

MATHICSE Technical Report

Nr. 12.2014
February 2014



Convergence of quasi-optimal sparse grid approximation of Hilbert-valued functions: application to random elliptic PDEs

Fabio Nobile, Lorenzo Tamellini, Raúl Tempone

Convergence of quasi-optimal sparse grid approximation of Hilbert-valued functions: application to random elliptic PDEs

F. Nobile · L. Tamellini · R. Tempone

the date of receipt and acceptance should be inserted later

Abstract In this work we provide a convergence analysis for the quasi-optimal version of the Stochastic Sparse Grid Collocation method we had presented in our previous work “On the optimal polynomial approximation of Stochastic PDEs by Galerkin and Collocation methods” [6]. Here the construction of a sparse grid is recast into a knapsack problem: a profit is assigned to each hierarchical surplus and only the most profitable ones are added to the sparse grid. The convergence rate of the sparse grid approximation error with respect to the number of points in the grid is then shown to depend on weighted summability properties of the sequence of profits. This argument is very general and can be applied to sparse grids built with any uni-variate family of points, both nested and non-nested. As an example, we apply such quasi-optimal sparse grid to the solution of a particular elliptic PDE with stochastic diffusion coefficients, namely the “inclusions problem”: we detail the convergence estimate obtained in this case, using polynomial interpolation on either nested (Clenshaw–Curtis) or non-nested (Gauss–Legendre) abscissas, verify its sharpness numerically, and compare the performance of the resulting quasi-optimal grids with a few alternative sparse grids construction schemes recently proposed in literature.

Keywords Uncertainty Quantification · random PDEs · linear elliptic equations · multivariate polynomial approximation · best M-terms polynomial approximation · Smolyak approximation · Sparse grids · Stochastic Collocation method

F. Nobile · L. Tamellini
CSQI - MATHICSE, Ecole Polytechnique Fédérale Lausanne, Station 8, CH 1015, Lausanne, Switzerland. E-mail: lorenzo.tamellini@epfl.ch, tel.: +41216934234, fax: +41216930545.

R. Tempone
Applied Mathematics and Computational Science, 4700, King Abdullah University of Science and Technology, Thuwal, 23955-6900, Kingdom of Saudi Arabia.

Mathematics Subject Classification (2000) 41A10, 65C20, 65N12, 65N30, 65N35

1 Introduction

Sparse grid polynomial approximation has emerged as one of the most appealing methods for approximating high-dimensional functions. Indeed, although as simple to use as a sampling strategy, it can converge significantly faster if the function at hand presents some degree of differentiability. A number of “off-the-shelf” sparse grid packages can be found on the web¹, further enhancing the spread of this technique among practitioners.

Yet, the sparse grid technique experiences a dramatic performance deterioration as the number of random variables increases, i.e. it suffers from the so-called “curse of dimensionality effect”, to which Monte Carlo sampling methods are instead essentially immune. To avoid, or at least alleviate, this undesirable feature, a number of approaches have been recently proposed. Among others, we mention the anisotropic sparse grid technique [4, 26] and the a-posteriori adaptive strategy investigated in [9, 12, 19, 21, 23]; here we further investigate the quasi-optimal sparse grids method proposed in [6]. Such method, applied to elliptic PDEs with random coefficients, consists in reformulating the problem of the construction of a sparse grid as a knapsack problem as first proposed in [9, 19, 21], and estimating the profit of each sparse grid component (*hierarchical surplus*) using combined a-priori/a-posteriori information, i.e. providing a-priori estimates whose constants are numerically tuned (hence the name “quasi-optimal”), that have been observed numerically to be quite sharp. The goal of this work is to present a convergence theorem for such “knapsack” sparse grids in terms of the weighted τ -summability of the profits: in particular, we extend and improve the preliminary estimates presented in [31]. Our result is general and can accommodate both the case of nested and non-nested abscissas. We mention that two other related works have appeared more recently in the literature, [29] and [12], addressing alternative convergence estimates for the knapsack-type sparse grids. See also [26] for an older estimate for the convergence of anisotropic sparse grid approximations: in that work, the construction of the sparse grids is based on error contributions only rather than profits like here, making our current theoretical and numerical results better.

As a specific application we consider an elliptic PDE with random coefficients, namely the so-called “inclusions problem” already discussed in [4, 8], whose solution falls in the class of analytic functions in polyellipses, see e.g. [3, 8, 13]. We will derive an estimate of the profits of the hierarchical surpluses for this family of functions and prove that such profits satisfy suitable weighted summability properties. We will then deduce, using the above-mentioned Theorem, rigorous convergence results for the corresponding quasi-

¹ see e.g. <http://www.ians.uni-stuttgart.de/spinterp> or <http://dakota.sandia.gov>

optimal sparse grids approximation: in particular, we will show that it converges sub-exponentially with a rate comparable to that of the optimal (“best M -terms”) L^2 approximation in the case of nested points, and with half the rate for non-nested points, cf. [8, Theorem 16]. We will then verify numerically the sharpness of the estimates thus obtained, using Clenshaw–Curtis and Gauss–Legendre points as specific representatives of the two families of nested and non-nested points, and compare the performances of the quasi-optimal sparse grids with that of a few other sparse grid schemes recently proposed in the literature.

The rest of the work is organized as follows. Section 2 defines the general approximation problem and introduces the sparse grid methodology. The quasi-optimal sparse grid construction is explained in Section 3, while the general convergence result in terms of the weighted τ -summability of the profits is given in Section 4, see Theorem 1. Section 5 introduces the above-mentioned class of polyellipse-analytic Hilbert-valued functions and builds on the previous general Theorem to derive rigorous convergence estimates for their quasi-optimal sparse grid approximation with nested and non-nested collocation points, see Theorems 2 and 3. In particular, the Theorems are stated at the beginning of the Section, and the rest of the Section is devoted to their proof. Section 6 introduces the “inclusion problem” and shows some numerical results that confirm the effectiveness of the proposed quasi-optimal strategy and the sharpness of the proposed convergence estimates, while conclusions and final remarks are presented in Section 7.

2 Sparse grid polynomial approximation of Hilbert space-valued functions

We consider the problem of approximating a multivariate function $u(\mathbf{y}) : \Gamma \rightarrow V$, where Γ is an N -variate hypercube $\Gamma = \Gamma_1 \times \Gamma_2 \times \dots \times \Gamma_N$ (with $\Gamma_n \subseteq \mathbb{R}$ and N possibly infinite), and V is a Hilbert space. Furthermore, we assume that each Γ_n is endowed with a probability measure $\varrho_n(y_n)dy_n$, so that $\varrho(\mathbf{y})d\mathbf{y} = \prod_{n=1}^N \varrho_n(y_n)$ is a probability measure on Γ , and we restrict our attention to functions in the Bochner space $L^2_\varrho(\Gamma; V)$, where

$$L^2_\varrho(\Gamma; V) = \left\{ u : \Gamma \rightarrow V \text{ s.t. } \int_\Gamma \|u(\mathbf{y})\|_V^2 \varrho(\mathbf{y})d\mathbf{y} < \infty \right\}.$$

Observe that, since V and $L^2_\varrho(\Gamma)$ are Hilbert spaces, $L^2_\varrho(\Gamma; V)$ can be equivalently understood as the tensor space $V \otimes L^2_\varrho(\Gamma)$, defined as the completion of formal sums $v = \sum_{k=1}^{k'} \phi_k \psi_k$, with $\phi_k \in V$ and $\psi_k \in L^2_\varrho(\Gamma)$, with respect to the inner product

$$(v, \widehat{v})_{V \otimes L^2_\varrho(\Gamma)} = \sum_{k, \ell} (\phi_k, \widehat{\phi}_\ell)_V, (\psi_k, \widehat{\psi}_\ell)_{L^2_\varrho(\Gamma)}.$$

In particular, we aim at approximating $u(\mathbf{y})$ with global polynomials over Γ , which is a sound approach if u is a smooth function of \mathbf{y} . To introduce the

polynomial subspace of $V \otimes L_{\sigma}^2(\Gamma)$ in which we will build our approximate solution, it is convenient to use a multi-index notation². Let $w \in \mathbb{N}$ be an integer index indicating the level of approximation, and $\Lambda(w)$ a sequence of index sets in \mathbb{N}^N such that $\Lambda(0) = \{\mathbf{0}\}$, $\Lambda(w) \subseteq \Lambda(w+1)$ for $w \geq 0$ and $\mathbb{N}^N = \bigcup_{w \in \mathbb{N}} \Lambda(w)$. Denoting by $\mathbb{P}_{\Lambda(w)}(\Gamma)$ the multivariate polynomial space

$$\mathbb{P}_{\Lambda(w)}(\Gamma) = \text{span} \left\{ \prod_{n=1}^N y_n^{p_n}, \mathbf{p} \in \Lambda(w) \right\},$$

we will look for an approximation

$$u_w \in V \otimes \mathbb{P}_{\Lambda(w)}(\Gamma) = \left\{ \sum_j v_j q_j(\mathbf{y}), v_j \in V, q_j \in \mathbb{P}_{\Lambda(w)}(\Gamma) \right\}.$$

Clearly, the polynomial space $\mathbb{P}_{\Lambda(w)}(\Gamma)$ should be designed to have good approximation properties while keeping the number of degrees of freedom as low as possible. Although this is a problem-dependent choice, using the classical Tensor Product polynomial space $\mathbb{P}_{TP(w)}(\Gamma)$, with $TP(w) = \{\mathbf{q} \in \mathbb{N}^N : \max_n q_n \leq w\}$, is in general not a good choice, as its dimension grows exponentially fast with the number of random variables N , i.e. $\dim \mathbb{P}_{TP(w)}(\Gamma) = (1+w)^N$. Valid alternative choices that have been widely used in literature are: the Total Degree polynomial space $\mathbb{P}_{TD(w)}(\Gamma)$, $TD(w) = \{\mathbf{q} \in \mathbb{N}^N : \sum_n q_n \leq w\}$, see e.g. [20, 24], that contains indeed only $\binom{N+w}{N}$ monomials but has approximation properties similar to the Tensor Product space; or the Hyperbolic Cross polynomial space $\mathbb{P}_{HC(w)}(\Gamma)$, $HC(w) = \{\mathbf{q} \in \mathbb{N}^N : \prod_n (q_n + 1) \leq w + 1\}$, see e.g. [1, 25, 30]. One could also introduce anisotropy in the approximation, to enrich the polynomial space only in those variables y_n which contribute the most to the total variability of the solution, see e.g. [4]. Several methods can be used to compute the polynomial approximation u_w (projection, interpolation, regression): in this work we consider the Sparse Grid Approximation Method, that we briefly review in the rest of this Section.

2.1 Sparse grid Approximation Method

For a given level of approximation $w \geq 0$, the sparse grid approximation method (see e.g. [5, 9] and references therein) consists in evaluating the function u in a set of W points $\mathbf{y}_1, \dots, \mathbf{y}_W \in \Gamma$, and building a global polynomial approximation u_w (not necessarily interpolatory) in a suitable space $V \otimes \mathbb{P}_{\Lambda(w)}(\Gamma)$.

² Throughout the rest of this work, \mathbb{N} will denote the set of integer numbers including 0, and \mathbb{N}_+ that of integer numbers excluding 0. Moreover, $\mathbf{0}$ will denote the vector $(0, 0, \dots, 0) \in \mathbb{N}^N$, $\mathbf{1}$ the vector $(1, 1, \dots, 1) \in \mathbb{N}^N$, and \mathbf{e}_j the j -th canonical vector in \mathbb{R}^N , i.e. a vector whose components are all zero but the j -th, whose value is one. Finally, given two vectors $\mathbf{v}, \mathbf{w} \in \mathbb{N}^N$, $\mathbf{v} \leq \mathbf{w}$ if and only if $v_j \leq w_j$ for every $1 \leq j \leq N$.

