# On the optimal polynomial approximation of stochastic PDEs by Galerkin and Collocation methods

BECK, J.; NOBILE, F.; TAMELLINI, L.; TEMPONE, R.

# On the optimal polynomial approximation of stochastic PDEs by Galerkin and Collocation methods [*]

Joakim Beck[♭], Fabio Nobile[♯], Lorenzo Tamellini[♯], Raul Tempone[♭]

May 12, 2011

[♭] Applied Mathematics and Computational Science
4700 - King Abdullah University of Science and Technology
Thuwal 23955-6900, Kingdom of Saudi Arabia
joakim.back.09@ucl.ac.uk, raul.tempone@kaust.edu.sa

[♯] MOX– Modellistica e Calcolo Scientifico
Dipartimento di Matematica "F. Brioschi"
Politecnico di Milano
via Bonardi 9, 20133 Milano, Italy
fabio.nobile@polimi.it, lorenzo.tamellini@mail.polimi.it

**Keywords**: Uncertainty Quantification, PDEs with random data, elliptic equations, multivariate polynomial approximation, best $M$-terms approximation, Stochastic Galerkin methods, Smolyak approximation, Sparse grids, Stochastic Collocation methods.

**AMS Subject Classification**: 41A10, 65C20, 65N12, 65N35

## Abstract

In this work we focus on the numerical approximation of the solution $u$ of a linear elliptic PDE with stochastic coefficients. The problem is rewritten as a parametric PDE and the functional dependence of the solution on the parameters is approximated by multivariate polynomials. We first consider the Stochastic Galerkin method, and rely on sharp estimates for the decay of the Fourier coefficients of the spectral expansion of $u$ on an orthogonal

polynomial basis to build a sequence of polynomial subspaces that features better convergence properties, in terms of error versus number of degrees of freedom, than standard choices such as Total Degree or Tensor Product subspaces.

We consider then the Stochastic Collocation method, and use the previous estimates to introduce a new class of Sparse Grids, based on the idea of selecting a priori the most profitable hierarchical surpluses, that, again, features better convergence properties compared to standard Smolyak or tensor product grids. Numerical results show the effectiveness of the newly introduced polynomial spaces and sparse grids.

# 1 Introduction

Many works have been recently devoted to the analysis and the improvement of the Stochastic Galerkin and Collocation techniques for Uncertainty Quantification for PDEs with random input data. These methods are promising since they can exploit the possible regularity of the solution with respect to the stochastic parameters to achieve faster convergence than sampling methods like Monte Carlo.

Stochastic Galerkin and Collocation methods can be classified as *parametric* techniques, since both approximate $u$, the solution of the PDE as a linear combination of suitable deterministic basis functions in probability space, typically polynomials or piecewise polynomials. In this work we focus only on global multivariate polynomial approximations. Stochastic Galerkin is a projection technique over a set of orthogonal polynomials with respect to the probability measure at hand (see e.g. [1, 15, 19, 27, 30]), while Collocation is a sum of Lagrangian interpolants over the probability space (see e.g. [2, 12, 29]).

The comparison between performances of these methods is a matter of study (see e.g. [3, 10]). However, both suffer the so-called "Curse of Dimensionality": using naive projections/interpolations over tensor product polynomials spaces/grids leads to computational costs that grow exponentially fast with the number of random variables. Therefore the main requirement for these methods to be appealing is the capability of retaining good approximations of $u$ while keeping the computational cost as low as possible.

In a Stochastic Galerkin setting this requirement can be translated to the implementation of algorithms able to compute what is known as "best $M$-terms approximation". In other words, the method should be able to establish a-priori the set of the $M$ most fruitful multivariate orthogonal polynomials in the spectral approximation of $u$, and to compute only those terms.

Important contributions in the study of the best $M$-terms approximation have been given by Schwab and co-workers: estimates on the decay of the coefficients of the spectral expansion of $u$ have been proved e.g. in [5, 8, 7]. In this work we will reformulate and slightly generalize the result given in [8, Corollary 6.1], and show on few numerical examples that the sequence of polynomial sub-

spaces built upon those estimates ("TD with factorial correction" sets, TD-FC in the following) performs better than classical choices such as Total Degree or Tensor Product in terms of error versus the dimension of the polynomial space.

In a Stochastic Collocation setting, the construction of an optimal grid can be recast to a classical knapsack problem and relies on the notion of profit of each hierarchical surplus composing the sparse grid, as introduced e.g. in [6]. The "Best $M$-Terms" grid is then the one built with the set of the $M$ most profitable hierarchical surpluses. In this work we propose a heuristic estimate of the profit of each hierarchical surplus, and use it to build a quasi optimal sparse grid. The estimates of the profit are in turn based on the estimates of the decay of the spectral expansion of $u$. Numerical investigations show that these new grids perform better than standard Smolyak grids as well as grids constructed with the dimension adaptive approach developed in [14, 17]. A similar knapsack approach to the construction of generalized optimal sparse grids has been proposed also in [16]. Our contribution extends and details the procedure to the case of PDEs with stochastic coefficients, working with analytic functions instead of $H^r_{mix}$ ones, and using sharp estimates for the profits of the hierarchical surpluses.

The paper is organized as follows. Section 2 defines the elliptic model problem of interest and gives general regularity results of the solution $u$. In Section 3 we first address the general procedure that leads to the Stochastic Galerkin approximation of $u$; next we state the estimate for the decay of the spectral approximation of $u$ and explain how to build practically the TD-FC polynomial subspaces that stem from it. In Section 3.2 we consider some simple numerical tests where we can build explicitly the best $M$-terms approximation, and we compare it with the TD-FC and with some standard choices of polynomial subspaces. In Section 4 we recall the construction of a general sparse grid, motivate our heuristic estimate of the profit of each hierarchical surplus and explain how to construct in practice optimized sparse grids based on such estimates. Section 4.2 shows on some simple test cases the effectiveness of the method and the sharpness of our heuristic estimates. Finally 5 draws some conclusions.

## 2 Problem setting

Let $D$ be a convex bounded polygonal domain in $\mathbb{R}^d$ and $(\Omega, \mathcal{F}, P)$ be a complete probability space. Here $\Omega$ is the set of outcomes, $\mathcal{F} \subset 2^\Omega$ is the $\sigma$-algebra of events and $P : \mathcal{F} \to [0, 1]$ is a probability measure. Consider the stochastic linear elliptic boundary value problem:

**Strong Formulation.** *find a random function, $u : \overline{D} \times \Omega \to \mathbb{R}$, such that $P$-almost everywhere in $\Omega$, or in other words almost surely (a.s.), the following equation holds:*

$$\begin{cases} -\operatorname{div}(a(\mathbf{x}, \omega) \nabla u(\mathbf{x}, \omega)) = f(\mathbf{x}) & \mathbf{x} \in D, \\ u(\mathbf{x}, \omega) = 0 & \mathbf{x} \in \partial D. \end{cases} \tag{1}$$

3

*where the operators* div *and* $\nabla$ *imply differentiation with respect to the physical coordinate only.*

We make the following assumptions on the random diffusion coefficient:

**Assumption 2.1** (Coercivity and continuity). *$a(\mathbf{x}, \omega)$ is strictly positive and bounded with probability 1, i.e. there exist $a_{min} > 0$ and $a_{max} < \infty$ such that $P(a_{min} \leq a(\mathbf{x}, \omega) \leq a_{max}, \forall \mathbf{x} \in \overline{D}) = 1$.*

**Assumption 2.2** (Finite dimensional noise). *$a(\mathbf{x}, \omega)$ is parametrized by a set of $N$ independent and identically distributed uniform random variables in $(-1, 1)$, $\mathbf{y}(\omega) = [y_1(\omega), ..., y_N(\omega)]^T : \Omega \to \mathbb{R}^N$.*

Observe that the assumption that the random variables are uniform is not that restrictive. Indeed, we could assume that $a$ is parametrized by $N$ random variables $z_i$, $i = 1, \ldots, n$ and introduce a non linear map $y_i = \Theta(z_i)$ that transforms each of them into uniform random variables, following the well known theory on copulas, see e.g. [20].

We denote by $\Gamma_n = (-1, 1)$ the image set of the random variable $y_n$, and let $\Gamma = \Gamma_1 \times \ldots \times \Gamma_N$. After Assumption 2.2 the random vector $\mathbf{y}$ has a joint probability density function $\rho : \Gamma \to \mathbb{R}_+$ that factorizes as $\rho(\mathbf{y}) = \prod_{n=1}^N \rho_n(y_n)$, $\forall \mathbf{y} \in \Gamma$, with $\rho_n = \frac{1}{2}$. Moreover, the solution $u$ of (1) depends on the single realization $\omega \in \Omega$ only through the value taken by the random vector $\mathbf{y}$. We can therefore replace the probability space $(\Omega, \mathcal{F}, P)$ with $(\Gamma, B(\Gamma), \rho(\mathbf{y})d\mathbf{y})$, where $B(\Gamma)$ denotes the Borel $\sigma$-algebra on $\Gamma$ and $\rho(\mathbf{y})d\mathbf{y}$ is the distribution measure of the vector $\mathbf{y}$. We denote with $L_\rho^2(\Gamma)$ the space of square integrable functions on $\Gamma$ with respect to the measure $\frac{1}{2^N}d\mathbf{y}$. Note that in the case the original random variables are not uniform but with bounded support, and a mapping $\Theta$ is not available, one could still reduce the problem to the uniform case, at the price of bounding $\rho(\mathbf{y})$ with $\|\rho(\mathbf{y})\|_\infty$.

In the rest of the paper we will use the following notation: given a multi-index $\mathbf{i} \in \mathbb{N}^N$ and a vector $\mathbf{r} \in \mathbb{R}^N$, we define $|\mathbf{i}| = \sum_{n=1}^N i_n$, $\mathbf{i}! = \prod_{n=1}^N (i_n!)$ and $\mathbf{r}^\mathbf{i} = \prod_{n=1}^N r_n^{i_n}$. We can now state a regularity assumption on $a(\mathbf{x}, \mathbf{y})$:

**Assumption 2.3** (Stochastic regularity). *$a(\mathbf{x}, \mathbf{y})$ is infinitely many times differentiable with respect to $\mathbf{y}$ and $\exists \mathbf{r} \in \mathbb{R}_+^N$ s.t.*

$$\left\| \frac{\partial_\mathbf{i} a}{a}(\cdot, \mathbf{y}) \right\|_{L^\infty(D)} \leq \mathbf{r}^\mathbf{i} \quad \forall \mathbf{y} \in \Gamma,$$

*where $\mathbf{i}$ is a multi-index in $\mathbb{N}^N$, $\partial_\mathbf{i} a = \dfrac{\partial^{i_1 + \ldots + i_N} a}{\partial y_1^{i_1} \cdots \partial y_N^{i_N}}$, and $\mathbf{r}$ is independent of $\mathbf{y}$.*

**Example 2.1** (Stochastic regularity). *A common situation of interest is when $a(\mathbf{x}, \omega)$ is an infinitely dimensional random field, suitably expanded in series (e.g. by a Karhunen-Loève or Fourier expansion) either as a* linear *expansion of the*

*form $a = a_0 + \sum_{n=1}^{\infty} b_n(\mathbf{x}) y_n$ with $b_n \in L^{\infty}(D)$ and $a_{min} = a_0 - \sum_{n=1}^{\infty} \|b_n\|_{L^{\infty}(D)}$, or an exponential expansion of the form $a = a_0 + \exp\left(\sum_{n=1}^{\infty} b_n(\mathbf{x}) y_n\right)$. Then the infinite series is truncated up to $N$ terms, with $N$ large enough to take into account a sufficiently large amount of the total variability. Both expansions comply with Assumption 2.3 taking $r_n = \|b_n\|_{L^{\infty}(D)}/a_{min}$ and $r_n = \|b_n\|_{L^{\infty}(D)}$, respectively.*

Finally, we denote by $V = H_0^1(D)$ the space of square integrable functions in $D$ with square integrable distributional derivatives and with zero trace on the boundary, equipped with the gradient norm $\|v\|_V = \|\nabla v\|_{L^2(D)}$, $\forall v \in V$. Its dual space will be denoted by $V'$. We are now in the position to write a weak formulation of problem (1):

**Weak Formulation.** *Find $u \in V \otimes L_\rho^2(\Gamma)$ such that $\forall v \in V \otimes L_\rho^2(\Gamma)$*

$$\int_\Gamma \int_D a(\mathbf{x}, \mathbf{y}) \nabla u(\mathbf{x}, \mathbf{y}) \cdot \nabla v(\mathbf{x}, \mathbf{y}) \, \rho(\mathbf{y}) \, d\mathbf{x} \, d\mathbf{y} = \int_\Gamma \int_D f(\mathbf{x}) v(\mathbf{x}, \mathbf{y}) \, \rho(\mathbf{y}) \, d\mathbf{x} \, d\mathbf{y}. \quad (2)$$

Under Assumption 2.1, the Lax-Milgram lemma yields that there exists a unique solution to problem (2) for any $f \in V'$. Moreover, the following estimate holds:

$$\|u\|_{V \otimes L_\rho^2(\Gamma)} \leq \frac{\|f\|_{V'}}{a_{min}}.$$

The solution $u$ can also be thought as a function defined in $\Gamma$ with solution in $V$, $u : \Gamma \to V$ and, thanks to the previous result, we have $u \in L_\rho^2(\Gamma, V) = V \otimes L_\rho^2$. In what follows we will often use the notation $u(\mathbf{y}) := u(\cdot, \mathbf{y}) \in V$ if no confusion arises.

Concerning the regularity of the solution with respect to $\mathbf{y}$, the following result holds, which generalizes the result given in [8] for the special case $a = a_0 + \sum_{n=1}^{N} b_n(x) y_n$.

**Theorem 2.1.** *Let $a(\mathbf{x}, \mathbf{y})$ be a diffusion coefficient for equation (1) that satisfies Assumptions 2.1 - 2.3. Then the derivatives of $u$ can be bounded as*

$$\|\partial_{\mathbf{i}} u(\mathbf{y})\|_V \leq C_0 |\mathbf{i}|! \, \tilde{\mathbf{r}}^{\mathbf{i}} \quad \forall \mathbf{y} \in \Gamma.$$

*Here $C_0 = \dfrac{\|f\|_{V'}}{a_{min}}$ and $\tilde{\mathbf{r}} = \left(\dfrac{1}{\log 2}\right) \mathbf{r}$, with $\mathbf{r}$ as in Assumption 2.3.*

The proof is technical; we thus postpone it to the Appendix. A consequence of Theorem 2.1 is that $u$ is analytic in every $\mathbf{y} \in \Gamma$.

