), i = 1, \dots, d \right\}.$$

Furthermore, we define the solution and test function space

$$W := H_0^1(D) \otimes L_\rho^2(\Gamma) = \left\{ w(\boldsymbol{x}, \boldsymbol{\mu}) \in H_0^1(D) \otimes L_\rho^2(\Gamma) \middle| \|w(\boldsymbol{x}, \boldsymbol{\mu})\|_{L^2 \otimes L_\rho^2} < \infty \right\},$$

where $\|\cdot\|_{L^2 \otimes L_\rho^2}$ is the norm induced by the inner product $\langle\,\cdot\,,\,\cdot\,\rangle_{L^2 \otimes L_\rho^2}$. The variational form of (2.1) can be written as: find $u$ in $W = H_0^1(D) \otimes L_\rho^2(\Gamma)$ such that

$$\mathfrak{B}(u, w) = \mathfrak{F}(w), \ \forall\ w \in W, \tag{3.1}$$

where

$$\mathfrak{B}(u, w) := \langle \mathfrak{L}(\boldsymbol{x}, \boldsymbol{\mu}, u(\boldsymbol{x}, \boldsymbol{\mu})), w(\boldsymbol{x}, \boldsymbol{\mu})\rangle_{L^2 \otimes L_\rho^2}, \ \mathfrak{F}(w) := \langle f(\boldsymbol{x}), w(\boldsymbol{x}, \boldsymbol{\mu})\rangle_{L^2 \otimes L_\rho^2}.$$

Since $\mathfrak{L}$ is affinely dependent on the parameter $\boldsymbol{\mu} \in \Gamma$ (see (2.2)) then $\mathfrak{B}$ has the following form

$$\mathfrak{B}(u, w) = \sum_{i=1}^K \mathfrak{B}_i(u, w), \tag{3.2}$$

where the component bilinear forms $\mathfrak{B}_i(\cdot, \cdot)$ for $i \in \mathrm{N}^+$ are defined as

$$\mathfrak{B}_i(u, w) := \left\langle \Theta_{\mathfrak{L}}^{(i)}(\boldsymbol{\mu})\mathfrak{L}_i(\boldsymbol{x}, u(\boldsymbol{x}, \boldsymbol{\mu})), w(\boldsymbol{x}, \boldsymbol{\mu}) \right\rangle_{L^2 \otimes L_\rho^2}. \tag{3.3}$$

3.2 Discretization

A discrete version of (3.1) is obtained by introducing a finite dimensional subspace to approximate $W$. Specifically, we first denote the finite dimensional subspaces of the corresponding stochastic and physical spaces by

$$S_p = \text{span}\,\{\Phi_j(\boldsymbol{\mu})\}_{j=1}^{n_\mu} \subseteq L_\rho^2(\Gamma)\}, \ V_h = \text{span}\,\{v_s(\boldsymbol{x})\}_{s=1}^{n_x} \subseteq H_0^1(D)\},$$

where $\Phi_j(\boldsymbol{\mu})$ and $v_s(\boldsymbol{x})$ refer to basis functions. We next define a finite dimensional subspace of the overall solution (and test) function space $W$ by

$$W_h^p := V_h \otimes S_p := \text{span}\,\{v(\boldsymbol{x})\Phi(\boldsymbol{\mu})\,|\,v \in V_h, \Phi \in S_p\}.$$

The stochastic Galerkin method seeks an approximation $u^{\text{ap}}(\boldsymbol{x}, \boldsymbol{\mu}) \in W_h^p$ such that

$$\mathfrak{B}(u^{\text{ap}}, w) = \mathfrak{F}(w), \ \forall \ w \in W_h^p.$$

Suppose $u^{\text{ap}}(\boldsymbol{x}, \boldsymbol{\mu})$ is defined as

$$u^{\text{ap}}(\boldsymbol{x}, \boldsymbol{\mu}) := \sum_{s=1}^{n_x} \sum_{j=1}^{n_\mu} u_{sj}\Phi_j(\boldsymbol{\mu})v_s(\boldsymbol{x}). \tag{3.4}$$

Since $\mathfrak{L}$ is affinely dependent on the random inputs (see (2.2)), we substitute (3.4) into (3.3) to obtain

$$\mathfrak{B}_i(u^{\text{ap}}(\boldsymbol{x}, \boldsymbol{\mu}), w) = \sum_{s=1}^{n_x} \sum_{j=1}^{n_\mu} u_{sj}\Phi(\boldsymbol{\mu}) \left\langle \Theta_{\mathfrak{L}}^{(i)}(\boldsymbol{\mu})\mathfrak{L}_i v_s(\boldsymbol{x}), w(\boldsymbol{x}, \boldsymbol{\mu}) \right\rangle_{L^2 \otimes L_\rho^2}. \tag{3.5}$$

Combining (3.5) with (3.1)–(3.2), we obtain a linear system for the unknown coefficients $u_{sj}$:

$$\left( \sum_{i=1}^K \boldsymbol{G}_i \otimes \boldsymbol{A}_i \right) \bar{\boldsymbol{u}} = \boldsymbol{h} \otimes \boldsymbol{f}, \tag{3.6}$$

where $\{\boldsymbol{G}_i\}_{i=1}^K$ are matrices of size $n_\mu \times n_\mu$, and $\boldsymbol{h}$ is a column vector of length $n_\mu$. They are defined via

$$\boldsymbol{G}_i(j, k) = \langle \Theta_{\mathfrak{L}}^{(i)}(\boldsymbol{\mu})\Phi_j(\boldsymbol{\mu}), \Phi_k(\boldsymbol{\mu})\rangle_{L_\rho^2}, \ \boldsymbol{h}(i) = \langle \Phi_i(\boldsymbol{\mu}), 1\rangle_{L_\rho^2}. \tag{3.7}$$

The matrices $\boldsymbol{A}_i$ and the vector $\boldsymbol{f}$ in (3.6) are defined through

$$\boldsymbol{A}_i(s, t) = \langle \mathfrak{L}_i v_s, v_t\rangle_{L^2}, \ \boldsymbol{f}(s) = \langle f, v_s\rangle_{L^2}, \tag{3.8}$$

where $s = 1, \ldots n_x$ and $t = 1, \ldots, n_x$. The vector $\bar{\boldsymbol{u}}$ in (3.6) is a column vector of length $n_x \times n_\mu$, and is defined by

$$\bar{\boldsymbol{u}} = \begin{bmatrix} \boldsymbol{u}_1 \\ \vdots \\ \boldsymbol{u}_{n_\mu} \end{bmatrix}, \ \text{where } \boldsymbol{u}_j = \begin{bmatrix} u_{1j} \\ \vdots \\ u_{n_x j} \end{bmatrix}, \ j = 1, \ldots, n_\mu.$$

3.3 ANOVA decomposition

We define some notations before introducing the ANOVA decomposition. Let $\mathbb{T} = \{t_1, \ldots, t_{|\mathbb{T}|}\}$ be a subset of $\mathbb{U} = \{1, \ldots, N\}$, where $|\mathbb{T}|$ is the cardinality of $\mathbb{T}$. For special case where $\mathbb{T} = \emptyset$, we set $|\mathbb{T}|$ to 0. Otherwise, we assume that $t_1 < t_2 < \ldots < t_{|\mathbb{T}|}$. In addition, for $\mathbb{T} \neq \emptyset$, let $\boldsymbol{\mu}_{\mathbb{T}}$ denote the $|\mathbb{T}|$-vector that contains the components of the vector $\boldsymbol{\mu}$ indexed by $\mathbb{T}$, i.e., $\boldsymbol{\mu}_{\mathbb{T}} = (\mu_{t_1}, \ldots, \mu_{t_{|\mathbb{T}|}})$. Furthermore, let $\rho_{\mathbb{T}}(\boldsymbol{\mu}_{\mathbb{T}})$ and $\Gamma_{\mathbb{T}}$ denote the probability density function and the image corresponding to $\boldsymbol{\mu}_{\mathbb{T}}$ respectively, i.e.,

$$\rho_{\mathbb{T}}(\boldsymbol{\mu}_{\mathbb{T}}) = \rho_{t_1}(\mu_{t_1}) \cdots \rho_{t_{|\mathbb{T}|}}(\mu_{t_{|\mathbb{T}|}}), \ \Gamma_{\mathbb{T}} = \Gamma_{t_1} \times \cdots \times \Gamma_{t_{|\mathbb{T}|}}, \ |\mathbb{T}| > 0.$$

For a given cardinality $k = 0, 1, \ldots, N$, we define

$$\mathfrak{T}_k := \{\mathbb{T} | \mathbb{T} \subseteq \mathbb{U}, |\mathbb{T}| = k\}, \ \mathfrak{T}_k^* := \cup_{i=1,\ldots,k} \mathfrak{T}_i.$$

The representation of $u(\boldsymbol{x}, \boldsymbol{\mu})$ in a form

$$\begin{aligned} u(\boldsymbol{x}, \boldsymbol{\mu}) &= u_0(\boldsymbol{x}) + \sum_{\mathbb{T} \in \mathfrak{T}_N^*} u_{\mathbb{T}}(\boldsymbol{x}, \boldsymbol{\mu}_{\mathbb{T}}) \\ &= u_0(\boldsymbol{x}) + \sum_{\mathbb{T} \in \mathfrak{T}_1} u_{\mathbb{T}}(\boldsymbol{x}, \boldsymbol{\mu}_{\mathbb{T}}) + \ldots + \sum_{\mathbb{T} \in \mathfrak{T}_N} u_{\mathbb{T}}(\boldsymbol{x}, \boldsymbol{\mu}_{\mathbb{T}}), \end{aligned} \tag{3.9}$$

is a called an ANOVA decomposition if

$$u_0(\boldsymbol{x}) = \int_{\Gamma} \rho(\boldsymbol{\mu}) u(\boldsymbol{x}, \boldsymbol{\mu}) \mathrm{d}\boldsymbol{\mu}, \tag{3.10}$$

$$\int_{\Gamma_{t_k}} \rho_{t_k}(\mu_{t_k}) u_{\mathbb{T}}(\boldsymbol{x}, \boldsymbol{\mu}_{\mathbb{T}}) \mathrm{d}\mu_{t_k} = 0, \ t_k \in \mathbb{T}, \ |\mathbb{T}| > 0. \tag{3.11}$$

We call $u_{\mathbb{T}}(\boldsymbol{x}, \boldsymbol{\mu}_{\mathbb{T}})$ in (3.9) the $|\mathbb{T}|$-th order term or $|\mathbb{T}|$-th order component function, and call $u_0(\boldsymbol{x})$ the 0-th order term or 0-th order component function for special case.

In this work, we assume that the components of the random vector $\boldsymbol{\mu}$ are independent. It follows from (3.11) that the terms in (3.9) can be expressed as integrals of $u(\boldsymbol{x}, \boldsymbol{\mu})$. To illustrate this, we first show that if $\mathbb{M} \nsubseteq \mathbb{T}$, then

$$\int_{\Gamma_{\mathbb{T}^c}} \rho_{\mathbb{T}^c}(\boldsymbol{\mu}_{\mathbb{T}^c}) u_{\mathbb{M}}(\boldsymbol{x}, \boldsymbol{\mu}_{\mathbb{M}}) \mathrm{d}\boldsymbol{\mu}_{\mathbb{T}^c} = 0, \tag{3.12}$$

where $\mathbb{T}^c$ represents the complementary set of $\mathbb{T}$, i.e., $\mathbb{T}^c = \mathbb{U} \backslash \mathbb{T}$, and the universal set is given by $\mathbb{U} = \{1, \ldots, N\}$. In the rest of this paragraph, we prove (3.12). Since $\mathbb{M} \nsubseteq \mathbb{T}$, there exists an element $m_k \in \mathbb{M}$ such that $m_k \notin \mathbb{T}$; or in other words, there exists an element $m_k \in \mathbb{M}$ such that $m_k \in \mathbb{T}^c$. Letting $\mathbb{P} = \mathbb{T}^c \backslash \{m_k\}$, we then have

$$\begin{aligned} \int_{\Gamma_{\mathbb{T}^c}} \rho_{\mathbb{T}^c}(\boldsymbol{\mu}_{\mathbb{T}^c}) u_{\mathbb{M}}(\boldsymbol{x}, \boldsymbol{\mu}_{\mathbb{M}}) \mathrm{d}\boldsymbol{\mu}_{\mathbb{T}^c} &= \int_{\Gamma_{\mathbb{P}}} \int_{\Gamma_{m_k}} \rho_{m_k}(\mu_{m_k}) \rho_{\mathbb{P}}(\boldsymbol{\mu}_{\mathbb{P}}) u_{\mathbb{M}}(\boldsymbol{x}, \boldsymbol{\mu}_{\mathbb{M}}) \mathrm{d}\mu_{m_k} \mathrm{d}\boldsymbol{\mu}_{\mathbb{P}} \\ &= \int_{\Gamma_{\mathbb{P}}} \left( \rho_{\mathbb{P}}(\boldsymbol{\mu}_{\mathbb{P}}) \int_{\Gamma_{m_k}} \rho_{m_k}(\mu_{m_k}) u_{\mathbb{M}}(\boldsymbol{x}, \boldsymbol{\mu}_{\mathbb{M}}) \mathrm{d}\mu_{m_k} \right) \mathrm{d}\boldsymbol{\mu}_{\mathbb{P}}. \end{aligned}$$

According to (3.11), we have

$$\int_{\Gamma_{m_k}} \rho_{m_k}(\mu_{m_k}) u(\boldsymbol{x}, \boldsymbol{\mu}_{\mathbb{M}}) \mathrm{d}\mu_{m_k} = 0,$$

which gives (3.12).

If $\mathbb{M} \subseteq \mathbb{T}$, then $\mathbb{M} \cap \mathbb{T}^c = \emptyset$, and we have

$$\int_{\Gamma_{\mathbb{T}^c}} \rho_{\mathbb{T}^c}(\boldsymbol{\mu}_{\mathbb{T}^c}) u_{\mathbb{M}}(\boldsymbol{x}, \boldsymbol{\mu}_{\mathbb{M}}) \mathrm{d}\boldsymbol{\mu}_{\mathbb{T}^c} = u_{\mathbb{M}}(\boldsymbol{x}, \boldsymbol{\mu}_{\mathbb{M}}) \int_{\Gamma_{\mathbb{T}^c}} \rho_{\mathbb{T}^c}(\boldsymbol{\mu}_{\mathbb{T}^c}) \mathrm{d}\boldsymbol{\mu}_{\mathbb{T}^c} = u_{\mathbb{M}}(\boldsymbol{x}, \boldsymbol{\mu}_{\mathbb{M}}). \qquad (3.13)$$

By using (3.12) and (3.13), we obtain

$$\int_{\Gamma_{\mathbb{T}^c}} \rho_{\mathbb{T}^c}(\boldsymbol{\mu}_{\mathbb{T}^c}) u(\boldsymbol{x}, \boldsymbol{\mu}) \mathrm{d}\boldsymbol{\mu}_{\mathbb{T}^c} = u_{\mathbb{T}}(\boldsymbol{x}, \boldsymbol{\mu}_{\mathbb{T}}) + \sum_{\mathbb{M} \subset \mathbb{T}} u_{\mathbb{M}}(\boldsymbol{x}, \boldsymbol{\mu}_{\mathbb{M}}).$$

This formula provides a means to compute the ANOVA terms, as described [23,42]:

$$u_{\mathbb{T}}(\boldsymbol{x}, \boldsymbol{\mu}_{\mathbb{T}}) = \int_{\Gamma_{\mathbb{T}^c}} \rho_{\mathbb{T}^c}(\boldsymbol{\mu}_{\mathbb{T}^c}) u(\boldsymbol{x}, \boldsymbol{\mu}) \mathrm{d}\boldsymbol{\mu}_{\mathbb{T}^c} - \sum_{\mathbb{M} \subset \mathbb{T}} u_{\mathbb{M}}(\boldsymbol{x}, \boldsymbol{\mu}_{\mathbb{M}}).$$

An important property of the ANOVA decomposition is that all the terms in (3.9) are orthogonal, as follows from (3.11). To illustrate this, let us assume that $\mathbb{T} \neq \mathbb{M}$, which implies the existence of an element $t_k \in \mathbb{T}$ such that $t_k \notin \mathbb{M}$ (if this is not the case, then there exist an element $m_k \in \mathbb{M}$ such that $m_k \notin \mathbb{T}$, and the proof follows a similar line of reasoning). let $\mathbb{S} = \mathbb{U} \backslash \{t_k\}$ be the complementary set of $\{t_k\}$, then we have

$$\int_{\Gamma} \rho(\boldsymbol{\mu}) u(\boldsymbol{x}, \boldsymbol{\mu}_{\mathbb{T}}) u(\boldsymbol{x}, \boldsymbol{\mu}_{\mathbb{M}}) \mathrm{d}\boldsymbol{\mu} = \int_{\Gamma_{\mathbb{S}}} \int_{\Gamma_{t_k}} \rho_{t_k}(\mu_{t_k}) \rho_{\mathbb{S}}(\boldsymbol{\mu}_{\mathbb{S}}) u(\boldsymbol{x}, \boldsymbol{\mu}_{\mathbb{T}}) u(\boldsymbol{x}, \boldsymbol{\mu}_{\mathbb{M}}) \mathrm{d}\mu_{t_k} \mathrm{d}\boldsymbol{\mu}_{\mathbb{S}}$$

$$= \int_{\Gamma_{\mathbb{S}}} \left( \rho_{\mathbb{S}}(\boldsymbol{\mu}_{\mathbb{S}}) u(\boldsymbol{x}, \boldsymbol{\mu}_{\mathbb{M}}) \int_{\Gamma_{t_k}} \rho_{t_k}(\mu_{t_k}) u(\boldsymbol{x}, \boldsymbol{\mu}_{\mathbb{T}}) \mathrm{d}\mu_{t_k} \right) \mathrm{d}\boldsymbol{\mu}_{\mathbb{S}}.$$

According to (3.11), we have

$$\int_{\Gamma_{t_k}} \rho_{t_k}(\mu_{t_k}) u(\boldsymbol{x}, \boldsymbol{\mu}_{\mathbb{T}}) \mathrm{d}\mu_{t_k} = 0,$$

which implies

$$\int_{\Gamma} \rho(\boldsymbol{\mu}) u(\boldsymbol{x}, \boldsymbol{\mu}_{\mathbb{T}}) u(\boldsymbol{x}, \boldsymbol{\mu}_{\mathbb{M}}) \mathrm{d}\boldsymbol{\mu} = 0.$$

Due to the orthogonality of ANOVA terms, the variance of $u(\boldsymbol{x}, \boldsymbol{\mu})$ is the summation of the variances of all the decomposition terms:

$$\mathrm{V}[u] = \sum_{k=1}^{N} \sum_{\mathbb{T} \in \mathfrak{T}_k} \mathrm{V}[u_{\mathbb{T}}(\boldsymbol{x}, \boldsymbol{\mu}_{\mathbb{T}})], \qquad (3.14)$$

where

$$\mathrm{V}[u_{\mathbb{T}}(\boldsymbol{x}, \boldsymbol{\mu}_{\mathbb{T}})] = \int_{\Gamma} \rho(\boldsymbol{\mu}) u_{\mathbb{T}}^2(\boldsymbol{x}, \boldsymbol{\mu}_{\mathbb{T}}) \mathrm{d}\boldsymbol{\mu} = \int_{\Gamma_{\mathbb{T}}} \rho_{\mathbb{T}}(\boldsymbol{\mu}_{\mathbb{T}}) u_{\mathbb{T}}^2(\boldsymbol{x}, \boldsymbol{\mu}_{\mathbb{T}}) \mathrm{d}\boldsymbol{\mu}_{\mathbb{T}}. \qquad (3.15)$$

3.4 Adaptive ANOVA decomposition

Note that the $k$-th order term in (3.9), i.e., $\sum_{\mathbb{T} \in \mathfrak{T}_k} u_\mathbb{T}(\boldsymbol{x}, \boldsymbol{\mu}_\mathbb{T})$, has $\binom{N}{k}$ terms. For high-dimensional problems, the total number of terms in (3.9) can be prohibitively large. This motivates the development of an adaptive ANOVA expansion for such problems. The adaptive ANOVA approach is expected to be a more efficient way to approximate the exact solution since only part of low order terms in (3.9) is activated based on certain criteria [42].

To determine which terms to include in the ANOVA decomposition, we define sensitivity indices for each term as follows:

$$\mathcal{S}_\mathbb{T} = \frac{\|\mathrm{V}\,[u_\mathbb{T}]\,\|_{L^2}}{\sum_{\mathbb{T} \in \mathfrak{T}_N^*} \|\mathrm{V}\,[u_\mathbb{T}]\,\|_{L^2}},$$

where $\|\cdot\|_{L^2}$ denotes the $L^2$ function norm. It follows from equations (3.14)–(3.15) that

$$0 \le \mathcal{S}_\mathbb{T} \le 1, \;\; \sum_{\mathbb{T} \in \mathfrak{T}_N^*} \mathcal{S}_\mathbb{T} = 1.$$

The intuitive way to select the important terms in the ANOVA decomposition (3.9) is that find the terms such that $\mathcal{S}_\mathbb{T} \ge \mathtt{TOL}$, where $\mathtt{TOL}$ is a given tolerance. However, computing all possible terms is computationally expensive since $u(\boldsymbol{x}, \boldsymbol{\mu})$ is the solution of a PDE with random inputs. Instead, we construct the higher order component functions based on the lower order terms in the following way.

Let $\mathfrak{J}_k \subseteq \mathfrak{T}_k$ denote the sets of active indices for each order. Using these active indices, the solution $u(\boldsymbol{x}, \boldsymbol{\mu})$ can be approximated by

$$u(\boldsymbol{x}, \boldsymbol{\mu}) \approx u_0(\boldsymbol{x}) + \sum_{\mathbb{T} \in \mathfrak{J}_1} u_\mathbb{T}(\boldsymbol{x}, \boldsymbol{\mu}_\mathbb{T}) + \sum_{\mathbb{T} \in \mathfrak{J}_2} u_\mathbb{T}(\boldsymbol{x}, \boldsymbol{\mu}_\mathbb{T}) + \ldots.$$

For the first order terms, all the terms are retained, i.e., $\mathfrak{J}_1 = \mathfrak{T}_1$. Suppose that $\mathfrak{J}_k$ is given for $k \le N - 1$, and define the relative variance $\gamma_\mathbb{T}$ as

$$\gamma_\mathbb{T} := \frac{\|\mathrm{V}\,[u_\mathbb{T}]\,\|_{L^2}}{\sum_{\mathbb{T} \in \mathfrak{J}_k^*} \|\mathrm{V}\,[u_\mathbb{T}]\,\|_{L^2}}, \;\; \mathbb{T} \in \mathfrak{J}_k^*, \tag{3.16}$$

where $\mathfrak{J}_k^* := \mathfrak{J}_1 \cup \cdots \cup \mathfrak{J}_k$. Then, the index set of the next order can be constructed via

$$\mathfrak{J}_{k+1} := \{\mathbb{T} | \mathbb{T} \in \mathfrak{T}_{k+1}, \text{ and } \forall\, \mathbb{S} \subseteq \mathbb{T} \text{ with } |\mathbb{S}| = k \text{ satisfies } \mathbb{S} \in \tilde{\mathfrak{J}}_k\},$$

where

$$\tilde{\mathfrak{J}}_k := \{\mathbb{T} \in \mathfrak{J}_k | \gamma_\mathbb{T} \ge \mathtt{TOL}\}.$$

---

**Algorithm 1** Adaptive ANOVA decomposition [42]

---

**Input:** $u(\boldsymbol{x}, \boldsymbol{\mu})$ and $\mathtt{TOL}$, set $k = 1$ and $\mathfrak{J}_1 = \{\{1\}, \ldots, \{N\}\}$.
**while** $(k < N)$ and $\mathfrak{J}_k \ne \emptyset$ **do**
   Compute $\gamma_\mathbb{T}$ for $\mathbb{T} \in \mathfrak{J}_k$.
   Set $\tilde{\mathfrak{J}}_k := \{\mathbb{T} \in \mathfrak{J}_k | \gamma_\mathbb{T} \ge \mathtt{TOL}\}$.
   Set $\mathfrak{J}_{k+1} := \{\mathbb{T} | \mathbb{T} \in \mathfrak{T}_{k+1}, \text{ and } \forall\, \mathbb{S} \subseteq \mathbb{T} \text{ with } |\mathbb{S}| = k \text{ satisfies } \mathbb{S} \in \tilde{\mathfrak{J}}_k\}$.
   $k = k + 1$.
**end while**

---

Algorithm 1 presents the pseudo-code for the adaptive ANOVA decomposition. However, computing the relative variance $\gamma_{\mathbb{T}}$ efficiently by classical methods can be challenging, especially when dealing with the high-dimensional random inputs that arise in the context of PDEs with random inputs. To address this issue, we propose an adaptive ANOVA stochastic Galerkin method. In the following sections, we provide details of how to compute the relative variance $\gamma_{\mathbb{T}}$ in the adaptive ANOVA stochastic Galerkin method.

## 4 Adaptive ANOVA stochastic Galerkin method

In the last section, we introduced the ANOVA decomposition as a method to capture the important features of the solution. By representing the component functions of the ANOVA decomposition as the generalized polynomial chaos expansion, an effective surrogate model for the problem (2.1) can be constructed, which is essential for accelerating the solution evaluation process for time intensive problems. There are two widely used approaches for this purpose: the generalized polynomial chaos (gPC) expansion [38], and the polynomial dimensional decomposition (PDD) [31]. In this work, we apply the gPC expansion to present the component functions in (3.9). However, it is noteworthy that the PDD can also be used similarly.

### 4.1 Generalized polynomial chaos expansion of component functions

Let us commence with the definition of gPC basis functions for a single random variable. Suppose that $\mu_k$ is a random variable with probability density function $\rho_k(\mu_k)$, where $k = 1, \ldots, N$. The gPC basis functions are the orthogonal polynomials satisfying

$$\int_{\Gamma_{t_k}} \rho_k(\mu_k)\phi_i^{(k)}(\mu_k)\phi_j^{(k)}(\mu_k)\mathrm{d}\mu_k = \delta_{i,j}, \tag{4.1}$$

where $i$ and $j$ are non-negative integers, and $\delta_{i,j}$ is the Kronecker delta.

For $N$ dimensional random variables, let $\boldsymbol{i} = (i_1, \ldots, i_N) \in \mathrm{N}^N$ be a multi-index with the total degree $|\boldsymbol{i}| = i_1 + \cdots + i_N$. Note that in this work, we assume that the components of $\boldsymbol{\mu}$, i.e., $\mu_1, \ldots, \mu_N$ are mutually independent, and thus the $N$-variate gPC basis functions are the products of the univariate gPC polynomials, i.e.,

$$\Phi_{\boldsymbol{i}}(\boldsymbol{\mu}) = \phi_{i_1}^{(1)}(\mu_1)\cdots\phi_{i_N}^{(N)}(\mu_N).$$

It follows from (4.1) that

$$\int_{\Gamma} \rho(\boldsymbol{\mu})\Phi_{\boldsymbol{i}}(\boldsymbol{\mu})\Phi_{\boldsymbol{j}}(\boldsymbol{\mu})\mathrm{d}\boldsymbol{\mu} = \delta_{\boldsymbol{i},\boldsymbol{j}},$$

where $\delta_{\boldsymbol{i},\boldsymbol{j}} = \delta_{i_1,j_1}\cdots\delta_{i_N,j_N}$.

We now consider the generalized polynomial chaos expansion of the $|\mathbb{T}|$-th order component function $u_{\mathbb{T}}(\boldsymbol{x}, \boldsymbol{\mu}_{\mathbb{T}})$, which requires some notations. Let $\mathbb{T}^c = \{t_1^c, \ldots, t_{|\mathbb{T}^c|}^c\}$ be the complementary set of $\mathbb{T}$, i.e., $\mathbb{T}^c = \mathbb{U}\backslash\mathbb{T}$, where the universal set is given by $\mathbb{U} = \{1, \ldots, N\}$. For any set $\mathbb{T} \subseteq \mathbb{U}$ with $|\mathbb{T}| > 0$, the gPC basis function corresponding to the multi-index $\boldsymbol{i}_{\mathbb{T}}$ is given by

$$\Phi_{\boldsymbol{i}_{\mathbb{T}}}(\boldsymbol{\mu}_{\mathbb{T}}) = \phi_{\boldsymbol{i}(t_1)}^{(t_1)}(\mu_{t_1})\cdots\phi_{\boldsymbol{i}(t_{|\mathbb{T}|})}^{(t_{|\mathbb{T}|})}(\mu_{t_{|\mathbb{T}|}}),$$

where $\boldsymbol{i}_{\mathbb{T}}$ denote the multi-index that contains the components of the multi-index $\boldsymbol{i}$ indexed by $\mathbb{T}$, i.e., $\boldsymbol{i}_{\mathbb{T}} = (\boldsymbol{i}(t_1), \ldots, \boldsymbol{i}(t_{|\mathbb{T}_k|}))$. In additional, let $\mathfrak{M}_{\mathbb{T}}$ be the set of multi-indices defined by

$$\mathfrak{M}_{\mathbb{T}} := \{\boldsymbol{i} | \boldsymbol{i} \in \mathrm{N}^N, \boldsymbol{i}(t_1) \neq 0, \ldots, \boldsymbol{i}(t_{|\mathbb{T}|}) \neq 0, \boldsymbol{i}(t_1^c) = 0, \ldots, \boldsymbol{i}(t_{|\mathbb{T}^c|}^c) = 0\}, \ 0 < |\mathbb{T}| < N.$$

For the special case $|\mathbb{T}| = 0$, we set $\mathfrak{M}_\mathbb{T} = \{\boldsymbol{i} | \boldsymbol{i} \in \mathrm{N}^N, \boldsymbol{i}(1) = 0, \ldots, \boldsymbol{i}(N) = 0\}$, and for the special case $|\mathbb{T}| = N$, we set $\mathfrak{M}_\mathbb{T} := \{\boldsymbol{i} | \boldsymbol{i} \in \mathrm{N}^N, \boldsymbol{i}(1) \neq 0, \ldots, \boldsymbol{i}(N) \neq 0\}$. We can then state the following theorem:

**Theorem 4.1** *Given $\boldsymbol{x} \in D$ and $\mathbb{T} \subseteq \mathbb{U}$ with $|\mathbb{T}| > 0$, assuming that the $|\mathbb{T}|$-th order component function $u_\mathbb{T}(\boldsymbol{x}, \boldsymbol{\mu}_\mathbb{T})$ belongs to $L_\rho^2(\Gamma)$, then the generalized polynomial chaos expansion of $u_\mathbb{T}(\boldsymbol{x}, \boldsymbol{\mu}_\mathbb{T})$ can be expressed by*

$$u_\mathbb{T}(\boldsymbol{x}, \boldsymbol{\mu}_\mathbb{T}) = \sum_{\boldsymbol{i} \in \mathfrak{M}_\mathbb{T}} u_{\boldsymbol{i}}(\boldsymbol{x}) \Phi_{\boldsymbol{i}}(\boldsymbol{\mu}), \tag{4.2}$$

*where $u_{\boldsymbol{i}}(\boldsymbol{x})$ is the coefficient of $\Phi_{\boldsymbol{i}}(\boldsymbol{\mu})$ defined by*

$$u_{\boldsymbol{i}}(\boldsymbol{x}) = \int_\Gamma \rho(\boldsymbol{\mu}) u_\mathbb{T}(\boldsymbol{x}, \boldsymbol{\mu}_\mathbb{T}) \Phi_{\boldsymbol{i}}(\boldsymbol{\mu}) \mathrm{d}\boldsymbol{\mu}.$$

*Proof* Since $u_\mathbb{T}(\boldsymbol{x}, \boldsymbol{\mu}_\mathbb{T}) \in L_\rho^2(\Gamma)$ and $\{\Phi_{|\boldsymbol{i}|}\}_{|\boldsymbol{i}|=0}^\infty$ forms an complete orthonormal basis of $L_\rho^2(\Gamma)$, we can express the generalized polynomial chaos expansion of $u_\mathbb{T}(\boldsymbol{x}, \boldsymbol{\mu}_\mathbb{T})$ as

$$u_\mathbb{T}(\boldsymbol{x}, \boldsymbol{\mu}_\mathbb{T}) = \sum_{|\boldsymbol{i}|=0}^\infty u_{\boldsymbol{i}}(\boldsymbol{x}) \Phi_{\boldsymbol{i}}(\boldsymbol{\mu}), \tag{4.3}$$

where $u_{\boldsymbol{i}}(\boldsymbol{x})$ are the coefficients of the expansion given by

$$u_{\boldsymbol{i}}(\boldsymbol{x}) = \int_\Gamma \rho(\boldsymbol{\mu}) u_\mathbb{T}(\boldsymbol{x}, \boldsymbol{\mu}_\mathbb{T}) \Phi_{\boldsymbol{i}}(\boldsymbol{\mu}) \mathrm{d}\boldsymbol{\mu}.$$

To prove the theorem, it suffices to show that if $\boldsymbol{i} \notin \mathfrak{M}_\mathbb{T}$, which means that there exists $t_k \in \mathbb{T}$ such that $\boldsymbol{i}(t_k) = 0$ or there exists $t_k^c \in \mathbb{T}^c$ such that $\boldsymbol{i}(t_k^c) \neq 0$, then $u_{\boldsymbol{i}}(\boldsymbol{x}) = 0$.

If there exists $t_k \in \mathbb{T}$ such that $\boldsymbol{i}(t_k) = 0$, let $\mathbb{S} = \mathbb{U} \backslash \{t_k\}$ be the complementary set of $\{t_k\}$, and note that $\phi_0^{(t_k)}(\mu_{t_k}) = 1$. Then, we have

$$\int_\Gamma \rho(\boldsymbol{\mu}) u_\mathbb{T}(\boldsymbol{x}, \boldsymbol{\mu}_\mathbb{T}) \Phi_{\boldsymbol{i}}(\boldsymbol{\mu}) \mathrm{d}\boldsymbol{\mu} = \int_{\Gamma_\mathbb{S}} \int_{\Gamma_{t_k}} \rho_\mathbb{S}(\boldsymbol{\mu}_\mathbb{S}) \rho_{t_k}(\mu_{t_k}) \Phi_{\boldsymbol{i}_\mathbb{S}}(\boldsymbol{\mu}_\mathbb{S}) u_\mathbb{T}(\boldsymbol{x}, \boldsymbol{\mu}_\mathbb{T}) \mathrm{d}\mu_{t_k} \mathrm{d}\boldsymbol{\mu}_\mathbb{S}$$

$$= \int_{\Gamma_\mathbb{S}} \left( \rho_\mathbb{S}(\boldsymbol{\mu}_\mathbb{S}) \Phi_{\boldsymbol{i}_\mathbb{S}}(\boldsymbol{\mu}_\mathbb{S}) \int_{\Gamma_{t_k}} \rho_{t_k}(\mu_{t_k}) u_\mathbb{T}(\boldsymbol{x}, \boldsymbol{\mu}_\mathbb{T}) \mathrm{d}\mu_{t_k} \right) \mathrm{d}\boldsymbol{\mu}_\mathbb{S}.$$

Using (3.11), we obtain

$$\int_{\Gamma_{t_k}} \rho_{t_k}(\mu_{t_k}) u_\mathbb{T}(\boldsymbol{x}, \boldsymbol{\mu}_\mathbb{T}) \mathrm{d}\mu_{t_k} = 0,$$

which implies that

$$u_{\boldsymbol{i}}(\boldsymbol{x}) = \int_\Gamma \rho(\boldsymbol{\mu}) u_\mathbb{T}(\boldsymbol{x}, \boldsymbol{\mu}_\mathbb{T}) \Phi_{\boldsymbol{i}}(\boldsymbol{\mu}) \mathrm{d}\boldsymbol{\mu} = 0.$$

On the other hand, if there exists $t_k^c \in \mathbb{T}^c$ such that $\boldsymbol{i}(t_k^c) \neq 0$, then we have

$$\int_\Gamma \rho(\boldsymbol{\mu}) u_\mathbb{T}(\boldsymbol{x}, \boldsymbol{\mu}_\mathbb{T}) \Phi_{\boldsymbol{i}}(\boldsymbol{\mu}) \mathrm{d}\boldsymbol{\mu} = \int_{\Gamma_\mathbb{T}} \int_{\Gamma_{\mathbb{T}^c}} \rho_\mathbb{T}(\boldsymbol{\mu}_\mathbb{T}) \rho_{\mathbb{T}^c}(\boldsymbol{\mu}_{\mathbb{T}^c}) u_\mathbb{T}(\boldsymbol{x}, \boldsymbol{\mu}_\mathbb{T}) \Phi_{\boldsymbol{i}_\mathbb{T}}(\boldsymbol{\mu}_\mathbb{T}) \Phi_{\boldsymbol{i}_{\mathbb{T}^c}}(\boldsymbol{\mu}_{\mathbb{T}^c}) \mathrm{d}\boldsymbol{\mu}_\mathbb{T} \mathrm{d}\boldsymbol{\mu}_{\mathbb{T}^c}$$

$$= \int_{\Gamma_\mathbb{T}} \rho_\mathbb{T}(\boldsymbol{\mu}_\mathbb{T}) u_\mathbb{T}(\boldsymbol{x}, \boldsymbol{\mu}_\mathbb{T}) \Phi_{\boldsymbol{i}_\mathbb{T}}(\boldsymbol{\mu}_\mathbb{T}) \mathrm{d}\boldsymbol{\mu}_\mathbb{T} \int_{\Gamma_{\mathbb{T}^c}} \rho_{\mathbb{T}^c}(\boldsymbol{\mu}_{\mathbb{T}^c}) \Phi_{\boldsymbol{i}_{\mathbb{T}^c}}(\boldsymbol{\mu}_{\mathbb{T}^c}) \mathrm{d}\boldsymbol{\mu}_{\mathbb{T}^c}.$$

Note that since the gPC basis function corresponding to $(0, \ldots, 0)$ is 1, and $\boldsymbol{i}_{\mathbb{T}^c} \neq (0, \ldots, 0)$, we have

$$\int_{\Gamma_{\mathbb{T}^c}} \rho_{\mathbb{T}^c}(\boldsymbol{\mu}_{\mathbb{T}^c}) \Phi_{\boldsymbol{i}_{\mathbb{T}^c}}(\boldsymbol{\mu}_{\mathbb{T}^c}) \mathrm{d}\boldsymbol{\mu}_{\mathbb{T}^c} = 0,$$

and thus

$$u_{\boldsymbol{i}}(\boldsymbol{x}) = \int_{\Gamma} \rho(\boldsymbol{\mu}) u_{\mathbb{T}}(\boldsymbol{x}, \boldsymbol{\mu}_{\mathbb{T}}) \Phi_{\boldsymbol{i}}(\boldsymbol{\mu}) \mathrm{d}\boldsymbol{\mu} = 0.$$

$\square$

**Theorem 4.2** *Suppose that $u(\boldsymbol{x}, \boldsymbol{\mu})$ can be expressed as*

$$u(\boldsymbol{x}, \boldsymbol{\mu}) = u_0(\boldsymbol{x}) + \sum_{\mathbb{T} \in \mathfrak{T}_1} u_{\mathbb{T}}(\boldsymbol{x}, \boldsymbol{\mu}_{\mathbb{T}}) + \ldots + \sum_{\mathbb{T} \in \mathfrak{T}_N} u_{\mathbb{T}}(\boldsymbol{x}, \boldsymbol{\mu}_{\mathbb{T}}), \tag{4.4}$$

*where*

$$u_{\mathbb{T}}(\boldsymbol{x}, \boldsymbol{\mu}_{\mathbb{T}}) = \sum_{\boldsymbol{i} \in \mathfrak{M}_{\mathbb{T}}} u_{\boldsymbol{i}}(\boldsymbol{x}) \Phi_{\boldsymbol{i}}(\boldsymbol{\mu}), \ \mathbb{T} \in \mathfrak{T}_N^*. \tag{4.5}$$

*Then the right hand side of (4.4) is the ANOVA decomposition of $u(\boldsymbol{x}, \boldsymbol{\mu})$.*

*Proof* To complete the proof, we only need to show that the right hand side of (4.4) satisfies (3.10) and (3.11). Using (4.4) and (4.5), we have

$$\int_{\Gamma} \rho(\boldsymbol{\mu}) u(\boldsymbol{x}, \boldsymbol{\mu}) \mathrm{d}\boldsymbol{\mu} = \int_{\Gamma} \rho(\boldsymbol{\mu}) u_0(\boldsymbol{x}) \mathrm{d}\boldsymbol{\mu} + \sum_{\mathbb{T} \in \mathfrak{T}_N^*} \int_{\Gamma} \rho(\boldsymbol{\mu}) u_{\mathbb{T}}(\boldsymbol{x}, \boldsymbol{\mu}_{\mathbb{T}}) \mathrm{d}\boldsymbol{\mu}$$

$$= u_0(\boldsymbol{x}) + \sum_{\mathbb{T} \in \mathfrak{T}_N^*} \sum_{\boldsymbol{i} \in \mathfrak{M}_{\mathbb{T}}} u_{\boldsymbol{i}}(\boldsymbol{x}) \int_{\Gamma} \rho(\boldsymbol{\mu}) \Phi_{\boldsymbol{i}}(\boldsymbol{\mu}) \mathrm{d}\boldsymbol{\mu}.$$

Since $\mathbb{T} \neq \emptyset$ when $\mathbb{T} \in \mathfrak{T}_N^*$, we have

$$\boldsymbol{i} \neq (0, \ldots, 0), \ \boldsymbol{i} \in \mathfrak{M}_{\mathbb{T}}.$$

Thus,

$$\int_{\Gamma} \rho(\boldsymbol{\mu}) \Phi_{\boldsymbol{i}}(\boldsymbol{\mu}) \mathrm{d}\boldsymbol{\mu} = 0,$$

which implies that

$$\int_{\Gamma} \rho(\boldsymbol{\mu}) u(\boldsymbol{x}, \boldsymbol{\mu}) \mathrm{d}\boldsymbol{\mu} = u_0(\boldsymbol{x}).$$

To show that the right hand side of (4.4) satisfies (3.11), suppose that $t_k \in \mathbb{T}$, $k = 1, \ldots, t_{|\mathbb{T}|}$, and let $\mathbb{S} = \mathbb{U} \backslash \{t_k\}$ be the complementary set of $\{t_k\}$. Then, by (4.5), we have

$$\int_{\Gamma_{t_k}} \rho_{t_k}(\mu_{t_k}) u_{\mathbb{T}}(\boldsymbol{x}, \boldsymbol{\mu}_{\mathbb{T}}) \mathrm{d}\mu_{t_k} = \sum_{\boldsymbol{i} \in \mathfrak{M}_{\mathbb{T}}} u_{\boldsymbol{i}}(\boldsymbol{x}) \int_{\Gamma_{t_k}} \rho_{t_k}(\mu_{t_k}) \Phi_{\boldsymbol{i}}(\boldsymbol{\mu}) \mathrm{d}\mu_{t_k}$$

$$= \sum_{\boldsymbol{i} \in \mathfrak{M}_{\mathbb{T}}} u_{\boldsymbol{i}}(\boldsymbol{x}) \Phi_{\boldsymbol{i}_{\mathbb{S}}}(\boldsymbol{\mu}_{\mathbb{S}}) \int_{\Gamma_{t_k}} \rho_{t_k}(\mu_{t_k}) \phi_{\boldsymbol{i}(t_k)}^{(t_k)}(\mu_{t_k}) \mathrm{d}\mu_{t_k}$$

Recall that $\phi_0^{(t_k)}(\mu_{t_k}) = 1$ and $\boldsymbol{i}(t_k) \neq 0$, we have

$$\int_{\Gamma_{t_k}} \rho_{t_k}(\mu_{t_k}) \phi_{\boldsymbol{i}(t_k)}^{(t_k)}(\mu_{t_k}) \mathrm{d}\mu_{t_i} = 0,$$

which implies that

$$\int_{\Gamma_{t_k}} \rho_{t_k}(\mu_{t_k}) u_{\mathbb{T}}(\boldsymbol{x}, \boldsymbol{\mu}_{\mathbb{T}}) \mathrm{d}\mu_{t_k} = 0, \ k = 1, \ldots, t_{|\mathbb{T}|}.$$

$\square$

In practical computations, the expansion given in (4.2) must be truncated to a finite number of terms. Following the approach in [36], we retain the terms with the total degree up to $p$. This yields the approximation

$$u_{\mathbb{T}}(\boldsymbol{x}, \boldsymbol{\mu}_{\mathbb{T}}) \approx \sum_{\boldsymbol{i} \in \mathfrak{M}_{\mathbb{T}}^p} u_{\boldsymbol{i}}(\boldsymbol{x}) \Phi_{\boldsymbol{i}}(\boldsymbol{\mu}), \tag{4.6}$$

where $\mathfrak{M}_{\mathbb{T}}^p$ is the set of multi-indices defined by

$$\mathfrak{M}_{\mathbb{T}}^p := \{\boldsymbol{i} | \boldsymbol{i} \in \mathfrak{M}_{\mathbb{T}} \text{ and } |\boldsymbol{i}| \leq p\}.$$

By inserting (4.6) into (3.9), we obtain the following expansion of $u(\boldsymbol{x}, \boldsymbol{\mu})$ in terms of gPC basis functions:

$$u(\boldsymbol{x}, \boldsymbol{\mu}) \approx u_p(\boldsymbol{x}, \boldsymbol{\mu}) = u_0(\boldsymbol{x}) + \sum_{\mathbb{T} \in \mathfrak{T}_1} \sum_{\boldsymbol{i} \in \mathfrak{M}_{\mathbb{T}}^p} u_{\boldsymbol{i}}(\boldsymbol{x}) \Phi_{\boldsymbol{i}}(\boldsymbol{\mu}) + \ldots + \sum_{\mathbb{T} \in \mathfrak{T}_N} \sum_{\boldsymbol{i} \in \mathfrak{M}_{\mathbb{T}}^p} u_{\boldsymbol{i}}(\boldsymbol{x}) \Phi_{\boldsymbol{i}}(\boldsymbol{\mu}), \tag{4.7}$$

where $u_p(\boldsymbol{x}, \boldsymbol{\mu})$ is the polynomial approximation of $u(\boldsymbol{x}, \boldsymbol{\mu})$ with the total degree up to $p$. By employing (4.6), we can easily compute the variance of $u_{\mathbb{T}}(\boldsymbol{x}, \boldsymbol{\mu}_{\mathbb{T}})$ using the following expression:

$$\mathrm{V}[u_{\mathbb{T}}] \approx \sum_{\boldsymbol{i} \in \mathfrak{M}_{\mathbb{T}}^p} u_{\boldsymbol{i}}^2(\boldsymbol{x}), \ \mathbb{T} \in \mathfrak{T}_N^*. \tag{4.8}$$

Fig. 4.1 illustrates the multi-indices of the gPC basis functions corresponding to the component functions of each order in (4.7). It is worth noting that the number of terms in (4.7) is given by

$$\binom{N}{0}\binom{p}{0} + \cdots + \binom{N}{N}\binom{p}{N} = \binom{N+p}{N}, \tag{4.9}$$

which identical to the number of terms in the generalized polynomial chaos expansion with the total degree up to $p$. The equation (4.9) is commonly referred to as the Vandermonde's identity or the Vandermonde's convolution. Interested readers can find more information on this topic in the relevant literature, such as [2].

## 4.2 Adaptive ANOVA stochastic Galerkin method

Based on the adaptive ANOVA decomposition and the gPC expansion of component functions, we can develop an adaptive ANOVA stochastic Galerkin method for the problem (2.1). The idea is quite simple, namely, we select the basis functions of stochastic space based on the adaptive ANOVA decomposition, in which the relative variance is computed by (3.16) and (4.8).

To give the algorithm of this procedure, some notations are needed. We first collect the basis functions associated with the $k$-th order component function $u_{\mathbb{T}}(\boldsymbol{x}, \boldsymbol{\mu}_{\mathbb{T}})$ and denote the set of their multi-indices with the total degree up to $p$ as $\mathfrak{M}_k^p := \cup_{\mathbb{T} \in \mathfrak{J}_k} \mathfrak{M}_{\mathbb{T}}^p$. Moreover, let us denote the set of all multi-indices as $\mathfrak{M}_k^{p^\dagger} := \mathfrak{M}_\emptyset \cup \mathfrak{M}_k^{p^*}$, where $\mathfrak{M}_k^{p^*} := \cup_{\mathbb{T} \in \mathfrak{J}_k^*} \mathfrak{M}_{\mathbb{T}}^p$. With those notations, Algorithm 2 gives the pseudocode for the adaptive ANOVA stochastic Galerkin method.

In the adaptive ANOVA stochastic Galerkin method, we select the set of multi-indices adaptively based on the ANOVA decomposition. Specifically, only part of the multi-indices with the
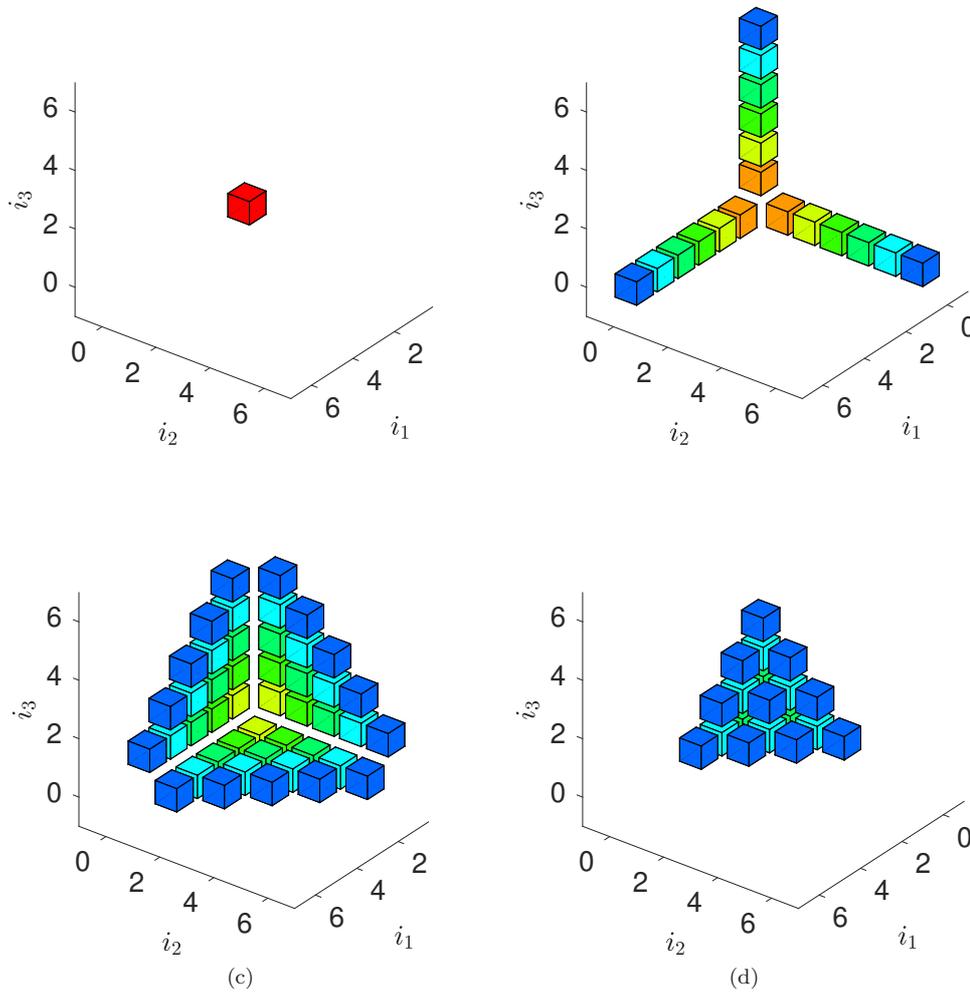
Fig. 4.1: Multi-indices of the gPC basis functions corresponding to the component functions of each order in 3 dimensions with the total degree up to 6, arranged according to the order of the component functions (from left to right): 0-th, first, second, and third order.

total degree up to $p$ will be retained, resulting in a much lower computational cost compared to the standard stochastic Galerkin method. It is worth noting that if the tolerance TOL is chosen small enough, all multi-indices will be selected, and the adaptive ANOVA stochastic Galerkin method will become equivalent to the standard stochastic Galerkin method.

## 5 Numerical results

In this section, we will explore two problems: a diffusion problem and a Helmholtz problem. All the results presented here are obtained using MATLAB R2015b on a desktop with a 2.90GHz

---

**Algorithm 2** Adaptive ANOVA stochastic Galerkin method

---

**Input:** The gPC order $p$ and the tolerance TOL in ANOVA decomposition.
Set $k = 1$, $\mathfrak{J}_0 = \emptyset$, $\mathfrak{J}_1 = \{\{1\}, \dots, \{N\}\}$ and compute $\boldsymbol{A}_i$ and $\boldsymbol{f}$ defined in (3.8).
**while** $(k < N)$ and $\mathfrak{J}_k \neq \emptyset$ **do**

　Generate the multi-indices $\mathfrak{M}_k^{p^\dagger}$ and compute $\boldsymbol{G}_i$ and $\boldsymbol{h}$ defined in (3.7).
　Solve the linear system (3.6) and compute $\gamma_{\mathbb{T}}$ for $\mathbb{T} \in \mathfrak{J}_k$ by (3.16) and (4.8).
　Set $\tilde{\mathfrak{J}}_k := \{\mathbb{T} \in \mathfrak{J}_k | \gamma_{\mathbb{T}} \geq$ TOL$\}$.
　Set $\mathfrak{J}_{k+1} := \{\mathbb{T} | \mathbb{T} \in \mathfrak{T}_{k+1}$, and $\forall\, \mathbb{S} \subseteq \mathbb{T}$ with $|\mathbb{S}| = k$ satisfies $\mathbb{S} \in \tilde{\mathfrak{J}}_k\}$.
　$k = k + 1$.
**end while**
**Return:** the approximation $u^{\mathrm{ap}}(\boldsymbol{x}, \boldsymbol{\mu})$, the mean function $u_0(\boldsymbol{x})$ and the variance function
$\sum_{\boldsymbol{i} \in \mathfrak{M}_{\mathbb{T}}^{p^*}} u_{\boldsymbol{i}}^2(\boldsymbol{x})$

---

Intel Core i7-10700 CPU. The CPU time reported in this paper correspond to the total time required to solve the linear systems in the respective procedures.

To assess the accuracy, we define the mean errors and the variance errors as follows:

$$\mathrm{E}_{\mathrm{ERR}} = \frac{\|\mathrm{E}\left[u_{\mathrm{REF}}(\boldsymbol{x}, \boldsymbol{\mu})\right] - \mathrm{E}\left[u^{\mathrm{ap}}(\boldsymbol{x}, \boldsymbol{\mu})\right]\|_{L^2}}{\|\mathrm{E}\left[u_{\mathrm{REF}}(\boldsymbol{x}, \boldsymbol{\mu})\right]\|_{L^2}},$$

$$\mathrm{V}_{\mathrm{ERR}} = \frac{\|\mathrm{V}\left[u_{\mathrm{REF}}(\boldsymbol{x}, \boldsymbol{\mu})\right] - \mathrm{V}\left[u^{\mathrm{ap}}(\boldsymbol{x}, \boldsymbol{\mu})\right]\|_{L^2}}{\|\mathrm{V}\left[u_{\mathrm{REF}}(\boldsymbol{x}, \boldsymbol{\mu})\right]\|_{L^2}}.$$

Here, $u_{\mathrm{REF}}(\boldsymbol{x}, \boldsymbol{\mu})$ is the reference solution, and $u^{\mathrm{ap}}(\boldsymbol{x}, \boldsymbol{\mu})$ is the approximate solution.

## 5.1 Test problem 1

In this problem, we investigate the diffusion equation with random inputs, given by

$$-\nabla \cdot (a(\boldsymbol{x}, \boldsymbol{\mu}) \nabla u(\boldsymbol{x}, \boldsymbol{\mu})) = 1 \quad \mathrm{in} \quad D \times \Gamma,$$
$$u(\boldsymbol{x}, \boldsymbol{\mu}) = 0 \quad \mathrm{on} \quad \partial D \times \Gamma,$$

where $D = [0, 1] \times [0, 1]$ is the spatial domain, and $\partial D$ represents the boundary of $D$. The diffusion coefficient $a(\boldsymbol{x}, \boldsymbol{\mu})$ is modeled as a truncated Karhunen–Loève (KL) expansion [15,9] of a random field with a mean function $a_0(\boldsymbol{x})$, a standard deviation $\sigma = 1/4$, and the covariance function $\mathrm{Cov}(\boldsymbol{x}, \boldsymbol{y})$ given by

$$\mathrm{Cov}(\boldsymbol{x}, \boldsymbol{y}) = \sigma^2 \exp\left(-\frac{|x_1 - y_1|}{c} - \frac{|x_2 - y_2|}{c}\right),$$

where $\boldsymbol{x} = [x_1, x_2]^T$, $\boldsymbol{y} = [y_1, y_2]^T$ and $c = 1/4$ is the correlation length. The KL expansion takes the form

$$a(\boldsymbol{x}, \boldsymbol{\mu}) = a_0(\boldsymbol{x}) + \sum_{i=1}^N a_i(\boldsymbol{x}) \mu_i = a_0(\boldsymbol{x}) + \sum_{i=1}^N \sqrt{\lambda_i} c_i(\boldsymbol{x}) \mu_i, \tag{5.1}$$

where $a_0(\boldsymbol{x}) = 1$, $\{\lambda_i, c_i(\boldsymbol{x})\}_{i=1}^N$ are the eigenpairs of $\mathrm{Cov}(\boldsymbol{x}, \boldsymbol{y})$, $\{\mu_i\}_{i=1}^N$ are uncorrelated random variables, and $N$ is the number of KL modes retained.

For this test problem, we assume that the random variables $\{\mu_i\}_{i=1}^N$ are independent and uniformly distributed within the range $[-1, 1]$. The parameters of $a(\boldsymbol{x}, \boldsymbol{\mu})$ are set as shown in

Table 5.1: Parameters of the diffusion coefficient $a(\boldsymbol{x}, \boldsymbol{\mu})$ in (5.1) and $K$ in (2.2).

| Case | $N$ | $K$ |
|------|-----|-----|
| I | 10 | 11 |
| II | 50 | 51 |

Table 5.1. In the physical domain, the meshgrid is set to $33 \times 33$ (i.e., $n_x = 33$). In the stochastic space, the total degree of gPC in the adaptive ANOVA stochastic Galerkin (AASG) method is specified as $p = 5$. The linear systems arising from both the standard and the adaptive ANOVA stochastic Galerkin methods are solved using the preconditioned conjugate gradients (CG) method, with a tolerance of $10^{-8}$ and the mean based preconditioner [26].

In this test problem, we compare the the adaptive ANOVA stochastic Galerkin (AASG) method with the anchored ANOVA stochastic collocation (AASC) method [23,42] and the Monte Carlo method (MCM). For the AASC method, we follow the method described in [23] with the relative mean $\tilde{\gamma}_{\mathbb{T}}$, i.e.,

$$\tilde{\gamma}_{\mathbb{T}} := \frac{\|\mathrm{E}\,[u_{\mathbb{T}}]\,\|_{L^2}}{\|\mathrm{E}\,[u_0]\,\|_{L^2}}, \ \mathbb{T} \in \mathfrak{J}_k^*,$$

as the criterion for selecting important terms within the ANOVA decomposition. Additionally, the mean value of $\boldsymbol{\mu}$ is used as the anchor point. In the AASC method, we adopt tensor style Gaussian quadrature points with a grid level of 5 as the collocation points, resulting in a total of $6^{|\mathbb{T}|}$ collocation (quadrature) points for each $\mathbb{T} \in \mathfrak{J}_k^*$. Furthermore, for the AASC method, the variance function is computed following the method proposed in [30]. For both the MCM and the AASC method, the linear systems are solved using the MATLAB backslash solver.

*5.1.1 Case I: a 10 dimensional diffusion problem*

We consider the AASG method with decreasing tolerances $\mathtt{TOL} = \{10^{-1}, 10^{-3}, 10^{-5}, 10^{-7}, 10^{-9}\}$ to demonstrate its effectiveness and efficiency. To access the accuracy, we obtain the reference solution $u_{\mathtt{REF}}(\boldsymbol{x}, \boldsymbol{\mu})$ using the standard stochastic Galerkin method with the total degree of up to $p = 7$.

Table 5.2: Performance of the AASG method for test problem 1 with $N = 10$.

| TOL | $|\mathfrak{J}_1|$ | $|\tilde{\mathfrak{J}}_1|$ | $|\mathfrak{J}_2|$ | $|\tilde{\mathfrak{J}}_2|$ | $|\mathfrak{J}_3|$ | $|\tilde{\mathfrak{J}}_3|$ | $|\mathfrak{J}_4|$ | $|\tilde{\mathfrak{J}}_4|$ | $|\mathfrak{J}_5|$ | $|\tilde{\mathfrak{J}}_5|$ | $k$ | $|\mathfrak{M}_k^{5\dagger}|$ | CPU time |
|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| $10^{-1}$ | 10 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 51 | 0.07 |
| $10^{-3}$ | 10 | 10 | 45 | 2 | 0 | 0 | 0 | 0 | 0 | 0 | 2 | 501 | 0.87 |
| $10^{-5}$ | 10 | 10 | 45 | 37 | 70 | 0 | 0 | 0 | 0 | 0 | 3 | 1201 | 3.35 |
| $10^{-7}$ | 10 | 10 | 45 | 45 | 120 | 75 | 60 | 0 | 0 | 0 | 4 | 2001 | 9.38 |
| $10^{-9}$ | 10 | 10 | 45 | 45 | 120 | 120 | 210 | 127 | 70 | 0 | 5 | 2821 | 18.71 |

Table 5.2 presents the number of active indices for each order of the ANOVA decomposition and the total number of selected gPC basis functions in the stochastic space. Furthermore, we report the computational time required for solving all linear systems that arise during the

execution of the while loop in Algorithm 2. It can be observed that the number of selected gPC basis functions increases as the tolerance TOL decreases, and therefore, the accuracy can be improved by reducing TOL. For a 10 dimensional problem, the number of gPC basis functions with the total degree up to 5 is $C_{15}^5 = 3003$. From the table, it can be seen that when TOL $= 10^{-9}$, almost all the gPC basis functions are selected in the AASG method. Further decreasing the tolerance in the AASG method results in the selection of all gPC basis functions with the total degree up to 5, making the AASG method equivalent to the standard stochastic Galerkin method.



(a) TOL $= 10^{-3}$.     (b) TOL $= 10^{-5}$.     (c) TOL $= 10^{-7}$.
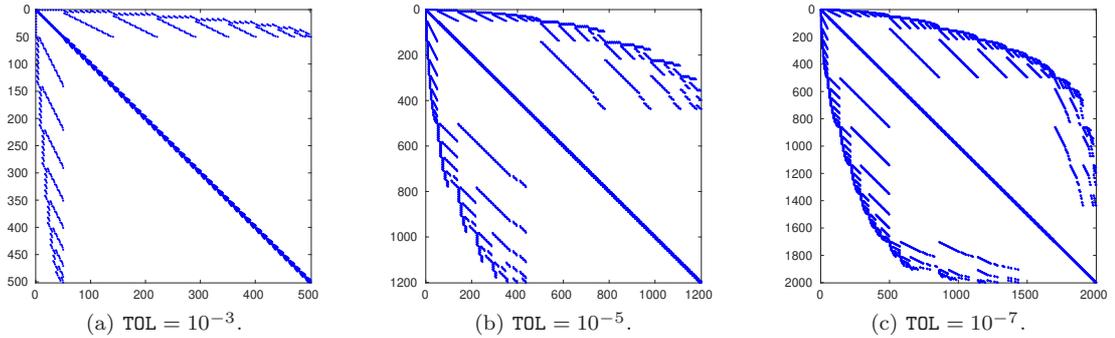
Fig. 5.1: Matrix block-structure (each block has dimension $n_x \times n_x$) for test problem 1 with $N = 10$.

Fig. 5.1 displays the block structure of the coefficient matrix of the resulting linear system. Each point in the figure represents a block of dimension $n_x \times n_x$. Furthermore, each nonzero block of the coefficient matrix has the same sparsity pattern as the corresponding deterministic problem. Therefore, the coefficient matrix is extremely large and sparse, and the resulting linear system should be solved using iterative methods.

Fig. 5.2 investigates the accuracy achieved by the three methods, presenting errors concerning both CPU time and stochastic degrees of freedom (DOF). For clarity, CPU time denotes the total time required to solve all the linear systems within the respective procedures. For the AASG method, stochastic degrees of freedom encompass the cumulative count of gPC basis functions generated during the execution of the while loop in Algorithm 2, while for the MCM and the AASC method, it corresponds to the total number of sample points used. The results indicate the notable efficiency of both the AASG method and the AASC method in comparison to the MCM, as they are hundreds of times faster in terms of CPU time and stochastic degrees of freedom. Furthermore, the AASG method outperforms the AASC method in terms of CPU time, and it is also evident that the AASG method provides a higher accuracy per stochastic degree of freedom compared to the AASC method.

### 5.1.2 Case II: a 50 dimensional diffusion problem

We consider the AASG method with decreasing tolerances TOL $= \{10^{-1}, 10^{-2}, 10^{-3}, 10^{-4}, 10^{-5}\}$ to demonstrate its effectiveness and efficiency. To access the accuracy of the AASG method and the MCM, we obtain the reference solution $u_{\text{REF}}(\boldsymbol{x}, \boldsymbol{\mu})$ using the AASG method with a tolerance of TOL $= 10^{-6}$.

Table 5.3 presents the number of active indices for each order of ANOVA decomposition and the total number of selected gPC basis functions in the stochastic space. Furthermore, we report

(a) Mean errors w.r.t. CPU time.

(b) Variance errors w.r.t. CPU time.

(c) Mean errors w.r.t. stochastic DOF.
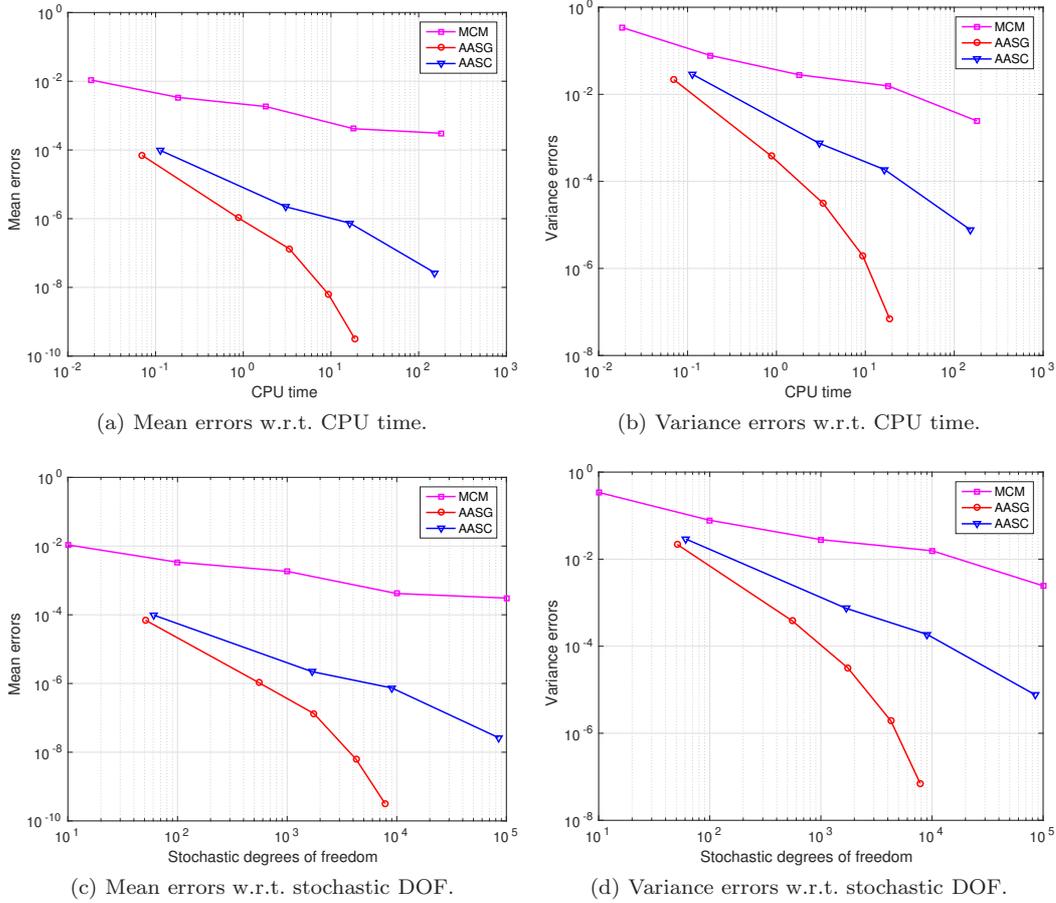
(d) Variance errors w.r.t. stochastic DOF.

Fig. 5.2: Comparison of errors with respect to CPU times and stochastic degrees of freedom for test problem 1 with $N = 10$, where both the total degree of gPC in the AASG method and the grid level in the AASC method are set to 5.

the computational time required for solving all linear systems that arise during the execution of the while loop in Algorithm 2. It can be observed that the number of selected gPC basis functions increases as the tolerance TOL decreases, and therefore, the accuracy can be improved by reducing TOL. It is worth noting that the number of gPC basis functions with the total degree up to $p = 5$ is $C_{55}^5 = 3478761$, which renders the standard stochastic Galerkin method practically infeasible for solving this problem in reasonable time.

Fig. 5.3 displays the block structure of the coefficient matrix of the resulting linear system. Each point in the figure represents a block of dimension $n_x \times n_x$. Furthermore, each nonzero block of the coefficient matrix has the same sparsity pattern as the corresponding deterministic problem. It is evident from the figure that the coefficient matrix in this case exhibits a sparser pattern than that of Case I.

The plot displayed in Fig. 5.4 compares the CPU times and the stochastic degrees of freedom required by the three methods in relation to mean and variance errors. The results clearly indicate that both the AASG method and the AASC method are significantly more efficient than
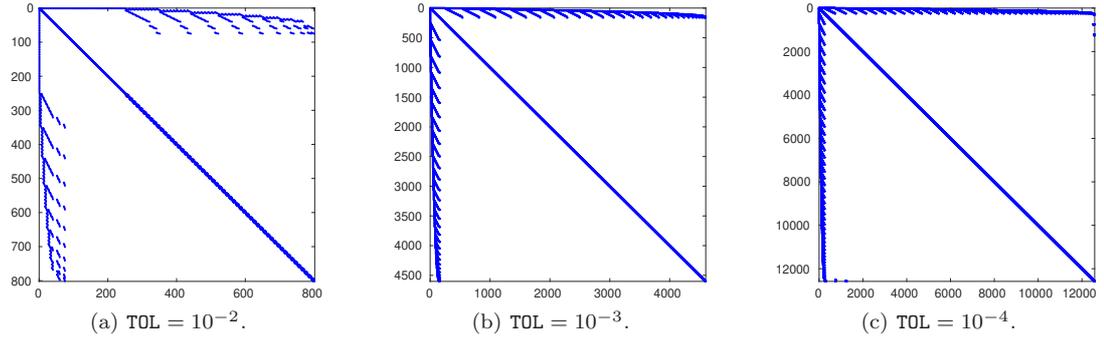
(a) TOL $= 10^{-2}$.　　　　　(b) TOL $= 10^{-3}$.　　　　　(c) TOL $= 10^{-4}$.

Fig. 5.3: Matrix block-structure (each block has dimension $n_x \times n_x$) for test problem 1 with $N = 50$.



(a) Mean errors w.r.t. CPU time.　　　　　(b) Variance errors w.r.t. CPU time.

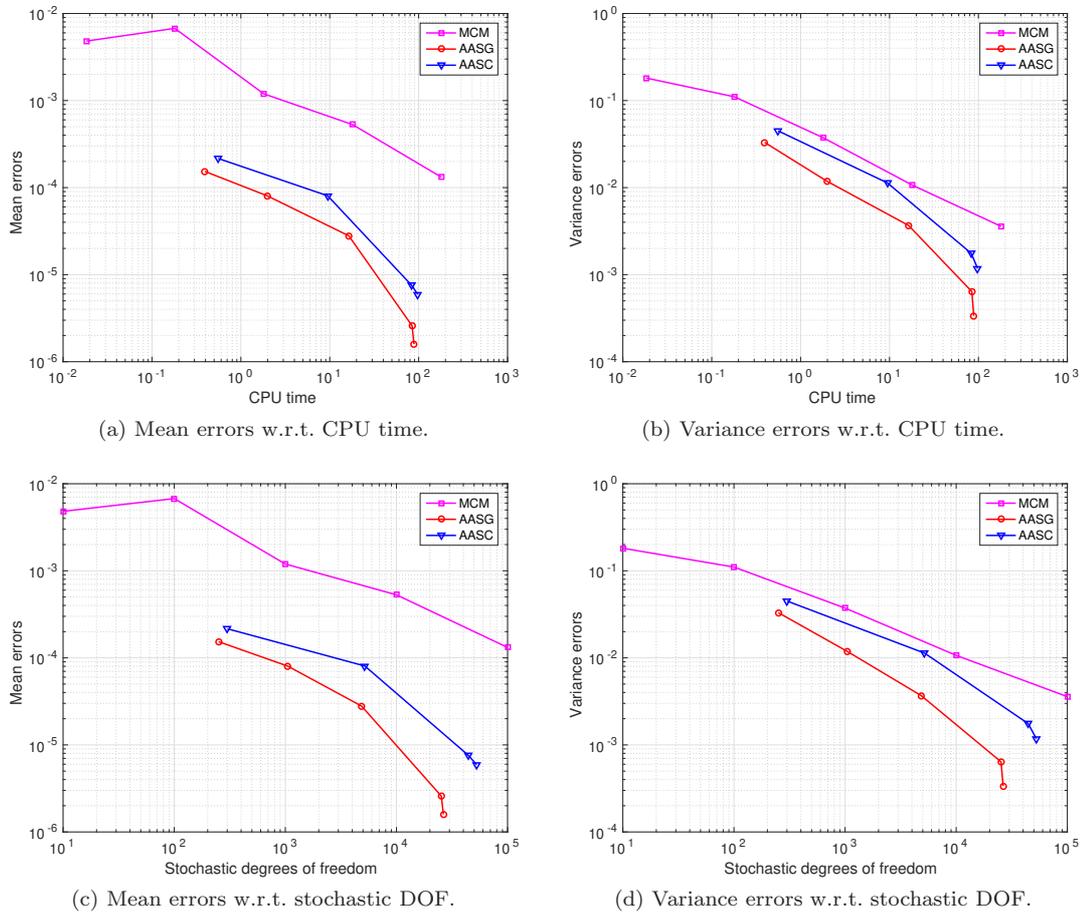(c) Mean errors w.r.t. stochastic DOF.　　　　　(d) Variance errors w.r.t. stochastic DOF.

Fig. 5.4: Comparison of errors with respect to CPU times and stochastic degrees of freedom for test problem 1 with $N = 50$, where both the total degree of gPC in the AASG method and the grid level in the AASC method are set to 5.

Table 5.3: Performance of the AASG method for test problem 1 with $N = 50$.

| TOL | $|\mathfrak{J}_1|$ | $|\tilde{\mathfrak{J}}_1|$ | $|\mathfrak{J}_2|$ | $|\tilde{\mathfrak{J}}_2|$ | $|\mathfrak{J}_3|$ | $|\tilde{\mathfrak{J}}_3|$ | $|\mathfrak{J}_4|$ | $|\tilde{\mathfrak{J}}_4|$ | $k$ | $|\mathfrak{M}_k^{5\dagger}|$ | CPU time |
|---|---|---|---|---|---|---|---|---|---|---|---|
| $10^{-1}$ | 50 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 251 | 0.39 |
| $10^{-2}$ | 50 | 11 | 55 | 0 | 0 | 0 | 0 | 0 | 2 | 801 | 1.99 |
| $10^{-3}$ | 50 | 30 | 435 | 0 | 0 | 0 | 0 | 0 | 2 | 4601 | 16.37 |
| $10^{-4}$ | 50 | 50 | 1225 | 15 | 8 | 0 | 0 | 0 | 3 | 12581 | 84.39 |
| $10^{-5}$ | 50 | 50 | 1225 | 83 | 120 | 0 | 0 | 0 | 3 | 13701 | 88.55 |
| $10^{-6}$ | 50 | 50 | 1225 | 377 | 1537 | 15 | 1 | 0 | 4 | 27876 | 284.41 |

the MCM, with improvements observed in both CPU time and stochastic degrees of freedom. Additionally, the AASG method outperforms the AASC method in terms of CPU time, and it is also evident that the AASG method provides a higher accuracy per stochastic degree of freedom compared to the AASC method.

5.2 Test problem 2

In this test problem, we consider the stochastic Helmholtz problem given by

$$\nabla^2 u + a^2(\boldsymbol{x}, \boldsymbol{\mu})u = f(\boldsymbol{x}) \quad \text{in} \quad D \times \Gamma,$$

with Sommerfeld radiation boundary condition. Here, $D = [0,1]^2$ is the domain of interest and the Helmholtz coefficient $a(\boldsymbol{x}, \boldsymbol{\mu})$ is a truncated KL expansion of a random field with a mean function $a_0(\boldsymbol{x})$, a standard deviation $\sigma = 2\pi$, and the covariance function $\text{Cov}(\boldsymbol{x}, \boldsymbol{y})$ given by

$$\text{Cov}(\boldsymbol{x}, \boldsymbol{y}) = \sigma^2 \exp\left(-\frac{|x_1 - y_1|}{c} - \frac{|x_2 - y_2|}{c}\right),$$

where $\boldsymbol{x} = [x_1, x_2]^T$, $\boldsymbol{y} = [y_1, y_2]^T$ and $c = 1$ is the correlation length. Note that the KL expansion takes the form

$$a(\boldsymbol{x}, \boldsymbol{\mu}) = a_0(\boldsymbol{x}) + \sum_{i=1}^{N} a_i(\boldsymbol{x})\mu_i = a_0(\boldsymbol{x}) + \sum_{i=1}^{N} \sqrt{\lambda_i} c_i(\boldsymbol{x})\mu_i, \quad (5.2)$$

where $a_0(\boldsymbol{x}) = 4 \cdot (2\pi)$, $\{\lambda_i, c_i(\boldsymbol{x})\}_{i=1}^{N}$ are the eigenpairs of $\text{Cov}(\boldsymbol{x}, \boldsymbol{y})$, $\{\mu_i\}_{i=1}^{N}$ are uncorrelated random variables, and $N$ is the number of KL modes retained. The Gaussian point source at the center of the domain is used as the source term, i.e.,

$$f(\boldsymbol{x}) = e^{-(8 \cdot 4)^2((x_1 - 0.5)^2 + (x_2 - 0.5)^2)}$$

For this test problem, we assume that the random variables $\{\mu_i\}_{i=1}^{N}$ are independent and uniformly distributed within the range $[-1, 1]$. The parameters of $a(\boldsymbol{x}, \boldsymbol{\mu})$ are set as shown in Table 5.4. We use the perfectly matched layers (PML) to simulate the Sommerfeld condition [4], and generate the matrices $\{\boldsymbol{A}_i\}_{i=1}^{K}$ using the codes associated with [22]. In the physical domain, the meshgrid is set to $33 \times 33$ (i.e., $n_x = 33$). In the stochastic space, the total degree of the gPC basis functions in the AASG method is set to $p = 6$. The linear systems arising from both the standard stochastic Galerkin method and the AASG method are solved using the preconditioned bi-conjugate gradient stabilized (Bi-CGSTAB) method, with a tolerance of $10^{-8}$ and the mean based preconditioner [26]. Furthermore, for the MCM, the linear systems is solved using the MATLAB backslash solver.

Table 5.4: Parameters of the Helmholtz coefficient $a(\boldsymbol{x}, \boldsymbol{\mu})$ in (5.2) and $K$ in (2.2).

| Case | $N$ | $K$ |
|------|-----|-----|
| I | 4 | 25 |
| II | 10 | 121 |

### 5.2.1 Case I: a 4 dimensional Helmholtz problem

We consider the AASG method with decreasing tolerances $\texttt{TOL} = \{10^{-1}, 10^{-2}, 10^{-3}, 10^{-4}, 10^{-5}\}$ to demonstrate its effectiveness and efficiency. To access the accuracy of the AASG method and the MCM, we obtain the reference solution $u_{\texttt{REF}}(\boldsymbol{x}, \boldsymbol{\mu})$ using the standard stochastic Galerkin method with the total degree of up to $p = 8$.

Table 5.5: Performance of the AASG method for test problem 2 with $N = 4$.

| TOL | $|\mathfrak{J}_1|$ | $|\tilde{\mathfrak{J}}_1|$ | $|\mathfrak{J}_2|$ | $|\tilde{\mathfrak{J}}_2|$ | $|\mathfrak{J}_3|$ | $|\tilde{\mathfrak{J}}_3|$ | $|\mathfrak{J}_4|$ | $|\tilde{\mathfrak{J}}_4|$ | $k$ | $|\mathfrak{M}_k^{6\dagger}|$ | CPU time |
|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| $10^{-1}$ | 4 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 25 | 0.22 |
| $10^{-2}$ | 4 | 4 | 6 | 3 | 0 | 0 | 0 | 0 | 2 | 115 | 1.52 |
| $10^{-3}$ | 4 | 4 | 6 | 5 | 2 | 1 | 0 | 0 | 3 | 155 | 3.62 |
| $10^{-4}$ | 4 | 4 | 6 | 6 | 4 | 3 | 0 | 0 | 3 | 195 | 4.46 |
| $10^{-5}$ | 4 | 4 | 6 | 6 | 4 | 4 | 1 | 1 | 4 | 210 | 7.95 |

Table 5.5 presents the number of active indices for each order of ANOVA decomposition and the total number of selected gPC basis functions in the stochastic space. Furthermore, we report the computational time required for solving all linear systems that arise during the execution of the while loop in Algorithm 2. It can be observed that the number of selected gPC basis functions increases as the tolerance $\texttt{TOL}$ decreases. For a 4 dimensional problem, there are $C_{10}^6 = 210$ gPC basis functions with the total degree up to 6. From the table, it can be seen that when $\texttt{TOL} = 10^{-5}$, all the gPC basis functions are selected in the AASG method, making the AASG method equivalent to the standard stochastic Galerkin method.

Fig. 5.5 illustrates the block structure of the coefficient matrix of the resulting linear system. Each point in the figure represents a block of dimension $n_x \times n_x$. Moreover, each nonzero block of the coefficient matrix has the same sparsity pattern as the corresponding deterministic problem. Although the coefficient matrix of the Helmholtz equation is much denser than that of the diffusion equation, it is still very sparse and thus should be solved by iterative methods.

Fig. 5.6 investigates the accuracy achieved by the two methods, presenting errors concerning both CPU time and stochastic degrees of freedom. For clarity, CPU time denotes the total time required to solve all the linear systems within the respective procedures. For the AASG method, stochastic degrees of freedom encompass the cumulative count of gPC basis functions generated during the execution of the while loop in Algorithm 2, while for the MCM, it corresponds to the total number of sample points used. The results indicate the notable efficiency of the AASG method in comparison to the MCM, as it is dozens of times faster in terms of CPU time and hundreds of times faster in terms of stochastic degrees of freedom. This distinction from the diffusion problem is noteworthy, as the performance of CPU time does not seem to align with
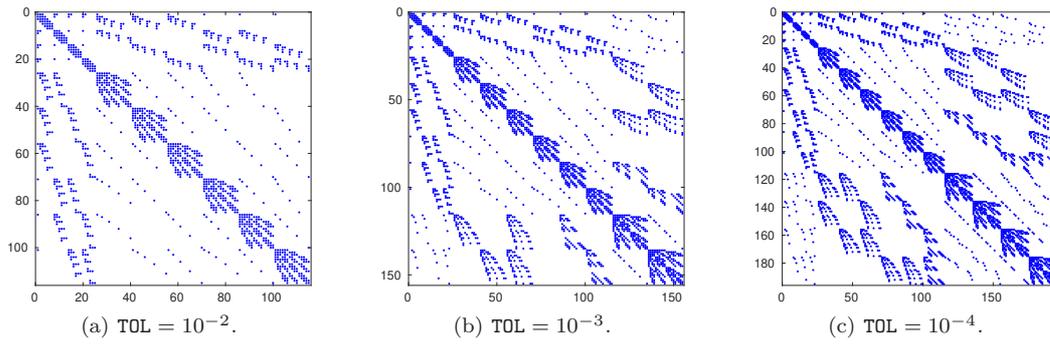
(a) TOL $= 10^{-2}$.              (b) TOL $= 10^{-3}$.              (c) TOL $= 10^{-4}$.

Fig. 5.5: Matrix block-structure (each block has dimension $n_x \times n_x$) for test problem 2 with $N = 4$.



(a) Mean errors w.r.t. CPU time.              (b) Variance errors w.r.t. CPU time.
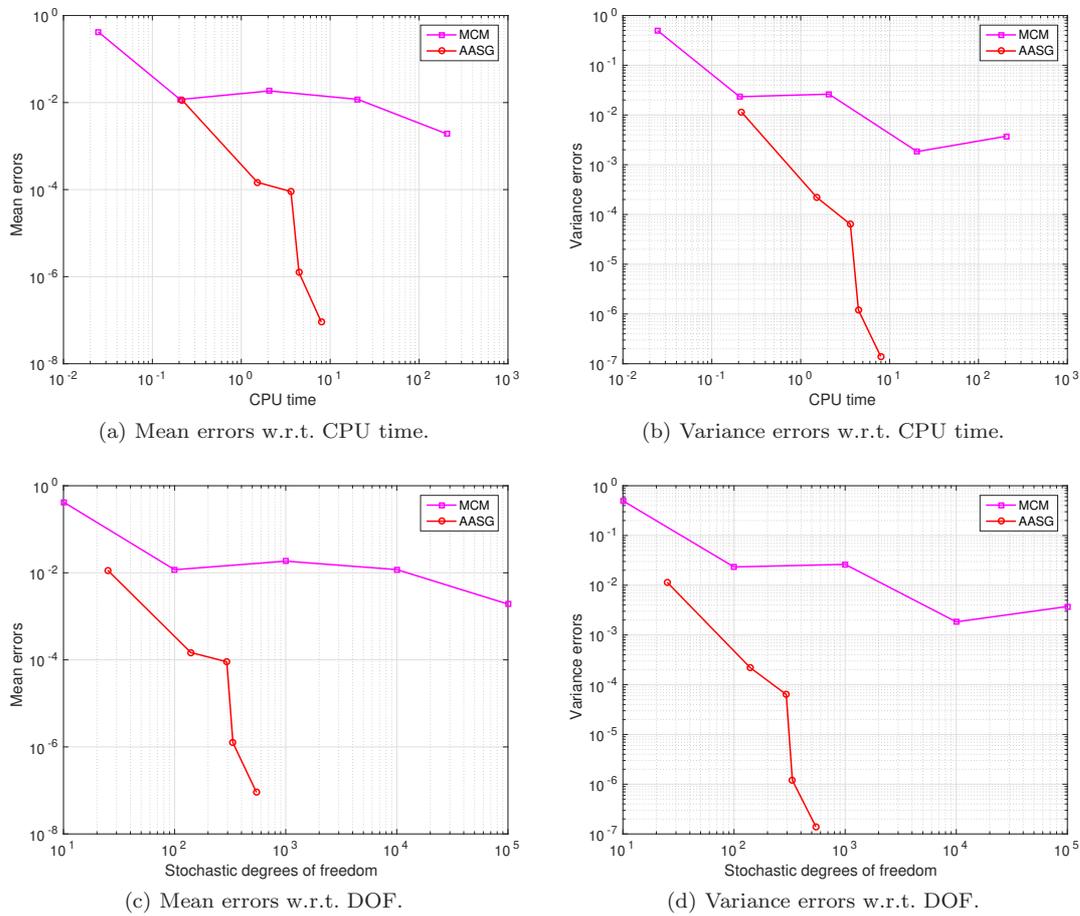
(c) Mean errors w.r.t. DOF.              (d) Variance errors w.r.t. DOF.

Fig. 5.6: Comparison of errors with respect to CPU times and stochastic degrees of freedom for test problem 2 with $N = 4$, where the total degree of gPC in the AASG method is set to 6.

that of stochastic degrees of freedom. This phenomenon can be attributed to the fact that the number of matrix-vector products computed in each iteration for solving the linear system (3.6) is proportional to $(N+1)^2$ for Helmholtz problems, while it is proportional to $N+1$ for diffusion problems.

### 5.2.2 Case II: a 10 dimensional Helmholtz problem

We consider the AASG method with decreasing tolerances $\texttt{TOL} = \{10^{-1}, 10^{-2}, 10^{-3}, 10^{-4}, 10^{-5}\}$ to demonstrate its effectiveness and efficiency. To access the accuracy of the AASG method and the MCM, we obtain the reference solution $u_{\texttt{REF}}(\boldsymbol{x}, \boldsymbol{\mu})$ using the standard stochastic Galerkin method with the total degree of up to $p = 8$.

Table 5.6: Performance of the AASG method for test problem 2 with $N = 10$.

| TOL | $|\mathfrak{J}_1|$ | $|\tilde{\mathfrak{J}}_1|$ | $|\mathfrak{J}_2|$ | $|\tilde{\mathfrak{J}}_2|$ | $|\mathfrak{J}_3|$ | $|\tilde{\mathfrak{J}}_3|$ | $|\mathfrak{J}_4|$ | $|\tilde{\mathfrak{J}}_4|$ | $k$ | $|\mathfrak{M}_k^{6^\dagger}|$ | CPU time |
|---|---|---|---|---|---|---|---|---|---|---|---|
| $10^{-1}$ | 10 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 61 | 0.72 |
| $10^{-2}$ | 10 | 9 | 36 | 6 | 0 | 0 | 0 | 0 | 2 | 601 | 14.03 |
| $10^{-3}$ | 10 | 10 | 45 | 14 | 7 | 3 | 0 | 0 | 3 | 876 | 41.77 |
| $10^{-4}$ | 10 | 10 | 45 | 32 | 55 | 21 | 0 | 0 | 3 | 1836 | 82.92 |
| $10^{-5}$ | 10 | 10 | 45 | 43 | 105 | 59 | 41 | 22 | 4 | 3451 | 281.31 |

Table 5.6 presents the number of active indices for each order of ANOVA decomposition and the total number of selected gPC basis functions in the stochastic space. Furthermore, we report the computational time required for solving all linear systems that arise during the execution of the while loop in Algorithm 2. It can be observed that the number of selected gPC basis functions increases as the tolerance $\texttt{TOL}$ decreases, and therefore, the accuracy can be improved by reducing $\texttt{TOL}$.



(a) $\texttt{TOL} = 10^{-2}$.          (b) $\texttt{TOL} = 10^{-3}$.          (c) $\texttt{TOL} = 10^{-4}$.
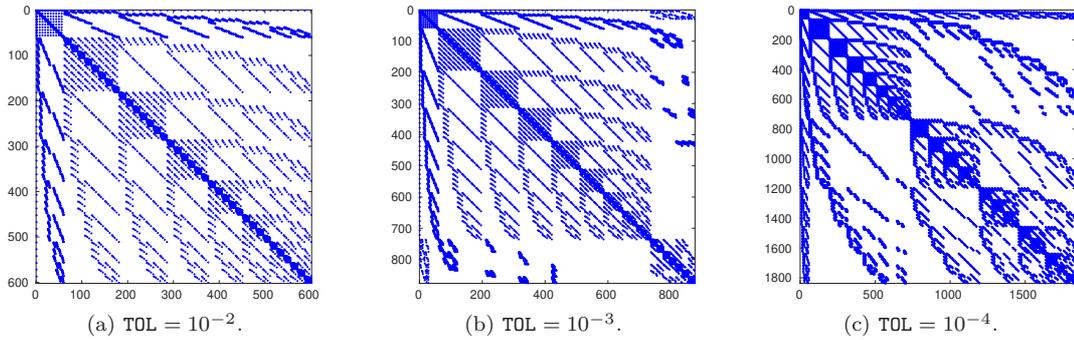
Fig. 5.7: Matrix block-structure (each block has dimension $n_x \times n_x$) for test problem 2 with $N = 10$.

Fig. 5.7 illustrates the block structure of the coefficient matrix of the resulting linear system. Each point in the figure represents a block of dimension $n_x \times n_x$. Moreover, each nonzero block of

the coefficient matrix has the same sparsity pattern as the corresponding deterministic problem. It can be observed that the coefficient matrix of the Helmholtz equation is much denser than that of the diffusion equation, which makes it more time consuming to solve than the diffusion problem.
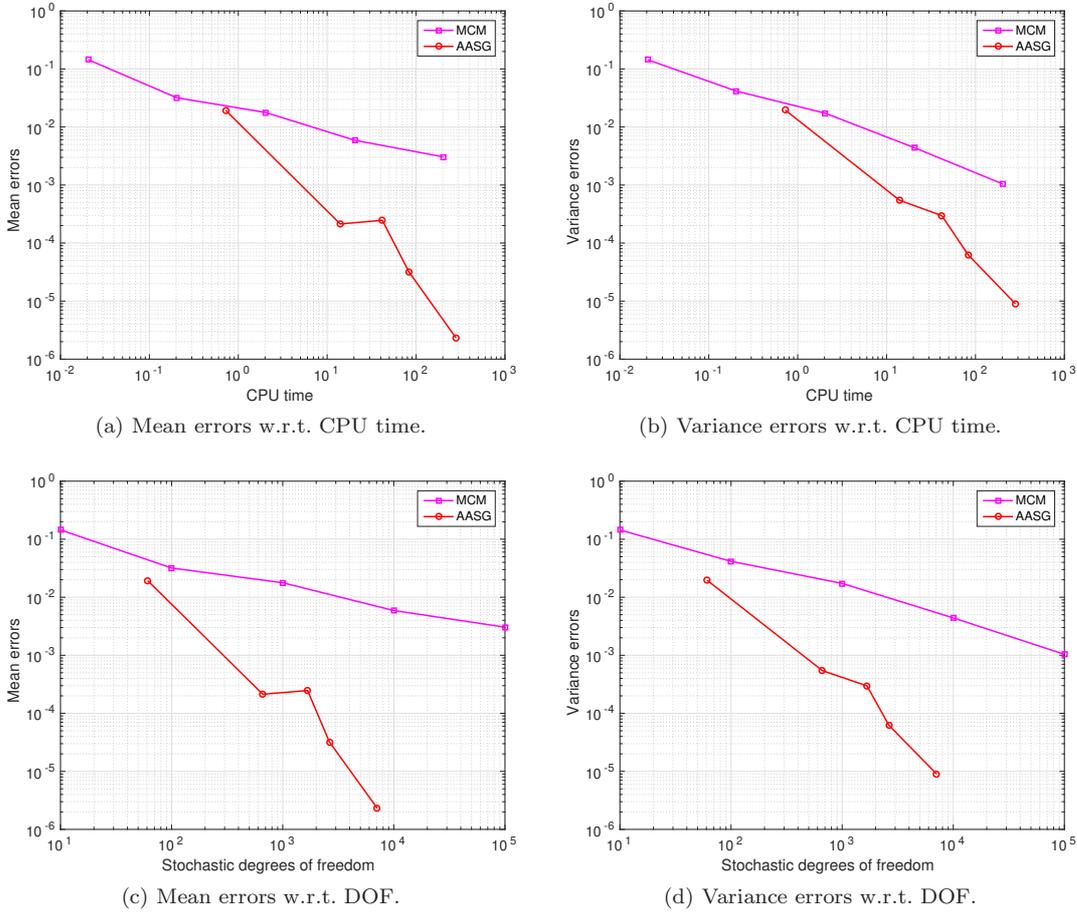


(a) Mean errors w.r.t. CPU time.

(b) Variance errors w.r.t. CPU time.

(c) Mean errors w.r.t. DOF.

(d) Variance errors w.r.t. DOF.

Fig. 5.8: Comparison of errors with respect to CPU times and stochastic degrees of freedom for test problem 2 with $N = 10$, where the total degree of gPC in the AASG method is set to 6.

Fig. 5.8 investigates the accuracy achieved by the two methods, presenting errors concerning both CPU time and stochastic degrees of freedom. The results indicate the notable efficiency of the AASG method in comparison to the MCM, as it is dozens of times faster in terms of CPU time and hundreds of times faster in terms of stochastic degrees of freedom. This distinction from the diffusion problem is noteworthy, as the performance of CPU time does not seem to align with that of stochastic degrees of freedom. This phenomenon can be attributed to the fact that the number of matrix-vector products computed in each iteration for solving the linear system (3.6) is proportional to $(N + 1)^2$ for Helmholtz problems, in contrast to $N + 1$ for diffusion problems.

# 6 Conclusion

In this work, we investigate the generalized polynomial chaos (gPC) expansion of component functions for the ANOVA decomposition, and present a concise form of the gPC expansion for each component function. With this formulation, we propose an adaptive ANOVA stochastic Galerkin method for solving partial differential equations with random inputs. The proposed method effectively selects basis functions in the stochastic space, enabling significant reduction in the dimension of the stochastic approximation space. Compared with anchored ANOVA methods, the proposed approach avoids the difficulty for selecting proper anchor points, which are crucial for achieving efficient approximations in the context of anchored ANOVA methods. Numerical simulations are conducted to demonstrate the effectiveness and efficiency of the proposed method. While our current focus is on selecting the basis in the stochastic space, future work will explore techniques for reducing computational costs in the physical space.

# References

1. Agarwal, N., Aluru, N.R.: A domain adaptive stochastic collocation approach for analysis of MEMS under uncertainties. Journal of Computational Physics **228**(20), 7662–7688 (2009). DOI 10.1016/j.jcp.2009.07.014
2. Askey, R.: Orthogonal polynomials and special functions. SIAM (1975)
3. Babuška, I., Tempone, R.l., Zouraris, G.E.: Galerkin finite element approximations of stochastic elliptic partial differential equations. SIAM Journal on Numerical Analysis **42**(2), 800–825 (2004). DOI 10.1137/S0036142902418680
4. Berenger, J.P.: A perfectly matched layer for the absorption of electromagnetic waves. Journal of Computational Physics **114**(2), 185–200 (1994). DOI 10.1006/jcph.1994.1159
5. Caflisch, R.E.: Monte Carlo and quasi-Monte Carlo methods. Acta Numerica **7**, 1–49 (1998)
6. Cheng, M., Hou, T.Y., Zhang, Z.: A dynamically bi-orthogonal method for time-dependent stochastic partial differential equations I: Derivation and algorithms. Journal of Computational Physics **242**, 843–868 (2013). DOI 10.1016/j.jcp.2013.02.033
7. Cheng, M., Hou, T.Y., Zhang, Z.: A dynamically bi-orthogonal method for time-dependent stochastic partial differential equations II: Adaptivity and generalizations. Journal of Computational Physics **242**, 753–776 (2013). DOI 0.1016/j.jcp.2013.02.020
8. Cho, H., Elman, H.C.: An adaptive reduced basis collocation method based on PCM ANOVA decomposition for anisotropic stochastic PDEs. International Journal for Uncertainty Quantification **8**(3) (2018). DOI 10.1615/Int.J.UncertaintyQuantification.2018024436
9. Elman, H., Furnival, D.: Solving the stochastic steady-state diffusion problem using multigrid. IMA Journal of Numerical Analysis **27**(4), 675–688 (2007). DOI 10.1093/imanum/drm006
10. Elman, H., Liao, Q.: Reduced basis collocation methods for partial differential equations with random coefficients. SIAM/ASA Journal on Uncertainty Quantification **1**, 192–217 (2013). DOI 10.1137/120881841
11. Elman, H.C., Ernst, O.G., O'Leary, D.P., Stewart, M.: Efficient iterative algorithms for the stochastic finite element method with application to acoustic scattering. Computer Methods in Applied Mechanics and Engineering **194**, 1037–1055 (2005). DOI 10.1016/j.cma.2004.06.028
12. Feng, X., Lin, J., Lorton, C.: An efficient numerical method for acoustic wave scattering in random media. SIAM/ASA Journal on Uncertainty Quantification **3**(1), 790–822 (2015). DOI 10.1137/140958232
13. Fishman, G.: Monte Carlo: concepts, algorithms, and applications. Springer Science & Business Media (2013)
14. Gao, Z., Hesthaven, J.S.: On ANOVA expansions and strategies for choosing the anchor point. Applied Mathematics and Computation **217**(7), 3274–3285 (2010). DOI 10.1016/j.amc.2010.08.061
15. Ghanem, R.G., Spanos, P.D.: Stochastic finite elements: a spectral approach. Courier Corporation (2003)
16. Guo, L., Narayan, A., Zhou, T.: Constructing least-squares polynomial approximations. SIAM Review **62**(2), 483–508 (2020). DOI 10.1137/18M1234151
17. Jakeman, J.D., Narayan, A., Zhou, T.: A generalized sampling and preconditioning scheme for sparse approximation of polynomial chaos expansions. SIAM Journal on Scientific Computing **39**(3), A1114–A1144 (2017). DOI 10.1137/16M1063885
18. Kämmerer, L., Potts, D., Taubert, F.: The uniform sparse FFT with application to PDEs with random coefficients. Sampling Theory, Signal Processing, and Data Analysis **20**(19) (2021). DOI 10.1007/s43670-022-00037-3

19. Lee, K., Elman, H.C.: A preconditioned low-rank projection method with a rank-reduction scheme for stochastic partial differential equations. SIAM Journal on Scientific Computing **39**(5), S828–S850 (2017). DOI 10.1137/16M1075582

20. Lee, K., Elman, H.C., Sousedik, B.: A low-rank solver for the Navier–Stokes equations with uncertain viscosity. SIAM/ASA Journal on Uncertainty Quantification **7**(4), 1275–1300 (2019). DOI 10.1137/17M1151912

21. Liao, Q., Lin, G.: Reduced basis ANOVA methods for partial differential equations with high-dimensional random inputs. Journal of Computational Physics **317**, 148–164 (2016). DOI 10.1016/j.jcp.2016.04.029

22. Liu, F., Ying, L.: Additive sweeping preconditioner for the Helmholtz equation. Multiscale Modeling & Simulation **14**(2), 799–822 (2016). DOI 10.1137/15M1017144

23. Ma, X., Zabaras, N.: An adaptive high-dimensional stochastic model representation technique for the solution of stochastic partial differential equations. Journal of Computational Physics **229**(10), 3884–3915 (2010). DOI 10.1016/j.jcp.2010.01.033

24. Musharbash, E., Nobile, F., Zhou, T.: Error analysis of the dynamically orthogonal approximation of time dependent random PDEs. SIAM Journal on Scientific Computing **37**(2), A776–A810 (2015). DOI 10.1137/140967787

25. Potts, D., Schmischke, M.: Approximation of high-dimensional periodic functions with Fourier-based methods. SIAM Journal on Numerical Analysis **59**(5), 2393–2429 (2021). DOI 10.1137/20M1354921

26. Powell, C.E., Elman, H.C.: Block-diagonal preconditioning for spectral stochastic finite element systems. IMA Journal of Numerical Analysis **29**(2), 350–375 (2009). DOI 10.1093/imanum/drn014

27. Powell, C.E., Silvester, D., Simoncini, V.: An efficient reduced basis solver for stochastic Galerkin matrix equations. SIAM Journal on Scientific Computing **39**(1), A141–A163 (2017). DOI 10.1137/15M1032399

28. Sobol', I.M.: Theorems and examples on high dimensional model representation. Reliability Engineering and System Safety **79**(2), 187–193 (2003). DOI 10.1016/S0951-8320(02)00229-6

29. Sudret, B.: Global sensitivity analysis using polynomial chaos expansions. Reliability Engineering and System Safety **93**(7), 964–979 (2008). DOI 10.1016/j.ress.2007.04.002

30. Tang, K., Congedo, P.M., Abgrall, R.: Sensitivity analysis using anchored ANOVA expansion and high-order moments computation. International Journal for Numerical Methods in Engineering **102**(9), 1554–1584 (2015). DOI 10.1002/nme.4856

31. Tang, K., Congedo, P.M., Abgrall, R.: Adaptive surrogate modeling by ANOVA and sparse polynomial dimensional decomposition for global sensitivity analysis in fluid simulation. Journal of Computational Physics **314**(1), 557–589 (2016). DOI 10.1016/j.jcp.2016.03.026

32. Tang, T., Zhou, T.: Convergence analysis for stochastic collocation methods to scalar hyperbolic equations with a random wave speed. Commun. Comput. Phys **8**(1), 226–248 (2010). DOI 10.4208/cicp.060109.130110a

33. Wan, X., Karniadakis, G.E.: An adaptive multi-element generalized polynomial chaos method for stochastic differential equations. Journal of Computational Physics **209**(2), 617–642 (2005). DOI 10.1016/j.jcp.2005.03.023

34. Wang, X.: On the approximation error in high dimensional model representation. In: 2008 Winter Simulation Conference, pp. 453–462. IEEE (2008). DOI 10.1109/WSC.2008.4736100

35. Williamson, K., Cho, H., Sousedík, B.: Application of adaptive ANOVA and reduced basis methods to the stochastic Stokes-Brinkman problem. Computational Geosciences **25**(3), 1191–1213 (2021). DOI 10.1007/s10596-021-10048-z

36. Xiu, D.: Numerical methods for stochastic computations: a spectral method approach. Princeton University Press (2010)

37. Xiu, D., Hesthaven, J.: High-order collocation methods for differential equations with random inputs. SIAM Journal on Scientific Computing **27**, 1118–1139 (2005). DOI 10.1137/040615201

38. Xiu, D., Karniadakis, G.E.: Modeling uncertainty in steady state diffusion problems via generalized polynomial chaos. Computer Methods in Applied Mechanics and Engineering **191**(43), 4927–4948 (2002). DOI 10.1016/S0045-7825(02)00421-8

39. Xiu, D., Karniadakis, G.E.: The Wiener-Askey polynomial chaos for stochastic differential equations. SIAM Journal on Scientific Computing **24**(2), 619–644 (2002). DOI 10.1137/S1064827501387826

40. Xiu, D., Karniadakis, G.E.: Modeling uncertainty in flow simulations via generalized polynomial chaos. Journal of Computational Physics **187**(1), 137–167 (2003). DOI doi:10.1016/S0021-9991(03)00092-5

41. Yan, L., Zhou, T.: Adaptive multi-fidelity polynomial chaos approach to Bayesian inference in inverse problems. Journal of Computational Physics **381**, 110–128 (2019). DOI 10.1016/j.jcp.2018.12.025

42. Yang, X., Choi, M., Lin, G., Karniadakis, G.E.: Adaptive ANOVA decomposition of stochastic incompressible and compressible flows. Journal of Computational Physics **231**(4), 1587–1614 (2012)