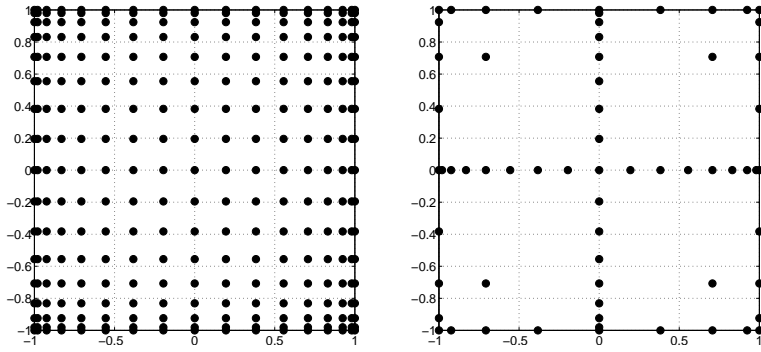


Fig. 1: Left: Tensor grid. Right: Sparse grid.

For each direction y_n we introduce a sequence of one-dimensional polynomial Lagrangian interpolant operators of increasing degree, indexed by $i_n \geq 1$:

$$\forall i_n \geq 1, \quad \mathcal{U}_n^{m(i_n)} : C^0(\Gamma_n) \rightarrow \mathbb{P}_{m(i_n)-1}(\Gamma_n),$$

where $m(i_n)$ is the number of collocation points used to build the interpolant at level i_n and $\mathbb{P}_q(\Gamma_n)$ is the set of polynomials in y_n of degree at most q . We require the level-to-nodes function $m : \mathbb{N} \rightarrow \mathbb{N}$ to satisfy the following assumptions:

$$m(0) = 0, \quad m(1) = 1, \quad m(i_n) < m(i_n + 1), \quad i_n \geq 1.$$

In addition, let $\mathcal{U}_n^0[f] = 0, \forall f \in C^0(\Gamma_n)$. Next, we introduce the difference operators $\Delta_n^{m(i_n)} = \mathcal{U}_n^{m(i_n)} - \mathcal{U}_n^{m(i_n-1)}$, and consider a sequence of index sets $\mathcal{I}(w) \subset \mathbb{N}_+^N$ such that $\mathcal{I}(w) \subset \mathcal{I}(w+1)$ and $\mathcal{I}(0) = \{\mathbf{1}\}$. We define the sparse grid approximation of $u : \Gamma \rightarrow V$ at level w as

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

As pointed out in [19], it is desirable that the sum (1) has some telescopic properties. To ensure this, we have to impose some additional constraints on \mathcal{I} . Following [19] we say that a set \mathcal{I} is *admissible*³ if

$$\forall \mathbf{i} \in \mathcal{I}, \quad \mathbf{i} - \mathbf{e}_j \in \mathcal{I} \text{ for } 1 \leq j \leq N \text{ such that } i_j > 1. \quad (2)$$

We refer to this property as *admissibility condition*, or *ADM* in short. Given a multi-index set \mathcal{I} , we will denote by \mathcal{I}^{ADM} the smallest admissible set such that $\mathcal{I} \subset \mathcal{I}^{ADM}$. The set of all evaluation points needed by (1) is called a *sparse grid*, and we denote its cardinality by $W_{\mathcal{I}(w),m}$. Note that (1) is indeed

³ Also known as *lower sets* or *downward closed set*, see e.g. [14].

equivalent to a linear combination of tensor grid interpolations each of which uses only “few” interpolation points (see e.g. [34]):

$$\mathcal{S}_{\mathcal{I}(w)}^m[u](\mathbf{y}) = \sum_{\mathbf{i} \in \mathcal{I}(w)^{ADM}} c_{\mathbf{i}} \bigotimes_{n=1}^N \mathcal{U}_n^{m(i_n)}[u](\mathbf{y}), \quad c_{\mathbf{i}} = \sum_{\substack{\mathbf{j} \in \{0,1\}^N \\ (\mathbf{i}+\mathbf{j}) \in \mathcal{I}(w)^{ADM}}} (-1)^{|\mathbf{j}|}. \quad (3)$$

Observe that many of the coefficients $c_{\mathbf{i}}$ in (3) may be zero: in particular $c_{\mathbf{i}}$ is zero whenever $\mathbf{i} + \mathbf{1} \in \mathcal{I}(w)^{ADM}$. To any sparse grid one can associate a corresponding quadrature formula $\mathcal{Q}_{\mathcal{I}(w)}^m[\cdot]$,

$$\forall f \in C^0(\Gamma), \quad \int_{\Gamma} f(\mathbf{y}) \varrho(\mathbf{y}) d\mathbf{y} \approx \int_{\Gamma} \mathcal{S}_{\mathcal{I}(w)}^m[f] \varrho(\mathbf{y}) d\mathbf{y} \approx \mathcal{Q}_{\mathcal{I}(w)}^m[f] = \sum_{j=1}^{W_{\mathcal{I}(w)}^m} f(\mathbf{y}_j) \beta_j,$$

for suitable $\beta_j \in \mathbb{R}$. In particular, given $g \in V'$, where V' is the dual space of V , the expected value of $\langle g, u \rangle$ can be computed as

$$\mathbb{E}[\langle g, u \rangle] \approx \mathcal{Q}_{\mathcal{I}(w)}^m[\langle g, u \rangle] = \sum_{j=1}^{W_{\mathcal{I}(w)}^m} \langle g, u(\mathbf{y}_j) \rangle \beta_j.$$

The sequence of sets $\mathcal{I}(w)$, the level-to-nodes function m and the family of collocation points to be used at each level characterize the sparse approximation operator $\mathcal{S}_{\mathcal{I}(w)}^m$ introduced in (1). The choice of $\mathcal{I}(w)$ will be the subject of the next Section; anisotropic sets $\mathcal{I}(w)$, that enrich the approximation in specific directions of the parameter space Γ have been studied in [4, 26]. As for the family of points, they should be chosen according to the probability measure $\varrho(\mathbf{y}) d\mathbf{y}$ on Γ for an optimal performance, e.g. the Gauss–Legendre points for the Uniform measure, and the Gauss–Hermite points for the Gaussian measure (see e.g. [33]). For good uniform approximations on $\Gamma = [-1, 1]^N$, Clenshaw–Curtis points are also a good choice. In the following, we will refer to *nested* points when the set of points used to build the operator $\mathcal{U}_n^{m(i_n)}$ is a subset of the points of the operator $\mathcal{U}_n^{m(i_{n+1})}$, and *non-nested* points otherwise. It is well-known that nested quadrature formulae have a lower degree of exactness than Gaussian quadrature formulae when approximating integrals of functions of one variable; however, the accuracy of Clenshaw–Curtis points is similar to that of Gauss–Legendre points (cf. e.g. [32]), and nestedness allows for significant savings in the sparse grids construction. This distinction will also play a central role in the following sections.

We finally point out (see also [4]) that given any polynomial space $\mathbb{P}_{\Lambda}(\Gamma)$, one can always find a sparse grid that delivers approximations in that space, simply by taking $m(i) = i$ and $\mathcal{I} = \{\mathbf{i} \in \mathbb{N}_+^N : \mathbf{i} - \mathbf{1} \in \Lambda\}$. Conversely, given a sparse grid approximation $\mathcal{S}_{\mathcal{I}(w)}^m$, the underlying polynomial space is $\mathbb{P}_{\Lambda}(\Gamma)$ with $\Lambda = \{\mathbf{q} \in \mathbb{N}^N : \mathbf{q} \leq m(\mathbf{i}) - \mathbf{1} \text{ for some } \mathbf{i} \in \mathcal{I}\}$.

3 Quasi-Optimal sparse grid construction

We now summarize and slightly generalize a procedure for quasi-optimal sparse grid construction, using the approach introduced in our previous work [6]; see also [9, 19, 21]. We begin by introducing the concept of work (or computational cost) associated to the construction of a given sparse grid approximation of u . Assuming that the computation of the points of the sparse grid itself is negligible, the main cost of the sparse grid approximation is then given by the evaluation of the target function u at each point of the sparse grid, so that the required work is proportional to the total number of points used, $W_{\mathcal{I}(w),m}$. For notational convenience, we thus use the same symbol $W_{\mathcal{I}(w),m}$ to denote the work of the sparse grid approximation of u . Next, for each multi-index $\mathbf{i} \in \mathcal{I}(w)$ we introduce the *hierarchical surplus* operator

$$\Delta^{m(\mathbf{i})} = \bigotimes_{n=1}^N \Delta^{m(i_n)},$$

so that the sparse grid approximation (1) can actually be seen as a sum of hierarchical surplus operators applied to u . The quasi-optimal sparse grid relies on the concept of *profit* of a hierarchical surplus: to this end, we associate to each hierarchical surplus an *error contribution* and a *work contribution*, i.e. the contribution to the total error (respectively cost) that can be ascribed to a specific hierarchical surplus composing a sparse grid.

Associated to the hierarchical surplus $\Delta^{m(\mathbf{i})}$, we first introduce the quantity

$$\delta E(\mathbf{i}) = \left\| \Delta^{m(\mathbf{i})}[u] \right\|_{V \otimes L^2_\varrho(\Gamma)}, \quad (4)$$

and observe that, for any multi-index set \mathcal{J} such that $\mathbf{i} \notin \mathcal{J}$ and $\mathcal{J}, \{\mathcal{J} \cup \mathbf{i}\}$ are admissible according to condition (2), we have

$$(u - \mathcal{S}_{\{\mathcal{J} \cup \mathbf{i}\}}^m[u]) - (u - \mathcal{S}_{\mathcal{J}}^m[u]) = \mathcal{S}_{\{\mathcal{J} \cup \mathbf{i}\}}^m[u] - \mathcal{S}_{\mathcal{J}}^m[u] = \Delta^{m(\mathbf{i})}[u],$$

so that $\delta E(\mathbf{i})$ can be considered as a good indicator of the error reduction due to the addition of $\Delta^{m(\mathbf{i})}$ to any sparse grid approximation of u (in other words, $\delta E(\mathbf{i})$ is independent of \mathcal{J}). We can then naturally define as *error contribution* any upper bound $\Delta E(\mathbf{i})$ for $\delta E(\mathbf{i})$,

$$\delta E(\mathbf{i}) \leq \Delta E(\mathbf{i}),$$

Defining a work contribution $\Delta W(\mathbf{i})$ is instead a more delicate issue. Indeed, one could define the quantity $\delta W(\mathbf{i}) = W_{\{\mathcal{J} \cup \mathbf{i}\},m} - W_{\mathcal{J},m}$ as the work contribution of the hierarchical surplus $\Delta^{m(\mathbf{i})}$; however, such quantity $\delta W(\mathbf{i})$ does in general depend on the starting set \mathcal{J} unless nested points are used (see Example 1 below). Therefore, in the case of nested points we can safely define an “exact” *work contribution* $\Delta W(\mathbf{i})$ as

$$\Delta W(\mathbf{i}) = \prod_{n=1}^N (m(i_n) - m(i_n - 1)), \quad (5)$$

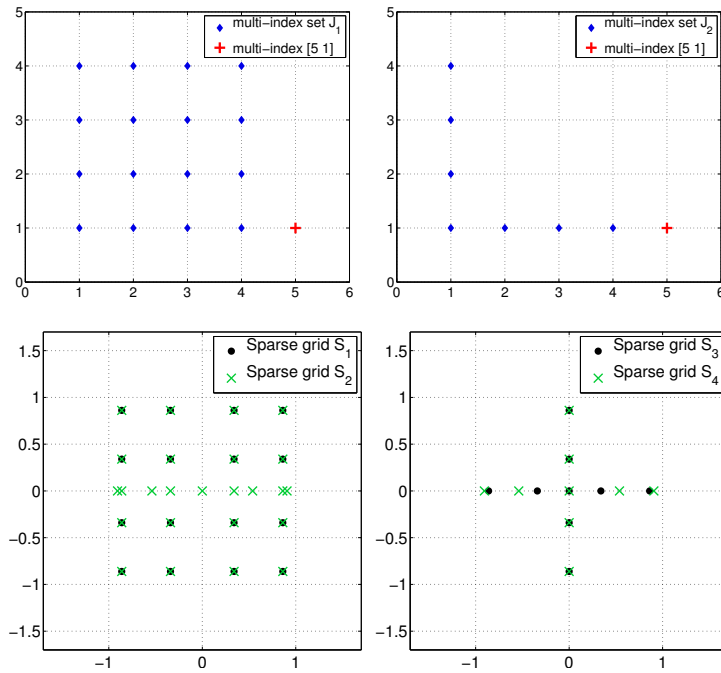


Fig. 2: The top row shows the two different multi-index sets considered in Example 1, i.e. \mathcal{J}_1 (left) and \mathcal{J}_2 (right), as well as the multi-index \mathbf{i} (red cross). The bottom row shows the resulting sparse grids $\mathcal{S}_1, \mathcal{S}_2, \mathcal{S}_3, \mathcal{S}_4$, corresponding to the sets $\mathcal{J}_1, \{\mathcal{J}_1 \cup \mathbf{i}\}, \mathcal{J}_2, \{\mathcal{J}_2 \cup \mathbf{i}\}$. The sparse grid \mathcal{S}_2 has 9 points more than \mathcal{S}_1 (left), while \mathcal{S}_3 and \mathcal{S}_4 have the same number of points (right).

which satisfies $W_{\mathcal{I}(w),m} = \sum_{\mathbf{i} \in \mathcal{I}(w)} \Delta W(\mathbf{i})$. On the other hand, in the case of non-nested points we consider the following definition of work contribution:

$$\Delta W(\mathbf{i}) = \prod_{n=1}^N m(i_n), \quad (6)$$

i.e. the cost of the tensor grid associated to \mathbf{i} , so that $W_{\mathcal{I}(w),m} \leq \sum_{\mathbf{i} \in \mathcal{I}(w)} \Delta W(\mathbf{i})$. This work contribution estimate is reasonable if one builds the (non nested) sparse grid “incrementally”, i.e. starting from $\mathcal{I} = \{\mathbf{1}\}$, adding one multi-index $\mathbf{i} \in \mathbb{N}_+^N$ to \mathcal{I} at a time and immediately evaluating the function u on the corresponding tensor grid $\bigotimes_{n=1}^N \mathcal{U}_n^{m(i_n)}$, cf. equation (3). By doing this, one does not exploit the fact that many tensor grids in the final formula (3) are multiplied by zero coefficients, and therefore $W_{\mathcal{I}(w),m} \leq \sum_{\mathbf{i} \in \mathcal{I}(w)} \Delta W(\mathbf{i})$.

Example 1 To show that $\delta W(\mathbf{i})$ is not uniquely defined when non-nested points are used, we take as an example the case of sparse grids built over $\Gamma = [-1, 1]^2$ using Gauss–Legendre points. We set $\mathbf{i} = (1, 5)$ and consider the multi-index sets \mathcal{J}_1 and \mathcal{J}_2 , shown in the top-left and top-right plots of Figure 2. We then

consider four different sparse grids: $\mathcal{S}_1, \mathcal{S}_2, \mathcal{S}_3, \mathcal{S}_4$, built using the multi-index sets $\mathcal{J}_1, \{\mathcal{J}_1 \cup \mathbf{i}\}, \mathcal{J}_2$, and $\{\mathcal{J}_2 \cup \mathbf{i}\}$ respectively.

Comparing \mathcal{S}_1 and \mathcal{S}_2 (bottom-left plot of Figure 2), we can see that adding \mathbf{i} to \mathcal{J}_1 results in a sparse grid with 9 additional points (i.e. $\delta W(\mathbf{i}) = 9$). Conversely, the comparison of \mathcal{S}_3 and \mathcal{S}_4 (bottom-right plot of Figure 2) shows that adding \mathbf{i} to \mathcal{J}_2 does not change the number of points of the sparse grid (i.e. $\delta W(\mathbf{i}) = 0$), since 4 new points are added but 4 points are no longer present.

Next, we introduce the *estimated profit* of a hierarchical surplus,

$$P(\mathbf{i}) = \frac{\Delta E(\mathbf{i})}{\Delta W(\mathbf{i})}, \quad (7)$$

and the sequence of decreasingly-ordered profits, $\{P_j^{ord}\}_{j \in \mathbb{N}_+}$

$$P_j^{ord} \geq P_{j+1}^{ord}.$$

It is also convenient to introduce a function that assigns the corresponding multi-index to the j -th ordered profit: we will denote such function as $\mathbf{i}(j)$, i.e. $P_j^{ord} = P(\mathbf{i}(j))$. Incidentally, note that as soon as two hierarchical surpluses have equal estimated profit, the map $\mathbf{i}(j)$ is not unique: in this case, any criterion to select a specific sequence can be used.

We can now define a quasi-optimal sparse grid at level w of approximation as the sparse grid including in the sum (1) only the w hierarchical surpluses with the highest profit (in the spirit of a knapsack problem), possibly made admissible according to condition (2):

$$\mathcal{I}(w) = \mathcal{J}(w)^{ADM}, \quad \mathcal{J}(w) = \{\mathbf{i}(1), \mathbf{i}(2), \dots, \mathbf{i}(w)\}. \quad (8)$$

Equivalently, one can automatically enforce the admissibility of the multi-index set by introducing the auxiliary profits (see also [11])

$$P^*(\mathbf{i}) = \max_{\mathbf{j} \geq \mathbf{i}} P(\mathbf{j}), \quad (9)$$

considering the decreasingly-ordered sequence $\{P_j^{*,ord}\}_{j \in \mathbb{N}_+}$

$$P_j^{*,ord} \geq P_{j+1}^{*,ord}, \quad (10)$$

and the new ordering $\mathbf{i}^*(j)$ such that $P_j^{*,ord} = P^*(\mathbf{i}^*(j))$ so that (8) can be rewritten as

$$\mathcal{I}(w) = \{\mathbf{i}^*(1), \mathbf{i}^*(2), \dots, \mathbf{i}^*(w)\}. \quad (11)$$

Definition 1 A sequence of profits $\{P(\mathbf{i})\}_{\mathbf{i} \in \mathbb{N}_+^N}$ is *monotone* if

$$\forall \mathbf{i} \in \mathbb{N}_+^N, \quad P^*(\mathbf{i}) = P(\mathbf{i}).$$

Notice that for a monotone sequence $\{P(\mathbf{i})\}_{\mathbf{i} \in \mathbb{N}_+^N}$, the set $\mathcal{J}(w) = \{\mathbf{i}(1), \mathbf{i}(2), \dots, \mathbf{i}(w)\}$ is always admissible (or downward closed).

The idea of constructing a sparse grid based on the profit of each hierarchical surplus (or other suitable “optimality indicators”) has been first proposed in a series of works [9, 19, 21], see also [23]. In particular, a possible approach could be an adaptive “greedy-type” algorithm in which the set \mathcal{I} is constructed iteratively: given a set $\mathcal{I}^{(k)}$ at the k -th iteration, one looks at the “neighbor” (or “margin”) $\mathcal{M}^{(k)}$ of $\mathcal{I}^{(k)}$ and adds to the set $\mathcal{I}^{(k)}$ the most profitable hierarchical surplus in $\mathcal{M}^{(k)}$,

$$\mathcal{I}^{(k+1)} = \mathcal{I}^{(k)} \cup \{\mathbf{i}\}, \quad \mathbf{i} = \operatorname{argmax}_{\mathbf{j} \in \mathcal{M}^{(k)}} P^*(\mathbf{j}),$$

see e.g. [19, 23] for a similar algorithm based however on optimality indicators other than profits. Clearly, this methodology implicitly assumes some kind of decay of the profits, or, equivalently, that the next most profitable multi-index always belongs to the margin of the current set $\mathcal{I}^{(k)}$. Moreover, the exploration of the margin $\mathcal{M}^{(k)}$ can be expensive in high dimensions. Therefore, in the context of elliptic PDEs with random coefficients, we proposed in [6] to add hierarchical surpluses based on a-priori error and work contribution estimates with numerically tuned constants (hybrid “a-priori”/“a-posteriori” estimates), that we observed numerically to be quite sharp and thus effective in reducing the sparse grid construction cost. Note that an analogous fully “a-priori” approach has been considered in [21, 9] in the context of wavelet-type approximation of high-dimensional deterministic PDEs.

4 Convergence estimate for the quasi-optimal sparse grid approximation

We now state and prove a convergence result of the quasi-optimal sparse grid approximation, built according to (11); see [31] for earlier versions and [29, 12] for an alternative estimate derived in the context of elliptic PDEs with random diffusion coefficients. We first recall a technical result, the so-called Stechkin Lemma, see e.g. [15] for a proof.

Lemma 1 (Stechkin) *Let $0 \leq p \leq q$, and let $\{a_j\}_{j \in \mathbb{N}_+}$ be a positive decreasing sequence. Then, for any $M > 0$*

$$\left(\sum_{j > M} (a_j)^q \right)^{1/q} \leq M^{-\frac{1}{p} + \frac{1}{q}} \left(\sum_{j \in \mathbb{N}_+} (a_j)^p \right)^{1/p}.$$

Theorem 1 (Quasi-optimal sparse grid convergence)

If the auxiliary profits (9) satisfy the weighted summability condition

$$\left(\sum_{\mathbf{i} \in \mathbb{N}_+^N} P^*(\mathbf{i})^\tau \Delta W(\mathbf{i}) \right)^{1/\tau} = C_P(\tau) < \infty \quad (12)$$

for some $0 < \tau \leq 1$, then

$$\left\| u - \mathcal{S}_{\mathcal{I}(w)}^m[u] \right\|_{V \otimes L^2_\varrho(\Gamma)} \leq W_{\mathcal{I}(w), m}^{1-1/\tau} C_P(\tau).$$

where $\mathcal{I}(w)$ is given by (11), and ΔW is given by (5) for grids with nested points and by (6) for grids with non nested points.

Proof We start by introducing the following auxiliary sequences:

- $\{\Delta W_j\}_{j \in \mathbb{N}_+}$ is the sequence of work contributions arranged using the same order as the sequence of the auxiliary profits (10). Note that this sequence will not be ordered in general.
- $\{Q_j\}_{j \in \mathbb{N}_+}$ is the sequence of the sum of the first j work contributions, i.e.

$$Q_0 = 0, \quad Q_j = \sum_{k=1}^j \Delta W_k.$$

- $\{\Delta E_j\}_{j \in \mathbb{N}_+}$ is the sequence of error contributions arranged using the same order as the sequence of the auxiliary profits (10). Again, this sequence will not be ordered in general.
- $\{\Delta \tilde{E}_k\}_{k \in \mathbb{N}_+}$ is a modification of the error contributions sequence $\{\Delta E_j\}_{j \in \mathbb{N}_+}$ just introduced, in which each ΔE_j is repeated a number of times equal to the corresponding work contribution. More precisely

$$\{\Delta \tilde{E}_k\}_{k \in \mathbb{N}_+} = \left\{ \underbrace{\Delta E_1, \Delta E_1, \dots}_{\Delta W_1 \text{ times}}, \underbrace{\Delta E_2, \Delta E_2, \dots}_{\Delta W_2 \text{ times}}, \dots \right\}$$

i.e. $\Delta \tilde{E}_{Q_{j-1}+s} = \Delta E_j$, $s = 1, \dots, \Delta W_j$.

- $\{\tilde{P}_k\}_{k \in \mathbb{N}_+}$ is the analogously modified sequence of auxiliary profits,

$$\{\tilde{P}_k\}_{k \in \mathbb{N}_+} = \left\{ \underbrace{\frac{\Delta E_1}{\Delta W_1}, \frac{\Delta E_1}{\Delta W_1}, \dots}_{\Delta W_1 \text{ times}}, \underbrace{\frac{\Delta E_2}{\Delta W_2}, \frac{\Delta E_2}{\Delta W_2}, \dots}_{\Delta W_2 \text{ times}}, \dots \right\} \quad (13)$$

i.e. $\tilde{P}_{Q_{j-1}+s} = P_j^{*, \text{ord}}$, $s = 1, \dots, \Delta W_j$.

For a generic sparse grid the following error decomposition holds:

$$\begin{aligned} \left\| u - \mathcal{S}_{\mathcal{I}(w)}^m[u] \right\|_{V \otimes L^2_\varrho(\Gamma)} &= \left\| \sum_{\mathbf{i} \notin \mathcal{I}(w)} \Delta^{m(\mathbf{i})}[u] \right\|_{V \otimes L^2_\varrho(\Gamma)} \\ &\leq \sum_{\mathbf{i} \notin \mathcal{I}(w)} \left\| \Delta^{m(\mathbf{i})}[u] \right\|_{V \otimes L^2_\varrho(\Gamma)} \leq \sum_{j > w} \Delta E_j. \end{aligned}$$

Next we recast the previous sum of error contributions in terms of the auxiliary sequence \tilde{P}_k in (13), i.e.

$$\sum_{j > w} \Delta E_j = \sum_{j > w} \sum_{s=1}^{\Delta W_j} \frac{\Delta \tilde{E}_{Q_{j-1}+s}}{\Delta W_j} = \sum_{k > Q_w} \tilde{P}_k.$$

We now apply Lemma 1 with $q = 1$ and $p = \tau$, obtaining

$$\left\| u - \mathcal{S}_{\mathcal{I}(w)}^m[u] \right\|_{V \otimes L^2_\rho(\Gamma)} \leq \sum_{k > Q_w} \tilde{P}_k \leq Q_w^{-1/\tau+1} \left(\sum_{k > 0} \tilde{P}_k^\tau \right)^{1/\tau},$$

and observe that

$$\left(\sum_{k > 0} \tilde{P}_k^\tau \right)^{1/\tau} = \left(\sum_{j > 0} \left(P_j^{*,ord} \right)^\tau \Delta W_j \right)^{1/\tau} = \left(\sum_{\mathbf{i} \in \mathbb{N}_+^N} P^*(\mathbf{i})^\tau \Delta W(\mathbf{i}) \right)^{1/\tau},$$

and

$$Q_w = \sum_{k=1}^w \Delta W_k \geq W_{\mathcal{I}(w),m},$$

due to the definitions of $\Delta W(\mathbf{i})$ in (5), (6). Then,

$$\left\| u - \mathcal{S}_{\mathcal{I}(w)}^m[u] \right\|_{V \otimes L^2_\rho(\Gamma)} \leq Q_w^{-1/\tau+1} \left(\sum_{\mathbf{i} \in \mathbb{N}_+^N} P^*(\mathbf{i})^\tau \Delta W(\mathbf{i}) \right)^{1/\tau} \leq C_P(\tau) W_{\mathcal{I}(w),m}^{1-1/\tau}.$$

which concludes the proof. \square

Remark 1 An alternative approach would consist in sorting the hierarchical surpluses according to error contribution estimates $\Delta E(\mathbf{i})$ rather than according to profits. Following the lines of the Theorem above, one could derive the following convergence estimate for the resulting sparse grid,

$$\left\| u - \mathcal{S}_{\mathcal{I}(w)}^m[u] \right\|_{V \otimes L^2_\rho(\Gamma)} \leq w^{1-1/\tau} C_E(\tau), \quad C_E(\tau) = \left(\sum_{\mathbf{i} \in \mathbb{N}_+^N} \Delta E(\mathbf{i})^\tau \right)^{1/\tau}.$$

However, recasting such estimate in terms of computational cost would not be possible without assumptions on the shape of the optimal set $\mathcal{I}(w)$.

5 Quasi-optimal sparse grid approximation of analytic functions on polyellipses

In this Section we apply the convergence Theorem 1 to a particular class of Hilbert-valued functions, which contains in particular the solution of some linear elliptic equations with random diffusion coefficients, as will be shown in Section 6. More precisely, we consider the class of functions $u : \Gamma = [1, 1]^N \rightarrow V$ that satisfies the following assumption

Assumption A1 (“Polyellipse analyticity”) *The function $u : \Gamma \rightarrow V$, $\Gamma = [-1, 1]^N$, admits a complex continuation $u^* : \mathbb{C}^N \rightarrow V$ that is a V -valued holomorphic function in any polyellipse*

$$\mathcal{E}_{\delta_1, \dots, \delta_N} = \prod_{n=1}^N \mathcal{E}_{n, \delta_n},$$

where $\mathcal{E}_{n, \delta_n}$ denotes the so-called Bernstein ellipse

$$\mathcal{E}_{n, \delta_n} = \left\{ z_n \in \mathbb{C} : \Re(z) \leq \frac{\delta_n + \delta_n^{-1}}{2} \cos \phi, \Im(z) \leq \frac{\delta_n - \delta_n^{-1}}{2} \sin \phi, \phi \in [0, 2\pi) \right\},$$

with $1 < \delta_n < \delta_n^*$, $n = 1, 2, \dots, N$. Moreover $\sup_{\mathbf{z} \in \mathcal{E}_{\delta_1, \dots, \delta_N}} \|u^*(\mathbf{z})\|_V \leq B_u$, and $B_u = B_u(\delta_1, \delta_2, \dots, \delta_N) \rightarrow \infty$ as $\delta_n \rightarrow \delta_n^*$, $n = 1, \dots, N$.

As already mentioned in the Introduction, sparse polynomial approximations are particularly suitable for this kind of functions, see e.g. [8]. We now state the convergence results for the quasi-optimal sparse grid approximation of functions satisfying Assumption A1, and devote the rest of this Section to their proof.

Definition 2 For a given family of collocation points, let $\mathbb{M}_n^{m(i_n)}$ be the norm of the interpolation operator $\mathcal{U}_n^{m(i_n)} : C^0(\Gamma_n) \rightarrow L^2_\varrho(\Gamma_n)$,

$$\mathbb{M}_n^{m(i_n)} = \sup_{\|f\|_{L^\infty(\Gamma_n)}=1} \left\| \mathcal{U}_n^{m(i_n)}[f] \right\|_{L^2_\varrho(\Gamma_n)}.$$

Definition 3 Given a level-to-nodes function $m(i_n)$, we define

$$d(i_n) = m(i_n) - m(i_n - 1).$$

Assumption A2 *For nested collocation points, there exists a constant $C_M > 0$ such that*

$$\frac{\mathbb{M}_n^{m(i_n)}}{d(i_n)} \leq C_M, \quad \forall i_n \in \mathbb{N},$$

while for non-nested collocation points, there exists a constant $C_M > 0$ such that

$$\frac{\mathbb{M}_n^{m(i_n)}}{m(i_n)} \leq C_M, \quad \forall i_n \in \mathbb{N}.$$

Assumption A3 *There exists a constant $C_m > 0$ such that there holds*

$$\frac{d(i_n + 1)}{d(i_n)} \leq C_m, \quad \forall i_n \in \mathbb{N}.$$

Remark 2 Assumption A3 allows for an exponential increase of the number of points from one level to another, of the type $m(i_n) \sim \gamma^{i_n}$. On the other hand, Assumption A2 requires that the operator norm $\mathbb{M}_n^{m(i_n)}$ grows slower than $d(i_n)$ for nested points and slower than $m(i_n)$ for non-nested points. Therefore, an exponential growth $m(i_n) \sim \gamma^{i_n}$ can accommodate even a polynomial growth of $\mathbb{M}_n^{m(i_n)}$, whereas a linear growth $m(i_n) \sim \alpha i_n$ requires a uniform boundedness of $\mathbb{M}_n^{m(i_n)}$.

Theorem 2 (nested quasi-optimal sparse grid convergence)

Let u be a function satisfying assumption A1. If the collocation points used are nested and satisfy Assumptions A2 and A3, then the quasi-optimal sparse grid approximation built according to the profits

$$P(\mathbf{i}) = C_E \prod_{n=1}^N \frac{e^{-g_n m(i_n-1)} \mathbb{M}_n^{m(i_n)}}{d(i_n)}, \quad g_n < \log \delta_n^*, \quad (14)$$

where C_E is a constant depending exponentially on N that will be specified in Lemma 6, converges as

$$\left\| u - \mathcal{S}_{\mathcal{I}(w)}^m[u] \right\|_{V \otimes L^2(\Gamma)} \leq \inf_{0 < \tau < 1} \mathcal{C}(N, \tau) W_{\mathcal{I}(w), m}^{1-1/\tau},$$

where

$$\mathcal{C}(N, \tau) = \left(C_E^\tau (C_M^\tau \widehat{C}_m)^N \prod_{n=1}^N \frac{e^{\tau g_n}}{e^{\tau g_n} - 1} \right)^{1/\tau}, \quad \widehat{C}_m = \max\{1, C_m\}.$$

Moreover, letting $g_m = \sqrt[N]{\prod_{n=1}^N g_n}$ denote the geometric mean of g_1, \dots, g_N , and assuming without loss of generality that $g_1 \leq g_2 \leq \dots \leq g_N$, there exist some constants $\alpha_L, \beta_L, C_{\log}$ that will be specified in Lemmas 4 and 5 such that

$$\left\| u - \mathcal{S}_{\mathcal{I}(w)}^m[u] \right\|_{V \otimes L^2(\Gamma)} \leq C_E C_M^N \exp \left(\left(C_{\log} - \frac{g_m \delta}{\widehat{C}_m e} \right) N \sqrt[N]{W_{\mathcal{I}(w), m}} \right), \quad (15)$$

that holds for $0 < \delta < 1 - \frac{1}{e}$ and for

$$W_{\mathcal{I}(w), m} > \left(\frac{g_N e \widehat{C}_m}{g_m (\alpha_L - \delta \beta_L)} \right)^N. \quad (16)$$

Theorem 3 (non-nested quasi-optimal sparse grid convergence)

Let u be a function satisfying assumption A1. If the collocation points used are non-nested and satisfy Assumptions A2 and A3, then the quasi-optimal sparse grids approximation built according to

$$P(\mathbf{i}) = C_E \prod_{n=1}^N \frac{e^{-g_n m(i_n-1)} \mathbb{M}_n^{m(i_n)}}{m(i_n)}, \quad g_n < \log \delta_n^*, \quad (17)$$

where C_E a constant depending exponentially on N that will be specified in Lemma 6, converges as

$$\left\| u - \mathcal{S}_{\mathcal{I}(w)}^m[u] \right\|_{V \otimes L^2_\varrho(\Gamma)} \leq \inf_{0 < \tau < 1} \mathcal{C}(N, \tau) W_{\mathcal{I}(w), m}^{1-1/\tau},$$

where

$$\mathcal{C}(N, \tau) = \left((C_E C_M^N)^\tau \prod_{n=1}^N \left(\widehat{C}_m \frac{e^{\tau g_n}}{e^{\tau g_n} - 1} + \frac{2}{\tau g_n} \frac{e^{\tau g_n/2}}{e^{\tau g_n/2} - 1} \right) \right)^{1/\tau}, \quad \widehat{C}_m = \max\{1, C_m\}.$$

Moreover, letting $g_m = \sqrt[N]{\prod_{n=1}^N g_n}$ denote the geometric mean of g_1, \dots, g_N , and assuming without loss of generality that $g_1 \leq g_2 \leq \dots \leq g_N$, there exist some constants $\alpha_L, \beta_L, C_{\log}$ that will be specified in Lemmas 4 and 5 such that

$$\left\| u - \mathcal{S}_{\mathcal{I}(w)}^m[u] \right\|_{V \otimes L^2_\varrho(\Gamma)} \leq C_E C_M^N \exp \left((2C_{\log} - \mathcal{K}g_m) N \sqrt[2N]{W_{\mathcal{I}(w), m}} \right), \quad (18)$$

with $\mathcal{K} = (\delta + 1 - \log 2)/(2\sqrt{e})$ that holds for $0 < \delta < 1 - \frac{1}{e}$ and for

$$W_{\mathcal{I}(w), m} > \max \left\{ \left(\frac{4\sqrt{e}g_n}{g_m(\alpha_L - \delta\beta_L)} \right)^{2N}, \left(\frac{\widehat{C}_m e^{3/2} g_n}{2g_m} \right)^{2N} \right\}. \quad (19)$$

The convergence estimates (15) and (18) just provided share the same structure of the result obtained for the optimal (“best M -terms”) L^2 approximation of u in [8, Theorem 16]. In particular:

1. they show that the convergence of the quasi-optimal sparse grid approximation is essentially sub-exponential, with a rate comparable to that of the optimal L^2 approximation in the case of nested points, and with half the rate for non-nested points. Such difference can be ascribed to the fact that the sparse grid construction on non-nested collocation points is not as efficient as its nested counterpart, and that the non-nested work contribution estimate is actually pessimistic.
2. Such convergence rate is obtained only after a sufficiently large amount of work (collocation points in this case, polynomial terms added in the expansion in the case of the optimal L^2 approximation) has been performed.
3. Both the convergence rate and the minimal work depend in a non trivial way on the choice of the parameters δ and C_{\log} .

A detailed discussion on the interplay between δ and C_{\log} can be found in [8, Remarks 17,19]; here we only mention that the two expressions in (19) are almost equivalent: for example, if the quadrature points are such that $\widehat{C}_m = 2$ (see Section 6.1), the first term is larger than the second one if $\frac{4}{e} \geq \alpha_L - \delta\beta_L$, i.e. if δ is larger than approximately 0.25, which is well inside the range of

feasible values of δ , cf. Lemma 4. In the numerical results Section we will show that a convergence estimate of the form

$$\left\| u - \mathcal{S}_{\mathcal{I}(w)}^m[u] \right\|_{V \otimes L^2_\varrho(\Gamma)} \leq \alpha \exp\left(-\beta N \sqrt[N]{W_{\mathcal{I}(w),m}}\right),$$

for nested points, and

$$\left\| u - \mathcal{S}_{\mathcal{I}(w)}^m[u] \right\|_{V \otimes L^2_\varrho(\Gamma)} \leq \alpha \exp\left(-\beta N \sqrt[2N]{W_{\mathcal{I}(w),m}}\right),$$

for non-nested points, with $\alpha \in \mathbb{R}$ and $\beta \in \mathbb{R}_+$, captures well the behavior of the computational error.

5.1 Preliminary results

We start by introducing the Chebyshev expansion of u (see e.g. [33]) and estimate the decay of its coefficients. To this end, we introduce the Chebyshev polynomials of the first kind $\Psi_q(t)$ on $[-1, 1]$, that are defined as the unique polynomials satisfying

$$\Psi_q(\cos \vartheta) = \cos(q\vartheta), \quad 0 \leq \vartheta \leq \pi, \quad q \in \mathbb{N}.$$

As a consequence, $|\Psi_q(t)| \leq 1$ on $[-1, 1]$, with $\Psi_q(1) = 1$ and $\Psi_q(-1) = (-1)^q$; moreover, they are orthogonal with respect to the weight $\rho_C(t) = 1/\sqrt{1-t^2}$:

$$\int_{-1}^1 \Psi_q(t) \Psi_\kappa(t) \rho_C(t) dt = \begin{cases} 0 & \kappa \neq q \\ \pi & \kappa = q = 0 \\ \pi/2 & \kappa = q \neq 0. \end{cases}$$

Lemma 2 *Let $\Psi_{q_n}(y_n)$ be the family of Chebyshev polynomials of the first kind on $\Gamma_n = [-1, 1]$, and let*

$$\Psi_{\mathbf{q}}(\mathbf{y}) = \prod_{n=1}^N \Psi_{q_n}(y_n), \quad \mathbf{q} = (q_1, q_2, \dots, q_N) \in \mathbb{N}^N.$$

be the generic N -variate Chebyshev polynomial. If the function u satisfies Assumption A1, then it admits a Chebyshev expansion

$$u(\mathbf{x}, \mathbf{y}) = \sum_{\mathbf{q} \in \mathbb{N}^N} u_{\mathbf{q}}(\mathbf{x}) \Psi_{\mathbf{q}}(\mathbf{y}),$$

with coefficients

$$u_{\mathbf{q}}(\mathbf{x}) = \frac{1}{\int_{\Gamma} \Psi_{\mathbf{q}}^2(\mathbf{y}) \varrho_C(\mathbf{y}) d\mathbf{y}} \int_{\Gamma} u(\mathbf{x}, \mathbf{y}) \Psi_{\mathbf{q}}(\mathbf{y}) \varrho_C(\mathbf{y}) d\mathbf{y},$$

which converges in $C^0(\Gamma, V)$, and whose coefficients $u_{\mathbf{q}}(\mathbf{x})$ are such that

$$\|u_{\mathbf{q}}\|_V \leq C_{Cheb} \prod_{n=1}^N e^{-g_n q_n}, \quad g_n = \log \delta_n \quad (20)$$

with $1 < \delta_n < \delta_n^*$, $C_{Cheb} = 2^{|\mathbf{q}|_0} B_u(\delta_1, \dots, \delta_N)$ where $\|\mathbf{q}\|_0$ denotes the number of non-zero elements of \mathbf{q} and $B_u(\delta_1, \dots, \delta_N)$ as in Assumption A1.

Proof The proof is a straightforward extension to the N -dimensional case of the argument in [16, Chapter 7, Theorem 8.1]; see also [3]. \square

Remark 3 An analogous bound could be proved for the decay of the coefficients of the Legendre expansion of u ,

$$\|u_{\mathbf{q}}\|_V \leq B_u(\delta_1, \dots, \delta_N) \prod_{n=1}^N \frac{r_n^{q_n}}{\tau_n \delta_n} \left(\sqrt{1 - r_n^2} + \mathcal{O}\left(\frac{1}{q_n^{1/3}}\right) \right),$$

with τ_n to be chosen in $(0, 1)$ and

$$r_n = r_n(\tau_n, \delta_n) = \frac{1}{1 + \delta_n(1 - \tau_n) + \sqrt{\delta_n^2(1 - \tau_n)^2 + 2\delta_n(1 - \tau_n)}},$$

see [2, 22]. Hence, the same analysis presented in this work could still be performed using the Legendre expansion of u instead of the Chebyshev one.

Remark 4 Lemma 2 and Remark 3 state that the convergence of the coefficients of both Chebyshev and Legendre expansions is essentially exponential with respect to the degree of approximation of each parameter. Our numerical experience shows that such bound is actually sharp, at least for the inclusion problem that will be discussed in Section 6; see also [8] for more examples on the Legendre expansion.

We close this Section with some technical lemmas which will be needed in the following analysis.

Lemma 3 *If $\mathcal{U}_n^{m(i_n)}$ is built over Gaussian abscissas, then $\mathbb{M}_n^{m(i_n)} = 1$.*

Proof Let $\mathcal{Q}_n^{m(i_n)}$ be the quadrature rule built over the same abscissas used for $\mathcal{U}_n^{m(i_n)}$,

$$\mathcal{Q}_n^{m(i_n)}[f] = \sum_{j=1}^{m(i_n)} f(\alpha_j) \omega_j.$$

Observe that since the considered abscissas are Gaussian, $\mathcal{Q}_n^{m(i_n)}$ is exact for polynomials of degree $2m(i_n) - 1$, and in particular $\sum_{j=1}^{m(i_n)} \omega_j = 1$; furthermore, the quadrature weights ω_j are positive. Next, observe that $\left(\mathcal{U}_n^{m(i_n)}[f(t)]\right)^2$

is a polynomial of degree $2(m(i_n) - 1)$: therefore, using the fact that $\mathcal{U}_n^{m(i_n)}$ is a Lagrangian interpolant, we have

$$\begin{aligned} \int_{\Gamma_n} \left(\mathcal{U}_n^{m(i_n)}[f] \right)^2 \varrho_n(t) dt &= \mathcal{Q} \left[\left(\mathcal{U}_n^{m(i_n)}[f] \right)^2 \right] = \sum_{j=1}^{m(i_n)} \left(\mathcal{U}_n^{m(i_n)}[f](\alpha_j) \right)^2 \omega_j \\ &= \sum_{j=1}^{m(i_n)} f^2(\alpha_j) \omega_j \leq \|f^2\|_{L^\infty(\Gamma_n)}, \end{aligned}$$

and this finishes the proof. \square

Lemma 4 For $0 < \epsilon < \frac{e-1}{e} = \epsilon_{max} \approx 0.63$, there holds

$$\frac{1}{1 - e^{-x}} \leq \frac{(1 - \epsilon)e}{x}, \quad 0 < x \leq x_{cr}(\epsilon).$$

Moreover, the function $x_{cr}(\epsilon)$ is concave and can be bounded from below as

$$\alpha_L - \beta_L \epsilon \leq x_{cr}(\epsilon), \quad 0 < \epsilon < \epsilon_{max}$$

with $\alpha_L \approx 2.49$ and $\beta_L = \alpha_L / \epsilon_{max}$.

Proof See [8, Lemma 13]. \square

Lemma 5 Given any $C_{log} \in (0, 1/e]$, and denoting by \bar{t} the largest root of the equation $\log t = C_{log} t$, then $\forall K > 0$ there holds

$$M \leq e^{C_{log} K} \sqrt[K]{M} \quad \forall M > \bar{t}^K, \quad \forall K > 0,$$

If $C_{log} = 1/e$, the bound holds for any $M > 0$.

Proof See [8, Lemma 14]. \square

5.2 Estimates of hierarchical surplus error contributions

Lemma 6 If u satisfies Assumption A1, then for each hierarchical surplus operator $\Delta^{m(i)}$ there holds

$$\delta E(\mathbf{i}) \leq \Delta E(\mathbf{i}) = C_E e^{-\sum_{n=1}^N g_n m(i_n-1)} \prod_{n=1}^N \mathbb{M}_n^{m(i_n)}.$$

with g_n as in Lemma 2 and $C_E = 2^N C_{Cheb} \prod_{n=1}^N \frac{1}{1 - e^{-g_n}}$.

Proof Let us consider again the Chebyshev expansion of u . From the definition (4) of $\delta E(\mathbf{i})$ we have

$$\begin{aligned}\delta E(\mathbf{i}) &= \left\| \Delta^{m(\mathbf{i})}[u] \right\|_{V \otimes L^2_{\mathbf{q}}(\Gamma)} = \left\| \Delta^{m(\mathbf{i})} \left[\sum_{\mathbf{q} \in \mathbb{N}^N} u_{\mathbf{q}} \Psi_{\mathbf{q}} \right] \right\|_{V \otimes L^2_{\mathbf{q}}(\Gamma)} \\ &= \left\| \sum_{\mathbf{q} \in \mathbb{N}^N} u_{\mathbf{q}} \Delta^{m(\mathbf{i})}[\Psi_{\mathbf{q}}] \right\|_{V \otimes L^2_{\mathbf{q}}(\Gamma)}\end{aligned}$$

Observe now that by construction of hierarchical surplus there holds $\Delta^{m(\mathbf{i})}[\Psi_{\mathbf{q}}] = 0$ for all Chebyshev polynomials $\Psi_{\mathbf{q}}$ such that $\exists n : q_n < m(i_n - 1)$. Therefore, the previous sum reduces to the multi-index set $\mathbf{q} \geq m(\mathbf{i} - \mathbf{1})$, and we have

$$\begin{aligned}\delta E(\mathbf{i}) &= \left\| \sum_{\mathbf{q} \geq m(\mathbf{i} - \mathbf{1})} u_{\mathbf{q}} \Delta^{m(\mathbf{i})}[\Psi_{\mathbf{q}}] \right\|_{V \otimes L^2_{\mathbf{q}}(\Gamma)} \\ &\leq \sum_{\mathbf{q} \geq m(\mathbf{i} - \mathbf{1})} \|u_{\mathbf{q}}\|_V \left\| \Delta^{m(\mathbf{i})}[\Psi_{\mathbf{q}}] \right\|_{L^2_{\mathbf{q}}(\Gamma)}.\end{aligned}$$

Next, using the definition of $\Delta^{m(\mathbf{i})}$ we bound

$$\begin{aligned}\left\| \Delta^{m(\mathbf{i})}[\Psi_{\mathbf{q}}] \right\|_{L^2_{\mathbf{q}}(\Gamma)} &= \prod_{n=1}^N \left\| \Delta^{m(i_n)}[\Psi_{q_n}] \right\|_{L^2_{q_n}(\Gamma_n)} \\ &\leq \prod_{n=1}^N 2 \mathbb{M}_n^{m(i_n)} \|\Psi_{q_n}\|_{L^\infty(\Gamma_n)} = \prod_{n=1}^N 2 \mathbb{M}_n^{m(i_n)}.\end{aligned}$$

Recalling estimate (20) for the decay of the Chebyshev coefficients of u , one obtains

$$\begin{aligned}\delta E(\mathbf{i}) &\leq \sum_{\mathbf{q} \geq m(\mathbf{i} - \mathbf{1})} \|u_{\mathbf{q}}\|_V \prod_{n=1}^N 2 \mathbb{M}_n^{m(i_n)} \leq 2^N \sum_{\mathbf{q} \geq m(\mathbf{i} - \mathbf{1})} C_{Cheb} \prod_{n=1}^N e^{-g_n q_n} \mathbb{M}_n^{m(i_n)} \\ &\leq 2^N C_{Cheb} \prod_{n=1}^N \mathbb{M}_n^{m(i_n)} \sum_{q_n \geq m(i_n - 1)} e^{-g_n q_n} \leq 2^N C_{Cheb} \prod_{n=1}^N \mathbb{M}_n^{m(i_n)} \frac{e^{-g_n m(i_n - 1)}}{1 - e^{-g_n}}.\end{aligned}$$

□

Remark 5 This bound had been already proposed without proof in [6], using the norm of the interpolation operator $\mathcal{U}_n^{m(i_n)} : C^0(\Gamma_n) \rightarrow L^\infty(\Gamma_n)$, i.e. the standard Lebesgue constant associated to $\mathcal{U}_n^{m(i_n)}$, instead of $\mathbb{M}_n^{m(i_n)}$.

5.3 Convergence result: nested case

We now focus on the case of nested sequences of collocation points. Observe that the profits (14) are derived by the profit definition (7) combining the work contribution (5) and Lemma 6.

Lemma 7 *Under Assumption A2, the auxiliary profits*

$$P^b(\mathbf{i}) = C_E C_M^N \prod_{n=1}^N e^{-g_n m(i_n-1)}$$

are such that

$$P(\mathbf{i}) \leq P^b(\mathbf{i}), \quad \forall \mathbf{i} \in \mathbb{N}_+^N, \quad (21)$$

where $P(\mathbf{i})$ are the profits in (14). Moreover, the sequence $\{P^b(\mathbf{i})\}_{\mathbf{i} \in \mathbb{N}_+^N}$ is monotone according to Definition 1, i.e.

$$P^{b,*}(\mathbf{i}) = \max_{\mathbf{j} \geq \mathbf{i}} P^b(\mathbf{j}) = P^b(\mathbf{i}), \quad \forall \mathbf{i} \in \mathbb{N}_+^N$$

and, under Assumption A3, it satisfies the weighted τ -summability condition (12) for every $0 < \tau < 1$. In particular, there holds

$$\sum_{\mathbf{i} \in \mathbb{N}_+^N} P^b(\mathbf{i})^\tau \Delta W(\mathbf{i}) \leq C_E^\tau (C_M^\tau \widehat{C}_m)^N \prod_{n=1}^N \frac{e^{\tau g_n}}{e^{\tau g_n} - 1},$$

with C_E as in Lemma 6 and $\widehat{C}_m = \max\{1, C_m\}$.

Proof Inequality (21) follows from Assumption A2, while the fact that the sequence $\{P^b(\mathbf{i})\}_{\mathbf{i} \in \mathbb{N}_+^N}$ is monotone is a straightforward consequence of its definition. As for the summability property, we start by observing that we can actually write the weighted sum $\sum_{\mathbf{i} \in \mathbb{N}_+^N} P^b(\mathbf{i})^\tau \Delta W(\mathbf{i})$ as a product of series depending on i_n only,

$$\sum_{\mathbf{i} \in \mathbb{N}_+^N} C_E^\tau \prod_{n=1}^N \left[\left(C_M e^{-g_n m(i_n-1)} \right)^\tau d(i_n) \right] = C_E^\tau C_M^{\tau N} \prod_{n=1}^N \sum_{i_n=1}^{\infty} \left(e^{-g_n m(i_n-1)} \right)^\tau d(i_n), \quad (22)$$

so that we only need to study the summability of

$$\mathbb{S}_n = \sum_{i_n=1}^{\infty} \left(e^{-g_n m(i_n-1)} \right)^\tau d(i_n), \quad n = 1, \dots, N.$$

We begin by taking out of the sum the term for $i_n = 1$ and using Assumption A3

$$\mathbb{S}_n = 1 + \sum_{i_n=2}^{\infty} e^{-\tau g_n m(i_n-1)} d(i_n) \leq 1 + C_m \sum_{i_n=2}^{\infty} e^{-\tau g_n m(i_n-1)} d(i_n - 1). \quad (23)$$

Next, observe that

$$\begin{aligned} e^{-\tau g_n m(i_n-1)} d(i_n - 1) &= e^{-\tau g_n m(i_n-1)} \left(m(i_n - 1) - m(i_n - 2) \right) \\ &\leq \sum_{j_n=m(i_n-2)+1}^{m(i_n-1)} e^{-j_n \tau g_n}, \end{aligned}$$

so that

$$\sum_{i_n=2}^{\infty} e^{-\tau g_n m(i_n-1)} d(i_n-1) \leq \sum_{j_n=1}^{\infty} e^{-\tau g_n j_n}. \quad (24)$$

Therefore, going back to (23) we obtain

$$\mathbb{S}_n \leq 1 + C_m \sum_{i_n=1}^{\infty} e^{-\tau g_n i_n} \leq \max\{1, C_m\} \sum_{i_n=0}^{\infty} e^{-\tau g_n i_n} = \widehat{C}_m \frac{e^{\tau g_n}}{e^{\tau g_n} - 1}.$$

and the proof is concluded by substituting this bound in (22). \square

We are now ready to give the full proof of Theorem 2.

Proof (of Theorem 2) In the case of nested collocation points, the quasi-optimal sparse grid is built using the profits (14). We first observe that, due to Lemma 7, there holds

$$P^*(\mathbf{i}) = \max_{\mathbf{j} \geq \mathbf{i}} P(\mathbf{i}) \leq \max_{\mathbf{j} \geq \mathbf{i}} P^b(\mathbf{i}) = P^b(\mathbf{i}).$$

Therefore, the profits $P^*(\mathbf{i})$ have (at least) the same τ -summability properties than $P^b(\mathbf{i})$, and thus from Theorem 1 we have

$$\begin{aligned} \left\| u - \mathcal{S}_{\mathcal{I}(w)}^m[u] \right\|_{V \otimes L^2_{\mathfrak{q}}(\Gamma)} &\leq W_{\mathcal{I}(w), m}^{1-1/\tau} \left(\sum_{\mathbf{i} \in \mathbb{N}_+^N} P^*(\mathbf{i})^\tau \Delta W(\mathbf{i}) \right)^{1/\tau} \\ &\leq W_{\mathcal{I}(w), m}^{1-1/\tau} \left(\sum_{\mathbf{i} \in \mathbb{N}_+^N} P^b(\mathbf{i})^\tau \Delta W(\mathbf{i}) \right)^{1/\tau}. \end{aligned}$$

Now, due to Lemma 7, the profits P^b satisfy the weighted τ -summability condition for every τ in $(0, 1)$, hence we can use any τ in the range $0 < \tau < 1$ to obtain a valid bound for the sparse grid error. Thus, we can choose the smallest bound, by minimizing the error estimate over τ : to this end, we follow closely the argument in [8, Theorem 16]. We have

$$\left\| u - \mathcal{S}_{\mathcal{I}(w)}^m[u] \right\|_{V \otimes L^2_{\mathfrak{q}}(\Gamma)} \leq W_{\mathcal{I}(w), m}^{1-1/\tau} \left(C_E^\tau (C_M^\tau \widehat{C}_m)^N \prod_{n=1}^N \frac{e^{\tau g_n}}{e^{\tau g_n} - 1} \right)^{1/\tau}, \quad (25)$$

and we want to minimize the right-hand side with respect to τ . We do not solve this minimization problem exactly, but rather take $\tau = e\mathcal{K}$, with $\mathcal{K}^N = \frac{\widehat{C}_m^N}{W_{\mathcal{I}(w), m} \prod_{n=1}^N g_n}$. The motivation for this choice is the following: if τ is small,

we can approximate $\frac{e^{\tau g_n}}{e^{\tau g_n} - 1} \approx \frac{1}{\tau g_n}$ and rewrite the right-hand side of (25) as

$$\begin{aligned} W_{\mathcal{I}(w),m}^{1-1/\tau} \left(C_E^\tau (C_M^\tau \widehat{C}_m)^N \prod_{n=1}^N \frac{e^{\tau g_n}}{e^{\tau g_n} - 1} \right)^{1/\tau} &\approx W_{\mathcal{I}(w),m}^{1-1/\tau} \left(C_E^\tau (C_M^\tau \widehat{C}_m)^N \frac{1}{\tau^N \prod_{n=1}^N g_n} \right)^{1/\tau} \\ &= W_{\mathcal{I}(w),m} C_E C_M^N \left(\frac{\widehat{C}_m^N}{\tau^N W_{\mathcal{I}(w),m} \prod_{n=1}^N g_n} \right)^{1/\tau} \\ &= W_{\mathcal{I}(w),m} C_E C_M^N \left(\frac{\mathcal{K}^N}{\tau^N} \right)^{1/\tau}. \end{aligned}$$

Therefore, we enforce

$$\begin{aligned} 0 &= \frac{d}{d\tau} \left(\frac{\mathcal{K}^N}{\tau^N} \right)^{1/\tau} = \frac{d}{d\tau} \exp \left(-\frac{1}{\tau} \log \frac{\tau^N}{\mathcal{K}^N} \right) \\ &= \exp \left(-\frac{1}{\tau} \log \frac{\tau^N}{\mathcal{K}^N} \right) \left[\frac{1}{\tau^2} \log \frac{\tau^N}{\mathcal{K}^N} - \frac{1}{\tau} \frac{N}{\tau} \right] \end{aligned}$$

that is

$$0 = \frac{N}{\tau^2} \left[\log \frac{\tau}{\mathcal{K}} - 1 \right], \quad (26)$$

resulting in $\tau = e\mathcal{K}$. We now insert this choice of τ in the original bound (25) obtaining

$$\left\| u - \mathcal{S}_{\mathcal{I}(w)}^m[u] \right\|_{V \otimes L^2_\epsilon(\Gamma)} \leq W_{\mathcal{I}(w),m} C_E C_M^N \left(\frac{\widehat{C}_m^N}{W_{\mathcal{I}(w),m}} \prod_{n=1}^N \frac{e^{g_n e\mathcal{K}}}{e^{g_n e\mathcal{K}} - 1} \right)^{1/\tau}.$$

Next, we bound each of the factors $e^{e g_n \mathcal{K}} / (e^{e g_n \mathcal{K}} - 1)$ by Lemma 4 (with $x = e g_n \mathcal{K}$), obtaining

$$\begin{aligned} \left\| u - \mathcal{S}_{\mathcal{I}(w)}^m[u] \right\|_{V \otimes L^2_\epsilon(\Gamma)} &\leq W_{\mathcal{I}(w),m} C_E C_M^N \left(\frac{\widehat{C}_m^N}{W_{\mathcal{I}(w),m}} \prod_{n=1}^N (1 - \epsilon_n) \frac{e g_m \sqrt[N]{W_{\mathcal{I}(w),m}}}}{e g_n \widehat{C}_m} \right)^{1/(e\mathcal{K})}, \\ &= W_{\mathcal{I}(w),m} C_E C_M^N \left(\prod_{n=1}^N (1 - \epsilon_n) \right)^{1/(e\mathcal{K})}. \end{aligned} \quad (27)$$

Note that the latter equation holds true for ϵ_n and $W_{\mathcal{I}(w),m}$ satisfying $e g_n \mathcal{K} \leq x_{cr}(\epsilon_n)$ (cf. again Lemma 4), which in turn is satisfied if we choose ϵ_n using the lower bound in Lemma 4, i.e.

$$e g_n \mathcal{K} = \frac{e g_n \widehat{C}_m}{\sqrt[N]{W_{\mathcal{I}(w),m} g_m}} = \alpha_L - \beta_L \epsilon_n \Rightarrow \epsilon_n = \left(\alpha_L - \frac{g_n e \widehat{C}_m}{g_m \sqrt[N]{W_{\mathcal{I}(w),m}}} \right) \frac{1}{\beta_L}.$$

Moreover, we also have to enforce $\epsilon_n > 0$ to ensure convergence of estimate (27), thus obtaining a constraint on $W_{\mathcal{I}(w),m}$. Namely, taken any $0 < \delta < \epsilon_{max}$ we require $\epsilon_n > \delta$, which implies

$$\delta < \left(\alpha_L - \frac{g_n e \widehat{C}_m}{g_m \sqrt[N]{W_{\mathcal{I}(w),m}}} \right) \frac{1}{\beta_L} \quad \Rightarrow \quad W_{\mathcal{I}(w),m} > \left(\frac{g_n e \widehat{C}_m}{g_m (\alpha_L - \delta \beta_L)} \right)^N.$$

Since we have assumed that the coefficients g_n are ordered increasingly, this condition has to be checked for $n = N$ only, hence (16). With this choice of $\epsilon_{W_{\mathcal{I}(w),m},n}$, equation (27) further simplifies to

$$\begin{aligned} \left\| u - \mathcal{S}_{\mathcal{I}(w)}^m[u] \right\|_{V \otimes L_2^2(\Gamma)} &\leq W_{\mathcal{I}(w),m} C_E C_M^N \exp \left(\sum_{n=1}^N \log(1 - \epsilon_n) \right)^{1/(e\mathcal{K})} \\ &\leq W_{\mathcal{I}(w),m} C_E C_M^N \exp \left(- \sum_{n=1}^N \epsilon_n \right)^{1/(e\mathcal{K})} \\ &\leq W_{\mathcal{I}(w),m} C_E C_M^N \exp \left(- \frac{N\delta}{e\mathcal{K}} \right) \end{aligned}$$

and the final result follows by recalling the definition of \mathcal{K} and using Lemma 5. \square

5.4 Convergence result: non-nested case

We now focus on the case of non-nested sequences of collocation points. Observe that the profits (17) are derived by the profit definition (7) combining the work contribution (6) and Lemma 6.

Lemma 8 *Under Assumption A2, the auxiliary profits*

$$P^b(\mathbf{i}) = C_E C_M^N \prod_{n=1}^N e^{-g_n m(i_n-1)},$$

are such that

$$P(\mathbf{i}) \leq P^b(\mathbf{i}), \quad \forall \mathbf{i} \in \mathbb{N}_+^N, \quad (28)$$

where $P(\mathbf{i})$ are the profits (17). Moreover, the sequence $\{P^b(\mathbf{i})\}_{\mathbf{i} \in \mathbb{N}_+^N}$ is monotone according to Definition 1, i.e.

$$P^{b,*}(\mathbf{i}) = \max_{\mathbf{j} \geq \mathbf{i}} P^b(\mathbf{j}) = P^b(\mathbf{i}), \quad \forall \mathbf{i} \in \mathbb{N}_+^N$$

and, under Assumption A3, it satisfies the weighted τ -summability condition (12) for every $0 < \tau < 1$. In particular, there holds

$$\sum_{\mathbf{i} \in \mathbb{N}_+^N} P^b(\mathbf{i})^\tau \Delta W(\mathbf{i}) \leq (C_E C_M^N)^\tau \prod_{n=1}^N \left(\widehat{C}_m \frac{e^{\tau g_n}}{e^{\tau g_n} - 1} + \frac{2}{\tau g_n e} \frac{e^{\tau g_n/2}}{e^{\tau g_n/2} - 1} \right).$$

Proof Inequality (28) follows from Assumption A2, while the fact that $\{P^b(\mathbf{i})\}_{\mathbf{i} \in \mathbb{N}_+^N}$ is monotone is a straightforward consequence of its definition. As for the summability property, we proceed as in the proof of Lemma 7, and rewrite

$$\begin{aligned} \sum_{\mathbf{i} \in \mathbb{N}_+^N} P^b(\mathbf{i})^\tau \Delta W(\mathbf{i}) &= \sum_{\mathbf{i} \in \mathbb{N}_+^N} C_E^\tau \prod_{n=1}^N \left[\left(C_M e^{-g_n m(i_n-1)} \right)^\tau m(i_n) \right] \\ &= C_E^\tau C_M^{\tau N} \prod_{n=1}^N \sum_{i_n=1}^{\infty} \left(e^{-g_n m(i_n-1)} \right)^\tau m(i_n), \end{aligned} \quad (29)$$

to study the summability of

$$\mathbb{S}_n = \sum_{i_n=1}^{\infty} \left(e^{-g_n m(i_n-1)} \right)^\tau m(i_n), \quad n = 1, \dots, N.$$

We split the sum as

$$\begin{aligned} \mathbb{S}_n &= 1 + \sum_{i_n=2}^{\infty} e^{-\tau g_n m(i_n-1)} m(i_n) \\ &= 1 + \sum_{i_n=2}^{\infty} e^{-\tau g_n m(i_n-1)} d(i_n) + \sum_{i_n=2}^{\infty} e^{-\tau g_n m(i_n-1)} m(i_n - 1) \end{aligned}$$

and we consider the two sums separately. The first one can be bounded as in Lemma 7, equations (23)-(24),

$$\sum_{i_n=2}^{\infty} e^{-\tau g_n m(i_n-1)} d(i_n) \leq C_m \sum_{i_n=1}^{\infty} e^{-\tau g_n i_n},$$

while the second one can be bounded as

$$\sum_{i_n=2}^{\infty} e^{-\tau g_n m(i_n-1)} m(i_n - 1) \leq \frac{2}{\tau g_n e} \sum_{i_n=2}^{\infty} e^{-\tau g_n m(i_n-1)/2}.$$

exploiting the elementary fact that for every $x > 0$ and for every $\epsilon > 0$, there holds $x \leq \frac{1}{\epsilon e} e^{\epsilon x}$. Combining the two bounds, we obtain

$$\begin{aligned} \mathbb{S}_n &\leq 1 + C_m \sum_{i_n=1}^{\infty} e^{-\tau g_n i_n} + \frac{2}{\tau g_n e} \sum_{i_n=2}^{\infty} e^{-\tau g_n m(i_n-1)/2} \\ &\leq \max\{1, C_m\} \sum_{i_n=0}^{\infty} e^{-\tau g_n i_n} + \frac{2}{\tau g_n e} \sum_{i_n=0}^{\infty} e^{-\tau g_n i_n/2} \\ &\leq \widehat{C}_m \frac{e^{\tau g_n}}{e^{\tau g_n} - 1} + \frac{2}{\tau g_n e} \frac{e^{\tau g_n/2}}{e^{\tau g_n/2} - 1}, \end{aligned}$$

and the proof is concluded by substituting this bound into (29). \square

We are now ready to give the full proof of Theorem 3.

Proof (of Theorem 3) In the case of non-nested collocation points, the quasi-optimal sparse grid is built using the profits (17). The proof is analogous to that of Theorem 2. Due to Lemma 8, there holds

$$P^*(\mathbf{i}) = \max_{\mathbf{j} \geq \mathbf{i}} P(\mathbf{i}) \leq \max_{\mathbf{j} \geq \mathbf{i}} P^b(\mathbf{i}) = P^b(\mathbf{i}),$$

so that from Theorem 1 and Lemma 8 we have

$$\begin{aligned} \left\| u - \mathcal{S}_{\mathcal{I}(w)}^m[u] \right\|_{V \otimes L^2_g(\Gamma)} &\leq W_{\mathcal{I}(w),m}^{1-1/\tau} \left(\sum_{\mathbf{i} \in \mathbb{N}_+^N} P^b(\mathbf{i})^\tau \Delta W_{\mathbf{j}} \right)^{1/\tau} \\ &\leq W_{\mathcal{I}(w),m}^{1-1/\tau} \left((C_E C_M^N)^\tau \prod_{n=1}^N \left(\hat{C}_m \frac{e^{\tau g_n}}{e^{\tau g_n} - 1} + \frac{2}{\tau g_n e} \frac{e^{\tau g_n/2}}{e^{\tau g_n/2} - 1} \right) \right)^{1/\tau}, \end{aligned} \quad (30)$$

to be minimized with respect to τ . Next, we suppose τ to be small, so that

$$\hat{C}_m \frac{e^{\tau g_n}}{e^{\tau g_n} - 1} + \frac{2}{\tau g_n e} \frac{e^{\tau g_n/2}}{e^{\tau g_n/2} - 1} \approx \frac{\hat{C}_m}{\tau g_n} + \frac{4}{(\tau g_n)^2 e} \approx \frac{4}{(\tau g_n)^2 e},$$

and

$$\left\| u - \mathcal{S}_{\mathcal{I}(w)}^m[u] \right\|_{V \otimes L^2_g(\Gamma)} \approx W_{\mathcal{I}(w),m}^{1-1/\tau} C_E C_M^N \left(\frac{4^N}{e^N \tau^{2N} \prod_{n=1}^N g_n^2} \right)^{1/\tau} = W_{\mathcal{I}(w),m} C_E C_M^N \left(\frac{\mathcal{K}^N}{\tau^{2N}} \right)^{1/\tau},$$

with $\mathcal{K}^N = \frac{4^N}{e^N W_{\mathcal{I}(w),m} g_m^{2N}}$. With some calculus analogous to (26), we then obtain

$$\frac{d}{d\tau} \left(\frac{\mathcal{K}^N}{\tau^{2N}} \right)^{1/\tau} = 0 \Leftrightarrow \log \frac{\tau^2}{\mathcal{K}} = 2 \Leftrightarrow \tau = e\sqrt{\mathcal{K}}.$$

We now go back to the original bound (30) and apply Lemma 4, obtaining

$$\begin{aligned} \left\| u - \mathcal{S}_{\mathcal{I}(w)}^m[u] \right\|_{V \otimes L^2_g(\Gamma)} &\leq W_{\mathcal{I}(w),m} C_E C_M^N \left(\frac{1}{W_{\mathcal{I}(w),m}} \prod_{n=1}^N \left(\hat{C}_m \frac{(1-\epsilon_n)e}{\tau g_n} + \frac{2}{\tau g_n e} \frac{(1-\epsilon_n)e}{\tau g_n/2} \right) \right)^{1/\tau}, \\ &= W_{\mathcal{I}(w),m} C_E C_M^N \left(\frac{1}{W_{\mathcal{I}(w),m}} \prod_{n=1}^N \frac{1-\epsilon_n}{\tau g_n} \left(\hat{C}_m e + \frac{4}{\tau g_n} \right) \right)^{1/\tau} \end{aligned} \quad (31)$$

that holds for $\tau g_n/2 \leq x_{cr}(\epsilon_n)$. This condition can be satisfied using the lower bound in Lemma 4, i.e. by choosing ϵ_n such that

$$\frac{e\sqrt{\mathcal{K}}g_n}{2} = \frac{4g_n\sqrt{e}}{g_m^{2N}\sqrt{W_{\mathcal{I}(w),m}}} = \alpha_L - \beta_L \epsilon_n \Rightarrow \epsilon_n = \left(\alpha_L - \frac{2g_n\sqrt{e}}{g_m^{2N}\sqrt{W_{\mathcal{I}(w),m}}} \right) \frac{1}{\beta_L}.$$

Note also that we will need $\epsilon_n > 0$ to ensure convergence of the estimate; namely, taken any $0 < \delta < \epsilon_{max}$ we require $\epsilon_n > \delta$, which implies

$$\delta < \left(\alpha_L - \frac{4g_n\sqrt{e}}{g_m^{2N}\sqrt{W_{\mathcal{I}(w),m}}} \right) \frac{1}{\beta_L} \quad \Rightarrow \quad W_{\mathcal{I}(w),m} > \left(\frac{4\sqrt{e}g_n}{g_m(\alpha_L - \delta\beta_L)} \right)^{2N}.$$

Moreover, under the additional assumption that $\widehat{C}_m e \leq 4/(\tau g_n)$, i.e.

$$\widehat{C}_m e \leq \frac{4}{e\sqrt{\mathcal{K}}g_n} = \frac{2g_m^{2N}\sqrt{W_{\mathcal{I}(w),m}}}{g_n\sqrt{e}} \quad \Rightarrow \quad W_{\mathcal{I}(w),m} > \left(\frac{\widehat{C}_m e^{3/2}g_n}{2g_m} \right)^{2N},$$

we can bound the term $(\widehat{C}_m e + 4/(\tau g_n))$ in (31) with $8/(\tau g_n)$, so that (31) can be rewritten as

$$\left\| u - \mathcal{S}_{\mathcal{I}(w)}^m[u] \right\|_{V \otimes L^2_\xi(\Gamma)} \leq W_{\mathcal{I}(w),m} C_E C_M^N \left(\frac{8^N}{W_{\mathcal{I}(w),m}} \prod_{n=1}^N \frac{1 - \epsilon_n}{\tau^2 g_n^2} \right)^{1/\tau}$$

which can be further simplified by inserting the nearly optimal value of τ previously computed:

$$\begin{aligned} \left\| u - \mathcal{S}_{\mathcal{I}(w)}^m[u] \right\|_{V \otimes L^2_\xi(\Gamma)} &\leq W_{\mathcal{I}(w),m} C_E C_M^N \left(\frac{8^N}{W_{\mathcal{I}(w),m}} \frac{1}{e^{2N}\mathcal{K}^N} \prod_{n=1}^N \frac{1 - \epsilon_n}{g_n^2} \right)^{1/\tau} \\ &= W_{\mathcal{I}(w),m} C_E C_M^N \left(\frac{8^N}{W_{\mathcal{I}(w),m}} \frac{e^N W_{\mathcal{I}(w),m} g_m^{2N}}{e^{2N} 4^N} \prod_{n=1}^N \frac{1 - \epsilon_n}{g_n^2} \right)^{1/\tau} \\ &= W_{\mathcal{I}(w),m} C_E C_M^N \left(\left(\frac{2}{e} \right)^N \prod_{n=1}^N (1 - \epsilon_n) \right)^{1/\tau} \\ &\leq W_{\mathcal{I}(w),m} C_E C_M^N \left(\left(\frac{2}{e} \right)^N e^{-N\delta} \right)^{1/\tau} \\ &= W_{\mathcal{I}(w),m} C_E C_M^N \exp \left(-\frac{N(\delta + 1 - \log 2)}{e\sqrt{\mathcal{K}}} \right) \\ &= W_{\mathcal{I}(w),m} C_E C_M^N \exp \left(-\frac{N(\delta + 1 - \log 2)g_m^{2N}\sqrt{W_{\mathcal{I}(w),m}}}{2\sqrt{e}} \right). \end{aligned}$$

The proof is then concluded by using Lemma 5. \square

6 Application to a diffusion problem with random inclusions

In this Section we show how the solution u of a certain class of elliptic PDEs with stochastic coefficients (namely, the so-called ‘‘inclusions problem’’ already examined in [4, 8]) satisfies the polyellipse analyticity condition A1; we will

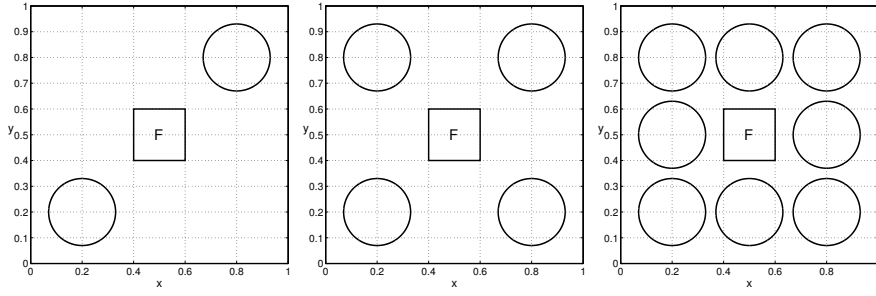


Fig. 3: Domains for the inclusions problem with 2, 4 and 8 circular inclusions.

then apply the previous Theorems 2 and 3 to establish the convergence of the quasi-optimal sparse grid approximation of u , using both nested and non-nested points, and verify numerically such convergence results.

Let D be a convex polygonal domain in \mathbb{R}^2 , and let \mathbf{y} be an N -variate random vector whose components y_1, \dots, y_N are independent uniform random variables over $\Gamma_i = [y_{min}, y_{max}]$. The support of the random vector \mathbf{y} is therefore the hypercube $\Gamma = \Gamma_1 \times \Gamma_2 \times \dots \times \Gamma_N$, the joint probability density function of \mathbf{y} is $\varrho(\mathbf{y}) = \prod_{n=1}^N \varrho_n(y_n) = \prod_{n=1}^N \frac{1}{y_{max} - y_{min}}$, and $(\Gamma, B(\Gamma), \varrho(\mathbf{y})d\mathbf{y})$ is a probability space, $B(\Gamma)$ being the Borel σ -algebra on Γ . We consider the stochastic elliptic problem

Problem 1 Find a real-valued function $u : \bar{D} \times \Gamma \rightarrow \mathbb{R}$, such that $\varrho(\mathbf{y})d\mathbf{y}$ -almost everywhere there holds:

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

where the operators div and ∇ imply differentiation with respect to the physical coordinate only, and the diffusion coefficient is:

$$a(\mathbf{x}, \mathbf{y}) = a_0 + \sum_{n=1}^N \gamma_n \chi_n(\mathbf{x}) y_n. \quad (32)$$

Here $\chi_n(\mathbf{x})$ are the indicator functions of the disjoint circular sub-domains $D_n \subset D = [0, 1]^2$ as in Figure 3, and a_0, γ_n are real coefficients such that $a(\mathbf{x}, \mathbf{y})$ is strictly positive and bounded, i.e. there exist two positive constants $0 < a_{min} < a_{max} < \infty$ such that

$$0 < a_{min} \leq a(\mathbf{x}, \mathbf{y}) \leq a_{max} < \infty, \quad (33)$$

$\varrho(\mathbf{y})d\mathbf{y}$ -almost surely, $\forall \mathbf{x} \in D$.

Next, let $V = H_0^1(D)$ be 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 $\|v\|_V = \|\nabla v\|_{L^2(D)}$. Using Lax–Milgram’s Lemma it is straightforward to show that Problem 1 is $\varrho(\mathbf{y})d\mathbf{y}$ -almost everywhere well-posed in V , due to the boundedness assumption (33). Similarly, it is easy to see that $u \in L^2_\varrho(\Gamma) \otimes V$, see e.g. [4, 6, 8].

Remark 6 In this work we do not address the discretization of the solution u in the physical variable \mathbf{x} . In this respect all results obtained here apply also to a discrete solution u_h , obtained by introducing e.g. a finite element discretization over a triangulation \mathcal{T}_h of the physical domain D , and a finite element space of piecewise continuous polynomials on \mathcal{T}_h , $V_h(D) \subset H^1_0(D)$.

We shall begin by reparametrizing the diffusion coefficient in terms of new random variables distributed over $[-1, 1]$. For the sake of notation, we will still denote the new variables as y_i , i.e. $y_i \sim \mathcal{U}(-1, 1)$. The new diffusion coefficient will be therefore:

$$a(\mathbf{x}, \mathbf{y}) = a_0 + \sum_{n=1}^N \gamma_n \chi_n(\mathbf{x}) \left(\frac{y_n + 1}{2} (y_{\max} - y_{\min}) + y_{\min} \right). \quad (34)$$

Lemma 9 *The complex continuation u^* of the solution u corresponding to a diffusion coefficient (34) is analytic in the region*

$$\Sigma = \Sigma_1 \times \Sigma_2 \times \dots \times \Sigma_N, \quad \Sigma_n = \{z_n \in \mathbb{C} : \Re(z_n) \geq T_n\},$$

with $-1 \geq T_n > T_n^* = \frac{2a_0 + \gamma_n(y_{\max} + y_{\min})}{\gamma_n(y_{\min} - y_{\max})}$. