

ON THE CONVERGENCE OF THE STOCHASTIC GALERKIN METHOD FOR RANDOM ELLIPTIC PARTIAL DIFFERENTIAL EQUATIONS *

ANTJE MUGLER¹ AND HANS-JÖRG STARKLOFF²

Abstract. In this article we consider elliptic partial differential equations with random coefficients and/or random forcing terms. In the current treatment of such problems by stochastic Galerkin methods it is standard to assume that the random diffusion coefficient is bounded by positive deterministic constants or modeled as lognormal random field. In contrast, we make the significantly weaker assumption that the non-negative random coefficients can be bounded strictly away from zero and infinity by random variables only and may have distributions different from a lognormal one. We show that in this case the standard stochastic Galerkin approach does not necessarily produce a sequence of approximate solutions that converges in the natural norm to the exact solution even in the case of a lognormal coefficient. By using weighted test function spaces we develop an alternative stochastic Galerkin approach and prove that the associated sequence of approximate solutions converges to the exact solution in the natural norm. Hereby, ideas for the case of lognormal coefficient fields from earlier work of Galvis, Sarkis and Gittelson are used and generalized to the case of positive random coefficient fields with basically arbitrary distributions.

Mathematics Subject Classification. 33C45, 35R60, 65N12, 65N30.

Received June 4, 2012. Revised October 9, 2012.

Published online July 9, 2013.

1. INTRODUCTION

The stochastic Galerkin method developed by Ghanem and Spanos in [22] and refined by many researchers (see *e.g.* the books of LeMaître and Knio [27] and Xiu [40], the Acta Numerica article of Schwab and Gittelson [35] and the references therein) is a powerful tool for the numerical solution of different types of random equations. In many applications like *e.g.* subsurface flows in porous media input parameters cannot be modeled as deterministic functions but depend also on a random component. In this setting, it is desirable to quantify the uncertainty of the modeled system. To this end a variant of the stochastic Galerkin method, based on Wiener-Hermite or generalized polynomial chaos expansions, is mainly used for computing the approximate solution of

Keywords and phrases. Equations with random data, stochastic Galerkin method, generalized polynomial chaos, spectral methods.

* *This work was supported by the Deutsche Forschungsgemeinschaft Priority Programme 1324.*

¹ Mathematisches Institut, Brandenburgische Technische Universität Cottbus, 03013 Cottbus, Germany.
mugler@math.tu-cottbus.de

² Fachgruppe Mathematik, Westsächsische Hochschule Zwickau, 08056 Zwickau, Germany.
hans.joerg.starkloff@fh-zwickau.de

the underlying random equation. The concept of polynomial chaos was introduced by Wiener in [39] and the concept of generalized polynomial chaos expansions was introduced by Xiu and Karniadakis in [38, 42]. Some fundamental properties of such expansions have been investigated in [13, 14].

In this work we focus on elliptic partial differential equations of the form

$$-\nabla \cdot (\kappa \nabla u) = f$$

with random coefficient κ and random force term f together with appropriate boundary conditions. If the random coefficient κ can be bounded strictly away from zero and infinity by deterministic constants, *i.e.*, there are constants $\underline{\kappa}, \bar{\kappa} > 0$, such that

$$0 < \underline{\kappa} \leq \kappa(x, \omega) \leq \bar{\kappa} < \infty \quad \text{a. e. and a. s.},$$

then using the Lax–Milgram theorem it can be shown, that there exists a unique solution of a corresponding stochastic variational problem (see *e.g.* Babuška *et al.* [2–4, 11], Schwab *et al.* [5, 6, 10, 15, 36] or the work of Matthies and Keese [31]). In addition, it can be shown that the standard stochastic Galerkin approach described in these works yields approximate solutions which converge to the exact solution in the natural norm.

However, in some applications where the coefficient describes a physical property which has to be positive (*e.g.* the permeability of a medium) the random coefficient is modeled for instance as a lognormal random field (see *e.g.* [16, 20, 21, 34, 46]) which cannot be bounded by deterministic constants. In this case one cannot use the Lax–Milgram theorem directly to prove the existence of a unique weak solution. Consequently, in recent works Galvis and Sarkis [17] and Gittelsohn [23] prove the existence result using alternative techniques, provide a stochastic Petrov–Galerkin formulation, *i.e.* weighted test functions, to produce a sequence of stochastic Petrov–Galerkin approximations and establish also convergence results for this sequence of numerical solutions. Thereby Gittelsohn investigate among others weak problems posed on the energy space or on spaces with modified measures and proved their well-posedness. He shows that the error of his Galerkin solution can be estimated by a best approximation error in a stronger norm and establish error bounds. Galvis and Sarkis pose a weak problem on different spaces for solution and test functions and prove the existence of a unique weak solution *via* inf-sup-techniques using a White noise framework. They introduce a stochastic Petrov–Galerkin approach by weighting the test functions in order to obtain a well-posed finite-dimensional problem and deduce *a priori* error estimates for the approximation error by using different sequences of weights. In a more recent work [18] they study also spatial and stochastic regularity of the solution. However, the question of convergence (or non-convergence) of corresponding standard stochastic Galerkin approximations, *i.e.* using non-weighted test functions, to the exact solution is not addressed by Gittelsohn or Galvis and Sarkis.

In this work we extend the results of Galvis and Sarkis and Gittelsohn, respectively, and consider more general random coefficients κ , for which we only assume some natural regularity for their realizations and that there are random variables $\kappa_{\min}, \kappa_{\max} > 0$ a. s. satisfying

$$0 < \kappa_{\min}(\omega) \leq \kappa(x, \omega) \leq \kappa_{\max}(\omega) < \infty \quad \text{a. e. and a. s.}$$

Thus, we are not limited only to lognormal random coefficients and allow coefficients with distributions different from a lognormal one and do not use special properties of the normal distribution in the proofs. So, we provide an abstract setting which hopefully make this model more applicable to problems in practice. We investigate the existence of a unique weak solution for elliptic problems with random coefficients of this type as well as the convergence of the standard stochastic Galerkin approximations. Furthermore we construct an example where the standard stochastic Galerkin approximations do not converge in the natural norm to the exact solution. In order to overcome this convergence problem we introduce an alternative stochastic Galerkin approach which uses a modified (a weighted) finite-dimensional test function space corresponding to a stochastic Petrov–Galerkin approach based mainly on the idea of the stochastic Petrov–Galerkin method published by Galvis and Sarkis in [17]. We extend their results for the lognormal case to our more general setting and prove also the convergence of the associated stochastic Petrov–Galerkin approximations to the exact solution in the natural norm.

The paper is organized as follows. In Section 2, we describe our model problem and present two different stochastic variational formulations with standard unweighted and weighted (bi)linear forms, respectively. We discuss the existence of a unique solution for each case. In Section 3, we apply the stochastic Galerkin method to both formulations. We arrive at two discrete stochastic variational problems which can be interpreted as two different stochastic Galerkin approaches for the variational formulation with unweighted (bi)linear forms. One approach corresponds to the standard stochastic Galerkin ansatz, and the other one corresponds to a stochastic Petrov–Galerkin ansatz where the test function space is weighted by a suitable random variable. For the stochastic Petrov–Galerkin method we establish also convergence results w. r. t. the natural norm. Examples, demonstrating that the standard stochastic Galerkin method is not stable and their stochastic Galerkin solutions may not converge in the natural norm, are given in Section 4. Here, the non-convergence is explicitly proved and illustrated numerically. In contrast, we show that our new stochastic Petrov–Galerkin approach yields a convergent sequence of approximate solutions.

2. STOCHASTIC VARIATIONAL PROBLEM FORMULATIONS

Let $D \subset \mathbb{R}^d$, $d \in \mathbb{N}$, be a bounded Lipschitz domain and $(\Omega, \mathfrak{A}, \mathbf{P})$ a probability space consisting of a sample space Ω , a σ -algebra \mathfrak{A} of subsets of Ω and a probability measure \mathbf{P} defined on the σ -algebra \mathfrak{A} . We consider the following boundary-value problem

$$\begin{aligned} -\nabla \cdot (\kappa(x, \omega) \nabla u(x, \omega)) &= f(x, \omega) & x \in D, \omega \in \Omega, \\ u(x, \omega) &= 0 & x \in \partial D, \omega \in \Omega, \end{aligned} \quad (2.1)$$

for a given random coefficient κ and a random forcing term f . In contrast to many articles (compare *e.g.* with Babuška *et al.* [2–4, 11] or Schwab *et al.* [5, 6, 10, 15, 36]) we do not assume that the coefficient κ can be uniformly bounded by deterministic constants. We only assume that each realization of the coefficient κ can be strictly bounded away from zero and infinity.

Assumption 2.1. The coefficient $\kappa : D \times \Omega \rightarrow \mathbb{R}$ is a strongly measurable random variable with values in $L^\infty(D)$ and there exist real-valued random variables κ_{\min} and κ_{\max} on $(\Omega, \mathfrak{A}, \mathbf{P})$ such that

$$0 < \kappa_{\min}(\omega) \leq \kappa(x, \omega) \leq \kappa_{\max}(\omega) < \infty \quad \text{a. e. and a. s.} \quad (2.2)$$

W. l. o. g. we assume that κ_{\min} and κ_{\max} are measurable with respect to the σ -algebra $\sigma(\kappa)$ which is generated by the random variable κ .

Remark 2.2. A random variable η with values in a Banach space \mathcal{B} is called strongly measurable iff there exists a sequence of countably-valued random variables in \mathcal{B} converging almost surely in \mathcal{B} to η (see *e.g.* [25]). We assume strongly measurable coefficients since the space $L^\infty(D)$ is a non-separable Banach space which causes some problems for probabilistic investigations. In most applications coefficients have additional regularity, for example, coefficients are assumed to be continuous or to belong to another separable subspace of $L^\infty(D)$.

Assumption 2.1 is satisfied for random coefficients with constant bounds as well as lognormal random fields under weak conditions but other $L^\infty(D)$ -valued random variables can be treated as well. Assumption 2.1 is not very restrictive: For any non-negative strongly measurable random variable κ in $L^\infty(D)$ where $1/\kappa$ belongs almost surely to $L^\infty(D)$, there exist random bounds $0 < \kappa_{\min} < \kappa_{\max}$ as in (2.2). Though Assumption 2.1 is fairly general the pathwise weak problem, *i.e.* a variational formulation associated with the spatial variables, is a well-posed problem. To see this, we define the corresponding pathwise bilinear form $b(\cdot, \cdot; \omega) : H^1(D) \times H^1(D) \rightarrow \mathbb{R}$ by

$$b(u, v; \omega) = \int_D \kappa(x, \omega) \nabla u(x) \cdot \nabla v(x) \, dx$$

for $\omega \in \Omega$ and we denote by $\langle g, v \rangle_{H^{-1}, \dot{H}^1}$ the duality pairing between $g \in H^{-1}(D)$ and $v \in \dot{H}^1(D)$. Now, assuming that f is a random variable with values in $H^{-1}(D)$, we consider a pathwise weak formulation of the boundary-value problem (2.1).

Problem 2.3 (Pathwise weak formulation). For a given random variable f with values in $H^{-1}(D)$ find a random variable \tilde{u} with values in $\dot{H}^1(D)$, such that

$$b(\tilde{u}(\omega), v; \omega) = \langle f(\omega), v \rangle_{H^{-1}, \dot{H}^1} \quad \text{for all } v \in \dot{H}^1(D) \tag{2.3}$$

holds almost surely.

Since almost every realization of the coefficient κ is bounded by Assumption 2.1 and f is a random variable with values in $H^{-1}(D)$, the Lax–Milgram theorem (see e.g. [8], Thm. 2.7.7) tells us that there exists a mapping $\tilde{u} : \Omega \rightarrow \dot{H}^1(D)$, $\omega \mapsto \tilde{u}(\omega)$ which satisfies the variational equation (2.3) almost surely. Furthermore, the Lax–Milgram theorem provides the estimate

$$\|\tilde{u}(\omega)\|_{\dot{H}^1(D)} \leq C \frac{\|f(\omega)\|_{H^{-1}(D)}}{\kappa_{\min}(\omega)} \quad \text{a. s.} \tag{2.4}$$

where $C > 0$ is a suitable constant independent of $\omega \in \Omega$. Moreover, due to Assumption 2.1 this mapping \tilde{u} can be chosen as a random variable which is measurable w.r.t. the σ -algebra $\sigma(\kappa, f)$ generated by the random variables κ and f . Thus, the mapping \tilde{u} is the unique solution of the pathwise weak formulation, i.e. of Problem 2.3. For further details we refer to [32].

In analogy to purely deterministic boundary-value problems we wish to formulate a variational problem in both spatial and random variables associated with (2.1). This requires a suitable Hilbert space of random variables with values in the Sobolev space $\dot{H}^1(D)$. Now, moments of random variables with values in appropriate functional spaces come into play. As it turns out it is not sufficient to deal with expectations with respect to the basic probability measure \mathbf{P} , but suitable weightings are required. Given a real-valued measurable random function $\varrho : D \times \Omega \rightarrow \mathbb{R}$ with $\varrho(\cdot, \cdot) > 0$ a.e. and a.s. which satisfies Assumption 2.1 we introduce the weighted space

$$\dot{U}_\varrho := \left\{ \eta : \Omega \rightarrow \dot{H}^1(D) \text{ measurable} : \|\eta\|_{\dot{U}_\varrho} < \infty \right\}$$

with

$$\|\eta\|_{\dot{U}_\varrho} := \left(\mathbf{E}_{\mathbf{P}} \left(\int_D |\nabla \eta(x)|^2 \varrho(x) \, dx \right) \right)^{1/2} = \left(\int_\Omega \int_D |\nabla \eta(x, \omega)|^2 \varrho(x, \omega) \, dx \, d\mathbf{P}(\omega) \right)^{1/2}$$

which is a Hilbert space together with the inner product

$$(\eta_1, \eta_2)_{\dot{U}_\varrho} := \mathbf{E}_{\mathbf{P}} \left(\int_D \nabla \eta_1(x) \cdot \nabla \eta_2(x) \varrho(x) \, dx \right).$$

If the weight function ϱ is independent of the spatial argument, i.e., if it is an almost surely positive real-valued random variable $\varrho > 0$, then the weighted spaces

$$U_\varrho^m := L^2(\Omega, \mathfrak{A}, \varrho d\mathbf{P}; H^m(D)) := \left\{ \eta : \Omega \rightarrow H^m(D) \text{ measurable} : \mathbf{E}_{\mathbf{P}} \left(\|\eta\|_{H^m(D)}^2 \varrho \right) < \infty \right\}, \quad m \in \mathbb{Z},$$

and

$$\dot{U}_\varrho^m := L^2(\Omega, \mathfrak{A}, \varrho d\mathbf{P}; \dot{H}^m(D)) := \left\{ \eta : \Omega \rightarrow \dot{H}^m(D) \text{ measurable} : \mathbf{E}_{\mathbf{P}} \left(\|\eta\|_{\dot{H}^m(D)}^2 \varrho \right) < \infty \right\}, \quad m \in \mathbb{N}_0$$

equipped with the norm $\|\eta\|_{U_\varrho^m} := \left(\mathbf{E}_\mathbf{P} \left(\|\eta\|_{H^m(D)}^2 \varrho\right)\right)^{1/2} = \left(\mathbf{E}_\mathbf{P} \left(\left(\sum_{|\alpha|\leq m} \|\partial^\alpha \eta\|_{L^2(D)}^2\right) \varrho\right)\right)^{1/2}$ and the inner product $(\eta_1, \eta_2)_{U_\varrho^m} := \mathbf{E}_\mathbf{P} \left((\eta_1, \eta_2)_{H^m(D)} \varrho\right)$ are also Hilbert spaces. Note that for $\eta \in \mathring{U}_\varrho^1$, *i.e.* for random variables η with values in $\mathring{H}^1(D)$, the seminorm

$$|\eta|_{U_\varrho^1} := \left(\mathbf{E}_\mathbf{P} \left(|\eta|_{H^1(D)}^2\right)\right)^{1/2} = \left(\mathbf{E}_\mathbf{P} \left(\varrho \int_D |\nabla \eta(x)|^2 dx\right)\right)^{1/2}$$

is equivalent to the norm $\|\cdot\|_{U_\varrho^1}$. If the space $L^2(\Omega, \mathfrak{A}, \varrho d\mathbf{P}) := \{\eta : \Omega \rightarrow \mathbb{R} \text{ measurable: } \mathbf{E}_\mathbf{P} (|\eta|^2 \varrho) < \infty\}$ is separable, then there exist isomorphisms to the corresponding tensor product spaces, *i.e.*

$$U_\varrho^m \cong H^m(D) \otimes L^2(\Omega, \mathfrak{A}, \varrho d\mathbf{P}), \quad m \in \mathbb{Z},$$

and

$$\mathring{U}_\varrho^m \cong \mathring{H}^m(D) \otimes L^2(\Omega, \mathfrak{A}, \varrho d\mathbf{P}), \quad m \in \mathbb{N}_0.$$

Furthermore, the dual space of \mathring{U}_ϱ^m can be identified with the space $U_{\varrho^{-1}}^{-m}$. For convenience we denote the unweighted spaces, *i.e.* $\varrho \equiv 1$, by U^m and \mathring{U}^m . On occasion we may replace the probability measure \mathbf{P} in the definition of U^m and \mathring{U}^m by another probability measure \mathbf{Q} and write $U_\mathbf{Q}^m$ and $\mathring{U}_\mathbf{Q}^m$, respectively, in this case.

We define the bilinear form

$$a(u, v) := \mathbf{E}_\mathbf{P} \left(\int_D \kappa(x) \nabla u(x) \cdot \nabla v(x) dx \right) = \int_\Omega \int_D \kappa(x, \omega) \nabla u(x, \omega) \cdot \nabla v(x, \omega) dx d\mathbf{P}(\omega),$$

and the linear form

$$\ell(v) := \mathbf{E}_\mathbf{P} \left(\langle f, v \rangle_{H^{-1}, \mathring{H}^1} \right) = \int_\Omega \langle f(\omega), v(\omega) \rangle_{H^{-1}, \mathring{H}^1} d\mathbf{P}(\omega).$$

Given Assumption 2.1 the weak problem posed in unweighted spaces, *i.e.*, for a given $f \in U^{-1}$ find $u \in \mathring{U}^1$ such that there holds

$$a(u, v) = \ell(v) \quad \text{for all } v \in \mathring{U}^1 \tag{2.5}$$

is in general *not* well-posed. For example, it may happen that there does not even exist a solution in \mathring{U}^1 for a given $f \in U^{-1}$ (see [32]). Note, that this stands in marked contrast to the case when the coefficient κ can be strictly bounded away from zero and infinity. Thus, we have to use an alternative stochastic variational formulation and arrive at a first weak problem posed on the associated energy space as in [2, 23].

Problem 2.4 (Unweighted weak formulation). For a given f find $\hat{u} \in \mathring{U}_\kappa$, such that

$$a(\hat{u}, v) = \ell(v) \quad \text{for all } v \in \mathring{U}_\kappa.$$

Thereby the term *unweighted* weak formulation refers to the fact that standard *unweighted* (bi)linear forms are used although the solution and test function spaces are weighted spaces. Afterwards we introduce also a *weighted* weak formulation, see Problem 2.6, where the (bi)linear forms are weighted.

To prove existence of a solution to Problem 2.4 we employ weighted spaces with the lower bound κ_{\min} in Assumption 2.1. We note that such a lower bound is not uniquely defined and some freedom in its choice is sometimes useful. Due to (2.2) there exist several real-valued random variables $\mu > 0$ a. s. satisfying

$$0 < c\mu(\omega) \leq \kappa(x, \omega) \quad \text{a. e. and a. s.} \tag{2.6}$$

for a constant $c > 0$. The random variable κ_{\min} is only a special case satisfying (2.6) with constant $c = 1$. W. l. o. g. we can assume that the random variable μ is also $\sigma(\kappa)$ -measurable.

Theorem 2.5. *If Assumption 2.1 holds and μ satisfies (2.6), then for any $f \in U_{\mu^{-1}}^{-1}$ there exists a unique $\hat{u} \in \mathring{U}_{\kappa}$ solving Problem 2.4.*

Proof. Since for $f \in U_{\mu^{-1}}^{-1}$ we have $\ell \in U_{\mu^{-1}}^{-1} \subset (\mathring{U}_{\kappa})^*$, the result follows immediately from Riesz' representation theorem. \square

Therefore Problem 2.4 is well-posed in the given weighted spaces. However, we are interested in a solution \hat{u} in the space \mathring{U}^1 rather than the space \mathring{U}_{κ} . For deterministic lower bounds the space \mathring{U}_{κ} is continuously embedded in \mathring{U}^1 . Unfortunately, this is in general not true if the coefficient κ can be bounded from below only by a random variable, *e.g.* μ . Clearly, if we assume a forcing term $f \in U_{\mu^{-2}}^{-1}$, then our discussion of the pathwise variational problem and the corresponding estimate (2.4) above tell us that the weak solution \hat{u} belongs also to the space \mathring{U}^1 , see also Corollary 2.8. In general, however, we cannot estimate the \mathring{U}^1 -norm in terms of the \mathring{U}_{κ} -norm.

A possible remedy is to recast Problem 2.4 above. By multiplying the pathwise variational equation (2.3) with the reciprocal of the random variable μ we obtain a random coefficient κ/μ which can now be strictly bounded away from zero by the constant c in inequality (2.6). The corresponding bilinear form a_{μ} and linear form ℓ_{μ} are the weighted versions of the bilinear form a and linear form ℓ , respectively, *i.e.*

$$a_{\mu}(u, v) := \mathbf{E}_{\mathbf{P}} \left(\frac{1}{\mu} \int_D \kappa(x) \nabla u(x) \cdot \nabla v(x) \, dx \right) = \int_{\Omega} \int_D \frac{\kappa(x, \omega)}{\mu(\omega)} \nabla u(x, \omega) \cdot \nabla v(x, \omega) \, dx \, d\mathbf{P}(\omega)$$

and

$$\ell_{\mu}(v) := \mathbf{E}_{\mathbf{P}} \left(\frac{1}{\mu} \langle f, v \rangle_{H^{-1}, \mathring{H}^1} \right) = \int_{\Omega} \frac{1}{\mu(\omega)} \langle f(\omega), v(\omega) \rangle_{H^{-1}, \mathring{H}^1} \, d\mathbf{P}(\omega).$$

Problem 2.6 (Weighted weak formulation). For a given f find $\bar{u} \in \mathring{U}_{\mu}^{\kappa}$, such that there holds

$$a_{\mu}(\bar{u}, v) = \ell_{\mu}(v) \quad \text{for all } v \in \mathring{U}_{\mu}^{\kappa}.$$

Theorem 2.7. *If Assumption 2.1 holds and μ satisfies (2.6), then for any $f \in U_{\mu^{-2}}^{-1}$ there exists a unique $\bar{u} \in \mathring{U}_{\mu}^{\kappa}$ solving Problem 2.6.*

Proof. Since for any $f \in U_{\mu^{-2}}^{-1}$ there holds $\ell_{\mu} \in U^{-1} \subset (\mathring{U}_{\mu}^{\kappa})^*$ the result follows immediately from Riesz' representation theorem. \square

Now, the space $\mathring{U}_{\mu}^{\kappa}$ is continuously embedded in \mathring{U}^1 . Thus, there exists a constant $C > 0$ such that

$$\|u\|_{U^1} \leq C \|u\|_{\mathring{U}_{\mu}^{\kappa}} \quad \text{for all } u \in \mathring{U}_{\mu}^{\kappa}.$$

Corollary 2.8. *Under Assumption 2.1 together with (2.6) for every $f \in U_{\mu^{-1}}^{-1} \cap U_{\mu^{-2}}^{-1}$ the weak formulations Problems 2.4 and 2.6 have unique solutions \hat{u} and \bar{u} in \mathring{U}^1 , respectively, which are almost surely equal to the pathwise solution \tilde{u} of the pathwise weak formulation Problem 2.3, *i.e.*,*

$$\hat{u} = \bar{u} = \tilde{u} \quad \text{a. s.}$$

as random variables in $\mathring{H}^1(D)$. Moreover, both weak solutions \hat{u} and \bar{u} are $\sigma(\kappa, f)$ -measurable.

Proof. Analogously to Theorem 3.3 in [32], we can show that both weak solutions \hat{u} and \bar{u} solve the pathwise weak formulation Problem 2.3 almost surely, the solution \tilde{u} of which is unique and $\sigma(\kappa, f)$ -measurable. \square

Therefore, we assume throughout the rest of the paper that $f \in U_{\mu^{-1}}^{-1} \cap U_{\mu^{-2}}^{-1}$ and do not distinguish between the three weak solutions, but refer to *the* weak solution denoted by \hat{u} .

Remark 2.9. The results in this article hold also for other boundary conditions as long as the bilinear forms a and a_μ form inner products on the corresponding solution spaces.

3. STOCHASTIC GALERKIN METHOD

The stochastic Galerkin method or stochastic finite element method (SFEM) is a very popular approach for the numerical solution of random equations. It is based on a stochastic variational formulation and uses sequences of finite-dimensional solution and test function spaces in order to obtain approximate solutions. However, to get reliable results a rigorous convergence analysis must be carried out. In the next section we will construct examples showing that the standard stochastic Galerkin approach applied to our model problem (2.1) is not stable and may not lead to a sequence of approximate solutions which converge in the natural norm to the exact solution. Therefore, we present in this section an alternative stochastic Galerkin approach and prove its stability.

The definition of the finite-dimensional solution and test function spaces requires a discretization in the spatial as well as in the stochastic dimension. Therefore, we parameterize the random input data by finitely many or countably many real-valued random variables which are referred to as basic random variables. For example, the popular Karhunen–Loève expansion of random fields (see *e.g.* [28]) yields such a representation. As observed, for example, by Karniadakis, Xiu and their co-authors in [29, 41–45] the choice/distribution of the basic random variables affects the approximation quality of the stochastic Galerkin solution. We mention that not only the distribution of the basic random variables is crucial here.

In the following we assume that for our model problem (2.1) a finite or infinite-dimensional random vector $\xi := (\xi_i)_{i \in I_\xi}$ with values in $\mathbb{R}^{|I_\xi|}$ is given such that κ and f are both measurable with respect to $\sigma(\xi)$. Hence, due to the Doob–Dynkin lemma (*e.g.* in [26]) there exist measurable functions $\kappa_\xi : \mathbb{R}^{|I_\xi|} \rightarrow L^\infty(D)$ and $f_\xi : \mathbb{R}^{|I_\xi|} \rightarrow H^{-1}(D)$ satisfying $\kappa = \kappa_\xi(\xi)$ in $L^\infty(D)$ and $f = f_\xi(\xi)$ in $H^{-1}(D)$ almost surely. Because of the assumptions on f we obtain a weak solution \hat{u} in \mathring{U}^1 . By Corollary 2.8 the weak solution \hat{u} is $\sigma(\kappa, f)$ -measurable and thus it is also $\sigma(\xi)$ -measurable. Consequently there exists a measurable function $\hat{u}_\xi : \mathbb{R}^{|I_\xi|} \rightarrow \mathring{H}^1(D)$, such that $\hat{u} = \hat{u}_\xi(\xi)$ a.s. in $\mathring{H}^1(D)$. In our setting all random variables of interest are $\sigma(\xi)$ -measurable, thus, we suppose w.l.o.g. that $\mathfrak{A} = \sigma(\xi)$ holds. We then identify the space $U^1 = L^2(\Omega, \sigma(\xi), \mathbf{P}; \mathring{H}^1(D))$ as appropriate solution space.

Moreover, we assume that the random vector ξ has the distribution $F_\xi^{\mathbf{P}}$ and each of the basic random variables ξ_i , $i \in I_\xi$, possesses finite moments of arbitrary order, *i.e.* $\mathbf{E}_{\mathbf{P}}|\xi_i|^n < \infty$ for all $n \in \mathbb{N}$, and has a continuous distribution function $F_{\xi_i}^{\mathbf{P}}$. Here, the distributions $F_\xi^{\mathbf{P}}$ and $F_{\xi_i}^{\mathbf{P}}$ are understood as the distributions of the random variables ξ and ξ_i w.r. t. the probability measure \mathbf{P} , *i.e.*, these are the distributions of ξ and ξ_i as random variables (with values in the spaces $\mathbb{R}^{|I_\xi|}$ and \mathbb{R} , respectively) on the probability space $(\Omega, \mathfrak{A}, \mathbf{P})$. On occasion we will consider these random variables also on probability spaces with other probability measures and thus their distributions change accordingly. For this reason we attach the corresponding probability measure to the distribution symbol of a random variable.

We choose the finite-dimensional solution space

$$U_{N,K,hp} := U_{hp} \otimes U_{N,K} \subset \mathring{H}^1(D) \otimes L^2(\Omega, \sigma(\xi), \mathbf{P}) \simeq L^2(\Omega, \sigma(\xi), \mathbf{P}; \mathring{H}^1(D)),$$

where U_{hp} is a finite-dimensional subspace of $\mathring{H}^1(D)$ obtained by a suitable version (h -, p - or hp -) of the deterministic finite element method and $U_{N,K}$ is a finite-dimensional subspace of $L^2(\Omega, \sigma(\xi_1, \dots, \xi_K), \mathbf{P}) \subset L^2(\Omega, \sigma(\xi), \mathbf{P})$ where $\{1, \dots, K\} \subseteq I_\xi$. Since we want to use generalized polynomial chaos (see *e.g.* [41, 42]), *i.e.*

polynomials in the underlying basic random variables $\xi_i, i \in I_\xi$, we construct the finite dimensional space $U_{N,K}$ as follows

$$U_{N,K} := \text{span} \left\{ \xi^\alpha := \prod_{i \in I_\xi} \xi_i^{\alpha_i}, \alpha \in \Lambda_{N,K} \right\},$$

where $\Lambda_{N,K}$ denotes a finite index set

$$\Lambda_{N,K} \subset \Lambda := \left\{ \alpha \in \mathbb{N}_0^{|I_\xi|} : \alpha \text{ has only finitely many non-zero entries} \right\}$$

such that the sequence of these finite-dimensional subsets $(\Lambda_{N,K})_{N,K}$ satisfies $\bigcup_{N,K} \Lambda_{N,K} = \Lambda$. For example, we may choose $\Lambda_{N,K} = J_1$, where

$$J_1 := \{ \alpha \in \Lambda : \alpha_i = 0 \ \forall i > K, |\alpha| \leq N \}, \quad |\alpha| := \sum_{i \in I_\xi} \alpha_i,$$

selects (complete) multivariate polynomials with bounded total degree, or $\Lambda_{N,K} = J_2$, where

$$J_2 := \{ \alpha \in \Lambda : \alpha_i = 0 \ \forall i > K, \alpha_i \leq N \ \forall 1 \leq i \leq K \}$$

selects (tensor product) multivariate polynomials with uniformly bounded degree of the individual basic random variables. Then, the discrete versions of the unweighted and the weighted weak formulation introduced in Section 2 read as follows.

Problem 3.1 (Discrete unweighted weak formulation). For a given f find $\hat{u}_{N,K,hp} \in U_{N,K,hp}$, such that there holds

$$a(\hat{u}_{N,K,hp}, v) = \ell(v) \quad \text{for all } v \in U_{N,K,hp}.$$

Problem 3.2 (Discrete weighted weak formulation). For a given f find $\bar{u}_{N,K,hp} \in U_{N,K,hp}$, such that there holds

$$a_\mu(\bar{u}_{N,K,hp}, v) = \ell_\mu(v) \quad \text{for all } v \in U_{N,K,hp}.$$

So Problem 3.1 is exactly that discrete weak problem which arises if we apply the stochastic Galerkin method to the weak problem posed on the unweighted spaces (2.5) no matter if this is well-posed or not. Alternatively, this means, Problem 3.1 corresponds to the discrete weak problem which results from applying the stochastic Galerkin method to the standard weak problem arising in the case of a coefficient κ which can be bounded by constants. Therefore, we call the approach yielding the discrete unweighted weak problem, *i.e.* Problem 3.1, the *standard* stochastic Galerkin approach. Thus, the solution $\hat{u}_{N,K,hp}$ is called the *standard* stochastic Galerkin solution or *standard* stochastic Galerkin approximation, respectively.

We then make the following observations.

Theorem 3.3. *Let Assumption 2.1 be fulfilled and $f \in U_{\mu^{-2}}^{-1} \cap U_{\mu^{-1}}^{-1}$.*

- (i) *If $\kappa_{\max} \in L^r(\Omega, \mathfrak{A}, \mathbf{P})$ for some $r > 1$, then there exists a unique stochastic Galerkin solution $\hat{u}_{N,K,hp}$ in $U_{N,K,hp}$ solving Problem 3.1. Moreover, $\hat{u}_{N,K,hp}$ is the best approximation in $U_{N,K,hp}$ of the weak solution \hat{u} with respect to the norm of \hat{U}_κ , *i.e.* $\|\hat{u} - \hat{u}_{N,K,hp}\|_{\hat{U}_\kappa} = \inf_{z \in U_{N,K,hp}} \|\hat{u} - z\|_{\hat{U}_\kappa}$.*
- (ii) *If $\kappa_{\max}/\mu \in L^r(\Omega, \mathfrak{A}, \mathbf{P})$ for some $r > 1$, then there exists a unique stochastic Galerkin solution $\bar{u}_{N,K,hp}$ in $U_{N,K,hp}$ solving Problem 3.2. Moreover, $\bar{u}_{N,K,hp}$ is the best approximation in $U_{N,K,hp}$ of the weak solution \hat{u} with respect to the norm of $\hat{U}_{\frac{\kappa}{\mu}}$, *i.e.* $\|\hat{u} - \bar{u}_{N,K,hp}\|_{\hat{U}_{\frac{\kappa}{\mu}}} = \inf_{z \in U_{N,K,hp}} \|\hat{u} - z\|_{\hat{U}_{\frac{\kappa}{\mu}}}$.*

Proof. In order to prove this result we show that the finite-dimensional subspace $U_{N,K,hp}$ is a subset of \hat{U}_κ and \hat{U}_μ , respectively. For any $v = \xi^\alpha \varphi$, $\alpha \in \Lambda$ and $\varphi \in U_{hp}$ this can be deduced for (i) from the integrability condition of κ_{\max} by

$$\begin{aligned} \|v\|_{\hat{U}_\kappa}^2 &= a(v, v) = \int_{\Omega} \int_D \kappa(x, \omega) |\nabla v(x, \omega)|^2 dx d\mathbf{P}(\omega) \leq \int_{\Omega} \int_D \kappa_{\max}(\omega) \xi(\omega)^{2\alpha} |\nabla \varphi(x)|^2 dx d\mathbf{P}(\omega) \\ &= |\varphi|_{H^1(D)}^2 \mathbf{E}\mathbf{P}(\kappa_{\max} \xi^{2\alpha}) \leq |\varphi|_{H^1(D)}^2 (\mathbf{E}\mathbf{P} \kappa_{\max}^r)^{1/r} (\mathbf{E}\mathbf{P} \xi^{2s\alpha})^{1/s} < \infty \end{aligned}$$

and for (ii) from the integrability condition on κ_{\max}/μ by

$$\begin{aligned} \|v\|_{\hat{U}_\mu}^2 &= a_\mu(v, v) = \int_{\Omega} \int_D \frac{\kappa(x, \omega)}{\mu(\omega)} |\nabla v(x, \omega)|^2 dx d\mathbf{P}(\omega) \leq \int_{\Omega} \int_D \frac{\kappa_{\max}(\omega)}{\mu(\omega)} \xi(\omega)^{2\alpha} |\nabla \varphi(x)|^2 dx d\mathbf{P}(\omega) \\ &= |\varphi|_{H^1(D)}^2 \mathbf{E}\mathbf{P} \left(\frac{\kappa_{\max}}{\mu} \xi^{2\alpha} \right) \leq |\varphi|_{H^1(D)}^2 \left(\mathbf{E}\mathbf{P} \left(\frac{\kappa_{\max}}{\mu} \right)^r \right)^{1/r} (\mathbf{E}\mathbf{P} \xi^{2s\alpha})^{1/s} < \infty, \end{aligned}$$

respectively, where $1 < s < \infty$ with $\frac{1}{r} + \frac{1}{s} = 1$ in both cases. Thus, there exist unique best approximations, *i.e.* orthogonal projections $\hat{u}_{N,K,hp} := \Pi_{U_{N,K,hp}} \hat{u}$ onto $U_{N,K,hp}$ in \hat{U}_κ and $\bar{u}_{N,K,hp} := \Pi_{U_{N,K,hp}} \hat{u}$ onto $U_{N,K,hp}$ in \hat{U}_μ satisfying

$$a(\hat{u} - \hat{u}_{N,K,hp}, v) = 0 \quad \text{and} \quad a_\mu(\hat{u} - \bar{u}_{N,K,hp}, v) = 0, \quad \text{respectively,}$$

for all $v \in U_{N,K,hp}$. In other words, since $a(\hat{u}, v) = \ell(v)$ and $a_\mu(\bar{u}, v) = \ell_\mu(v)$, we have arrived at

$$a(\hat{u}_{N,K,hp}, v) = \ell(v) \quad \text{and} \quad a_\mu(\bar{u}_{N,K,hp}, v) = \ell_\mu(v), \quad \text{respectively,}$$

for all $v \in U_{N,K,hp}$. □

Notably, the discrete version of the weighted weak problem can also be interpreted as a discrete version of the unweighted variational formulation by using a so-called stochastic Petrov–Galerkin ansatz where the finite-dimensional test function space differs from the finite-dimensional solution space. If we define the weighted finite-dimensional test function space

$$V_{N,K,hp} := \left\{ \frac{u}{\mu} : u \in U_{N,K,hp} \right\},$$

we see that the problem of finding $\bar{u}_{N,K,hp} \in U_{N,K,hp}$ which satisfies

$$a_\mu(\bar{u}_{N,K,hp}, v) = \ell_\mu(v) \quad \text{for all } v \in U_{N,K,hp}$$

is equivalent to finding $\bar{u}_{N,K,hp} \in U_{N,K,hp}$ which satisfies

$$a(\bar{u}_{N,K,hp}, v) = \ell(v) \quad \text{for all } v \in V_{N,K,hp}.$$

This means, the weighting $1/\mu$ of the (bi)linear forms is transferred to the finite-dimensional test function space. Thus, we arrive at the discrete weighted weak problem using the unweighted bilinear form a and linear form ℓ (instead of a_μ and ℓ_μ) on the finite-dimensional solution space $U_{N,K,hp}$ and the finite-dimensional test function space $V_{N,K,hp}$ which equals the space $U_{N,K,hp}$ weighted by $1/\mu$.

Altogether, this implies, that we have identified two different stochastic Galerkin approaches for the very same variational equation: We are looking for a solution in $U_{N,K,hp}$ which satisfies the equation

$$a(u, v) = \ell(v)$$

for all $v \in U_{N,K,hp}$ or for all $v \in V_{N,K,hp}$, respectively. In general this defines two different Galerkin projections and hence two different solutions in the same finite-dimensional space.

If we know in addition that the finite-dimensional solution spaces $(U_{N,K,hp})_{N,K,hp}$ lie dense in the space \mathring{U}_κ or \mathring{U}_μ , respectively, we can conclude that the approximation errors tend to zero, *i.e.*

$$\|\hat{u} - \hat{u}_{N,K,hp}\|_{\mathring{U}_\kappa} \rightarrow 0 \quad \text{and} \quad \|\hat{u} - \bar{u}_{N,K,hp}\|_{\mathring{U}_\mu} \rightarrow 0.$$

In many cases of practical interest, however, we are not interested in the convergence in these weighted norms but in the approximation error in the unweighted, the U^1 -norm which therefore we call also the *natural norm*. Here the weighted weak formulation is advantageous due to the continuous embedding of \mathring{U}_μ in \mathring{U}^1 which allows us to deduce

$$\|\hat{u} - \bar{u}_{N,K,hp}\|_{U^1} \leq C \|\hat{u} - \bar{u}_{N,K,hp}\|_{\mathring{U}_\mu} = C \inf_{z \in U_{N,K,hp}} \|\hat{u} - z\|_{\mathring{U}_\mu} \tag{3.1}$$

together with an upper estimate of the form

$$\|\hat{u} - \bar{u}_{N,K,hp}\|_{U^1} \leq C \inf_{z \in U_{N,K,hp}} \|\hat{u} - z\|_{U^1_{\frac{\kappa_{\max}}{\mu}}}. \tag{3.2}$$

Hence, we can estimate the approximation error of the stochastic Galerkin projection $\bar{u}_{N,K,hp}$ in U^1 by a best approximation error in a stronger norm related to a weighting by the real-valued random variable κ_{\max}/μ . Due to the assumption $\kappa_{\max}/\mu \in L^r(\Omega, \mathfrak{A}, \mathbf{P})$ for some $r > 1$ the expectation $c_{\frac{\kappa_{\max}}{\mu}} := \mathbf{E}_{\mathbf{P}}(\kappa_{\max}/\mu)$ is finite and the measure \mathbf{Q} defined by

$$d\mathbf{Q} := \frac{1}{c_{\frac{\kappa_{\max}}{\mu}}} \frac{\kappa_{\max}}{\mu} d\mathbf{P}$$

is a probability measure. In the following we consider the function spaces $U_{\mathbf{Q}}^m$ and $\mathring{U}_{\mathbf{Q}}^m$ instead of $U_{\frac{\kappa_{\max}}{\mu}}^m$ and $\mathring{U}_{\frac{\kappa_{\max}}{\mu}}^m$, $m \in \mathbb{Z}$, which coincide as sets with $U_{\frac{\kappa_{\max}}{\mu}}^m$ and $\mathring{U}_{\frac{\kappa_{\max}}{\mu}}^m$ but are much easier to handle due to the corresponding probability space $(\Omega, \mathfrak{A}, \mathbf{Q})$ at hand where w.l.o.g. we have assumed $\mathfrak{A} = \sigma(\xi)$.

Corollary 3.4. *If $f \in U_{\mu^{-2}}^{-1} \cap U_{\mu^{-1}}^{-1}$, $\kappa_{\max}/\mu \in L^r(\Omega, \mathfrak{A}, \mathbf{P})$ for some $r > 1$ and $\hat{u} \in \mathring{U}_{\mathbf{Q}}^1$, then there holds*

$$\|\hat{u} - \bar{u}_{N,K,hp}\|_{U^1} \leq C \inf_{z \in U_{N,K,hp}} \|\hat{u} - z\|_{U_{\mathbf{Q}}^1} \tag{3.3}$$

for a constant $C > 0$ independent of N, K, h, p and \hat{u} .

Proof. This follows immediately from the assumptions and the estimates (3.1) and (3.2). □

Remark 3.5. If the coefficient κ can be bounded from below by a positive constant, then the finite-dimensional test function space $V_{N,K,hp}$ coincides with the solution space $U_{N,K,hp}$ and thus weighting is not necessary. If, additionally, the coefficient κ can also be bounded from above by a constant, then the result of Corollary 3.4 is equivalent to the deterministic theory.

3.1. Convergence analysis

Next, we proceed with a closer look at the estimate of the approximation error (3.3) and identify and analyze different sources of discretization errors. To this end, we introduce three different orthogonal projections and

denote by

$$\begin{aligned} \Pi_{\dot{U}_{\mathbf{Q},N,K,hp}^1} &: \dot{U}_{\mathbf{Q}}^1 \rightarrow U_{N,K,hp} \\ &\text{the orthogonal projection onto } U_{N,K,hp}, \\ \Pi_{\dot{U}_{\mathbf{Q},N,K}^1} &: \dot{U}_{\mathbf{Q}}^1 \rightarrow \dot{H}^1(D) \otimes U_{N,K} \\ &\text{the orthogonal projection onto } \dot{H}^1(D) \otimes U_{N,K}, \\ \Pi_{\dot{U}_{\mathbf{Q},K}^1} &: \dot{U}_{\mathbf{Q}}^1 \rightarrow \dot{H}^1(D) \otimes L^2(\Omega, \sigma(\xi_1, \dots, \xi_K), \mathbf{Q}) \\ &\text{the orthogonal projection onto} \\ &\dot{H}^1(D) \otimes L^2(\Omega, \sigma(\xi_1, \dots, \xi_K), \mathbf{Q}) \simeq L^2(\Omega, \sigma(\xi_1, \dots, \xi_K), \mathbf{Q}; \dot{H}^1(D)), \end{aligned}$$

respectively. All projections are carried out w. r. t. the inner product associated with the $U_{\mathbf{Q}}^1$ -norm. Assuming $\hat{u} \in \dot{U}_{\mathbf{Q}}^1$ we can estimate the approximation error of the stochastic Petrov–Galerkin approximation $\bar{u}_{N,K,hp}$ to the exact solution by using (3.3) for $z = \Pi_{\dot{U}_{\mathbf{Q},N,K,hp}^1} \hat{u}$ as follows

$$\|\hat{u} - \bar{u}_{N,K,hp}\|_{U^1} \leq C \left[\|\hat{u} - \Pi_{\dot{U}_{\mathbf{Q},K}^1} \hat{u}\|_{U_{\mathbf{Q}}^1} + \|\Pi_{\dot{U}_{\mathbf{Q},K}^1} \hat{u} - \Pi_{\dot{U}_{\mathbf{Q},N,K}^1} \hat{u}\|_{U_{\mathbf{Q}}^1} + \|\Pi_{\dot{U}_{\mathbf{Q},N,K}^1} \hat{u} - \Pi_{\dot{U}_{\mathbf{Q},N,K,hp}^1} \hat{u}\|_{U_{\mathbf{Q}}^1} \right]. \quad (3.4)$$

Hence this error has three components, namely an approximation error induced by discretizing in the spatial dimension

$$e_1 := \|\Pi_{\dot{U}_{\mathbf{Q},N,K}^1} \hat{u} - \Pi_{\dot{U}_{\mathbf{Q},N,K,hp}^1} \hat{u}\|_{U_{\mathbf{Q}}^1}$$

and two approximation errors arising from the discretization in the stochastic dimension: We consider only polynomials in the first K basic random variables

$$e_2 := \|\hat{u} - \Pi_{\dot{U}_{\mathbf{Q},K}^1} \hat{u}\|_{U_{\mathbf{Q}}^1}$$

and bound their polynomial degrees by N

$$e_3 := \|\Pi_{\dot{U}_{\mathbf{Q},K}^1} \hat{u} - \Pi_{\dot{U}_{\mathbf{Q},N,K}^1} \hat{u}\|_{U_{\mathbf{Q}}^1}.$$

Again, the first error term e_1 is the error arising from the spatial discretization of the Sobolev space $\dot{H}^1(D)$ by the finite-dimensional finite element space U_{hp} . Therefore, the spatial approximation error can be bounded using standard arguments from the theory of the deterministic finite element methods. Here, we employ an hp -version of the finite element method (see e.g. [9]). For example, under the assumptions in Section 5.8.3 in [9] and given that the domain D is obtained from the reference domain $\hat{D} := (-1, 1)^d$, $d = 1, 2, 3$, by a smooth, invertible and regular mapping $F : \hat{D} \rightarrow D$ (in the sense of Section 5.8.3 in [9]) there holds the following.

Corollary 3.6. *For $\hat{u} \in U_{\mathbf{Q}}^k \cap \dot{U}_{\mathbf{Q}}^1$, $k \geq 2$, there holds*

$$e_1 = \|\Pi_{\dot{U}_{\mathbf{Q},N,K}^1} \hat{u} - \Pi_{\dot{U}_{\mathbf{Q},N,K,hp}^1} \hat{u}\|_{U_{\mathbf{Q}}^1} \leq \tilde{C} h^{m-1} p^{-(k-1)} \|\hat{u}\|_{U_{\mathbf{Q}}^k}$$

where $m = \min\{k, p + 1\}$ and the constant $\tilde{C} > 0$ is independent of N , K , h , p and \hat{u} .

Proof. Based on the results in [9] (see e.g. Eq. (5.8.25)) for the deterministic case we find

$$\|\Pi_{\dot{U}_{\mathbf{Q},N,K}^1} \hat{u}(\omega) - \Pi_{\dot{U}_{\mathbf{Q},N,K,hp}^1} \hat{u}(\omega)\|_{H^1(D)} \leq C h^{m-1} p^{-(k-1)} \|\hat{u}(\omega)\|_{H^k(D)} \quad \text{a. s.}$$

with $m = \min\{k, p + 1\}$ and a suitable constant $C > 0$ independent of N , K , h , p , \hat{u} and $\omega \in \Omega$. Squaring the formula above and taking the expectation $\mathbf{E}_{\mathbf{Q}}$ w. r. t. the probability measure \mathbf{Q} yields

$$\mathbf{E}_{\mathbf{Q}} \left\| \Pi_{\dot{U}_{\mathbf{Q},N,K}^1} \hat{u} - \Pi_{\dot{U}_{\mathbf{Q},N,K,hp}^1} \hat{u} \right\|_{H^1(D)}^2 \leq C^2 h^{2(m-1)} p^{-2(k-1)} \mathbf{E}_{\mathbf{Q}} \left(\|\hat{u}\|_{H^k(D)}^2 \right) = C^2 h^{2(m-1)} p^{-2(k-1)} \|\hat{u}\|_{U_{\mathbf{Q}}^k}^2$$

with $m = \min\{k, p + 1\}$ and a suitable constant $C > 0$ independent of N , K , h , p and \hat{u} . □

Remark 3.7. If, in addition, the mapping F is affine we can prove, analogously to [9], the result

$$e_1 = \|\Pi_{\dot{U}_{\mathbf{Q},N,K}^1} \hat{u} - \Pi_{\dot{U}_{\mathbf{Q},N,K,h,p}^1} \hat{u}\|_{U_{\mathbf{Q}}^1} \leq Ch^{m-1}p^{-(k-1)} \left(\mathbf{E}_{\mathbf{Q}} |\hat{u}|_{H^{m;k}(D)}^2 \right)^{1/2}$$

with $m = \min\{k, p + 1\}$, a constant $C > 0$ and where $|\cdot|_{H^{m;k}(D)}$ denotes the seminorm

$$|u|_{H^{m;k}(D)} := \left(\sum_{j=m}^k \sum_{i=1}^d \left\| \frac{\partial^j}{\partial x_i^j} u \right\|_{L^2(D)}^2 \right)^{1/2}$$

where only pure derivatives in each spatial direction appear (cf. Canuto *et al.* [9], p. 318 and formula (5.8.12) on p. 314).

Clearly, additional smoothness of the weak solution in the spatial argument implies faster convergence. Note that the approximation error e_1 does not depend on the basic random variables.

The second error term usually occurs only if countably many basic random variables are used to represent the random input data. For finitely many basic random variables $\xi = (\xi_i)_{i \in I_\xi}$, *i.e.* $I_\xi = \{1, \dots, M\} \subset \mathbb{N}$, the finite-dimensional solution space $U_{N,K,h,p}$ is usually defined for $K = M$. Thus, the orthogonal projection of the weak solution \hat{u} onto the space $\dot{H}^1(D) \otimes L^2(\Omega, \sigma(\xi_1, \dots, \xi_K), \mathbf{Q})$ is the weak solution itself and the error e_2 equals zero. For a countable set of basic random variables it can be shown that the error e_2 converges to zero.

Corollary 3.8. *Let $I_\xi = \mathbb{N}$, then there holds for $\hat{u} \in \dot{U}_{\mathbf{Q}}^1$*

$$e_2 = \|\hat{u} - \Pi_{\dot{U}_{\mathbf{Q},K}^1} \hat{u}\|_{U_{\mathbf{Q}}^1} \rightarrow 0, \quad K \rightarrow \infty.$$

Proof. The orthogonal projection $\Pi_{\dot{U}_{\mathbf{Q},K}^1} v$ onto $L^2(\Omega, \sigma(\xi_1, \dots, \xi_K), \mathbf{Q}; \dot{H}^1(D))$ is equal to the conditional expectation of $v \in \dot{U}_{\mathbf{Q}}^1$ w.r.t. the σ -algebra spanned by the random variables ξ_1, \dots, ξ_K (see *e.g.* [37]). This means

$$\Pi_{\dot{U}_{\mathbf{Q},K}^1} v = \mathbf{E}_{\mathbf{Q}} \left(v \mid (\xi_1, \dots, \xi_K) \right).$$

Since

$$\sigma \left(\bigcup_{K \geq 1} \sigma(\xi_1, \dots, \xi_K) \right) = \sigma(\xi),$$

it follows (see *e.g.* [7])

$$\left\| v - \Pi_{\dot{U}_{\mathbf{Q},K}^1} v \right\|_{U_{\mathbf{Q}}^1} = \left\| v - \mathbf{E}_{\mathbf{Q}} \left(v \mid (\xi_1, \dots, \xi_K) \right) \right\|_{U_{\mathbf{Q}}^1} \rightarrow 0, \quad K \rightarrow \infty$$

for any $v \in L^2(\Omega, \sigma(\xi), \mathbf{Q}; \dot{H}^1(D))$. Due to $\hat{u} \in L^2(\Omega, \sigma(\xi), \mathbf{Q}; \dot{H}^1(D))$ this implies immediately

$$\|\hat{u} - \Pi_{\dot{U}_{\mathbf{Q},K}^1} \hat{u}\|_{U_{\mathbf{Q}}^1} \rightarrow 0, \quad K \rightarrow \infty$$

which completes the proof. □

Investigations of the third error term e_3 are based on properties of generalized polynomial chaos. In order to represent an arbitrary $\dot{H}^1(D)$ -valued random variable in $\dot{U}_{\mathbf{Q}}^1 = L^2(\Omega, \sigma(\xi), \mathbf{Q}; \dot{H}^1(D))$ in terms of a generalized polynomial chaos expansion the multivariate polynomials in the basic random variables ξ must be complete in $L^2(\Omega, \sigma(\xi), \mathbf{Q})$. Since $\kappa_{\max}/\mu \in L^r(\Omega, \mathfrak{A}, \mathbf{P})$ for some $r > 1$ there holds $\mathbf{E}_{\mathbf{Q}} |\xi_i|^n < \infty$ for all $n \in \mathbb{N}$ and

$i \in I_\xi$. Thus, there exists a set of multivariate orthonormal polynomials $\{p_\alpha(\xi), \alpha \in \Lambda\}$ w. r. t. the probability measure \mathbf{Q} , *i.e.*

$$\mathbf{E}_{\mathbf{Q}}(p_\alpha(\xi)p_\beta(\xi)) = \delta_{\alpha,\beta}, \quad \text{for } \alpha, \beta \in \Lambda.$$

Here, the distribution $F_\xi^{\mathbf{Q}}$ of the basic random variable $\xi = (\xi_i)_{i \in I_\xi}$ as random variable on the probability space $(\Omega, \mathfrak{A}, \mathbf{Q})$ is defined by

$$F_\xi^{\mathbf{Q}}(y) = \mathbf{Q}(\xi \leq y) = \mathbf{E}_{\mathbf{Q}} \mathbb{1}_{\{\xi \leq y\}} = \frac{1}{c \frac{\kappa_{\max}}{\mu}} \mathbf{E}_{\mathbf{P}} \left(\mathbb{1}_{\{\xi \leq y\}} \frac{\kappa_{\max}}{\mu} \right)$$

Since $\frac{\kappa_{\max}}{\mu}$ is $\sigma(\xi)$ -measurable, there exists a measurable function $t_{\frac{\kappa_{\max}}{\mu}} : \mathbb{R}^{|I_\xi|} \rightarrow \mathbb{R}$ with $t_{\frac{\kappa_{\max}}{\mu}}(\xi) = \frac{\kappa_{\max}}{\mu}$. Therefore, the distribution $F_\xi^{\mathbf{Q}}$ of ξ is determined by

$$F_\xi^{\mathbf{Q}}(dy) = \frac{1}{c \frac{\kappa_{\max}}{\mu}} t_{\frac{\kappa_{\max}}{\mu}}(y) F_\xi^{\mathbf{P}}(dy).$$

The completeness of the polynomials $\{p_\alpha(\xi), \alpha \in \Lambda\}$ in the space $L^2(\Omega, \sigma(\xi), \mathbf{Q})$ is then equivalent to the density of the polynomials $\{p_\alpha(y), \alpha \in \Lambda\}$ in $L^2(\mathbb{R}^{|I_\xi|}, \mathfrak{B}(\mathbb{R}^{|I_\xi|}), F_\xi^{\mathbf{Q}}(dy))$. Some necessary conditions to establish this property are discussed in [13]. For example, the polynomials form an orthonormal basis of $L^2(\Omega, \sigma(\xi), \mathbf{Q})$ if for any marginal distribution

$$F_{\xi_1, \dots, \xi_K}^{\mathbf{Q}}(y_1, \dots, y_K) = \int F_\xi^{\mathbf{Q}}(y_1, \dots, y_K, (dy_i)_{i>K}) = \int \frac{1}{c \frac{\kappa_{\max}}{\mu}} t_{\frac{\kappa_{\max}}{\mu}}(y) F_\xi^{\mathbf{P}}(y_1, \dots, y_K, (dy_i)_{i>K})$$

$K \in \mathbb{N}$ with $\{1, \dots, K\} \subseteq I_\xi$, there exists a constant $a > 0$ such that

$$\int \exp \left(a \sqrt{\sum_{i=1}^K y_i^2} \right) F_{\xi_1, \dots, \xi_K}^{\mathbf{Q}}(dy_1, \dots, dy_K) < \infty.$$

In summary, this condition is e.g. fulfilled if there exists a constant $a > 0$ satisfying

$$\int \exp \left(a \sqrt{\sum_{i \in I_\xi} y_i^2} \right) F_\xi^{\mathbf{Q}}(dy) = \mathbf{E}_{\mathbf{Q}} \exp \left(a \sqrt{\sum_{i \in I_\xi} \xi_i^2} \right) = \frac{1}{c \frac{\kappa_{\max}}{\mu}} \mathbf{E}_{\mathbf{P}} \left(\exp \left(a \sqrt{\sum_{i \in I_\xi} \xi_i^2} \right) t_{\frac{\kappa_{\max}}{\mu}}(\xi) \right) < \infty.$$

Assuming that the polynomials are dense in $L^2(\Omega, \sigma(\xi), \mathbf{Q})$ the weak solution $\hat{u} \in \mathring{U}_{\mathbf{Q}}^1$ possesses a generalized polynomial chaos expansion, *i.e.*

$$\hat{u} = \sum_{\alpha \in \Lambda} \hat{u}_\alpha p_\alpha(\xi),$$

since it is $\sigma(\xi)$ -measurable. The generalized polynomial chaos coefficients are given by

$$\hat{u}_\alpha = \mathbf{E}_{\mathbf{Q}}(\hat{u} p_\alpha(\xi)) = \mathbf{E}_{\mathbf{Q}}(\hat{u}_\xi(\xi) p_\alpha(\xi)) \in \mathring{H}^1(D), \quad \alpha \in \Lambda.$$

Notably, the orthogonal projections $\Pi_{\mathring{U}_{\mathbf{Q},K}^1} \hat{u}$ and $\Pi_{\mathring{U}_{\mathbf{Q},N,K}^1} \hat{u}$ admit also a generalized polynomial chaos representation, namely,

$$\Pi_{\mathring{U}_{\mathbf{Q},K}^1} \hat{u} = \sum_{d(\alpha) \leq K} \hat{u}_\alpha p_\alpha(\xi), \quad d(\alpha) := \max\{i \in \mathbb{N} : \alpha_i \neq 0\}$$

and

$$\Pi_{\mathring{U}_{\mathbf{Q},N,K}^1} \hat{u} = \sum_{\alpha \in \Lambda_{N,K}} \hat{u}_\alpha p_\alpha(\xi).$$

Summarizing these results we arrive at the following corollary.

Corollary 3.9. *If the polynomials $\{p_\alpha(\xi), \alpha \in \Lambda\}$ form an orthonormal basis of $L^2(\Omega, \sigma(\xi), \mathbf{Q})$, and $\hat{u} \in \hat{U}_{\mathbf{Q}}^1$, then the approximation error e_3 converges to zero, i.e.*

$$e_3 = \|\Pi_{\hat{U}_{\mathbf{Q},K}^1} \hat{u} - \Pi_{\hat{U}_{\mathbf{Q},N,K}^1} \hat{u}\|_{U_{\mathbf{Q}}^1} \rightarrow 0 \quad \text{when } N, K \rightarrow \infty.$$

In particular, if the polynomials $\{p_\alpha(\xi), \alpha \in \Lambda : d(\alpha) \leq K\}$ form an orthonormal basis of $L^2(\Omega, \sigma(\xi_1, \dots, \xi_K), \mathbf{Q})$ for any $K \in \mathbb{N}$ with $\{1, \dots, K\} \subseteq I_\xi$, we have for any K

$$e_3 = \|\Pi_{\hat{U}_{\mathbf{Q},K}^1} \hat{u} - \Pi_{\hat{U}_{\mathbf{Q},N,K}^1} \hat{u}\|_{U_{\mathbf{Q}}^1} \rightarrow 0 \quad \text{for } N \rightarrow \infty.$$

Remark 3.10. A necessary condition for the completeness of the polynomials $\{p_\alpha(\xi), \alpha \in \Lambda\}$ in $L^2(\Omega, \sigma(\xi), \mathbf{Q})$ is their completeness in the space $L^2(\Omega, \sigma(\xi), \mathbf{P})$.

Combining the results of the Corollaries 3.6, 3.8 and 3.9 we can now formulate the precise conditions for which the sequence of stochastic Petrov–Galerkin solutions $\bar{u}_{N,K,hp}$ converges to the weak solution \hat{u} in \hat{U}^1 .

Corollary 3.11. *For $\hat{u} \in U_{\mathbf{Q}}^k \cap \hat{U}_{\mathbf{Q}}^1$, $k \geq 2$, if the stochastic Petrov–Galerkin solutions $\bar{u}_{N,K,hp}$ exist and the orthonormal polynomials $\{p_\alpha(\xi), \alpha \in \Lambda\}$ are complete in $L^2(\Omega, \sigma(\xi), \mathbf{Q})$, then the sequence of stochastic Petrov–Galerkin solutions $\bar{u}_{N,K,hp}$ converge in \hat{U}^1 to the weak solution \hat{u} .*

Proof. The result follows immediately from the estimate (3.4) and the Corollaries 3.4, 3.6, 3.8 and 3.9. □

This means that for the stochastic Petrov–Galerkin solutions $\bar{u}_{N,K,hp}$ associated with the weighted test function space we can establish relatively weak conditions which guarantee convergence in U^1 . We would like to emphasize that for standard stochastic Galerkin solutions $\hat{u}_{N,K,hp}$ associated with unweighted test functions, we need stronger assumptions to prove convergence in U^1 . To illustrate this point we refer to Example 4.1 in Section 4 ahead where a convergent sequence of stochastic Petrov–Galerkin solutions $\bar{u}_{N,K,hp}$ is constructed whereas the sequence of standard stochastic Galerkin solutions $\hat{u}_{N,K,hp}$ fails to converge. In ongoing research we are also studying the basic random variables’ impact on the approximation quality and the rate of convergence of the stochastic Petrov–Galerkin solutions $\bar{u}_{N,K,hp}$.

3.2. Assembly of the stochastic Galerkin equations and computational aspects

The use of weighted test functions affects the assembly of the associated stochastic Galerkin system. In what follows we outline the differences to the unweighted case (see *e.g.* [12, 19, 27]). Specifically, we address the calculation of the stochastic Petrov–Galerkin solution $\bar{u}_{N,K,hp}$ and associated characteristics, for example, the expectation function $\mathbf{E}_{\mathbf{P}} \bar{u}_{N,K,hp}$ and second order moment function $\mathbf{E}_{\mathbf{P}} \bar{u}_{N,K,hp}^2$.

Let $\kappa_{\max}/\mu \in L^r(\Omega, \mathfrak{A}, \mathbf{P})$ for some $r > 1$. Assume $\{q_\alpha(\xi)\varphi_i, \alpha \in \Lambda_{N,K}, i \in I_{hp}\}$ is a basis of the space $U_{N,K,hp}$ where the sets $\{\varphi_i, i \in I_{hp}\}$ and $\{q_\alpha(\xi), \alpha \in \Lambda_{N,K}\}$ are bases of the spaces U_{hp} and $U_{N,K}$, respectively. Then a basis of the finite-dimensional test function $V_{N,K,hp}$ space is given by

$$\left\{ \frac{q_\alpha(\xi)}{\mu} \varphi_i, \alpha \in \Lambda_{N,K}, i \in I_{hp} \right\}.$$

Hence, the stochastic Petrov–Galerkin solution $\bar{u}_{N,K,hp} = \sum_{\substack{\alpha \in \Lambda_{N,K} \\ i \in I_{hp}}} u_{\alpha i} \varphi_i q_\alpha(\xi)$ satisfies the linear system of equations

$$A \underline{u} = \underline{f}$$

where

$$A := \left(a(q_\alpha \varphi_i, q_\beta \varphi_j / \mu) \right)_{\beta j, \alpha i}, \quad \underline{f} := \left(\ell(\varphi_j q_\beta / \mu) \right)_{\beta j} \quad \text{and} \quad \underline{u} := \left(u_{\alpha i} \right)_{\alpha i}.$$

Using the parameterized input data $\kappa_\xi(\xi)$ with $\kappa_\xi : \mathbb{R}^{|\mathcal{I}_\xi|} \rightarrow L^\infty(D)$ and $f_\xi(\xi)$ with $f_\xi : \mathbb{R}^{|\mathcal{I}_\xi|} \rightarrow H^{-1}(D)$ we obtain

$$\begin{aligned} a(\varphi_i q_\alpha, \varphi_j q_\beta / \mu) &= \mathbf{E}_\mathbf{P} \left(\left(\int_D [\kappa_\xi(\xi)](x) \nabla \varphi_i(x) \cdot \nabla \varphi_j(x) \, dx \right) q_\alpha(\xi) q_\beta(\xi) / \mu \right) \\ &= \int_D \mathbf{E}_\mathbf{P} \left(\frac{[\kappa_\xi(\xi)](x)}{\mu} q_\alpha(\xi) q_\beta(\xi) \right) \nabla \varphi_i(x) \cdot \nabla \varphi_j(x) \, dx \end{aligned}$$

and

$$\ell(\varphi_j q_\beta / \mu) = \mathbf{E}_\mathbf{P} \left(\langle f_\xi(\xi), \varphi_j \rangle_{H^{-1}, \dot{H}^1} q_\beta(\xi) / \mu \right) = \left\langle \mathbf{E}_\mathbf{P} \left(\frac{f_\xi(\xi)}{\mu} q_\beta(\xi) \right), \varphi_j \right\rangle_{H^{-1}, \dot{H}^1}$$

for any $i, j \in I_{hp}$, $\alpha, \beta \in \Lambda_{N,K}$.

It is important to choose suitable polynomials $\{q_\alpha(\xi), \alpha \in \Lambda_{N,K}\}$ so that the computation of the expectations in the linear system is actually feasible. To this end, we discuss two options.

First, it is possible to employ the polynomials $\{\tilde{p}_\alpha(\xi), \alpha \in \Lambda\}$ which form an orthonormal basis of $L^2(\Omega, \sigma(\xi), \mathbf{P})$. Thus, evaluation of the expected values

$$\mathbf{E}_\mathbf{P} \left(\kappa_\xi(\xi) \tilde{p}_\alpha(\xi) \tilde{p}_\beta(\xi) / \mu \right) \quad \text{and} \quad \mathbf{E}_\mathbf{P} \left(f_\xi(\xi) \tilde{p}_\beta(\xi) / \mu \right)$$

requires the calculation of the generalized polynomial chaos expansions of $\kappa_\xi(\xi)/\mu$ and $f_\xi(\xi)/\mu$ or their generalized polynomial chaos coefficients, respectively. The corresponding characteristics of the stochastic Petrov–Galerkin solution $\bar{u}_{N,K, hp}$, *i.e.* expected value and second order moment function

$$\mathbf{E}_\mathbf{P} \bar{u}_{N,K, hp} \quad \text{and} \quad \mathbf{E}_\mathbf{P} \bar{u}_{N,K, hp}^2$$

are determined as usual by

$$\mathbf{E}_\mathbf{P} \bar{u}_{N,K, hp} = \mathbf{E}_\mathbf{P} \left(\sum_{\substack{\alpha \in \Lambda_{N,K} \\ i \in I_{hp}}} u_{\alpha i} \varphi_i \tilde{p}_\alpha(\xi) \right) = \sum_{\substack{\alpha \in \Lambda_{N,K} \\ i \in I_{hp}}} u_{\alpha i} \varphi_i \underbrace{\mathbf{E}_\mathbf{P} (\tilde{p}_\alpha(\xi))}_{=\delta_{\alpha,0}} = \sum_{i \in I_{hp}} u_{0i} \varphi_i$$

and

$$\begin{aligned} \mathbf{E}_\mathbf{P} \bar{u}_{N,K, hp}^2 &= \mathbf{E}_\mathbf{P} \left(\sum_{\substack{\alpha, \beta \in \Lambda_{N,K} \\ i, j \in I_{hp}}} u_{\alpha i} u_{\beta j} \varphi_i \varphi_j \tilde{p}_\alpha(\xi) \tilde{p}_\beta(\xi) \right) \\ &= \sum_{\substack{\alpha, \beta \in \Lambda_{N,K} \\ i, j \in I_{hp}}} u_{\alpha i} \varphi_i u_{\beta j} \varphi_j \underbrace{\mathbf{E}_\mathbf{P} (\tilde{p}_\alpha(\xi) \tilde{p}_\beta(\xi))}_{=\delta_{\alpha, \beta}} = \sum_{\alpha \in \Lambda_{N,K}} \sum_{i, j \in I_{hp}} u_{\alpha i} u_{\alpha j} \varphi_i \varphi_j = \sum_{\alpha \in \Lambda_{N,K}} \left(\sum_{i \in I_{hp}} u_{\alpha i} \varphi_i \right)^2. \end{aligned}$$

A second option can be obtained if the expectation of the random variable μ^{-1} is finite and $\mathbf{E}_\mathbf{P} |\xi_i|^n \mu^{-1} < \infty$ holds for all $n \in \mathbb{N}$ and $i \in \mathcal{I}_\xi$. For example, this is satisfied if $H^{-1}(D) \subset U_{\mu^{-2}}^{-1}$. Then, *w.l.o.g.* we can assume that the expectation of μ^{-1} is equal to one. In this case the polynomials $\{\tilde{p}_\alpha(\xi), \alpha \in \Lambda\}$ which are orthonormal in $L^2(\Omega, \mathfrak{A}, \mu^{-1} d\mathbf{P})$ can serve as basis polynomials $q_\alpha(\xi)$. Then, the computation of the expected values $\mathbf{E}_\mathbf{P} \left(\kappa_\xi(\xi) \tilde{p}_\alpha(\xi) \tilde{p}_\beta(\xi) / \mu \right)$ and $\mathbf{E}_\mathbf{P} \left(f_\xi(\xi) \tilde{p}_\beta(\xi) / \mu \right)$ requires only the generalized polynomial chaos expansions of κ_ξ and f_ξ in the polynomials $\{\tilde{p}_\alpha(\xi), \alpha \in \Lambda\}$. However, evaluating the first and second order moment function of the stochastic Petrov–Galerkin solution is more complicated than in the first option since the polynomials $\{\tilde{p}_\alpha(\xi), \alpha \in \Lambda\}$ do not satisfy $\mathbf{E}_\mathbf{P} (\tilde{p}_\alpha(\xi) \tilde{p}_\beta(\xi)) = C_\alpha \delta_{\alpha, \beta}$. Thus, we have to compute

$$\mathbf{E}_\mathbf{P} \bar{u}_{N,K, hp} = \mathbf{E}_\mathbf{P} \left(\sum_{\substack{\alpha \in \Lambda_{N,K} \\ i \in I_{hp}}} u_{\alpha i} \varphi_i \tilde{p}_\alpha(\xi) \right) = \sum_{\substack{\alpha \in \Lambda_{N,K} \\ i \in I_{hp}}} u_{\alpha i} \varphi_i \mathbf{E}_\mathbf{P} \tilde{p}_\alpha(\xi)$$

and

$$\mathbf{E_P} \bar{u}_{N,K,hp}^2 = \mathbf{E_P} \left(\sum_{\substack{\alpha, \beta \in \Lambda_{N,K} \\ i, j \in I_{hp}}} u_{\alpha i} u_{\beta j} \varphi_i \varphi_j \tilde{p}_\alpha(\xi) \tilde{p}_\beta(\xi) \right) = \sum_{\substack{\alpha, \beta \in \Lambda_{N,K} \\ i, j \in I_{hp}}} u_{\alpha i} \varphi_i u_{\beta j} \varphi_j \mathbf{E_P} (\tilde{p}_\alpha(\xi) \tilde{p}_\beta(\xi)).$$

Suppose now, that the generalized polynomial chaos expansions of the polynomials $\tilde{p}_\alpha(\xi)$, $\alpha \in \Lambda_{N,K}$, are given in terms of sums of the polynomials $\{\bar{p}_\alpha(\xi), \alpha \in \Lambda\}$, *i.e.*, we know the generalized polynomial chaos coefficients $c_\beta^{(\alpha)}$ satisfying

$$\tilde{p}_\alpha(\xi) = \sum_{\beta \in \Lambda_{N,K}} c_\beta^{(\alpha)} \bar{p}_\beta(\xi).$$

These coefficients are sometimes called *connection coefficients* in the literature [24, 30, 33]. We then obtain the stochastic Petrov–Galerkin solution in terms of the polynomials $\{\bar{p}_\alpha(\xi), \alpha \in \Lambda\}$ by rearranging

$$\bar{u}_{N,K,hp} = \sum_{\substack{\alpha \in \Lambda_{N,K} \\ i \in I_{hp}}} u_{\alpha i} \varphi_i \tilde{p}_\alpha(\xi) = \sum_{\substack{\alpha \in \Lambda_{N,K} \\ i \in I_{hp}}} u_{\alpha i} \varphi_i \sum_{\beta \in \Lambda_{N,K}} c_\beta^{(\alpha)} \bar{p}_\beta(\xi) = \sum_{\beta \in \Lambda_{N,K}} \bar{p}_\beta(\xi) \left(\sum_{\substack{\alpha \in \Lambda_{N,K} \\ i \in I_{hp}}} c_\beta^{(\alpha)} u_{\alpha i} \varphi_i \right).$$

The expected value and second order moment function can be computed as follows

$$\mathbf{E_P} \bar{u}_{N,K,hp} = \sum_{\substack{\alpha \in \Lambda_{N,K} \\ i \in I_{hp}}} c_0^{(\alpha)} u_{\alpha i} \varphi_i \quad \text{and} \quad \mathbf{E_P} \bar{u}_{N,K,hp}^2 = \sum_{\beta \in \Lambda_{N,K}} \left(\sum_{\substack{\alpha \in \Lambda_{N,K} \\ i \in I_{hp}}} c_\beta^{(\alpha)} u_{\alpha i} \varphi_i \right)^2.$$

If both sets of polynomials are basis polynomials of the finite-dimensional subspace $U_{N,K}$, then the resulting stochastic Galerkin solutions coincide, of course. Depending on the problem, however, it can happen that the two options differ with respect to computational aspects and performance.

Example 3.12. Assume that for a given boundary-value problem the basic random variable is a single standard Gaussian random variable ξ and that the random variable $1/\mu = t_{\mu^{-1}}(\xi) = \exp(\frac{\varepsilon}{2}\xi^2)/\sigma$ with $\sigma = 1/\sqrt{1-\varepsilon}$ for $0 < \varepsilon < 1$ serves as an admissible weighting random variable. Then the basic random variable ξ has also a Gaussian distribution on the probability space $(\Omega, \mathfrak{A}, \mu^{-1}d\mathbf{P})$ and the corresponding density function $f_\xi^{\mu^{-1}d\mathbf{P}}$ is given by

$$\begin{aligned} f_\xi^{\mu^{-1}d\mathbf{P}}(y) &= t_{\mu^{-1}}(y) f_\xi^{\mathbf{P}}(y) = \frac{1}{\sqrt{2\pi\sigma}} \exp\left(\frac{\varepsilon}{2}y^2\right) \exp\left(-\frac{y^2}{2}\right) \\ &= \frac{1}{\sqrt{2\pi\sigma}} \exp\left(-\frac{y^2(1-\varepsilon)}{2}\right) = \frac{1}{\sqrt{2\pi\sigma}} \exp\left(-\frac{y^2}{2\sigma^2}\right), \quad y \in \mathbb{R}. \end{aligned}$$

The associated orthonormal polynomials $\{\tilde{p}_n(\xi), n \in \mathbb{N}_0\}$ in $L^2(\Omega, \sigma(\xi), \mu^{-1}d\mathbf{P})$ are the normalized probabilistic Hermite polynomials $\{\text{He}_n(\xi), n \in \mathbb{N}_0\}$ scaled in the argument, *i.e.*

$$\tilde{p}_n(\xi) = \frac{1}{\sqrt{n!}} \text{He}_n\left(\frac{\xi}{\sigma}\right), \quad n \in \mathbb{N}_0.$$

If we assemble the stochastic Galerkin system by the first option described above, the original Hermite polynomials

$$\bar{p}_n(\xi) = \frac{\text{He}_n(\xi)}{\sqrt{n!}}, \quad n = 0, \dots, N$$

are used to span the finite-dimensional subspace $U_{N,K}$ with $K = 1$. Therefore, we have to calculate for instance the polynomial chaos coefficients of $\kappa_\xi(\xi)/\mu$, *i.e.*

$$\mathbf{E_P} \left(\frac{\kappa_\xi(\xi)}{\mu} \frac{\text{He}_n(\xi)}{\sqrt{n!}} \right), \quad n \in \mathbb{N}_0.$$

Using the second proposed option we have to calculate the generalized polynomial chaos coefficients of κ and f w. r. t. the polynomials $\{\tilde{p}_n(\xi), n \in \mathbb{N}_0\}$. For the coefficient κ this requires evaluation of the polynomial chaos coefficients of $\kappa_\xi(\sigma\xi)$ by the change of variables $z = y/\sigma$

$$\begin{aligned} \mathbf{E}_{\mathbf{P}} \left(\frac{\kappa_\xi(\xi)}{\mu} \frac{\text{He}_n(\xi/\sigma)}{\sqrt{n!}} \right) &= \int_{\mathbb{R}} \kappa_\xi(y) \frac{\text{He}_n(y/\sigma)}{\sqrt{n!}} t_{\mu^{-1}}(y) f_\xi^{\mathbf{P}}(y) \, dy = \int_{\mathbb{R}} \kappa_\xi(y) \frac{\text{He}_n(y/\sigma)}{\sqrt{n!}} f_\xi^{\mu^{-1}\mathbf{dP}}(y) \, dy \\ &= \int_{\mathbb{R}} \kappa_\xi(\sigma z) \frac{\text{He}_n(z)}{\sqrt{n!}} f_\xi^{\mathbf{P}}(z) \, dz = \mathbf{E}_{\mathbf{P}} \left(\kappa_\xi(\sigma\xi) \frac{\text{He}_n(\xi)}{\sqrt{n!}} \right), \quad n \in \mathbb{N}_0. \end{aligned}$$

In order to obtain the first and second order moments of the approximate solution, we need the generalized polynomial chaos coefficients $c_i^{(n)}$ satisfying

$$\frac{\text{He}_n(\xi/\sigma)}{\sqrt{n!}} = \sum_{i=0}^n c_i^{(n)} \frac{\text{He}_i(\xi)}{\sqrt{i!}} \quad n \in \mathbb{N}_0.$$

Fortunately, in this case analytic expressions for these coefficients are available:

$$c_i^{(n)} = \begin{cases} \frac{\sqrt{n!}}{\sqrt{i!} \left(\frac{n-i}{2}\right)!} \sigma^{-n} \left(\frac{1-\sigma^2}{2}\right)^{(n-i)/2}, & n+i \text{ even,} \\ 0, & \text{otherwise,} \end{cases} \quad i = 0, \dots, n, \quad n \in \mathbb{N}_0.$$

If the multivariate orthonormal polynomials $\{\bar{p}_\alpha(\xi), \alpha \in \Lambda\}$ or $\{\tilde{p}_\alpha(\xi), \alpha \in \Lambda\}$ can be represented as tensor products of univariate orthonormal polynomials, then the assembly of the stochastic Galerkin system simplifies significantly. To answer this question, we need to investigate if the basic random variables $(\xi_i)_{i \in I_\xi}$ are independent as random variables on the probability spaces $(\Omega, \mathfrak{A}, \mathbf{P})$ or $(\Omega, \mathfrak{A}, \mu^{-1}\mathbf{dP})$, respectively. Let $t_{\mu^{-1}} : \mathbb{R}^{|I_\xi|} \rightarrow \mathbb{R}$ be a measurable function such that there holds

$$t_{\mu^{-1}}(\xi) = \mu^{-1}.$$

Then the distribution $F_\xi^{\mu^{-1}\mathbf{dP}}$ of ξ on $(\Omega, \mathfrak{A}, \mu^{-1}\mathbf{dP})$ is

$$F_\xi^{\mu^{-1}\mathbf{dP}}(dy) = t_{\mu^{-1}}(y) F_\xi^{\mathbf{P}}(dy).$$

Analogously, the distribution $F_{\xi_i}^{\mu^{-1}\mathbf{dP}}$ of each random variable $\xi_i, i \in I_\xi$, is given by

$$F_{\xi_i}^{\mu^{-1}\mathbf{dP}}(dy_i) = t_{\mu^{-1},i}(y_i) F_{\xi_i}^{\mathbf{P}}(dy_i),$$

where $t_{\mu^{-1},i} : \mathbb{R} \rightarrow \mathbb{R}$ is the measurable function satisfying

$$t_{\mu^{-1},i}(\xi_i) = \mathbf{E}_{\mathbf{P}}(\mu^{-1} | \xi_i), \quad i \in I_\xi.$$

Then, the independence of the basic random variables $(\xi_i)_{i \in I_\xi}$ on the respective probability spaces holds if and only if the associated distribution $F_\xi^{\mathbf{P}}$ or $F_\xi^{\mu^{-1}\mathbf{dP}}$ is the product of the distributions of the individual random variables $\xi_i, i \in I_\xi$, *i.e.*,

$$F_\xi^{\mathbf{P}}(dy) = \prod_{i \in I_\xi} F_{\xi_i}^{\mathbf{P}}(dy_i) \quad \text{or} \quad F_\xi^{\mu^{-1}\mathbf{dP}}(dy) = \prod_{i \in I_\xi} F_{\xi_i}^{\mu^{-1}\mathbf{dP}}(dy_i),$$

respectively. If, for instance, the random variables $(\xi_i)_{i \in I_\xi}$ are independent on the probability space $(\Omega, \mathfrak{A}, \mathbf{P})$ and the function $t_{\mu^{-1}}$ can be represented as product

$$t_{\mu^{-1}}(y) = \prod_{i \in I_\xi} t_{\mu^{-1},i}(y_i),$$

then, the random variables $(\xi_i)_{i \in I_\xi}$ are also independent on the probability space $(\Omega, \mathfrak{A}, \mu^{-1}\mathbf{dP})$.

4. EXAMPLES

In this section we present two examples of boundary-value problems which illustrate the different convergence behavior of the standard stochastic Galerkin approach and the stochastic Petrov–Galerkin approach, respectively. Because the difference is related to the stochastic aspects of the problems, we construct very simple one-dimensional boundary-value problems where the random coefficients depend on a single basic random variable. In both examples we will observe that

- the standard stochastic Galerkin approach which is usually used (see *e.g.* [27, 40]) fails to provide a sequence of approximate solutions converging to the exact solution in the natural norm,
- the stochastic Petrov–Galerkin approach with a suitable weighted test function space yields a sequence of approximate solutions converging to the exact solution in the natural norm.

Both results, non-convergence for the standard stochastic Galerkin approach and convergence for the stochastic Petrov–Galerkin approach, respectively, are illustrated numerically. In addition, we prove non-convergence of the standard stochastic Galerkin approach for Example 4.1 analytically.

In both examples we study a one-dimensional boundary-value problem of the type

$$-(\kappa(x, \omega)u'(x, \omega))' = f(x, \omega), \quad x \in (0, 1), \quad \omega \in \Omega,$$

with homogeneous Dirichlet boundary conditions

$$u(0, \omega) = u(1, \omega) = 0, \quad \omega \in \Omega,$$

or mixed boundary conditions

$$u(0, \omega) = 0 \quad \text{and} \quad \kappa(1, \omega)u'(1, \omega) = F(\omega), \quad \omega \in \Omega.$$

Though we do not consider mixed boundary conditions in our model problem (2.1) in Section 1, the theory developed in the preceding sections still applies. To eliminate the spatial discretization error as far as possible we use a single Gauss–Lobatto–Legendre spectral element of suitable degree p for the spatial discretization. Thus, we are able to observe the effect of the different discretizations of the stochastic space. In the experiments we always compare the corresponding relative approximation error in the U^1 -norm – the natural norm, *i.e.*,

$$\frac{\|\hat{u} - \hat{u}_{N,K,hp}\|_{U^1}}{\|\hat{u}\|_{U^1}} \quad \text{and} \quad \frac{\|\hat{u} - \bar{u}_{N,K,hp}\|_{U^1}}{\|\hat{u}\|_{U^1}},$$

respectively. Furthermore, since we consider a single basic random variable, we drop the index K in the notation of the finite-dimensional subspaces and stochastic Galerkin solutions for the sake of convenience.

Example 4.1. We consider the boundary-value problem on the domain $D = (0, 1)$ with random coefficient $\kappa(x, \omega) = \exp(\xi(\omega))$ and random forcing term $f(x, \omega) = \exp(\xi(\omega))|\xi(\omega) - 1|$ where $\xi \sim \mathcal{N}(0, 1)$ is a standard Gaussian random variable. This means, we are looking for the solution of

$$\begin{aligned} -(\exp(\xi(\omega))u'(x, \omega))' &= \exp(\xi(\omega))|\xi(\omega) - 1|, & x \in (0, 1), \quad \omega \in \Omega, \\ u(0, \omega) &= u(1, \omega) = 0, & \omega \in \Omega. \end{aligned}$$

Since the coefficient κ is strongly measurable in $L^\infty(D)$ and there holds

$$0 < \kappa_{\min} := \exp(\xi) \leq \kappa \leq \exp(\xi) =: \kappa_{\max} < \infty$$

we can choose $\mu = \kappa_{\min}$. Since $f \in U_{\exp(-2\xi)}^{-1} \cap U_{\exp(-\xi)}^{-1}$ we can also conclude that there exists a unique weak solution $\hat{u} \in \mathring{U}^1$ given by

$$\hat{u}(x, \omega) = [\hat{u}_\xi(\xi(\omega))](x) = \frac{1}{2}|\xi(\omega) - 1|(x - x^2), \quad x \in (0, 1), \quad \omega \in \Omega.$$

This is of course also a pathwise strong solution.

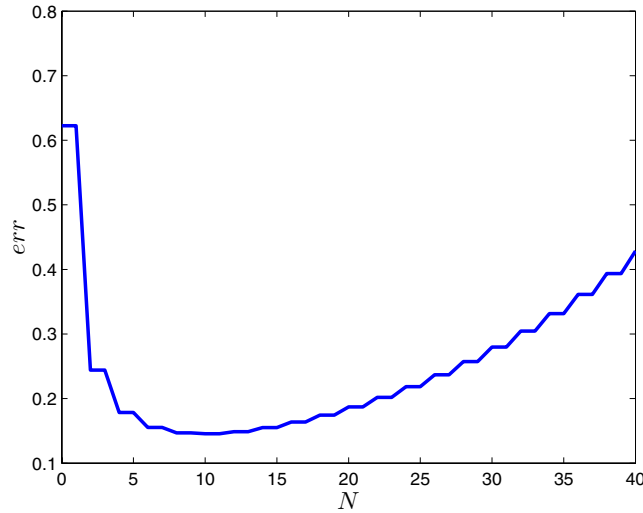


FIGURE 1. Relative approximation error err of the standard stochastic Galerkin solution $\hat{u}_{N,hp}$ for different orders N of polynomial chaos.

For the standard stochastic Galerkin approach we choose identical finite-dimensional solution and test function spaces $U_{N,hp} = U_{hp} \otimes U_N$ where

$$U_N = \text{span} \left\{ \frac{\text{He}_n(\xi)}{\sqrt{n!}}, n = 0, \dots, N \right\}$$

is spanned by the first $N + 1$ probabilistic Hermite polynomials $\{\text{He}_n(\xi), n \in \mathbb{N}_0\}$. The relative approximation error

$$\text{err} = \frac{\|\hat{u} - \hat{u}_{N,hp}\|_{U^1}}{\|\hat{u}\|_{U^1}}$$

of the corresponding standard stochastic Galerkin solutions $\hat{u}_{N,hp}$ is shown in Figure 1. Here, we use a single Gauss–Lobatto–Legendre spectral element of degree $p = 3$ for the spatial discretization. We clearly see that the approximation error does not converge to zero for increasing orders of the chaos polynomials. In fact, we can prove that the sequence of standard stochastic Galerkin solutions $\hat{u}_{N,hp}$ does not converge in \hat{U}^1 to the exact solution for $N \rightarrow \infty$. To see this, we note that due to the special structure of the boundary-value problem the standard stochastic Galerkin solution

$$\hat{u}_{N,hp}(x, \omega) = \hat{u}_{hp}(x)\hat{u}_N(\xi(\omega))$$

is the product of the deterministic function $\hat{u}_{hp}(x)$ which satisfies the equation

$$\int_D \hat{u}'_{hp}(x)v'_{hp}(x) dx = \int_D v_{hp}(x) dx \quad \text{for all } v_{hp} \in U_{hp} \tag{4.1}$$

and the random variable $\hat{u}_N(\xi)$ satisfying the equation

$$\mathbf{E}_{\mathbf{P}} \left(\exp(\xi)\hat{u}_N(\xi) \frac{\text{He}_m(\xi)}{\sqrt{m!}} \right) = \mathbf{E}_{\mathbf{P}} \left(\exp(\xi)|\xi - 1| \frac{\text{He}_m(\xi)}{\sqrt{m!}} \right) \quad \text{for all } m = 0, \dots, N. \tag{4.2}$$

We denote $\hat{u}_D(x) := \frac{1}{2}(x - x^2)$ and $\hat{u}_\Omega(\xi) := |\xi - 1|$ and obtain the following lower bound

$$\begin{aligned} \|\hat{u} - \hat{u}_{N,hp}\|_{U^1} &= \|\hat{u}_D \hat{u}_\Omega(\xi) - \hat{u}_{hp} \hat{u}_N(\xi)\|_{U^1} = \|\hat{u}_\Omega(\xi)(\hat{u}_D - \hat{u}_{hp}) + \hat{u}_{hp}(\hat{u}_\Omega(\xi) - \hat{u}_N(\xi))\|_{U^1} \\ &\geq \|\hat{u}_{hp}(\hat{u}_\Omega(\xi) - \hat{u}_N(\xi))\|_{U^1} - \|\hat{u}_\Omega(\xi)(\hat{u}_D - \hat{u}_{hp})\|_{U^1} \\ &= \|\hat{u}_{hp}\|_{H^1(D)} \|\hat{u}_\Omega(\xi) - \hat{u}_N(\xi)\|_{L^2(\Omega, \mathfrak{A}, \mathbf{P})} - \|\hat{u}_\Omega(\xi)\|_{L^2(\Omega, \mathfrak{A}, \mathbf{P})} \|\hat{u}_D - \hat{u}_{hp}\|_{H^1(D)}. \end{aligned}$$

In view of

$$\|\hat{u}_D - \hat{u}_{hp}\|_{H^1(D)} \rightarrow 0 \quad \text{and} \quad \|\hat{u}_\Omega(\xi)\|_{L^2(\Omega, \mathfrak{A}, \mathbf{P})} < \infty$$

the second error term $(\|\hat{u}_\Omega(\xi)\|_{L^2(\Omega, \mathfrak{A}, \mathbf{P})} \|\hat{u}_D - \hat{u}_{hp}\|_{H^1(D)})$ converges to zero and the norm of $\|\hat{u}_{hp}\|_{H^1(D)}$ converges to $\|\hat{u}_D\|_{H^1(D)}$. However, the stochastic part $\hat{u}_N(\xi)$ does not converge to $\hat{u}_\Omega(\xi)$ in $L^2(\Omega, \mathfrak{A}, \mathbf{P})$ as shown in detail in the Appendix A. In fact, there exists a subsequence $(\hat{u}_{N_k}(\xi))_{k \in \mathbb{N}_0}$ of $(\hat{u}_N(\xi))_{N \in \mathbb{N}_0}$ whose second order moments tend to infinity which implies that the error $\|\hat{u}_\Omega(\xi) - \hat{u}_{N_k}(\xi)\|_{L^2(\Omega, \mathfrak{A}, \mathbf{P})}$ goes to infinity, *i.e.*,

$$\|\hat{u}_\Omega(\xi) - \hat{u}_{N_k}(\xi)\|_{L^2(\Omega, \mathfrak{A}, \mathbf{P})} \rightarrow \infty, \quad k \rightarrow \infty.$$

Of course, this also implies that the corresponding approximation error *err* goes to infinity, *i.e.*,

$$\|\hat{u} - \hat{u}_{N_k, hp}\|_{U^1} \rightarrow \infty, \quad k, p \rightarrow \infty.$$

Thus, the sequence of standard stochastic Galerkin solutions $\hat{u}_{N, hp}$ does not converge to the exact solution \hat{u} in \hat{U}^1 .

Alternatively, we employ the proposed stochastic Petrov–Galerkin approach. W.l.o.g. we can normalize the random variable μ such that the expectation of its reciprocal μ^{-1} is one, *i.e.*, we take

$$\mu := \exp(1/2) \exp(\xi).$$

We define the weighted finite-dimensional test function space

$$V_{N, hp} = \{u \exp(-\xi), u \in U_{N, hp}\} = \text{span} \left\{ u_{hp} \frac{\text{He}_n(\xi)}{\sqrt{n!}} \exp(-\xi), u_{hp} \in U_{hp}, n = 0, \dots, N \right\}.$$

Since $\kappa_{\max} = \exp(-1/2)\mu$ the probability measure \mathbf{Q} coincides with the probability measure \mathbf{P} and thus $\hat{U}^1 = \hat{U}_{\mathbf{Q}}^1$. Hence, based on Corollary 3.4, we obtain for the corresponding stochastic Galerkin solution $\bar{u}_{N, hp}$ the estimate

$$\|\hat{u} - \bar{u}_{N, hp}\|_{U^1} \leq C \inf_{z \in U_{N, hp}} \|\hat{u} - z\|_{U^1}$$

where exactly the same norms appear on the left and right hand side of the inequality. Due to $\hat{u} \in U^k \cap \hat{U}^1$ for any $k \geq 1$ and the completeness of the Hermite polynomials in the space $L^2(\Omega, \sigma(\xi), \mathbf{P})$ the sequence of the stochastic Petrov–Galerkin solutions $\bar{u}_{N, hp}$ converges to the exact solution. This is illustrated by the approximation error *err* on the left side of Figure 2. Here we have again used a single Gauss–Lobatto–Legendre spectral element of degree $p = 3$ for the spatial discretization. Again, any solution $\bar{u}_{N, hp}$ is the product of a deterministic function \bar{u}_{hp} and a random variable $\bar{u}_N(\xi)$ satisfying analogous equations as (4.1) and (4.2) for \hat{u}_{hp} and $\hat{u}_N(\xi)$. Note that in this case, the stochastic part $\bar{u}_N(\xi)$ is the *exact* polynomial chaos approximation of order N of $\hat{u}_\Omega(\xi)$, hence we compute the best approximation with respect to the norm of the space $L^2(\Omega, \sigma(\xi), \mathbf{P})$. For example, for $N = 0$ the expectation of the stochastic Petrov–Galerkin solution $\bar{u}_{0, hp}$ is the exact expectation of \hat{u} as function of the spatial argument $x \in D$. On the right side of Figure 2 (using a single Gauss–Lobatto–Legendre spectral element of degree $p = 50$ for the spatial discretization) the second order moment function of the exact solution (blue-□ line)

$$\mathbf{E}_{\mathbf{P}} \hat{u}^2(x, \cdot) = \frac{1}{2}(x - x^2)^2, \quad x \in D,$$

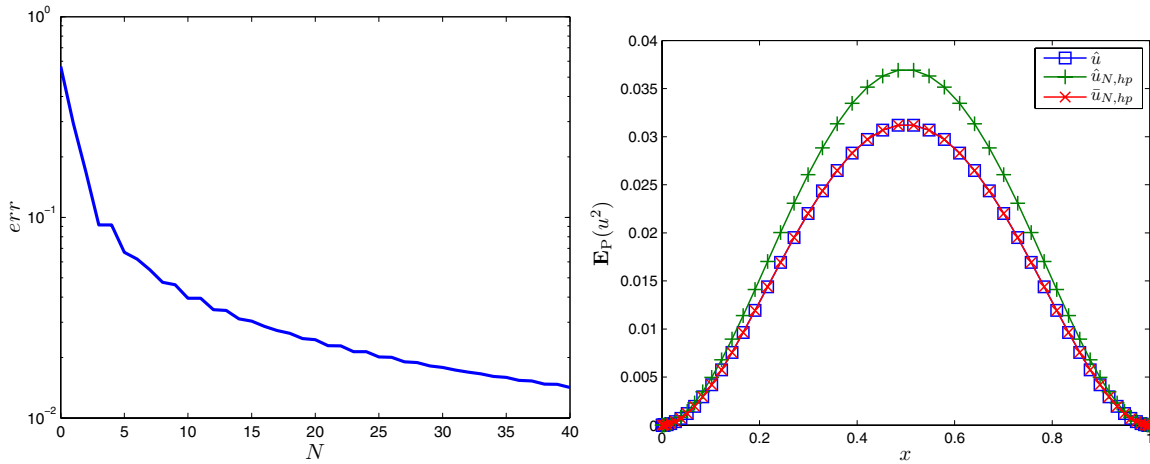


FIGURE 2. Left: Relative approximation error err of the stochastic Petrov–Galerkin solution $u_{N,hp}$ for different orders N of polynomial chaos. Right: Second order moment function of the exact solution \hat{u} and the standard stochastic Galerkin solution $\hat{u}_{N,hp}$ and stochastic Petrov–Galerkin solution $\bar{u}_{N,hp}$ for polynomial chaos of order $N = 40$.

is compared with the second order moment function of the stochastic Galerkin solutions $\hat{u}_{N,hp}$ (green-+ line) and $\bar{u}_{N,hp}$ (red-× line) for generalized polynomial chaos order $N = 40$. As expected, the second order moment function of the stochastic Petrov–Galerkin solution $\bar{u}_{N,hp}$ approximates well the second order moment function of the exact solution whereas the second order moment function of the standard stochastic Galerkin solution $\hat{u}_{N,hp}$ differs substantially.

Our second example is again a one-dimensional boundary-value problem but with a more sophisticated random coefficient, a deterministic forcing term and mixed boundary conditions.

Example 4.2. We study the following boundary-value problem

$$\begin{aligned} -(\exp(-\xi^2(\omega)x/10) u'(x, \omega))' &= f(x), & x \in (0, 1), \omega \in \Omega, \\ u(0, \omega) &= 0, & \omega \in \Omega, \\ \exp(-\xi^2(\omega)/10) u'(1, \omega) &= \exp(-6\xi^2(\omega)), & \omega \in \Omega, \end{aligned}$$

with random coefficient $\kappa(x, \omega) = \exp(-\xi^2(\omega)x/10)$, deterministic forcing term $f \in C[0, 1] \subset H^{-1}(D)$ and random boundary-value $F(\omega) = \exp(-6\xi^2(\omega))$ where the random variable ξ is again standard Gaussian. The coefficient κ depends on $\omega \in \Omega$ as well as on the spatial argument $x \in D$, is strongly measurable and bounded by the random variables

$$\kappa_{\min} := \exp(-\xi^2/10) \quad \text{and} \quad \kappa_{\max} := 1,$$

respectively, satisfying

$$0 < \kappa_{\min} \leq \kappa \leq \kappa_{\max} < \infty \quad \text{a. e. and a. s.}$$

The exact solution \hat{u} of the boundary-value problem is explicitly given by

$$\hat{u}(x, \omega) = \int_0^x \exp(\xi^2(\omega)y/10) \left(\exp(-6\xi^2(\omega)) + \int_y^1 f(z) dz \right) dy.$$

Alternatively, any random variable of the form $\exp(-\frac{\varepsilon}{10}\xi^2)$ with $\varepsilon \geq 1$ is also a lower bound for the coefficient κ . Since we want to work with random variables for which the reciprocal has expectation one and is integrable of

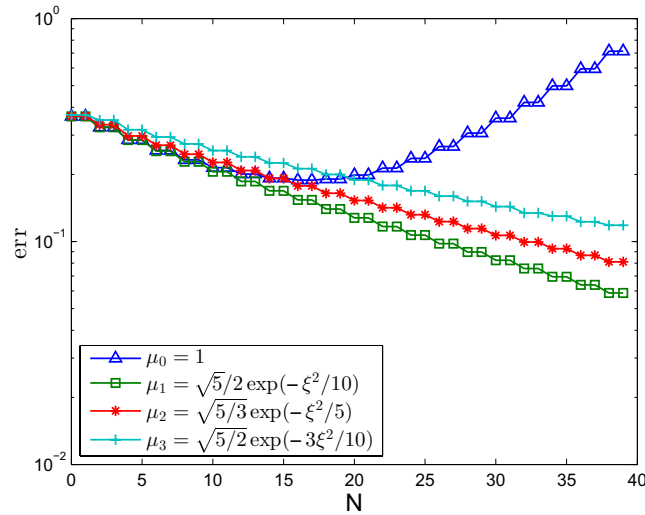


FIGURE 3. Relative approximation error err as function of the generalized polynomial chaos order N for different weights μ_ε , $\varepsilon = 0, 1, 2, 3$.

an order greater than one we consider

$$\mu_\varepsilon := \sqrt{\frac{5}{5-\varepsilon}} \exp\left(-\frac{\varepsilon}{10}\xi^2\right), \quad 1 \leq \varepsilon < 5.$$

Then, for $\varepsilon \leq 3$ the expectation

$$\mathbf{E}_{\mathbf{P}}\left(|\hat{u}|_{H^1}^2 \mu_\varepsilon^{-1}\right) = \sqrt{\frac{5-\varepsilon}{5}} \mathbf{E}_{\mathbf{P}}\left(|\hat{u}|_{H^1}^2 \exp\left(\frac{\varepsilon}{10}\xi^2\right)\right)$$

is also finite. Note that for $\varepsilon > 3$ the latter property is in general not satisfied. For $1 \leq \varepsilon \leq 3$ there exists a unique weak solution \hat{u} in U^1 , especially in U_{-1}^1 . Therefore, any μ_ε with $1 \leq \varepsilon \leq 3$ can serve as a suitable weighting random variable. This allows us to define the weighted test function spaces

$$V_{N,hp}^{(\varepsilon)} = \{u/\mu_\varepsilon, u \in U_{N,hp}\}$$

where

$$U_{N,hp} = U_{hp} \otimes U_N \quad \text{with} \quad U_N = \text{span}\{\xi^n, n = 0, \dots, N\}$$

is the finite-dimensional solution space. Now, any sequence of stochastic Petrov–Galerkin solutions $(\bar{u}_{N,hp}^{(\varepsilon)})_{N,hp}$ with $1 \leq \varepsilon \leq 3$ converges to the exact solution \hat{u} in U^1 , since there holds $\hat{u} \in U_{\frac{\kappa_{\max}}{\mu_\varepsilon}} \cap U_{\frac{\kappa_{\max}}{\mu_\varepsilon}}^k$ for any $k \geq 1$ and the corresponding polynomials are dense in $L^2(\Omega, \sigma(\xi), \kappa_{\max}/\mu_\varepsilon d\mathbf{P})$. For the constant forcing term $f \equiv 1$ we compute the relative approximation errors err of the stochastic Petrov–Galerkin solutions $\bar{u}_{N,hp}^{(\varepsilon)}$ for different values ε , namely $\varepsilon = 1, 2, 3$, using a single Gauss–Lobatto–Legendre spectral element of degree $p = 50$ for the spatial discretization. The results are presented in Figure 3. Obviously the error decays according to the theory developed in the preceding sections. For comparison purposes the relative approximation error of the standard stochastic Galerkin solution $\hat{u}_{N,hp}$ for the unweighted test function space, *i.e.* $\varepsilon = 0$, is also shown in Figure 3. After a period of decrease the approximation error increases suggesting that this sequence of standard stochastic Galerkin solutions does not converge to the exact solution \hat{u} in U^1 .

5. SUMMARY

In this work we have presented two possible stochastic variational formulations – an unweighted and a weighted one – for elliptic partial differential equations with random coefficients where the coefficients are bounded strictly away from zero and infinity by random variables. For each formulation we have established a solution theory and applied the stochastic Galerkin method for the numerical solution. It turns out that the stochastic Galerkin method employed on the unweighted weak formulation corresponds to the standard stochastic Galerkin approach whereas the stochastic Galerkin method employed on the weighted weak formulation can be interpreted as a stochastic Petrov–Galerkin approach since the test function space does no longer coincide with the solution space but corresponds to the solution space weighted by a suitable random variable. We have analyzed the convergence of these two different approximations to the exact solution. We have presented an example where the standard stochastic Galerkin approach produces approximations which do not converge to the exact solution in the natural norm. This does not mean, however, that this approach never yields convergent approximate solutions to the weak problem. But the convergence of the standard stochastic Galerkin approximation depends on the underlying boundary-value problem and must be proved separately. As possible remedy we have presented the stochastic Petrov–Galerkin approach with a modified test function space. We have established conditions for the convergence of the stochastic Petrov–Galerkin solutions in the natural norm and presented examples where this approach produces convergent approximations despite the failure of the standard stochastic Galerkin approach.

APPENDIX A. ON THE CONVERGENCE OF $\hat{u}_N(\xi)$ IN EXAMPLE 4.1

We consider the algebraic equation

$$\kappa u = f$$

for a given random coefficient κ with

$$\kappa(\omega) = \exp(\xi(\omega))$$

and random right side f with

$$f(\omega) = \exp(\xi(\omega))|\xi(\omega) - 1|$$

where the random variable ξ is standard Gaussian. Clearly, the exact solution is

$$u(\omega) = u_\xi(\xi(\omega)) = |\xi(\omega) - 1| \in L^2(\Omega, \sigma(\xi), \mathbf{P}).$$

As discussed in Example 4.1 of Section 4 the finite-dimensional solution and test function space of the standard stochastic Galerkin approach is the subspace

$$U_N = \text{span} \left\{ \frac{\text{He}_n(\xi)}{\sqrt{n!}}, n = 0, \dots, N \right\}$$

spanned by the probabilistic Hermite polynomials. Thus, the corresponding solution $u_N \in U_N$ solving

$$\mathbf{E}_{\mathbf{P}} \left(\exp(\xi) u_N \frac{\text{He}_m(\xi)}{\sqrt{m!}} \right) = \mathbf{E}_{\mathbf{P}} \left(\exp(\xi) |\xi - 1| \frac{\text{He}_m(\xi)}{\sqrt{m!}} \right) \quad \text{for all } m = 0, \dots, N \tag{A.1}$$

is equal to the stochastic part $\hat{u}_N(\xi)$ of the stochastic Galerkin solution $\hat{u}_{N,hp}$ of Example 4.1. Since $u_N \in U_N$ satisfies (A.1) there holds

$$\mathbf{E}_{\mathbf{P}} (\kappa u_N v) - \mathbf{E}_{\mathbf{P}} (f v) = \mathbf{E}_{\mathbf{P}} (\kappa (u_N - u) v) = 0 \quad \text{for all } v \in U_N.$$

Thus, the stochastic Galerkin solution u_N is the best approximation to the exact solution $u \in L^2(\Omega, \sigma(\xi), \kappa d\mathbf{P})$ in the subspace U_N . In other words, u_N is the orthogonal projection of u onto U_N in $L^2(\Omega, \sigma(\xi), \kappa d\mathbf{P})$, *i.e.* with

respect to the weighted inner product $\langle v_1, v_2 \rangle_\kappa := \mathbf{E}_\mathbf{P}(\kappa v_1 v_2)$. To prove that the solution u_N does not converge in quadratic mean to the exact solution u we consider the orthonormal basis $\{\tilde{p}_n(\xi), n \in \mathbb{N}_0\}$ of the weighted space $L^2(\Omega, \sigma(\xi), \exp(\xi)d\mathbf{P})$ which is given by

$$\tilde{p}_n(\xi) := \frac{e^{-1/4} \text{He}_n(\xi - 1)}{\sqrt{n!}}, \quad n \in \mathbb{N}_0.$$

Hence, the stochastic Galerkin solution u_N can be written as

$$u_N = \sum_{n=0}^N \langle u, \tilde{p}_n(\xi) \rangle_\kappa \tilde{p}_n(\xi).$$

Assume that the polynomial chaos coefficients $c_i^{(n)}, i = 0, \dots, n, n \in \mathbb{N}_0$, in the expansions

$$\tilde{p}_n(\xi) = \sum_{i=0}^n c_i^{(n)} \frac{\text{He}_i(\xi)}{\sqrt{i!}},$$

are known, then the stochastic Galerkin solution of order N with respect to the Hermite polynomials takes on the form

$$u_N = \sum_{n=0}^N \langle u, \tilde{p}_n(\xi) \rangle_\kappa \tilde{p}_n(\xi) = \sum_{n=0}^N \langle u, \tilde{p}_n(\xi) \rangle_\kappa \sum_{i=0}^n c_i^{(n)} \frac{\text{He}_i(\xi)}{\sqrt{i!}} = \sum_{i=0}^N \frac{\text{He}_i(\xi)}{\sqrt{i!}} \sum_{n=i}^N \langle u, \tilde{p}_n(\xi) \rangle_\kappa c_i^{(n)}.$$

If the stochastic Galerkin solutions u_N converge to the exact solution, their second order moments must converge to the exact solution's second order moment since

$$\| \|u\|_{L^2(\Omega, \mathfrak{A}, \mathbf{P})} - \|u_N\|_{L^2(\Omega, \mathfrak{A}, \mathbf{P})} \| \leq \|u - u_N\|_{L^2(\Omega, \mathfrak{A}, \mathbf{P})} \rightarrow 0, \quad N \rightarrow \infty,$$

i.e., there must hold

$$\|u_N\|_{L^2(\Omega, \mathfrak{A}, \mathbf{P})}^2 = \mathbf{E}_\mathbf{P}(u_N)^2 = \sum_{i=0}^N \left(\sum_{n=i}^N \langle u, \tilde{p}_n(\xi) \rangle_\kappa c_i^{(n)} \right)^2 \rightarrow \|u\|_{L^2(\Omega, \mathfrak{A}, \mathbf{P})}^2 = \mathbf{E}_\mathbf{P}u^2 < \infty, \quad N \rightarrow \infty.$$

However, we show that there exists a subsequence of stochastic Galerkin solutions $(u_N)_{N \in \mathbb{N}}$ where the second order moments diverge to infinity.

To see this, we evaluate the generalized polynomial chaos coefficients $\langle u, \tilde{p}_n(\xi) \rangle_\kappa$ for $n \in \mathbb{N}_0$ by the change of variables $y = x - 1$ yielding

$$\begin{aligned} \langle u, \tilde{p}_n(\xi) \rangle_\kappa &= \int_{-\infty}^{\infty} u_\xi(x) \frac{e^{-1/4} \text{He}_n(x - 1)}{\sqrt{n!}} \exp(x) \frac{\exp(-\frac{1}{2}x^2)}{\sqrt{2\pi}} dx \\ &= e^{1/4} \int_{-\infty}^{\infty} u_\xi(x) \frac{\text{He}_n(x - 1)}{\sqrt{n!}} \frac{\exp(-\frac{(x-1)^2}{2})}{\sqrt{2\pi}} dx \\ &= e^{1/4} \int_{-\infty}^{\infty} u_\xi(y + 1) \frac{\text{He}_n(y)}{\sqrt{n!}} \frac{\exp(-\frac{y^2}{2})}{\sqrt{2\pi}} dy = e^{1/4} \mathbf{E}_\mathbf{P} \left(u_\xi(\xi + 1) \frac{\text{He}_n(\xi)}{\sqrt{n!}} \right) \\ &= e^{1/4} \mathbf{E}_\mathbf{P} \left(|\xi| \frac{\text{He}_n(\xi)}{\sqrt{n!}} \right) = \begin{cases} e^{1/4} \frac{(-1)^{n/2+1} \sqrt{n!}}{\sqrt{2\pi} 2^{n/2-1} (\frac{n}{2})!(n-1)}, & n \text{ even,} \\ 0, & \text{otherwise.} \end{cases} \end{aligned}$$

Expanding the polynomials $\tilde{p}_n(\xi)$ in polynomial chaos leads us to the sums

$$\tilde{p}_n(\xi) = \frac{e^{-1/4}\text{He}_n(\xi - 1)}{\sqrt{n!}} = \sum_{i=0}^n c_i^{(n)} \frac{\text{He}_i(\xi)}{\sqrt{i!}}$$

where the coefficients are given by

$$c_i^{(n)} = e^{-1/4} \frac{(-1)^{n-i} \sqrt{n!}}{(n-i)! \sqrt{i!}}$$

due to Formula 22.12.8 in [1]. Combining these two results we obtain

$$\begin{aligned} \langle u, \tilde{p}_n(\xi) \rangle_{\kappa} c_i^{(n)} &= \begin{cases} e^{1/4} \frac{(-1)^{n/2+1} \sqrt{n!}}{\sqrt{2\pi} 2^{n/2-1} (\frac{n}{2})! (n-1)} e^{-1/4} \frac{(-1)^{n-i} \sqrt{n!}}{(n-i)! \sqrt{i!}}, & n \text{ even,} \\ 0, & \text{otherwise,} \end{cases} \\ &= \begin{cases} \frac{(-1)^{3n/2+1-i} n!}{\sqrt{2\pi} 2^{n/2-1} (\frac{n}{2})! (n-1) (n-i)! \sqrt{i!}}, & n \text{ even,} \\ 0, & \text{otherwise.} \end{cases} \end{aligned}$$

Now, we investigate the second order moments of the subsequence $(u_{2N})_{N \in \mathbb{N}_0}$ of the stochastic Galerkin solutions, that is,

$$\mathbf{E}_{\mathbf{P}}(u_{2N})^2 = \sum_{i=0}^{2N} \left(\sum_{n=i}^{2N} \langle u, \tilde{p}_n(\xi) \rangle_{\kappa} c_i^{(n)} \right)^2.$$

In order to show that this sequence goes to infinity for $N \rightarrow \infty$ we consider a single term in this sum, namely the summand for $i = 2N - 3$:

$$\begin{aligned} &\left(\sum_{n=2N-3}^{2N} \langle u, \tilde{p}_{2n}(\xi) \rangle_{\kappa} c_i^{(n)} \right)^2 \\ &= \left(\langle u, \tilde{p}_{2N-2}(\xi) \rangle_{\kappa} c_{2N-3}^{(2N-2)} + \langle u, \tilde{p}_{2N}(\xi) \rangle_{\kappa} c_{2N-3}^{(2N)} \right)^2 \\ &= \left(\frac{(-1)^{N+1} (2N-2)!}{\sqrt{2\pi} 2^{N-2} (N-1)! (2N-3) \sqrt{(2N-3)!}} + \frac{(-1)^N (2N)!}{\sqrt{2\pi} 2^{N-1} N! (2N-1) 3! \sqrt{(2N-3)!}} \right)^2 \\ &= \left(\frac{(2N-2)!}{\sqrt{2\pi} 2^{N-1} (N-1)! \sqrt{(2N-3)!}} \left(\frac{(-1)^{N+1} 2}{2N-3} + \frac{(-1)^N}{3} \right) \right)^2 \\ &= \frac{(2N-2)!^2}{2\pi 2^{2N-2} (N-1)!^2 (2N-3)!} \left(\frac{4}{(2N-3)^2} - \frac{4}{3(2N-3)} + \frac{1}{9} \right) \\ &= \frac{(2N-2)! 2^{-(2N-2)}}{\pi (N-1)!^2} \left(\frac{4(N-1)}{(2N-3)^2} - \frac{4(N-1)}{3(2N-3)} + \frac{N-1}{9} \right). \end{aligned}$$

Then, applying the Stirling's formula (see *e.g.* [1]) yields

$$\underbrace{\frac{(2N-2)! 2^{-(2N-2)}}{\pi (N-1)!^2}}_{\simeq \frac{1}{\sqrt{N}}} \underbrace{\left(\frac{4(N-1)}{(2N-3)^2} - \frac{4(N-1)}{3(2N-3)} + \frac{N-1}{9} \right)}_{\simeq N} \simeq \sqrt{N}$$

when $N \rightarrow \infty$ and thus we obtain

$$\begin{aligned} \mathbf{E}_{\mathbf{P}}(u_{2N})^2 &= \sum_{i=0}^{2N} \left(\sum_{n=i}^{2N} \langle u, \tilde{p}_n(\xi) \rangle_{\kappa} c_i^{(n)} \right)^2 \\ &\geq \frac{(2N-2)!2^{-(2N-2)}}{\pi(N-1)!^2} \left(\frac{4(N-1)}{(2N-3)^2} - \frac{4(N-1)}{3(2N-3)} + \frac{N-1}{9} \right) \rightarrow \infty \end{aligned}$$

for $N \rightarrow \infty$. This implies that the sequence of stochastic Galerkin solutions $(u_N)_{N \in \mathbb{N}_0}$ fails to converge in quadratic mean to the exact solution u .

REFERENCES

- [1] M. Abramowitz and I.A. Stegun, *Handbook of Mathematical Functions*. Dover Publications, Inc, New York (1965).
- [2] I. Babuška, F. Nobile and R. Tempone, A Stochastic Collocation Method for Elliptic Partial Differential Equations with Random Input Data. *SIAM J. Numer. Anal.* **45** (2007) 1005–1034.
- [3] I. Babuška, R. Tempone and G.E. Zouraris, Galerkin Finite Element Approximations of Stochastic Elliptic Partial Differential Equations. *SIAM J. Numer. Anal.* **42** (2004) 800–825.
- [4] I. 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** (2005) 1251–1294.
- [5] M. Bieri, R. Andreev and C. Schwab, Sparse tensor discretization of elliptic SPDEs. *SIAM J. Sci. Comput.* **31** (2009) 4281–4304.
- [6] M. Bieri and C. Schwab, Sparse high order FEM for elliptic sPDEs. *Comput. Methods Appl. Mech. Engrg.* **198** (2009) 1149–1170.
- [7] A. Bobrowski, *Functional Analysis for Probability and Stochastic Processes*. Cambridge University Press, Cambridge UK (2005).
- [8] S.C. Brenner and L.R. Scott, The Mathematical Theory of Finite Element Methods, 2nd ed. *Texts in Appl. Math.*, vol. 15. Springer-Verlag, New York (2002).
- [9] C. Canuto, M.Y. Hussaini, A. Quarteroni and T.A. Zang, *Spectral Methods: Fundamentals in Single Domains*. Springer-Verlag, Berlin Heidelberg (2006).
- [10] A. Cohen, R. DeVore and C. Schwab, Convergence Rates of Best N-term Galerkin Approximations for a Class of Elliptic sPDEs. *Foundations Comput. Math.* **10** (2010) 615–646.
- [11] M. K. Deb, I. Babuška and J.T. Oden, Solution of stochastic partial differential equations using Galerkin finite element techniques. *Comput. Methods Appl. Mech. Engrg.* **190** (2001) 6359–6372.
- [12] M. Eiermann, O.G. Ernst and E. Ullmann, Computational aspects of the stochastic finite element method. *Comput. Visualiz. Sci.* **10** (2007) 3–15.
- [13] O. Ernst, A. Mugler, E. Ullmann and H.J. Starkloff, On the convergence of generalized polynomial chaos. *ESAIM: M2AN* **46** (2012) 317–339.
- [14] R.V. Field and M. Grigoriu, *Convergence Properties of Polynomial Chaos Approximations for L^2 -Random Variables*, Sandia Report SAND2007-1262 (2007).
- [15] P. Frauenfelder, C. Schwab and R.A. Todor, Finite elements for elliptic problems with stochastic coefficients. *Comput. Methods Appl. Mech. Engrg.* **194** (2005) 205–228.
- [16] R.A. Freeze, A Stochastic-Conceptual Analysis of One-Dimensional Groundwater Flow in Nonuniform Homogeneous Media. *Water Resources Research* (1975) 725–741.
- [17] J. Galvis and M. Sarkis, Approximating infinity-dimensional stochastic Darcy’s Equations without uniform ellipticity. *SIAM J. Numer. Anal.* **47** (2009) 3624–3651.
- [18] J. Galvis and M. Sarkis, Regularity results for the ordinary product stochastic pressure equation, to appear in *SIAM J. Math. Anal.* (preprint 2011) 1–31.
- [19] R. Ghanem, Ingredients for a general purpose stochastic finite elements implementation. *Comput. Methods Appl. Mech. Engrg.* **168** (1999) 19–34.
- [20] R. Ghanem, Stochastic Finite Elements with Multiple Random Non-Gaussian Properties. *J. Engrg. Mech.* **125** (1999) 26–40.
- [21] R. Ghanem and S. Dham, Stochastic Finite Element Analysis for Multiphase Flow in Heterogeneous Porous Media. *Transport in Porous Media* **32** (1998) 239–262.
- [22] R. Ghanem and P.D. Spanos, *Stochastic Finite Elements: A Spectral Approach*. Springer-Verlag, New York (1991).
- [23] C.J. Gittelson, *Stochastic Galerkin discretization of the log-normal isotropic diffusion problem*. *Math. Models Methods Appl. Sci.* **20** (2010) 237–263.
- [24] E. Godoy, A. Ronveaux, A. Zarzo and I. Area, Minimal recurrence relations for connection coefficients between classical orthogonal polynomials: Continuous case. *J. Computat. Appl. Math.* **84** (1997) 257–275.
- [25] E. Hille and R.S. Phillips, *Functional Analysis and Semi-Groups*, Colloquium Publications. *Amer. Math. Soc.* **31** (1957).

- [26] O. Kallenberg, *Foundations of modern probability*. Springer-Verlag, Berlin (2001).
- [27] O.P. Le Maître and O.M. Knio, *Spectral Methods for Uncertainty Quantification. Scientific Computation: With Applications to Computational Fluid Dynamics*. Springer-Verlag (2010).
- [28] M. Loève, *Probability Theory II. 4th Edition*. Springer-Verlag, New York, Heidelberg, Berlin (1978).
- [29] D. Lucor, C. H Su and G.E. Karniadakis, Generalized polynomial chaos and random oscillators. *Int. J. Numer. Methods Engrg.* **60** (2004) 571–596.
- [30] P. Maroni and Z. Rocha, Connection coefficients between orthogonal polynomials and the canonical sequence: an approach based on symbolic computation. *Numer. Algor.* **47** (2008) 291–314.
- [31] H.G. Matthies and A. Keese, Galerkin methods for linear and nonlinear elliptic stochastic partial differential equations. *Comput. Methods Appl. Mech. Engrg.* **194** (2005) 1295–1331.
- [32] A. Mugler and H.J. Starkloff, On elliptic partial differential equations with random coefficients. *Stud. Univ. Babeş-Bolyai Math.* **56** (2011) 473–487.
- [33] A. Narayan and J.S. Hesthaven, Computation of connection coefficients and measure modifications for orthogonal polynomials. *BIT Numer. Math.* (2011).
- [34] W. Nowak, *Geostatistical Methods for the Identification of Flow and Transport Parameters in the Subsurface*, Ph.D. Thesis. Universität Stuttgart (2005).
- [35] C. Schwab and C.J. Gittelsohn, Sparse tensor discretizations of high-dimensional parametric and stochastic PDEs. *Acta Numer.* **20** (2011) 291–467.
- [36] R. A. Todor and C. Schwab, Convergence rates for sparse chaos approximations of elliptic problems with stochastic coefficients. *IMA J. Numer. Anal.* **27** (2007) 232–261.
- [37] H. Umegaki and A.T. Bharucha-Reid, Banach Space-Valued Random Variables and Tensor Products of Banach Spaces. *J. Math. Anal. Appl.* **31** (1970) 49–67.
- [38] X. Wan and G.E. Karniadakis, Beyond Wiener Askey Expansions: Handling Arbitrary PDFs. *J. Scient. Comput.* **27** (2006) 455–464.
- [39] N. Wiener, Homogeneous Chaos. *Amer. J. Math.* **60** (1938) 897–936.
- [40] D. Xiu, *Numerical methods for stochastic computations: A spectral method approach*. Princeton Univ. Press, Princeton and NJ (2010).
- [41] D. Xiu and G.E. Karniadakis, Modeling uncertainty in steady state diffusion problems via generalized polynomial chaos. *Comput. Methods Appl. Mech. Engrg.* **191** (2002) 4927–4948.
- [42] D. Xiu and G.E. Karniadakis, The Wiener-Askey Polynomial Chaos for Stochastic Differential Equations. *SIAM J. Sci. Comput.* **24** (2002) 619–644.
- [43] D. Xiu and G.E. Karniadakis, A new stochastic approach to transient heat conduction modeling with uncertainty. *Inter. J. Heat and Mass Transfer* **46** (2003) 4681–4693.
- [44] D. Xiu and G.E. Karniadakis, Modeling uncertainty in flow simulations via generalized polynomial chaos. *J. Comput. Phys.* **187** (2003) 137–167.
- [45] D. Xiu, D. Lucor, C. H Su and G.E. Karniadakis, *Performance Evaluation of Generalized Polynomial Chaos*, Computational Science – ICCS 2003, edited by P.M.A. Sloot, D. Abramson, A.V. Bogdanov, J.J. Dongarra, Albert Y. Zomaya and Y.E. Gorbachev, *Lect. Notes Comput. Sci.*, vol. 2660. Springer Verlag (2003).
- [46] D. Zhang, *Stochastic Methods for Flow in Porous Media. Coping with Uncertainties*. Academic Press, San Diego, CA (2002).