**Corollary 2.1.** *Under the hypotheses of Theorem 2.1, given $\varepsilon > 0$, for every $\mathbf{y}_0 \in \Gamma$ the Taylor series of $u$ converges in the disk*

$$\mathcal{D}(\mathbf{y}_0) = \left\{ \mathbf{y} \in \mathbb{R}^N : \tilde{\mathbf{r}} \cdot \mathrm{abs}\,(\mathbf{y} - \mathbf{y}_0) < 1 \right\}.$$

*where* $\mathrm{abs}\,(\mathbf{v}) = (|v_1|,\dots,|v_N|)^T$. *Therefore* $u : \Gamma \to V$ *is analytic and can be extended analytically to the set*

$$\Sigma = \left\{ \mathbf{y} \in \mathbb{R}^N : \exists\, \mathbf{y}_0 \in \Gamma \text{ s.t. } \tilde{\mathbf{r}} \cdot \mathrm{abs}\,(\mathbf{y} - \mathbf{y}_0) < 1 \right\}.$$

**Proof.**    Use first Theorem 2.1 to bound the norm of the Taylor expansion of $u(\mathbf{y})$ centered in $\mathbf{y}_0 \in \Gamma$ as

$$\left\| \sum_{k=0}^{\infty} \sum_{|\mathbf{j}|=k} \frac{\partial_{\mathbf{j}} u(\mathbf{y}_0)}{\mathbf{j}!} (\mathbf{y} - \mathbf{y}_0)^{\mathbf{j}} \right\|_V \leq \sum_{k=0}^{\infty} \sum_{|\mathbf{j}|=k} C_0 \tilde{\mathbf{r}}^{\mathbf{j}} \frac{|\mathbf{j}|!}{\mathbf{j}!} \mathrm{abs}\,(\mathbf{y} - \mathbf{y}_0)^{\mathbf{j}}.$$

Next exploit the generalized Newton binomial formula, that states that for $\alpha_1,\dots,\alpha_N \in \mathbb{R}_+$ and $k \in \mathbb{N}$ we have

$$\sum_{|\mathbf{j}|=k} \frac{k!}{\mathbf{j}!} \boldsymbol{\alpha}^{\mathbf{j}} = \left( \sum_{n=1}^{N} \alpha_n \right)^k,$$

to rewrite the bound on the norm of the Taylor series as

$$\left\| \sum_{k=0}^{\infty} \sum_{|\mathbf{j}|=k} \frac{\partial_{\mathbf{j}} u(\mathbf{y}_0)}{\mathbf{j}!} (\mathbf{y} - \mathbf{y}_0)^{\mathbf{j}} \right\|_V \leq C_0 \sum_{k=0}^{\infty} \left( \sum_{n=1}^{N} \tilde{r}_n |y_n - y_{0,n}| \right)^k.$$

Thus the Taylor series of $u$ converges to $u$ in the disk $\mathcal{D}(\mathbf{y}_0)$. Therefore $u$ is analytic and admits an analytic extension in $\Sigma$. $\qquad\square$

# 3   Stochastic Galerkin method

We now seek an approximation of the solution $u$ with respect to $\mathbf{y}$ by global polynomials.

As anticipated in the introduction, we remark that the choice of the polynomial space is critical when the number $N$ of input random variables is large, since the number of stochastic degrees of freedom might grow very quickly with $N$, even exponentially when isotropic tensor product polynomial spaces are used (see Table 1). This effect is known as the *curse of dimensionality.*

Several choices of polynomial spaces that mitigate this phenomenon have been proposed in the literature, see e.g. [3]. Each of these polynomial spaces is built as the span of a properly selected subset of a multivariate orthonormal polynomial basis $\{\mathcal{L}_p(\mathbf{y})\}_{p\in\mathbb{N}}$ for $L^2_\rho(\Gamma)$, to retain good approximating properties with only a finite number of basis functions.

Since $L^2_\rho(\Gamma) = \bigotimes_{n=1}^{N} L^2_{\rho_n}(\Gamma_n)$, the elements of an orthonormal basis can be built as products of orthonormal polynomials for each of the directions $y_n$, $\{L_{p_n}(y_n)\}_{p_n\in\mathbb{N}}$ ; we can thus index the multivariate orthonormal polynomials basis functions $\mathcal{L}_p(\mathbf{y})$ with multi-indices $\mathbf{p} = (p_1,\dots,p_N)$

$$\mathcal{L}_{\mathbf{p}}(\mathbf{y}) = \prod_{n=1}^{N} L_{p_n}(y_n).$$

Then, by construction, the set $\{\mathcal{L}_{\mathbf{p}}(\mathbf{y})\}_{\mathbf{p}\in\mathbb{N}^N}$ is a $\rho(\mathbf{y})d\mathbf{y}$-orthonormal basis in $L^2_\rho(\Gamma)$, i.e. such that $\int_\Gamma \mathcal{L}_{\mathbf{p}}(\mathbf{y})\mathcal{L}_{\mathbf{q}}(\mathbf{y})\rho(\mathbf{y})\,d\mathbf{y} = 1$ if $\mathbf{p} = \mathbf{q}$ and 0 otherwise.

Let now $w \in \mathbb{N}$ be an integer index indicating the level of approximation, and $\Lambda(w)$ a sequence of increasing index sets such that

$$\Lambda(0) = \{(0,\ldots,0)\}, \;\; \Lambda(w) \subseteq \Lambda(w+1) \subset \mathbb{N}^N \text{ for } w \geq 0, \;\; \mathbb{N}^N = \bigcup_{w\in\mathbb{N}} \Lambda(w). \quad (3)$$

Denoting by $\mathbb{P}_{\Lambda(w)}(\Gamma)$ the multivariate polynomial space

$$\mathbb{P}_{\Lambda(w)}(\Gamma) = span\{\mathcal{L}_{\mathbf{p}}(\mathbf{y}), \;\; \mathbf{p} \in \Lambda(w)\}, \quad (4)$$

the Stochastic Galerkin (SG) approximation consists in restricting the weak formulation (2) to the subspace $V \otimes \mathbb{P}_{\Lambda(w)}(\Gamma)$ and reads:

**Galerkin Formulation.** *Find $u_w \in V \otimes \mathbb{P}_{\Lambda(w)}(\Gamma)$ such that $\forall\, v_w \in V \otimes \mathbb{P}_{\Lambda(w)}(\Gamma)$*

$$\int_\Gamma \int_D a(\mathbf{x},\mathbf{y})\nabla u_w(\mathbf{x},\mathbf{y}) \cdot \nabla v_w(\mathbf{x},\mathbf{y})\,\rho(\mathbf{y})\,d\mathbf{x}\,d\mathbf{y} = \int_\Gamma \int_D f(\mathbf{x})v_w(\mathbf{x},\mathbf{y})\,\rho(\mathbf{y})\,d\mathbf{x}\,d\mathbf{y},$$
$$(5)$$

where, due to the orthonormality of $\{\mathcal{L}_{\mathbf{p}}(\mathbf{y})\}_{\mathbf{p}\in\Lambda(w)}$,

$$u_w(\mathbf{x},\mathbf{y}) = \sum_{\mathbf{p}\in\Lambda(w)} u_{\mathbf{p}}(\mathbf{x})\mathcal{L}_{\mathbf{p}}(\mathbf{y}), \;\; u_{\mathbf{p}}(\mathbf{x}) = \int_\Gamma u(\mathbf{x},\mathbf{y})\mathcal{L}_{\mathbf{p}}(\mathbf{y})\rho(\mathbf{y})d\mathbf{y} \;\forall \mathbf{p} \in \Lambda(w).$$
$$(6)$$

Commonly used spaces $\mathbb{P}_{\Lambda(w)}(\Gamma)$ are listed in Table 1; for further details, see [3] and references therein.

| | index set $\Lambda(w)$ | Dimension $|\Lambda(w)|$ |
|---|---|---|
| **Tensor product (TP)** | $\{\mathbf{p} \in \mathbb{N}^N : \max_{n=1\ldots,N} p_n \leq w\}$ | $(1+w)^N$ |
| **Total degree (TD)** | $\{\mathbf{p} \in \mathbb{N}^N : \sum_{n=1}^N p_n \leq w\}$ | $\binom{N+w}{N}$ |
| **Hyperbolic cross (HC)** | $\{\mathbf{p} \in \mathbb{N}^N : \prod_{n=1}^N (p_n+1) \leq w+1\}$ | $(w+1)(\log(w+1))^{N-1}$ |

Table 1: Examples of typical polynomial spaces. The result for HC is only an estimate; its proof can be found e.g. in [18].

One could also consider anisotropic versions of these spaces (see e.g. [1, 3, 21]) as in Table 2, where $\boldsymbol{\alpha} = (\alpha_1,\ldots,\alpha_N) \in \mathbb{R}^N_+$ is a vector of positive weights and $\alpha_{min} = \min_n \alpha_n$. We can interpret these weights as a measure of the importance of each random variable $y_n$ on the solution: the smaller the weight, the higher degree we allow in the corresponding variable.

The family of orthonormal monodimensional polynomials will of course depend on the measure of each $\Gamma_n$ (*Generalized Polynomial Chaos*). In the case of uniform random variables, one can use the well-known orthonormal Legendre polynomials; the $p$-th Legendre polynomial can be computed recursively (see e.g. [13]), or explicitly with the Rodrigues' formula:

$$L_{p_n}(t) = \frac{(-1)^n\sqrt{2p_n+1}}{2^{p_n} p_n!} \frac{d^{p_n}}{dt^{p_n}}\left((1-t^2)^{p_n}\right). \quad (7)$$

| Tensor product (TP) | $\Lambda(w) = \{\mathbf{p} \in \mathbb{N}^N : \max_{n=1\ldots,N} \alpha_n p_n \leq \alpha_{min} w\}$ |
|---|---|
| Total degree (TD) | $\Lambda(w) = \{\mathbf{p} \in \mathbb{N}^N : \sum_{n=1}^{N} \alpha_n p_n \leq \alpha_{min} w\}$ |
| Hyperbolic cross (HC) | $\Lambda(w) = \{\mathbf{p} \in \mathbb{N}^N : \prod_{n=1}^{N} (p_n + 1)^{\frac{\alpha_n}{\alpha_{min}}} \leq w + 1\}$ |

Table 2: Corresponding anisotropic version of the polynomial spaces on Table 1.

We recall Hermite polynomials for Gaussian measures and Laguerre polynomials for Exponential measures; see [30] for the general Askey scheme. Necessary conditions for the convergence of the Generalized Polynomial Chaos expansion can be found e.g. in [11].

Now let $\phi(\mathbf{x})$ be a basis function for the physical space $V$. Inserting $v_w = \phi(\mathbf{x})\mathcal{L}_{\mathbf{q}}(\mathbf{y})$ with $\mathbf{q} \in \Lambda(w)$ as test functions in the weak formulation (5) will result in a set of equations in weak form for the coefficients $u_{\mathbf{p}}(\mathbf{x})$ that will be generally coupled due to the term $a(\mathbf{x}, \mathbf{y})\mathcal{L}_{\mathbf{p}}(\mathbf{y})\mathcal{L}_{\mathbf{q}}(\mathbf{y})$ in the equation (5). See for instance the works [3, 22, 23] for further details on space discretization and on the numerical solution of such system of equations.

## 3.1 Quasi-optimal choice of polynomial spaces

A question that naturally arises in the context of Galerkin approximation concerns the best choice of the polynomial space to be used, to get maximum accuracy for a given dimension $M$ of the space (*best $M$-terms* approximation). In other words, we look for an index set $\mathcal{S}_M \subset \mathbb{N}^N$ with cardinality $M$ that minimizes the projection error

$$\|u - \sum_{\mathbf{p} \in \mathcal{S}_M} u_{\mathbf{p}}\mathcal{L}_{\mathbf{p}}\|^2_{V \otimes L^2_\rho(\Gamma)} = \sum_{\mathbf{p} \notin \mathcal{S}_M} \|u_{\mathbf{p}}\|^2_V, \tag{8}$$

where the equivalence is a consequence of Parseval's equality and the completeness of $\{\mathcal{L}_{\mathbf{p}}\}_{\mathbf{p} \in \Lambda(w)}$ in $L^2_\rho(\Gamma)$.

### 3.1.1 Abstract construction

The obvious solution to this problem is to take the set $\mathcal{S}_M$ that contains the $M$ coefficients $u_{\mathbf{p}}$ with largest norm. This solution of course is not constructive; what we need are sharp estimates of the decay of the coefficients $\|u_{\mathbf{p}}\|_V$, based only on computable quantities, to be used in the approximation of the set $\mathcal{S}_M$. Actually, assuming that an estimate of the type

$$\|u_{\mathbf{p}}\|_V \leq G(\mathbf{p}) \tag{9}$$

is available, one can define an index set $\Lambda_\epsilon$ by selecting all multi-indices $\mathbf{p}$ for which the *estimated decay* of the corresponding Legendre coefficient is above a fixed threshold $\epsilon \in \mathbb{R}_+$,

$$\Lambda_\epsilon = \{\mathbf{p} \in \mathbb{N}^N : G(\mathbf{p}) \geq \epsilon\},$$

or equivalently

$$\Lambda(w) = \left\{ \mathbf{p} \in \mathbb{N}^N : -\log G(\mathbf{p}) \leq w, \, w = \lceil -\log \epsilon \rceil \right\}. \tag{10}$$

If the sequence $\Lambda(w)$ covers $\mathbb{N}^N$ as $w$ goes to infinity, the corresponding $u_w$ will converge to $u$ and, if the bound $G(\mathbf{p})$ in (9) is sharp, $\Lambda(w)$ will be a "quasi optimal" approximation of the best $M$-terms approximation, where now $M$ denotes the cardinality of $\Lambda(w)$.

### 3.1.2   A preliminary example

Assume for a moment that $u$ factorizes, i.e. it can be written as a product of 1D analytic functions in the stochastic variables, $u(\mathbf{x}, \mathbf{y}) = f(\mathbf{x}) \prod_{n=1}^N v_n(y_n)$. If we denote with $v_{n,p_n}$ the Legendre coefficients of the factor $v_n$, i.e.

$$v_{n,p_n} = \int_{\Gamma_n} v_n(y_n) L_{p_n}(y_n) \rho_n(y_n) dy_n,$$

the Legendre coefficients of $u$ are given simply by

$$u_{\mathbf{p}}(\mathbf{x}) = f(\mathbf{x}) \prod_{n=1}^N v_{n,p_n}. \tag{11}$$

Now, from classical approximation theory ([9, 26]) it is well known that, if $v_n$ is analytic in $\Gamma_n$, the coefficient $v_{n,p_n}$ is exponentially decaying in $p_n$ with a certain rate $g_n$, $|v_{n,p_n}| \leq c(g_n) e^{-g_n p_n}$; as a consequence we easily obtain a sharp bound on the Legendre coefficients of $u$,

$$\|u_{\mathbf{p}}\|_V \leq \|f\|_V \, \mathrm{C} \, e^{-\sum_n g_n p_n}, \quad \mathrm{C} = \prod_{n=1}^N c(g_n). \tag{12}$$

Substituting this bound in (10), we get that a quasi optimal choice of polynomial sets for a separable function of the form (11) is the anisotropic TD sets sequence defined in Table 2 with weights $\alpha_n = g_n$.

### 3.1.3   General case

In the general case things deriving sharp estimates on the decay of $\|u_{\mathbf{p}}\|_V$ is a more delicate goal. Seminal works in this direction are [5, 8, 7], where estimates of the decay of the Legendre coefficients are provided. We consider here a slight generalization of the result in [8, Corollary 6.1] and show numerically that the polynomial sets built on these modified estimates behave closely to the true best $M$-terms  approximation.

Under Assumptions 2.1 - 2.3 it is possible to prove that the following estimate holds for the Legendre coefficients. Again, a similar result is given in [8] for the special case $a = a_0 + \sum_{n=1}^N b_n(x) y_n$.

**Proposition 3.1.** *Consider equation (1), suppose that the diffusion coefficient $a$ satisfies Assumptions 2.1 - 2.3, let $\mathbf{r}$ be as in Assumption 2.3 and $C_0$ be as in Theorem 2.1. Then the $V$-norm of the Legendre coefficients $u_{\mathbf{p}}$ can be bounded as*

$$\|u_{\mathbf{p}}\|_V \leq C_0 e^{-\sum_n g_n p_n} \frac{|\mathbf{p}|!}{\mathbf{p}!}, \quad g_n = -\log(\, r_n/(\sqrt{3}\log 2)\,). \tag{13}$$

**Proof.** We follow closely the proof in [8, Corollary 6.1]. We start from the definition of the Legendre coefficients (6) and the Rodrigues' formula for the Legendre polynomials (7). Integrating by parts and commuting the order of the $V$-norm and the integral over $\Gamma$, thanks to the properties of the Bochner integral, we have

$$\|u_{\mathbf{p}}\|_{V(D)} = \left\| \int_\Gamma u(\cdot,\mathbf{y}) \mathcal{L}_{\mathbf{p}}(\mathbf{y})\rho(\mathbf{y})d\mathbf{y} \right\|_V$$
$$\leq \frac{\prod_{n=1}^N \sqrt{2p_n+1}}{2^{|\mathbf{p}|}\mathbf{p}!} \int_\Gamma \left\| \partial_{\mathbf{p}}^{\mathbf{y}} u(\cdot,\mathbf{y}) \right\|_V \prod_{n=1}^N (1-y_n^2)^{p_n}\rho(\mathbf{y})d\mathbf{y}.$$

It has been shown in [8] that

$$I(\mathbf{p}) = \prod_{n=1}^N \sqrt{2p_n+1} \int_\Gamma \prod_{n=1}^N (1-y_n^2)^{p_n}\rho(\mathbf{y})d\mathbf{y} \leq \left(\frac{2}{\sqrt{3}}\right)^{|\mathbf{p}|}. \tag{14}$$

Thus we have

$$\|u_{\mathbf{p}}\|_{V(D)} \leq \max_{\mathbf{y}\in\Gamma} \left\| \partial_{\mathbf{p}}^{\mathbf{y}} u(\cdot,\mathbf{y}) \right\|_V I(\mathbf{p}) \frac{1}{2^{|\mathbf{p}|}\mathbf{p}!},$$

and the proof is completed using Theorem 2.1 to estimate $\max_{\mathbf{y}\in\Gamma}\left\| \partial_{\mathbf{p}}^{\mathbf{y}} u(\cdot,\mathbf{y}) \right\|_V$:

$$\|u_{\mathbf{p}}\|_{V(D)} \leq C_0|\mathbf{p}|! \left(\frac{1}{\log 2}\mathbf{r}\right)^{\mathbf{P}} \left(\frac{2}{\sqrt{3}}\right)^{|\mathbf{p}|} \frac{1}{2^{|\mathbf{p}|}\mathbf{p}!}$$
$$= C_0 \left(\frac{1}{\sqrt{3}\log 2}\mathbf{r}\right)^{\mathbf{P}} \frac{|\mathbf{p}|!}{\mathbf{p}!} = C_0 e^{\sum_n p_n \log\left(\frac{r_n}{\sqrt{3}\log 2}\right)} \frac{|\mathbf{p}|!}{\mathbf{p}!}. \tag{15}$$

$\square$

**Example 3.1.** *To motivate bound (13), assume that in the model problem (1) the forcing term is deterministic, $f = f(\mathbf{x})$, and the diffusion coefficient is constant in space, $a = a(\mathbf{y}) = 1 + \sum_{i=1}^N b_i y_i$, with $b_i > 0$. As explained in Remark 2.1, for such a diffusion coefficient Assumption 2.3 holds with $a_{min} = 1 - \sum_i b_i$ and $r_i = b_i/a_{min}$. Moreover, let us denote with $g \in V$ the solution of the auxiliary problem*

$$\begin{cases} \Delta g(\mathbf{x}) = f(\mathbf{x}) & \mathbf{x} \in D, \\ g(\mathbf{x}) = 0 & \mathbf{x} \in \partial D. \end{cases}$$

*Under these hypotheses we can derive an analytic expression for $u$ and its derivatives with respect to $\mathbf{y}$,*

$$u(\mathbf{x},\mathbf{y}) = g(\mathbf{x})\frac{1}{1+\sum_{i=1}^N b_i y_i}, \quad \partial_{\mathbf{p}}u(\mathbf{x},\mathbf{y}) = g(\mathbf{x})\frac{|\mathbf{p}|!\, \mathbf{b}^{\mathbf{P}}}{\left(1+\sum_{i=1}^N b_i y_i\right)^{|\mathbf{p}|+1}}. \tag{16}$$

*We can exploit this fact to compute explicitly a bound for the V-norm of the Legendre coefficients $u_{\mathbf{p}}$ of $u$. Actually, using again Rodrigues' formula (7) in the definition of $u_{\mathbf{p}}$, and integrating by parts, we obtain*

$$u_{\mathbf{p}}(\mathbf{x}) = \int_{\Gamma} u(\mathbf{x}, \mathbf{y}) \mathcal{L}_{\mathbf{p}}(\mathbf{y}) \rho(\mathbf{y}) d\mathbf{y}$$

$$\leq g(\mathbf{x}) \frac{\mathbf{b}^{\mathbf{p}}}{2^{|\mathbf{p}|}} \frac{|\mathbf{p}|!}{\mathbf{p}!} \frac{1}{(a_{min})^{|\mathbf{p}|+1}} \prod_{i=1}^{N} (-1)^n \sqrt{2p_n+1} \int_{[-1,1]} (1-y_n^2)^{p_n} \frac{1}{2} dy_n.$$

*Finally we exploit bound (14), pass to the V-norm and use the fact that $\|g\|_V = \|f\|_{V'}$ to obtain*

$$\|u_{\mathbf{p}}\|_V \leq \frac{\|f\|_{V'}}{a_{min}} \frac{\mathbf{b}^{\mathbf{p}}}{(a_{min})^{|\mathbf{p}|}} \left(\frac{1}{\sqrt{3}}\right)^{|\mathbf{p}|} \frac{|\mathbf{p}|!}{\mathbf{p}!}.$$

*This can be recast using the definition of the rate $r_i = b_i/a_{min}$ and of the constant $C_0$ in Theorem 2.1 to*

$$\|u_{\mathbf{p}}\|_V \leq C_0 e^{-\sum_n g_n p_n} \frac{|\mathbf{p}|!}{\mathbf{p}!}, \quad g_n = -\log(r_n/\sqrt{3}),$$

*that is precisely the bound derived in Proposition 3.1, with a slight modification on the rate $g_i$. Numerical results in the next Section will again cover this particular example, showing that the bound proposed yields good approximating properties.*

**Remark 3.1.** *Since $u(\mathbf{x}, \cdot)$ is analytic in $\Gamma$ (see Corollary 2.1), it can be shown that $u$ always admits a converging Legendre expansion. In spite of this, the estimate (13) in the previous Proposition does not ensure that the norm of the coefficients $\|u_{\mathbf{p}}\|_V$ of the expansion is decaying for any value of the coefficients $r_n$ when $|\mathbf{p}| \to \infty$, nor that the Legendre series is convergent; sufficient conditions for this to be true are given in the next Preposition.*

*This is a clear indication that estimate (13) is not sharp. Other estimates derived using complex analysis arguments are available and always predict a decay of $\|u_{\mathbf{p}}\|_V$ for $|\mathbf{p}| \to \infty$ (see e.g. [7]). On the other hand, we have observed that the behaviour of the Legendre coefficients is well described by a bound of the type of (13), if the rates $g_n$ are estimated numerically rather than analytically. See Section 3.2 for numerical evidence on the quality of the bound proposed.*

For a given set $\Lambda$, let $\overline{w}$ be the index of the largest TD set included in $\Lambda$:

$$\overline{w} = \max\{\widetilde{w} \in \mathbb{N} : TD(\widetilde{w}) \subseteq \Lambda(w)\}.$$

The following Proposition holds:

**Proposition 3.2.** *Given an increasing sequence of index sets $\Lambda(w)$ with $\overline{w} \to \infty$, the estimate* (13) *in Proposition 3.1 implies that a sufficient condition for the Legendre series $u_w$ defined in* (6) *to converge uniformly to u is*

$$\sum_{i=1}^{N} r_n < \log 2. \tag{17}$$

**Proof.** It is enough to prove that if condition (17) holds then the sequence $u_w = \sum_{\mathbf{p} \in \Lambda(w)} u_{\mathbf{p}} \mathcal{L}_{\mathbf{p}}$ is Cauchy with respect to the norm $\|\cdot\|_{L^\infty(\Gamma;V)}$, for the sequence $\Lambda(w)$ considered. As a consequence $u_w$ converges uniformly to its limit $u$.

To prove that $u_w$ is Cauchy, let $w_1, w_2 \in \mathbb{N}$ such that $w_1 < w_2$. It holds

$$\left\| \sum_{\mathbf{p} \in \Lambda(w_2)} u_{\mathbf{p}}(\mathbf{x}) \mathcal{L}_{\mathbf{p}}(\mathbf{y}) - \sum_{\mathbf{p} \in \Lambda(w_1)} u_{\mathbf{p}}(\mathbf{x}) \mathcal{L}_{\mathbf{p}}(\mathbf{y}) \right\|_{L^\infty(\Gamma;V)} =$$

$$\left\| \sum_{\mathbf{p} \in \Lambda(w_2) \setminus \Lambda(w_1)} u_{\mathbf{p}}(\mathbf{x}) \mathcal{L}_{\mathbf{p}}(\mathbf{y}) \right\|_{L^\infty(\Gamma;V)} \leq \sum_{\mathbf{p} \in \Lambda(w_2) \setminus \Lambda(w_1)} \|u_{\mathbf{p}}(\mathbf{x}) \mathcal{L}_{\mathbf{p}}(\mathbf{y})\|_{L^\infty(\Gamma;V)} \leq$$

$$\sum_{\mathbf{p} \notin TD(\overline{w}_1)} \|u_{\mathbf{p}}(\mathbf{x}) \mathcal{L}_{\mathbf{p}}(\mathbf{y})\|_{L^\infty(\Gamma;V)} = \sum_{\mathbf{p} \notin TD(\overline{w}_1)} \|u_{\mathbf{p}}(\mathbf{x})\|_V \|\mathcal{L}_{\mathbf{p}}(\mathbf{y})\|_{L^\infty(\Gamma)}.$$

Now use estimate (13) in Proposition 3.1 to bound $\|u_{\mathbf{p}}(\mathbf{x})\|_V$. Furthermore note that the $L^\infty(\Gamma)$-norm of the orthonormal Legendre polynomials can be bounded as

$$\|\mathcal{L}_{\mathbf{p}}(\mathbf{y})\|_{L^\infty(\Gamma)} = \prod_{n=1}^{N} \sqrt{2p_n + 1} \leq \left(\sqrt{3}\right)^{|\mathbf{p}|} \quad \forall \mathbf{p} \in \mathbb{N}^N,$$

so that

$$\sum_{\mathbf{p} \notin TD(\overline{w}_1)} \|u_{\mathbf{p}}(\mathbf{x})\|_V \|\mathcal{L}_{\mathbf{p}}(\mathbf{y})\|_{L^\infty(\Gamma)} \leq C_0 \sum_{\mathbf{p} \notin TD(\overline{w}_1)} \left(\frac{1}{\log 2} \mathbf{r}\right)^{\mathbf{p}} \frac{|\mathbf{p}|!}{\mathbf{p}!}$$

$$= C_0 \sum_{|\mathbf{p}| \geq \overline{w}_1} \left(\frac{1}{\log 2} \mathbf{r}\right)^{\mathbf{p}} \frac{|\mathbf{p}|!}{\mathbf{p}!} = C_0 \sum_{s=\overline{w}_1}^{\infty} \left(\sum_{n=1}^{N} \frac{1}{\log 2} r_n\right)^s,$$

that tends to 0 if condition (17) holds, where we have exploited the generalized Newton binomial formula as in Corollary 2.1. $\square$

**Remark 3.2.** *Condition* (17) *in Proposition 3.2 can be weakened by improving bound* (14). *We recall the definition of $I(\mathbf{p}) = \prod_{n=1}^{N} \sqrt{2p_n + 1} \int_{\Gamma_n} (1 - y_n^2)^{p_n} \rho(y_n) dy_n$. Integrating p times by parts, one obtains*

$$\int_{-1}^{1} (1 - t^2)^p \frac{1}{2} dt = \frac{2^{2p}(p!)^2}{(2p + 1)!}.$$

*Using Stirling's approximation formula*

$$p! = \sqrt{2\pi p} \left(\frac{p}{e}\right)^p e^{\lambda_p}, \quad \frac{1}{12p + 1} \leq \lambda_p \leq \frac{1}{12p},$$

*one can then bound*

$$I(p) \leq \sqrt{\frac{\pi}{2}} \quad \Rightarrow \quad I(\mathbf{p}) \leq \left(\frac{\pi}{2}\right)^{N/2}.$$

*Note that this bound is sharp, even for small values of* $|\mathbf{p}|$. *Using this result rather then* (14) *in* (15) *we obtain*

$$\|u_{\mathbf{p}}\|_{V(D)} \leq C_0 \left(\frac{\pi}{2}\right)^{N/2} \left(\frac{1}{2\log 2}\mathbf{r}\right)^{\mathbf{p}} \frac{|\mathbf{p}|!}{\mathbf{p}!} \tag{18}$$

*and, as a consequence, condition* (17) *becomes*

$$\sum_{i=1}^{N} r_n < \frac{2\log 2}{\sqrt{3}}. \tag{19}$$

*Note however that this is only a little improvement, being* $\log 2 = 0.69$ *and* $2\log 2/\sqrt{3} = 0.80$; *moreover, since* $\pi/2 > 1$, *bound* (18) *does not imply that the Legendre coefficients of* $u$ *decay regardless of the number of random variables, which was the case for the initial estimate* (13); *therefore, condition* (19) *holds for fixed* $N$, *while* (17) *is independent of* $N$.

Following again the abstract procedure in Section 3.1.1, we substitute the estimate (13) in the general quasi optimal set expression (10). This results in the following expression for the quasi optimal polynomial sets for a general non factorizing $u$,

$$\Lambda(w) = \left\{\mathbf{p} \in \mathbb{N}^N : \sum_{n=1}^{N} g_n p_n - \log \frac{|\mathbf{p}|!}{\mathbf{p}!} \leq w\right\}. \tag{20}$$

We refer to these sets as TD-FC sets ("TD with factorial correction" sets). We can indeed interpret the factor $\log \frac{|\mathbf{p}|!}{\mathbf{p}!}$ appearing in (20) as a correction factor to the TD space to take into account the intrinsic coupling between directions in the stochastic space; observe that this correction is always isotropic.

As pointed out in Remark 3.1, the quantities $g_n$ appearing in (20) are better estimated numerically by a sequence of monovariate analyses: one could indeed increase the polynomial degree in one random variable at a time while keeping degree zero in all the others variables and estimate numerically the exponential rate of convergence. Observe that in such monovariate analyses the factorial term does not appear so the expected convergence rate is precisely $\sim e^{-g_n p_p}$. In the numerical results presented in the next section we have used this strategy, which seems to work particularly well.

**Remark 3.3.** *Observe that* $\Lambda(w)$ *actually depends on the number of input variables* $N$. *One can extend the definition of* $\Lambda(w)$ *also to the case where* $\mathbf{p}$ *is a*

*sequence of natural numbers ("infinite dimensional probability space") with only a finite number of non zero terms, provided the sequence $g_n \to 0$ as $n \to \infty$. This is an alternative way to work with random fields, without truncating them a priori to a certain level (see e.g. [8, 7, 21]).*

## 3.2   Numerical Tests

In this section we show the performance of the TD-FC sets (20) compared to the isotropic and anisotropic versions of TD sets defined in Tables 1 and 2, as well as the best M-term approximation. We consider the following elliptic problem in one physical dimension

$$\begin{cases} -(a(x,\mathbf{y})u(x,\mathbf{y})')' = 1 & x \in D = (0,1), \mathbf{y} \in \Gamma \\ u(0,\mathbf{y}) = u(1,\mathbf{y}) = 0, & \mathbf{y} \in \Gamma \end{cases} \tag{21}$$

with different choices of diffusion coefficient $a(x,\mathbf{y})$, for which Assumptions 2.1 - 2.3 hold. We focus on a linear functional $\psi : V \to \mathbb{R}$ of the solution, so that $\psi(u)$ is a scalar random variable, function of $\mathbf{y}$ only. In our examples, $\psi$ is defined as $\psi(v) = v(\frac{1}{2})$.

To obtain the best $M$-terms approximation we compute explicitly all the Legendre coefficients of $\psi(u)$ in a sufficiently large index set $\mathbb{U}$ evaluating the integrals $\psi_{\mathbf{p}} = \int_\Gamma \psi(u)L_{\mathbf{p}}(\mathbf{y})\rho(\mathbf{y})d\mathbf{y}$ with a high-level sparse grid as reference values. We order then the coefficients in decreasing order, according to their modulus, and take the first $M$ terms of the reordered sequence as the best $M$-terms approximation.

The rates $\mathbf{g}$ used to build the TD-FC space, as well as the anisotropic TD space, are computed numerically, with a sequence of 1D analyses. For each random variable $1 \le n \le N$, we consider the subset $\mathbb{U}_n = \{\mathbf{p} \in \mathbb{U} : \mathbf{p}_i = 0 \text{ if } i \ne n, \mathbf{p}_n = 0, 1, 2, \ldots\}$; according to (13), the decay of the Legendre coefficients for this particular choice of multi-indices is $|\psi_{\mathbf{p}}| \sim e^{-g_n p_n}$, and we can then estimate the rate $g_n$ via a linear interpolation on the quantities $\log |\psi_{\mathbf{p}}|, \mathbf{p} \in \mathbb{U}_n$.

**Test 1: space independent diffusion coefficient**

The first case we consider has two random variables $(y_1, y_2)$ and a diffusion coefficient $a(x, \mathbf{y}) = 1 + 0.1y_1 + 0.5y_2$; results are shown in Figures 1-2.

Figure 1(a) shows the Legendre coefficients ordered in lexicographic order, giving this peculiar sawtooth shape. The first tooth corresponds to multi-indices of the form $[0, k]$, the second one to $[1, k]$ and so on. We have also added to the plot the estimate (13) in Proposition 3.1 of the magnitude of the Legendre coefficients, which leads to the TD-FC sets (20), as well as the estimate (12) which leads to the anisotropic TD spaces as in Table 2, with $\boldsymbol{\alpha}_n = g_n$. The plot suggests that estimate (13) is quite sharp, whereas the estimate corresponding to the TD space underestimates considerably the Legendre coefficients. This result highlights the importance of the factorial term in (13). We expect, therefore, that

14

(a) Legendre coefficients in lexicographic order and their corresponding estimates based on either TD-FC or TD approximations.

(b) Convergence of different polynomial approximations, measured as $\|\psi(u) - \psi(u_w)\|_{L^2_\rho(\Gamma)}$ versus dimension of polynomial space.
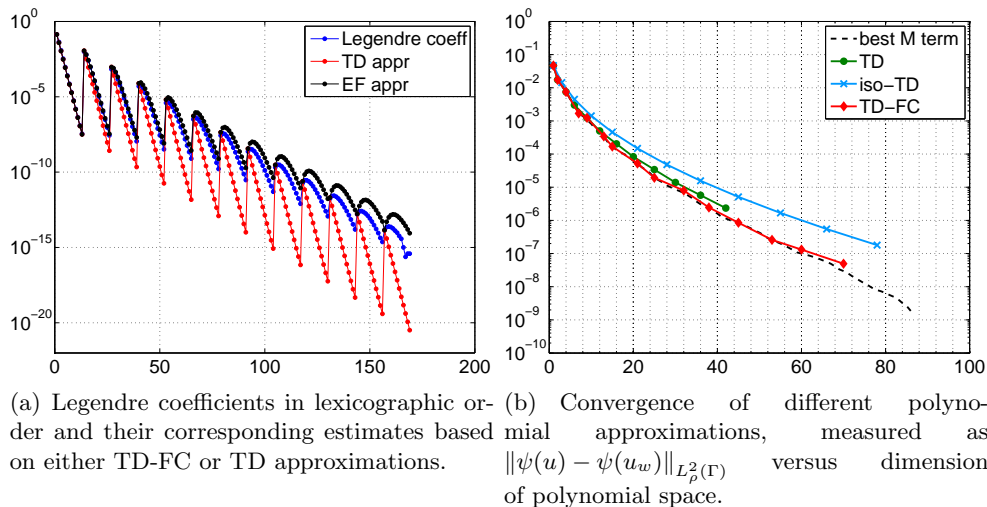
Figure 1: Results for $a(x, \mathbf{y}) = 1 + 0.1y_1 + 0.5y_2$. Here we have $g \simeq (2.49, 1.27)$, $\mathbb{U} = \text{TP}(12)$, Legendre coefficients computed with a standard Smolyak sparse grid of level 9, with Gauss-Legendre abscissae.

the TD-FC approximation performs better than the aniso-TD one. Moreover, we point out the non intuitive fact that the Legendre coefficients $\psi_{\mathbf{p}}$ *are not strictly decreasing in absolute value* when listed in the lexicographic order. As an example, $|\psi_{[5\,0]}| < |\psi_{[5\,1]}|$, and the same holds for all teeth but the first few.

Figure 1(b) shows convergence plots for the error in $L^2_\rho$-norm for the various polynomial spaces used versus the dimension of the polynomial space. As the TD-FC sequence is the only sequence that captures correctly the non decreasing behaviour of the Legendre coefficients in lexicographic order, the convergence of the TD-FC sequence in Figure 1(b) is the closest to the best $M$-terms approximation, even though the anisotropic TD space give good results as well. We also point out the poor performance of the standard isotropic TD space compared to both the anisotropic TD and the TD-FC spaces: this confirms the importance of using anisotropic spaces to reduce computational costs.

It is also useful to visualize the isolines of the Legendre coefficients of the expansion of $\psi(u)$ and to compare them with the isolines corresponding to estimates (13) for TD-FC sets, and (12) for iso and aniso TD sets, see Figure 2. The closer the matching of the sequence of sets with the true decay of the Legendre coefficients, the faster the $L^2$ convergence of the approximation for $\psi$ will be. The key property of the decay of the Legendre coefficients is the rounded shape of the isolines (see Figure 2(a)), properly caught only with the factorial term $\frac{|\mathbf{i}|!}{\mathbf{i}!}$ in the TD-FC set formula (Figure 2(b)). Also from these plots one can see the fact that the Legendre coefficients are not strictly decreasing in lexicographic order: actually close to the borders the isolines tend to bend "backward", so that for example the index [7, 1] belongs to a lower isoline than [7, 0]. However,

15

(a) Legendre coefficients isolines

(b) opt sets

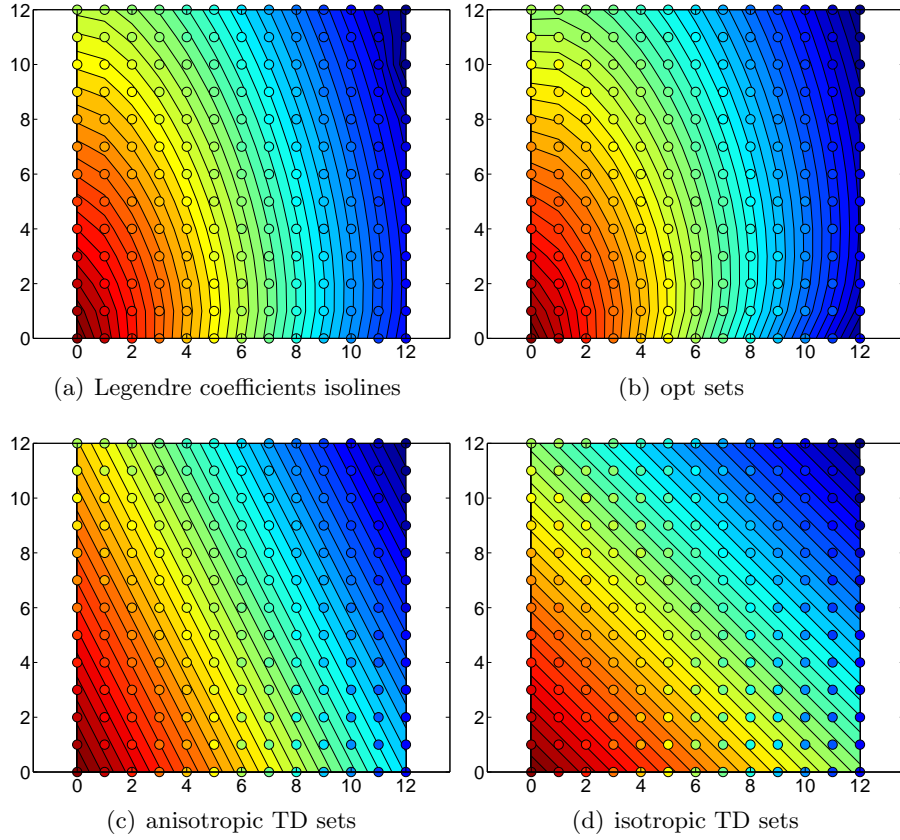(c) anisotropic TD sets

(d) isotropic TD sets

Figure 2: Isolines of estimated Legendre coefficients: a) true values, computed with high level sparse grids; b) estimate (13) leading to TD-FC sets; c) estimate (12) leading to aniso-TD sets with $\alpha_n = g_n$; d) estimate (12) with $\alpha_n = 1 \, \forall n = 1, \ldots, N$, leading to standard TD sets as in Table 1. In all plots, each dot represents a multi-index in $\mathbb{N}^2$, and it is coloured according to the size of the corresponding exact coefficient in the Legendre expansion for $\psi$; on the background the isolines.

as appears from results in Figure 1, approximating the isolines with "mean" straight lines as it is done in the anisotropic TD (Figure 2(c)) gives quite good results as well. On the other hand, using the wrong slopes for TD sets, like in isotropic TD sets (2(d)), will result in general in poor approximation properties.
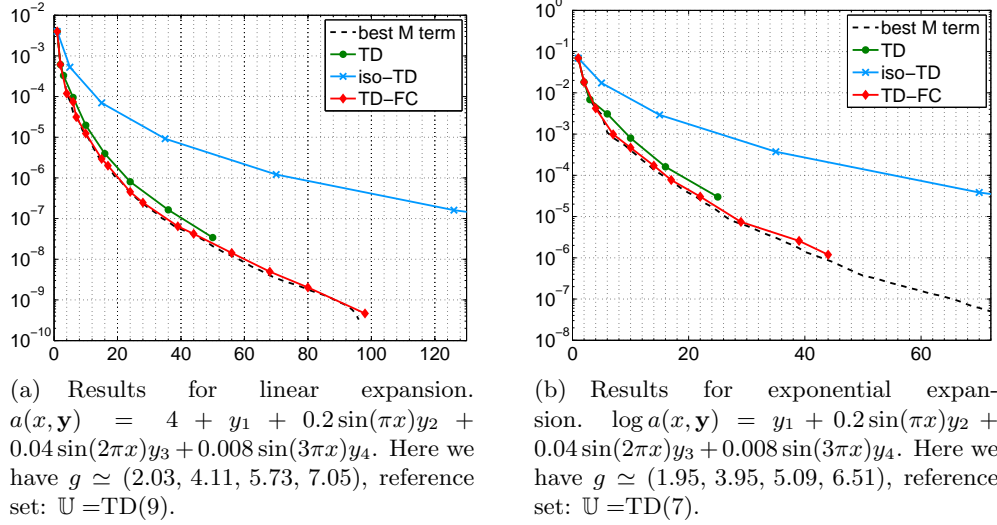
16

(a) Results for linear expansion. $a(x, \mathbf{y}) = 4 + y_1 + 0.2\sin(\pi x)y_2 + 0.04\sin(2\pi x)y_3 + 0.008\sin(3\pi x)y_4$. Here we have $g \simeq (2.03, 4.11, 5.73, 7.05)$, reference set: $\mathbb{U} = \text{TD}(9)$.

(b) Results for exponential expansion. $\log a(x, \mathbf{y}) = y_1 + 0.2\sin(\pi x)y_2 + 0.04\sin(2\pi x)y_3 + 0.008\sin(3\pi x)y_4$. Here we have $g \simeq (1.95, 3.95, 5.09, 6.51)$, reference set: $\mathbb{U} = \text{TD}(7)$.

Figure 3: Convergence of polynomial approximations for elliptic equation with the coefficient $a$ depending also on $x$. Convergence measured as $\|\psi(u) - \psi(u_w)\|_{L^2_\rho(\Gamma)}$ versus the dimension of polynomial space.

### Test 2: space dependent diffusion coefficient

We now consider the following two expansions:

- $a(x, \mathbf{y}) = 4 + y_1 + 0.2\sin(\pi x)y_2 + 0.04\sin(2\pi x)y_3 + 0.008\sin(3\pi x)y_4$,

- $\log a(x, \mathbf{y}) = y_1 + 0.2\sin(\pi x)y_2 + 0.04\sin(2\pi x)y_3 + 0.008\sin(3\pi x)y_4$.

and look at the functional $\psi(v) = v(0.7)$ (the functional $\psi(v) = v(1/2)$ is not suited for analysis in this case as, by symmetry, many of the Legendre coefficients are zero). Figure 3 shows the results, and again we see that the TD-FC approximation is the best performing, with anisotropic TD closely following and isotropic TD far worse.

### Test 3: factorizable diffusion coefficient

Let us now give an example on the case of a factorizable $u$, as in Section 3.1.2. We recall that Section 3.1.2 states that if we can express the solution $u(\mathbf{x}, \mathbf{y})$ as a product $u(\mathbf{x}, \mathbf{y}) = f(\mathbf{x}) \prod_n v_n(y_n)$ then the Legendre coefficients can be computed as a product of 1D Legendre coefficients, and thus the optimal estimate is (12), leading to aniso-TD sets, rather than estimate (13) leading to what we have called TD-FC sets. To support our thesis, we now consider $a(\mathbf{y}) = (1 + 0.6y_1)(1 + 0.6y_2)$, so that the solution of (21) is $u(x, \mathbf{y}) = \frac{x(1-x)}{2a(\mathbf{y})}$.

The convergence plots for of $\psi(u)$ are shown in Figure 4 and confirm that in this case TD is the optimal choice, and is very close to the best $M$-terms
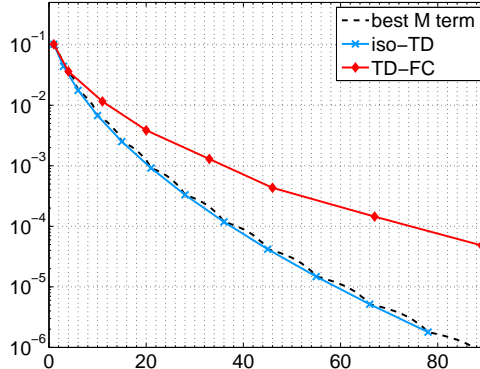
17

Figure 4: convergence of polynomial approximation of the elliptic equation (21) with the coefficient $a$ that factorizes with respect to $\mathbf{y}$, $a = (1+0.6y_1)(1+0.6y_2)$. $g \simeq (1.08,\ 1.08)$, $\mathbb{U} = \mathrm{TP}(12)$. Convergence measured as $\|\psi(u) - \psi(u_w)\|_{L^2_\rho(\Gamma)}$ versus number of Legendre coefficients (dimension of polynomial space).

approximation. Note that in this example the isotropic and anisotropic versions of TD coincide, since the two factors of $v$ are the same.

## 3.3 Alternative estimates for diffusion coefficients in exponential form

Let us consider again the model problem (21), with diffusion coefficient in exponential form $\log a(\mathbf{x}, \mathbf{y}) = \sum_{n=1}^{N} c_n y_n$. The solution is $u(x, \mathbf{y}) = \frac{x(1-x)}{2} \frac{1}{\prod_{n=1}^{N} e^{c_n y_n}}$, therefore the solution is in separable form with $v_n = e^{c_n y_n}$; as a consequence, following the arguments in Section 3.1.2 on factorizable functions, we have $\|u_{\mathbf{p}}\|_V = \|f(\mathbf{x})\|_V \prod_{n=1}^{N} |v_{n,p_n}|$, where $v_{n,p_n}$ indicates the $p_n$-th Legendre coefficients of $v_n$. In this case, however, we expect the decay of $v_{n,p_n}$ to be faster than exponential, since $v_n(y_n)$ is an entire function. Actually, the following lemma holds:

**Lemma 3.1.** *Given problem* (21) *with diffusion coefficient* $\log a(\mathbf{x}, \mathbf{y}) = \sum_{n=1}^{N} c_n y_n$, *the $V$-norm of the Legendre coefficients of $u$ can be bounded as*

$$\|u_{\mathbf{p}}\|_V \leq C_e \frac{e^{-\sum_{n=1}^{N} g_n p_n}}{\mathbf{p}!}, \tag{22}$$

*with* $g_n = -\log \frac{|c_n|}{\sqrt{3}}$ *and* $C_e = \|f\|_V\, e^{\sum_{n=1}^{N} |c_n|}$.

**Proof.** Since $\|u_{\mathbf{p}}\|_V = \|f(\mathbf{x})\|_V \prod_{n=1}^{N} |v_{n,p_n}|$ we only need to estimate $|v_{n,p_n}|$. Re-

calling the definition of $I(p)$ given in the proof of Proposition 3.1, one gets

$$|v_{n,p_n}| = \left| \int_{-1}^{1} L_{p_n}(y_n) v_n(y_n) \frac{dy_n}{2} \right| = \frac{\sqrt{2p_n+1}}{2^{p_n} p_n!} \left| \int_{-1}^{1} e^{-c_n y_n} \left( \frac{d}{dy} \right)^{p_n} (1 - y_n^2)^{p_n} \frac{dy_n}{2} \right| =$$

$$\sqrt{2p_n+1} \frac{|c_n|^{p_n}}{2^{p_n} p_n!} \int_{-1}^{1} e^{-c_n y_n} (1 - y_n^2)^{p_n} \frac{dy_n}{2} \le \frac{|c_n|^{p_n} e^{|c_n|}}{2^{p_n} p_n!} I(p) \le \frac{|c_n|^{p_n} e^{|c_n|}}{\sqrt{3}^{p_n} p_n!}.$$

The thesis follows setting $g_n = -\log \frac{|c_n|}{\sqrt{3}}$. $\qquad\qquad\square$

**Remark 3.4.** *Observe that in (22) the coefficient $u_{\mathbf{p}}$ will tend to zero as $|\mathbf{p}| \to \infty$ even when $g_n > \sqrt{3}$ for all $n = 1, \dots, N$.*

As a consequence, the abstract optimal space (10) becomes in this case

$$\Lambda(w) = \left\{ \sum_{n=1}^{N} p_n g_n + \sum_{n=1}^{N} \log(p_n!) \le w \right\}. \tag{23}$$

We refer to this set as anisotropic "factorial TD", or aniso-fTD in short.

We now guess that even in the more general case where $\log a(\mathbf{x}, \mathbf{y}) = \sum_{n=1}^{N} c_n(\mathbf{x}) y_n$ an estimate of the type of (22) for the Legendre coefficients of the solution holds, for some $g_n$, $n + 1, \dots, N$. We have tested this space on two cases

- $\log(a(x, \mathbf{y}) + 0.01) = 0.2 y_1 + 2 y_2$ (constant coefficients) ;

- $\log a(x, \mathbf{y}) = y_1 + 0.2 \sin(\pi x) y_2 + 0.04 \sin(2\pi x) y_3 + 0.008 \sin(3\pi x) y_4$ (sin expansion, this one is the same as in Test 2).

Again, the rates $g_n$ appearing in formula (23) can be estimated numerically with a least square approach. We will refer to these new rates as $\tilde{g}_n$ to stress the fact that they are different from the $g_n$ we use in TD and TD-FC spaces.

The corresponding results are shown in Figure 5, and show that actually fTD is competing with TD-FC .

## 4   Stochastic Collocation

The Stochastic Collocation (SC) Finite Element method consists in collocating problem (1) in a set of points $\{\mathbf{y}_j \in \Gamma, \quad j = 1, \dots, M_w\}$, i.e. computing the corresponding solutions $u(\cdot, \mathbf{y}_j)$ and building a global polynomial approximation $u_w$, not necessarily interpolatory, upon those evaluations: $u_w(\mathbf{x}, \mathbf{y}) = \sum_{j=1}^{M_w} u(\mathbf{x}, \mathbf{y}_j) \tilde{\psi}_j(\mathbf{y})$ for suitable multivariate polynomials $\{\tilde{\psi}_j\}_{j=1}^{M_w}$.

Building the set of evaluation points $\{\mathbf{y}_j\}$ as a cartesian product of monodimensional grids becomes quickly unfeasible, since the computational cost grows exponentially fast with the number of stochastic dimensions needed. We consider instead the so-called *sparse grid* procedure, originally introduced by Smolyak in [25] for high dimensional quadrature purposes; see also [4, 6] for polynomial interpolation. In the following we briefly review and generalize this construction.
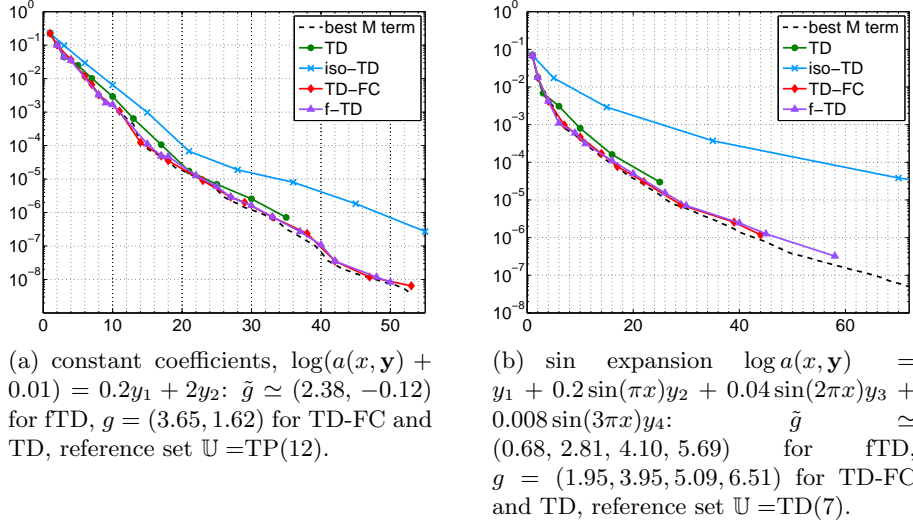
(a) constant coefficients, $\log(a(x, \mathbf{y}) + 0.01) = 0.2y_1 + 2y_2$: $\tilde{g} \simeq (2.38, -0.12)$ for fTD, $g = (3.65, 1.62)$ for TD-FC and TD, reference set $\mathbb{U} = TP(12)$.

(b) sin expansion $\log a(x, \mathbf{y}) = y_1 + 0.2 \sin(\pi x)y_2 + 0.04 \sin(2\pi x)y_3 + 0.008 \sin(3\pi x)y_4$: $\tilde{g} \simeq (0.68, 2.81, 4.10, 5.69)$ for fTD, $g = (1.95, 3.95, 5.09, 6.51)$ for TD-FC and TD, reference set $\mathbb{U} = TD(7)$.

Figure 5: convergence of polynomial spaces for elliptic equation with "shifted" exponential $a(x, \mathbf{y})$, using f-TD space.

For each direction $y_n$ we introduce a sequence of one dimensional polynomial interpolant operators of increasing order: $\mathcal{U}_n^{m(i)} : C^0(\Gamma_n) \to \mathbb{P}_{m(i)-1}(\Gamma_n)$. Here $i \geq 1$ denotes the level of approximation and $m(i)$ the number of collocation points used to build the interpolation at level $i$. As a consequence, $\mathcal{U}_n^{m(i)}[q] = q$ if $q$ is a polynomial of degree up to $m(i) - 1$. We require the function $m$ to satisfy the following assumptions: $m(0) = 0$, $m(1) = 1$ and $m(i) < m(i+1)$ for $i \geq 1$. In addition, let $\mathcal{U}_n^0[q] = 0$, $\forall q \in C^0(\Gamma_n)$.

Next we introduce the difference operators $\Delta_n^{m(i)} = \mathcal{U}_n^{m(i)} - \mathcal{U}_n^{m(i-1)}$, an integer value $w \geq 0$, multi-indices $\mathbf{i} \in \mathbb{N}_+^N$ and a sequence of index sets $\mathcal{I}(w)$ such that $\mathcal{I}(w) \subset \mathcal{I}(w+1)$ and $\mathcal{I}(0) = \{(1, 1, \ldots, 1)\}$. We define the sparse grid approximation of $u : \Gamma \to V$ at level $w$ as

$$u_w(\mathbf{y}) = \mathcal{S}_{\mathcal{I}(w)}^m[u](\mathbf{y}) = \sum_{\mathbf{i} \in \mathcal{I}(w)} \bigotimes_{n=1}^N \Delta_n^{m(i_n)}[u](\mathbf{y}). \tag{24}$$

As pointed out in [14], it is desiderable that the sum (24) has some telescopic properties. To ensure this we have to impose some additional constraints on $\mathcal{I}$. Following [14] we say that a set $\mathcal{I}$ is *admissible* if $\forall \mathbf{i} \in \mathcal{I}$

$$\mathbf{i} - \mathbf{e}_j \in \mathcal{I} \text{ for } 1 \leq j \leq N, i_j > 1. \tag{25}$$

We refer to this property as *admissibility condition*, or *ADM* in short. Given a set $\mathcal{I}$ we will denote by $\mathcal{I}^{ADM}$ the smallest set such that $\mathcal{I} \subset \mathcal{I}^{ADM}$ and $\mathcal{I}^{ADM}$ is admissible.

20

It is now possible to rewrite (24) in terms of linear combinations of tensor grids interpolations:

$$u_w(\mathbf{y}) = \sum_{\mathbf{i} \in \mathcal{I}(w)^{ADM}} c_i \bigotimes_{n=1}^N \mathcal{U}_n^{m(i_n)}[u](\mathbf{y}), \quad c_i = \sum_{\substack{\mathbf{j}=\{0,1\}^N: \\ \mathbf{i}+\mathbf{j} \in \mathcal{I}(w)^{ADM}}} (-1)^{|\mathbf{j}|}. \qquad (26)$$

Observe that many coefficients $c_i$ in (26) are zero. The set of all evaluation points needed is called *sparse grid* and denoted by $\mathcal{H}_{\mathcal{I}(w)}^m \subset \Gamma$ (see Figure 6). We also introduce the tensor notation

$$m(\mathbf{i}) = \prod_{n=1}^N m(i_n), \quad \Delta^{m(\mathbf{i})}[u] = \bigotimes_{n=1}^N \Delta^{m(i_n)}[u], \quad \mathcal{U}^{m(\mathbf{i})}[u] = \bigotimes_{n=1}^N \mathcal{U}^{m(i_n)}[u].$$

To fully characterize the sparse approximation operator $\mathcal{S}_{\mathcal{I}(w)}^m$ introduced in (24) one has to provide the sequence of sets $\mathcal{I}(w)$, the relation $m(i)$ between the level $i$ and the number of points in the corresponding one dimensional polynomial interpolation formula $\mathcal{U}^{m(i)}$, and the family of points to be used at each level, e.g. Clenshaw-Curtis or Gauss abscissae (see e.g. [28]).

In what follows we will consider Clenshaw-Curtis abscissae and the "doubling" rule $m(i) = db(i)$,

$$db(i) = \begin{cases} 0 \text{ if } i = 0 \\ 1 \text{ if } i = 1 \\ 2^{i-1} + 1, \text{ if } i > 1, \end{cases} \qquad (27)$$

which leads to nested grids. The classical Smolyak sparse grid (SM) uses $\mathcal{I}(w) = \{\mathbf{i} \in \mathbb{N}_+^N : |\mathbf{i} - \mathbf{1}| \leq w\}$. A quasi optimal choice of $\mathcal{I}(w)$ will be discussed in the next Section.

## 4.1   Quasi-optimal sparse grids

We now aim at constructing the quasi-optimal sparse grid for the Stochastic collocation method, i.e. we aim at choosing the best sequence of sets of indices. Let us define the error associated to a sparse grid as

$$E(\mathcal{S}_{\mathcal{I}(w)}^m) = \left\| u - \mathcal{S}_{\mathcal{I}(w)}^m[u] \right\|_{V \otimes L_\rho^2(\Gamma)},$$

and the work $W(\mathcal{S})$ as the number of evaluations needed, i.e.

$$W(\mathcal{S}_{\mathcal{I}(w)}^m) = |\mathcal{H}_{\mathcal{I}(w)}^m|.$$

Our goal is then to find the optimal set $\mathcal{S}$ that minimizes the error with a total work smaller or equal to a maximum work $W$, or alternatively the set that minimizes the work with an error smaller than or equal to a given threshold $\epsilon$.
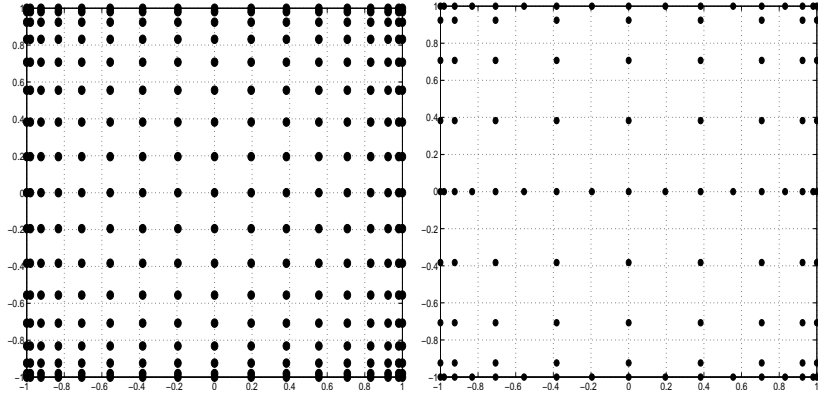
Figure 6: comparison between a tensor grid (left) and the TD-FC sparse grid (right) derived with the procedure explained in Section 4.1.

This is a classical knapsack problem and we adopt a greedy algorithm to solve it. To this end we define the error and work contribution of a multi-index $\mathbf{i}$. Let $\mathcal{J}$ be any set of indices such that $\mathbf{i} \notin \mathcal{J}$ and $\{\mathcal{J} \cup \mathbf{i}\}$ is admissible. Then the error contribution of $\mathbf{i}$ is

$$\Delta E(\mathbf{i}) = \left\| \mathcal{S}^m_{\{\mathcal{J} \cup \mathbf{i}\}}[u] - \mathcal{S}^m_{\mathcal{J}}[u] \right\|_{V \otimes L^2_\rho(\Gamma)} \tag{28}$$

and the work contribution is

$$\Delta W(\mathbf{i}) = |W(\mathcal{S}^m_{\{\mathcal{J} \cup \mathbf{i}\}}) - W(\mathcal{S}^m_{\mathcal{J}})|. \tag{29}$$

Observe that the error contribution defined in (28) is always independent of the set $\mathcal{J}$, since indeed

$$\Delta E(\mathbf{i}) = \left\| \sum_{\mathbf{j} \in \{\mathcal{J} \cup \mathbf{i}\}} \Delta^{m(\mathbf{j})}[u] - \sum_{\mathbf{j} \in \{\mathcal{J}\}} \Delta^{m(\mathbf{j})}[u] \right\|_{V \otimes L^2_\rho(\Gamma)} = \left\| \Delta^{m(\mathbf{i})}[u] \right\|_{V \otimes L^2_\rho(\Gamma)}. \tag{30}$$

On the other hand, the work contribution (29) will depend in general on the set $\mathcal{J}$, except in the case of nested abscissae, as for Clenshaw Curtis nodes, which is the case considered here. In this case indeed the evaluation of the extra term $\Delta^{m(\mathbf{i})}[u] = \bigotimes_{n=1}^N (\mathcal{U}^{m(i_n)} - \mathcal{U}^{m(i_n - 1)})[u]$ implies evaluations only in the extra points added at level $i_n$ in each direction, irrespectively of the set $\mathcal{J}$, provided that $\mathcal{J}$ is admissible.

Following [6, 14] we can now define the profit of an index $\mathbf{i}$ as

$$P(\mathbf{i}) = \frac{\Delta E(\mathbf{i})}{\Delta W(\mathbf{i})}$$

and identify the optimal sparse approximation operator $\mathcal{S}^*$ as the one using the set of most profitable indices, i.e. $\mathcal{I}^*(\epsilon) = \{\mathbf{i} \in \mathbb{N}^N_+ : P(\mathbf{i}) \geq \epsilon\}$.

22

To build the set $\mathcal{I}^*$ we rely on sharp estimates for both $\Delta E(\mathbf{i})$ and $\Delta W(\mathbf{i})$. Since, using Clenshaw-Curtis abscissae and the doubling rule $db(\cdot)$, we get nested grids, we can compute *exactly* $\Delta W(\mathbf{i})$ as

$$\Delta W(\mathbf{i}) = \prod_{n=1}^{N}(db(i_n) - db(i_n - 1)), \tag{31}$$

with $db(i_n)$ as in (27).

On the other hand, deriving a rigorous bound for $\Delta E(\mathbf{i})$ is not as easy. For instance, through numerical investigations on the model function $f(y_1, y_2) = \frac{1}{1+c_1 y_1 + c_2 y_2}$, one can conjecture the size of a generic $\Delta E(\mathbf{i})$ to be closely related to the norm of the corresponding Legendre coefficient $f_{m(\mathbf{i}-1)}$, with a correcting factor due to the interpolation operator norm. To be more precise, we conjecture the following estimate for $\Delta E(\mathbf{i})$, whenever $f$ is an analytic function:

$$\Delta E(\mathbf{i})[f] \lesssim \left\| f_{m(\mathbf{i}-1)} \right\|_V \prod_{n=1}^{N} \mathbb{L}_n^{m(i_n)}, \tag{32}$$

where $a \lesssim b$ means that there exists a constant $c$ independent of $\mathbf{i}$ such that $a \leq cb$ and $\mathbb{L}_n^{m(i)}$ is the Lebesgue constant for the interpolation operator $\mathcal{U}_n^{m(i)}$, defined as

$$\mathbb{L}_n^{m(i)} = \sup_{v \in C^0(\Gamma_n)} \frac{\left\| \mathcal{U}_n^{m(i)} v \right\|_{L^\infty(\Gamma_n)}}{\|v\|_{L^\infty(\Gamma_n)}}.$$

For Clenshaw-Curtis abscissae with doubling relation the Lebesgue constant can be shown to be

$$\mathbb{L}(db(i)) = \frac{2}{\pi} \log(db(i_n) + 1) + 1,$$

see e.g. [24] and references therein. Figure 7 shows the quality of estimate (32), and numerical results in next Section also confirm that such an estimate is accurate enough for our purposes.

Starting from (31) and (32), we can estimate the profit of each index, and estimate the sequence $\mathcal{S}_{\mathcal{I}^*(\epsilon)}$ of quasi-optimal grids with

$$\mathcal{I}^*(\epsilon) = \left\{ \mathbf{i} \in \mathbb{N}_+^N : \frac{C_0 \exp\left(-\sum_{n=1}^{N} db(i_n - 1)g_n\right) \frac{|db(\mathbf{i}-1)|!}{db(\mathbf{i}-1)!} \prod_{n=1}^{N} \mathbb{L}_n^{m(i_n)}}{\prod_{n=1}^{N}(db(i_n) - db(i_n - 1))} \geq \epsilon \right\}^{ADM} \tag{33}$$

with $\epsilon > 0 \in \mathbb{R}$. Equivalently, for $w = 0, 1, \dots$ we can define the sequence of sets

$$\mathcal{I}^*(w) = \left\{ \mathbf{i} \in \mathbb{N}_+^N : \sum_{i=n}^{N} db(i_n - 1)g_n - \log\frac{|db(\mathbf{i}-1)|!}{db(\mathbf{i}-1)!} - \sum_{n=1}^{N} \log\frac{\frac{2}{\pi}\log(db(i_n)+1)+1}{db(i_n) - db(i_n - 1)} \leq w \right\}^{ADM} \tag{34}$$

that will be used in (24) to build the quasi optimal sparse grids. We will refer to these "quasi best $M$-terms grids" as EW grids ("Error-Work" grids).

$$\text{(a) } f = \frac{1}{1+0.1y_1+0.1y_2} \qquad \text{(b) } f = \frac{1}{1+0.3y_1+0.3y_2} \qquad \text{(c) } f = \frac{1}{1+0.1y_1+0.5y_2}$$
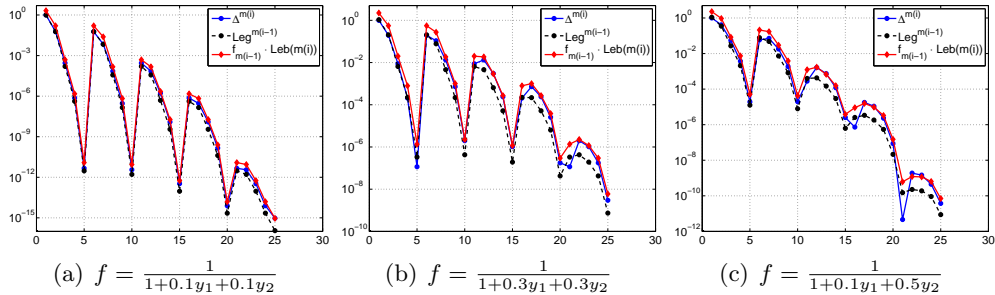
Figure 7: Numerical comparison between $\Delta E(\mathbf{i})$ and $|u_{m(\mathbf{i-1})}|$ for a scalar function $u$ of the form $f(y_1, y_2) = \frac{1}{1+c_1y_1+c_2y_2}$. Both the $\Delta E(\mathbf{i})$ for $\mathbf{i} \in TP(4)$ and the corresponding Legendre coefficients $|f_{m(\mathbf{i}-1)}|$ have been computed with a standard sparse grid SM(10).

**Remark 4.1.** *Observe that estimate (31) of the work $\Delta W(\mathbf{i})$ associated to a multi-index $\mathbf{i}$ is valid only if the underlying set of multi-indices is admissible. This is why in formulae (33) and (34) we have explicitly enforced the admissibility condition in the construction of the optimal set.*

## 4.2 Numerical tests on sparse grids

In this Section we consider the same problem as in Section 3.2 and use it to test the performance of the TD-FC grids derived above, comparing them with the classical SM grid and the best $M$-terms approximation.

To approximate the best $M$-terms we again consider a sufficiently large set $\mathbb{U}$ of multi-indices and for each of them we compute $\Delta W(\mathbf{i})$, $\Delta E(\mathbf{i})$ and their profit $P(\mathbf{i})$. Next, we sort the multi-indices according to $P(\mathbf{i})$, modify the sequence to fulfil the *ADM* condition (25) and compute the sparse grids according to this sequence.

We remark that the procedure just described only leads to an approximation of the best $M$-terms solution. Indeed, on the one hand replacing the total error $E(\mathcal{S})$ with the sum $\sum_{\mathbf{i}} \Delta E(\mathbf{i})$ provides only an upper bound that could be pessimistic because of possible cancellations, since the details $\Delta^{m(\mathbf{i})}[u]$ are not mutually orthogonal, in general. On the other hand, the fact that the most profitable index may be not admissible suggests that the solution cannot be found with a greedy algorithm. Here the coefficients $g_n$ in (34) are estimated numerically as in Section 3.2.

We also compare our results with the dimension adaptive algorithm [14], in the implementation proposed in [17] and available at

`http://www.ians.uni-stuttgart.de/spinterp`

This is an adaptive algorithm that given a sparse grid $\mathcal{S}_{\mathcal{I}}$ explores all neighbour multi-indices and adds to $\mathcal{I}$ the most profitable one. The algorithm implemented
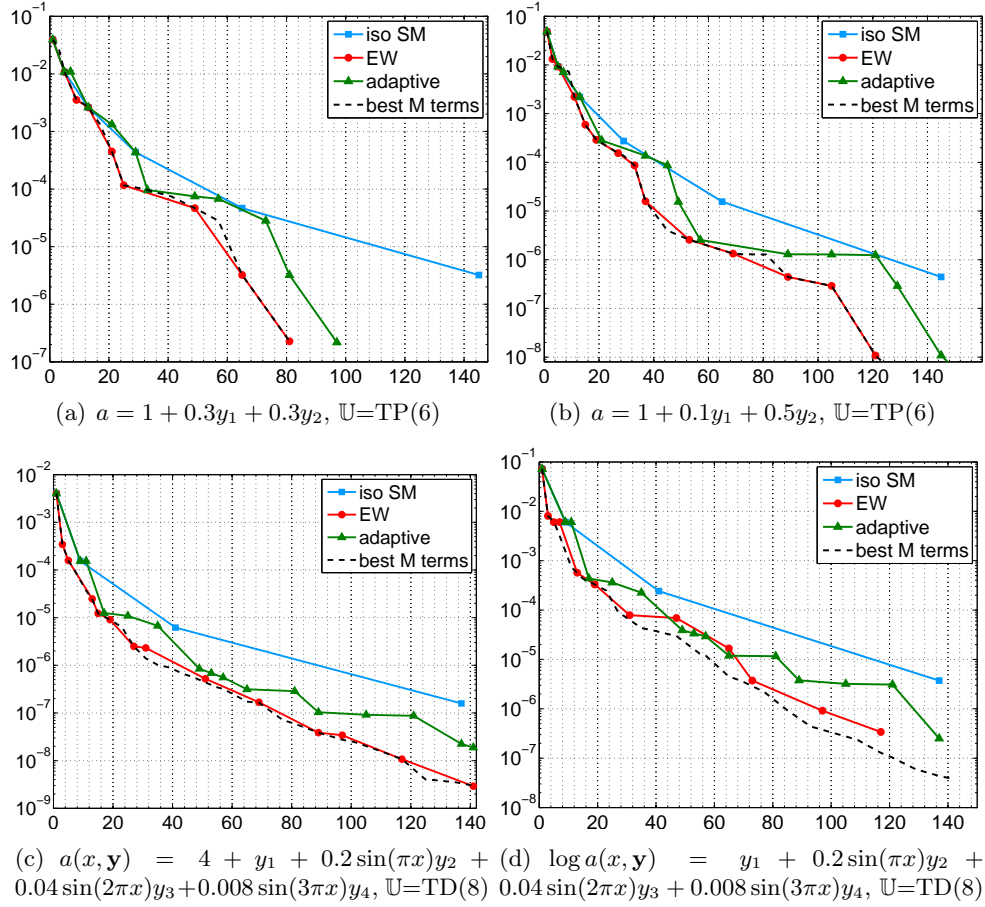
24

(a) $a = 1 + 0.3y_1 + 0.3y_2$, $\mathbb{U}$=TP(6)

(b) $a = 1 + 0.1y_1 + 0.5y_2$, $\mathbb{U}$=TP(6)

(c) $a(x, \mathbf{y}) = 4 + y_1 + 0.2\sin(\pi x)y_2 + 0.04\sin(2\pi x)y_3 + 0.008\sin(3\pi x)y_4$, $\mathbb{U}$=TD(8)

(d) $\log a(x, \mathbf{y}) = y_1 + 0.2\sin(\pi x)y_2 + 0.04\sin(2\pi x)y_3 + 0.008\sin(3\pi x)y_4$, $\mathbb{U}$=TD(8)

Figure 8: Results for TD-FC sparse grids compared with best $M$-terms , isotropic Smolyak and dimension adaptive algorithm. Convergence is measured as $\|\psi(u) - \psi(u_w)\|_{L^2(\Gamma)}$ versus number of evaluations (grid points).

in [17] has a tunable parameter $\widetilde{\omega}$ that allows one to move continuously from the classical Smolyak formula ($\widetilde{\omega} = 0$) to the fully adaptive algorithm ($\widetilde{\omega} = 1$). Following [17], in the present work we have set $\widetilde{\omega} = 0.9$, that numerically has been proved to be a good performing choice. The cost of this algorithm is the *total* number of evaluations needed, including also those necessary to explore all neighbours, to find the most profitable multi-index.

Figure 8 shows the convergence of the quantity $\|\psi(u) - \psi(u_w)\|_{L^2_\rho(\Gamma)}$ versus the number of grid points, for the different sparse grids considered. The $L^2_\rho$-norm has been computed with a high level isotropic Smolyak grid. The TD-FC grid is the best performing, even compared to the a-posteriori dimension adaptive algorithm [17], and the closest to the best $M$-terms grids sequence.

**Remark 4.2.** *A similar approach, based on estimates for $\Delta E$ and $\Delta W$ is possi-*

25

*ble also for the case of not nested grid points, as for the Gauss-Legendre quadrature points. However, in this case the estimate of $\Delta W$ is "path dependent" and any "path independent" estimate will be too pessimistic to build effective index sets.*

# 5 Conclusions

In this work we have proposed a new sequence of polynomial subspaces (TD-FC spaces in short) to be used in the solution of elliptic stochastic PDEs with Stochastic Galerkin method in the case of a solution that depends analytically on all random variables. The new polynomial spaces are based on sharp estimates of the decay of the Legendre coefficients.

The performances of TD-FC spaces have been assessed through some simple test cases. Here we have compared TD-FC with some standard choices of polynomial spaces and with the best $M$-terms approximation of the solution, that can be explicitly built for the examples considered. Results show that the TD-FC spaces perform better than the standard anisotropic TD ones, and are close to the best $M$-terms approximation a clear indication that our estimates of the decay of the Legendre coefficients are sharp. However, standard spaces may still have reasonable performances, if used in an appropriate anisotropic framework.

Using the estimate for the decay of the Legendre coefficients we have also defined a new class of sparse grids to be used in the context of Stochastic Collocation, relying on the concept of profit of each multi-index in the sparse grid. Again numerical tests show that these new sparse grids outperform the classical Smolyak construction and perform better than the a-posteriori dimension adaptive algorithm proposed in [14] (see also [17]). The reason is that our algorithm picks up the hierarchical surpluses based purely on a priori estimates, that turn out to be quite sharp, and does not have any extra cost to explore neighbor points as the algorithm in [17] does.

The new polynomial spaces and sparse grids proposed here are valid in the case of analytic dependence of the solution on the random variables. We point out, however, that the general strategy outlined in Sections 3.1 and 4.1 on how to build optimal polynomial spaces / sparse grids, is applicable to any problem. Of course, this strategy requires a sharp estimate of the decay of the coefficients of the spectral expansion of the solution on a orthonormal hierarchical basis (not necessarily polynomial). This step is highly problem dependent and should be analyze carefully in each situation, as we did here for a linear elliptic PDE with a stochastic coefficient dependent on uniformly distributed random variables.

# A  Proof of Theorem 2.1

Let us consider two sufficiently smooth $N$-dimensional functions $f(\mathbf{y}), g(\mathbf{y})$ : $\mathbb{R}^N \to \mathbb{R}$; an index $i \in \mathbb{N}$, $1 \leq i \leq N$; a set $\mathcal{S}$ of indices with cardinality $\mathscr{S}$; a multi-index $\mathbf{s} \in \mathbb{N}^N$. We use the following notation:

- $\partial_i f$ denotes the derivative of $f$ in the $i$-th direction: $\partial_i f = \frac{\partial}{\partial y_i} f$;

- $\partial_{\mathcal{S}} f$ denotes the $\mathscr{S}$-th order mixed derivative of $f$ with respect to all the directions included in $\mathcal{S}$. As an example, if $\mathcal{S} = \{1\,1\,2\,4\,4\,4\}$ then

$$\partial_{\mathcal{S}} f = \partial_{1\,1\,2\,4\,4\,4} f = \frac{\partial^6}{\partial_{y_1} \partial_{y_1} \partial_{y_2} \partial_{y_4} \partial_{y_4} \partial_{y_4}} f = \frac{\partial^6}{\partial_{y_1}^2 \partial_{y_2} \partial_{y_4}^3} f.$$

- $\mathbf{s}$ is the multi-index corresponding to the set $\mathcal{S}$ such that $\partial_{\mathcal{S}} f = \partial_{\mathbf{s}} f$. In the previous example $\mathbf{s} = [2\,1\,0\,3]$ is the multi-index corresponding to the set $\mathcal{S} = \{1\,1\,2\,4\,4\,4\}$.

**Lemma A.1** (generalized Leibniz rule). *Given a set of indices $\mathcal{K}$ with cardinality $\mathscr{K}$ and two functions $f, g : \mathbb{R}^N \to \mathbb{R}$, $f, g \in \mathcal{C}^{\mathscr{K}}(\mathbb{R}^N)$,*

$$\partial_{\mathcal{K}}(fg) = \sum_{\mathcal{S} \in \mathcal{P}(\mathcal{K})} \frac{\partial^{\mathscr{S}} f}{\prod_{i \in \mathcal{S}} \partial y_i} \frac{\partial^{\mathscr{K}-\mathscr{S}} g}{\prod_{i \notin \mathcal{S}} \partial y_i} = \sum_{\mathcal{S} \in \mathcal{P}(\mathcal{K})} \partial_{\mathcal{S}} f \, \partial_{\mathcal{K} \setminus \mathcal{S}} g, \qquad (35)$$

*where $\mathcal{P}(\mathcal{K})$ represents the power set of $\mathcal{K}$.*

**Lemma A.2.** *Let $a(\mathbf{x}, \mathbf{y})$ be a diffusion coefficient for equation (1) that satisfies Assumptions 2.1 - 2.3. Then the derivatives of $u$ can be bounded as*

$$\|\partial_{\mathbf{k}} u(\mathbf{y})\|_V \leq C_0 d_{|\mathbf{k}|} \mathbf{r}^{\mathbf{k}} \quad \forall \mathbf{y} \in \Gamma,$$

*where $C_0 = \frac{\|f\|_{V'}}{a_{min}}$, $\mathbf{r}$ as in Assumption 2.3, and $\{d_n\}_{n \in \mathbb{N}}$ is a sequence defined as:*

$$d_0 = 1, \, d_n = \sum_{i=0}^{n-1} \binom{n}{i} d_i. \qquad (36)$$

**Proof.**

We start by rewriting the statement using the correspondence between $\mathbf{k}$ and its equivalent set $\mathcal{K}$

$$\|\partial_{\mathbf{k}} u(\cdot, \mathbf{y})\|_V = \|\partial_{\mathcal{K}} \nabla u(\cdot, \mathbf{y})\|_{L^2(D)} \leq C_0 d_{\mathscr{K}} \mathbf{r}^{\mathbf{k}}, \quad \forall \mathbf{y} \in \Gamma.$$

We will first prove something closely related, namely

$$\left\| \sqrt{a(\cdot, \mathbf{y})} \partial_{\mathcal{K}} \nabla u(\cdot, \mathbf{y}) \right\|_{L^2(D)} \leq \frac{\|f\|_{V'}}{\sqrt{a_{min}}} d_{\mathscr{K}} \mathbf{r}^{\mathbf{k}} \quad \forall \mathbf{y} \in \Gamma, \qquad (37)$$

from which the previous inequality follows immediately. Let us start with a weak formulation of (1) in the physical space only, i.e.

*Find $u \in V \otimes L^2_\rho(\Gamma)$ such that for almost every $\mathbf{y} \in \Gamma$ it holds*

$$\int_D a(\mathbf{x}, \mathbf{y}) \nabla u(\mathbf{x}, \mathbf{y}) \cdot \nabla v(\mathbf{x}) d\mathbf{x} = \int_D f(\mathbf{x}) v(\mathbf{x}) d\mathbf{x} \quad \forall v \in V.(\mathbf{y}) \tag{38}$$

According to Lemma A.1, the $\partial_{\mathcal{K}}$ derivative of this weak formulation with respect to $\mathbf{y}$ is

$$\int_D \sum_{\mathcal{S} \in \mathcal{P}(\mathcal{K})} \partial_{\mathcal{S}} \nabla u(\mathbf{x}, \mathbf{y}) \partial_{\mathcal{K} \setminus \mathcal{S}} a(\mathbf{x}, \mathbf{y}) \nabla v(\mathbf{x}) d\mathbf{x} = 0,$$

and putting in evidence the $\partial_{\mathcal{K}} \nabla u$ term

$$\int_D a(\mathbf{x}, \mathbf{y}) \partial_{\mathcal{K}} \nabla u(\mathbf{x}, \mathbf{y}) \nabla v(\mathbf{x}) d\mathbf{x} = -\int_D \sum_{\mathcal{S} \in \mathcal{P}(\mathcal{K}), \mathcal{S} \neq \mathcal{K}} \partial_{\mathcal{S}} \nabla u(\mathbf{x}, \mathbf{y}) \partial_{\mathcal{K} \setminus \mathcal{S}} a(\mathbf{x}, \mathbf{y}) \nabla v(\mathbf{x}) d\mathbf{x}.$$

Next we choose $v = \partial_{\mathcal{K}} u$ and use Cauchy-Schwarz inequality on the right hand side:

$$\left\| \sqrt{a(\cdot, \mathbf{y})} \, \partial_{\mathcal{K}} \nabla u(\cdot, \mathbf{y}) \right\|^2_{L^2(D)} \leq$$

$$\sum_{\mathcal{S} \in \mathcal{P}(\mathcal{K}), \, \mathcal{S} \neq \mathcal{K}} \left\| \frac{\partial_{\mathcal{K} \setminus \mathcal{S}} a}{a}(\cdot, \mathbf{y}) \right\|_{L^\infty(D)} \left\| \sqrt{a(\cdot, \mathbf{y})} \partial_{\mathcal{S}} \nabla u(\cdot, \mathbf{y}) \right\|_{L^2(D)} \left\| \sqrt{a(\cdot, \mathbf{y})} \partial_{\mathcal{K}} \nabla u(\cdot, \mathbf{y}) \right\|_{L^2(D)}.$$

Now simplify $\left\| \sqrt{a(\cdot, \mathbf{y})} \partial_{\mathcal{K}} \nabla u(\cdot, \mathbf{y}) \right\|_{L^2(D)}$ on both sides and reorder the sum on the right hand side according to the cardinality of the subsets $\mathcal{S}$. From here on we omit the dependence of $a$ and $u$ on $\mathbf{x}, \mathbf{y}$, to have a lighter notation. We have

$$\left\| \sqrt{a} \, \partial_{\mathcal{K}} \nabla u \right\|_{L^2(D)} \leq \sum_{i=0}^{\mathcal{K}-1} \sum_{\mathcal{S} \in \mathcal{P}(\mathcal{K}), \mathcal{S} = i} \left\| \frac{\partial_{\mathcal{K} \setminus \mathcal{S}} a}{a} \right\|_{L^\infty(D)} \left\| \sqrt{a} \partial_{\mathcal{S}} \nabla u \right\|_{L^2(D)}. \tag{39}$$

We are finally in the position to prove (37). We will proceed by induction on (37), using (39) and Assumption 2.3 on the decay of $a$.

**Case $\mathcal{K} = 0$.** In this case (37) reads

$$\left\| \sqrt{a} \nabla u \right\|_{L^2(D)} \leq \frac{\|f\|_{V'}}{\sqrt{a_{min}}} d_0,$$

which is true setting $d_0 = 1$.

**Case $\mathcal{K} = 1$.** If $\mathcal{K} = \{j\}$, $1 \leq j \leq N$, (37) reads

$$\left\| \sqrt{a} \partial_j \nabla u \right\|_{L^2(D)} \leq \frac{\|f\|_{V'}}{\sqrt{a_{min}}} d_1 r_j = \frac{\|f\|_{V'}}{\sqrt{a_{min}}} r_j \binom{1}{0} d_0 = \frac{\|f\|_{V'}}{\sqrt{a_{min}}} r_j d_0 = \frac{\|f\|_{V'}}{\sqrt{a_{min}}} r_j.$$

To prove this, consider (39). Using Assumption 2.3 and the result for case $\mathcal{K} = 0$ one has precisely

$$\left\| \sqrt{a} \partial_j \nabla u \right\|_{L^2(D)} \leq \left\| \frac{\partial_j a}{a} \right\|_{L^\infty(D)} \left\| \sqrt{a} \nabla u \right\|_{L^2(D)} \leq r_j \frac{\|f\|_{V'}}{\sqrt{a_{min}}}.$$

28

**General $\mathscr{K}$.** Consider now a general $\mathcal{K}$, and suppose (37) holds for any set $\mathcal{S}$ with cardinality $\mathscr{K} - 1$. Use this induction hypothesis and again Assumption 2.3 on (39), denoting with $\mathbf{s}*$ the multi-index corresponding to the set $\mathcal{K} \setminus \mathcal{S}$. This yields

$$\left\| \sqrt{a} \partial_{\mathcal{K}} \nabla u \right\|_{L^2(D)} \leq \sum_{i=0}^{\mathscr{K}-1} \sum_{\mathcal{S} \in \mathcal{P}(\mathcal{K}), \mathscr{S}=i} \mathbf{r}^{\mathbf{s}*} \frac{\|f\|_{V'}}{\sqrt{a_{min}}} d_{\mathscr{S}} \mathbf{r}^{\mathbf{s}}.$$

Next note that:

$$\mathbf{r}^{\mathbf{s}*} \mathbf{r}^{\mathbf{s}} = \prod_{j \in \mathcal{K} \setminus \mathcal{S}} r_j \prod_{j \in \mathcal{S}} r_j = \prod_{j \in \mathcal{K}} r_j = \mathbf{r}^{\mathbf{k}},$$

and that the number of subsets $\mathcal{S}$ with cardinality $i$ is $\begin{pmatrix} \mathscr{K} \\ i \end{pmatrix}$. Then

$$\left\| \sqrt{a} \partial_{\mathcal{K}} \nabla u \right\|_{L^2(D)} \leq \frac{\|f\|_{V'}}{\sqrt{a_{min}}} \mathbf{r}^{\mathbf{k}} \sum_{i=0}^{\mathscr{K}-1} \begin{pmatrix} \mathscr{K} \\ i \end{pmatrix} d_i = \frac{\|f\|_{V'}}{\sqrt{a_{min}}} \mathbf{r}^{\mathbf{k}} d_{\mathscr{K}}$$

which proves the result.

$\square$

**Lemma A.3.** *The sequence $\{d_n\}_{n \in \mathbb{N}}$ defined in (36) can bounded as*

$$d_n \leq \left( \frac{1}{\log 2} \right)^n n! \tag{40}$$

**Proof.** From definition (36) we have

$$d_n = \sum_{i=0}^{n-1} \begin{pmatrix} n \\ i \end{pmatrix} d_i = \sum_{i=0}^{n-1} \frac{n!}{i!(n-i)!} d_i.$$

Let $f_n = \dfrac{d_n}{n!}$; the recurrency relation then becomes

$$f_n = \sum_{i=0}^{n-1} \frac{f_i}{(n-i)!}, \quad f_0 = f_1 = 1. \tag{41}$$

We now show by induction that $f_n \leq C\alpha^n$, with $C, \alpha \in \mathbb{R}$. Enforcing $1 = f_0 \leq C$ and $1 = f_1 \leq C\alpha$ results in $C \geq 1$ and $\alpha \geq 1$. Next, we reorder the sum in (41) and exploit the inductive hypothesis:

$$f_n = \sum_{i=0}^{n-1} \frac{f_{n-1-i}}{(1+i)!} \leq \sum_{i=0}^{n-1} \frac{C\alpha^{n-1-i}}{(1+i)!} = C\alpha^n \sum_{i=0}^{n-1} \frac{\alpha^{-(1+i)}}{(1+i)!} = C\alpha^n \left( e^{\frac{1}{\alpha}} - 1 \right) \leq C\alpha^n,$$

where the last inequality holds true provided we choose $e^{\frac{1}{\alpha}} - 1 \leq 1$. Therefore we take $\alpha = (\log 2)^{-1}$ and $C = 1$, yielding $f_n \leq (\log 2)^{-n}$ and $d_n \leq (\log 2)^{-n} n!$

$\square$

**Theorem 1.** *Let $a(\mathbf{x}, \mathbf{y})$ be a diffusion coefficient for equation (1) that satisfies Assumptions 2.1 - 2.3. Then the derivatives of $u$ can be bounded as*

$$\|\partial_{\mathbf{i}} u(\mathbf{y})\|_V \le C_0 |\mathbf{i}|! \, \tilde{\mathbf{r}}^{\mathbf{i}} \quad \forall \mathbf{y} \in \Gamma.$$

*Here $C_0 = \dfrac{\|f\|_{V'}}{a_{min}}$ and $\tilde{\mathbf{r}} = \left(\dfrac{1}{\log 2}\right) \mathbf{r}$, with $\mathbf{r}$ as in Assumption 2.3.*

**Proof.**  Combine Lemma A.2 and A.3. □

# References

[1] I. M. Babuška, R. Tempone, and G. E. Zouraris.  Galerkin finite element approximations of stochastic elliptic partial differential equations. *SIAM J. Numer. Anal.*, 42(2):800–825, 2004.

[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(3):1005–1034, 2007.

[3] J. Bäck, F. Nobile, L. Tamellini, and R. Tempone.  Stochastic spectral Galerkin and collocation methods for PDEs with random coefficients:  a numerical comparison. In J.S. Hesthaven and E.M. Ronquist, editors, *Spectral and High Order Methods for Partial Differential Equations*, volume 76 of *Lecture Notes in Computational Science and Engineering*, pages 43–62. Springer, 2011. Selected papers from the ICOSAHOM '09 conference, June 22-26, Trondheim, Norway.

[4] V. Barthelmann, E. Novak, and K. Ritter.  High dimensional polynomial interpolation on sparse grids. *Adv. Comput. Math.*, 12(4):273–288, 2000.

[5] M. Bieri, R. Andreev, and C. Schwab. Sparse tensor discretization of elliptic sPDEs. SAM-Report 2009-07, Seminar für Angewandte Mathematik, ETH, Zurich, 2009.

[6] H.J Bungartz and M. Griebel.  Sparse grids. *Acta Numer.*, 13:147–269, 2004.

[7] A. Cohen, R. DeVore, and C. Schwab. Analytic regularity and polynomial approximation of parametric and stochastic elliptic PDEs.  SAM-Report 2010-03, Seminar für Angewandte Mathematik, ETH, Zurich, 2010.

[8] A. Cohen, R. DeVore, and C. Schwab.  Convergence rates of best $n$-term Galerkin approximations for a class of elliptic sPDEs. *Foundations of Computational Mathematics*, 10:615–646, 2010. 10.1007/s10208-010-9072-2.

[9] P.J. Davis. *Interpolation and approximation.* Dover Publications Inc., New York, 1975. Republication, with minor corrections, of the 1963 original, with a new preface and bibliography.

[10] H. C. Elman, C. W. Miller, E. T. Phipps, and R. S. Tuminaro. Assessment of Collocation and Galerkin approaches to linear diffusion equations with random data. *International Journal for Uncertainty Quantification*, 1(1):19–33, 2011.

[11] O. G. Ernst, A. Mugler, H.-J. Starkloff, and E. Ullmann. On the convergence of generalized polynomial chaos expansions. Submitted, 2010.

[12] B. Ganapathysubramanian and N. Zabaras. Sparse grid collocation schemes for stochastic natural convection problems. *Journal of Computational Physics*, 225(1):652–685, 2007.

[13] W. Gautschi. *Orthogonal Polynomials: Computation and Approximation.* Oxford University Press, Oxford, 2004.

[14] T. Gerstner and M. Griebel. Dimension-adaptive tensor-product quadrature. *Computing*, 71(1):65–87, 2003.

[15] R. G. Ghanem and P. D. Spanos. *Stochastic Finite Elements: a Spectral Approach.* Springer–Verlag, New York, 1991.

[16] M. Griebel and S. Knapek. Optimized general sparse grid approximation spaces for operator equations. *Math. Comp.*, 78(268):2223–2257, 2009.

[17] A. Klimke. *Uncertainty modeling using fuzzy arithmetic and sparse grids.* PhD thesis, Universität Stuttgart, Shaker Verlag, Aachen, 2006.

[18] C. Lubich. *From quantum to classical molecular dynamics: reduced models and numerical analysis.* Zurich lectures in advanced mathematics. European Mathematical Society, 2008.

[19] H. G. Matthies and A. Keese. Galerkin methods for linear and nonlinear elliptic stochastic partial differential equations. *Comput. Methods Appl. Mech. Engrg.*, 194(12-16):1295–1331, 2005.

[20] Roger B. Nelsen. *An introduction to copulas.* Springer Series in Statistics. Springer, New York, second edition, 2006.

[21] F. Nobile, R. Tempone, and C.G. Webster. An anisotropic sparse grid stochastic collocation method for partial differential equations with random input data. *SIAM J. Numer. Anal.*, 46(5):2411–2442, 2008.

[22] M.F. Pellissetti and R.G. Ghanem. Iterative solution of systems of linear equations arising in the context of stochastic finite elements. *Adv. Eng. Software*, 31:607–616, 2000.

[23] C.E. Powell and H.C. Elman. Block-diagonal preconditioning for spectral stochastic finite-element systems. *IMA J. Numer. Anal.*, 29(2):350–375, 2009.

[24] A. Quarteroni, R. Sacco, and F. Saleri. *Numerical mathematics*, volume 37 of *Texts in Applied Mathematics*. Springer-Verlag, Berlin, second edition, 2007.

[25] S.A. Smolyak. Quadrature and interpolation formulas for tensor products of certain classes of functions. *Dokl. Akad. Nauk SSSR*, 4:240–243, 1963.

[26] G. Szegö. *Orthogonal polynomials*. Colloquium Publications - American Mathematical Society. American Mathematical Society, 1939.

[27] R. A. Todor and C. Schwab. Convergence rates for sparse chaos approximations of elliptic problems with stochastic coefficients. *IMA J Numer Anal*, 27(2):232–261, 2007.

[28] Lloyd N. Trefethen. Is Gauss quadrature better than Clenshaw-Curtis? *SIAM Rev.*, 50(1):67–87, 2008.

[29] D. Xiu and J.S. Hesthaven. High-order collocation methods for differential equations with random inputs. *SIAM J. Sci. Comput.*, 27(3):1118–1139, 2005.

[30] D. Xiu and G.E. Karniadakis. The Wiener-Askey polynomial chaos for stochastic differential equations. *SIAM J. Sci. Comput.*, 24(2):619–644, 2002.

# MOX Technical Reports, last issues

## Dipartimento di Matematica "F. Brioschi",
## Politecnico di Milano, Via Bonardi 9 - 20133 Milano (Italy)

**23/2011**  BECK, J.; NOBILE, F.; TAMELLINI, L.; TEMPONE, R.
*On the optimal polynomial approximation of stochastic PDEs by Galerkin and Collocation methods*

**22/2011**  AZZIMONTI, L.; IEVA, F.; PAGANONI, A.M.
*Nonlinear nonparametric mixed-effects models for unsupervised classification*

**21/2011**  AMBROSI, D.; PEZZUTO, S.
*Active stress vs. active strain in mechanobiology: constitutive issues*

**20/2011**  ANTONIETTI, P.F.; HOUSTON, P.
*Preconditioning high–order Discontinuous Galerkin discretizations of elliptic problems*

**19/2011**  PASSERINI, T.; SANGALLI, L.; VANTINI, S.; PICCINELLI, M.; BACIGALUPPI, S.; ANTIGA, L.; BOCCARDI, E.; SECCHI, P.; VENEZIANI, A.
*An Integrated Statistical Investigation of the Internal Carotid Arteries hosting Cerebral Aneurysms*

**18/2011**  BLANCO, P.; GERVASIO, P.; QUARTERONI, A.
*Extended variational formulation for heterogeneous partial differential equations*

**16/2011**  MESIN, L; AMBROSI, D.
*Spiral waves on a contractile tissue*

**17/2011**  QUARTERONI, A.; ROZZA, G.; MANZONI, A.
*Certified Reduced Basis Approximation for Parametrized Partial Differential Equations and Applications*

**15/2011**  ARGIENTO, R.; GUGLIELMI, A.; SORIANO J.
*A semiparametric Bayesian generalized linear mixed model for the reliability of Kevlar fibres*

**14/2011**  ANTONIETTI, P.F.; MAZZIERI, I.; QUARTERONI, A.; RAPETTI, F.
*Non-Conforming High Order Approximations for the Elastic Wave Equation*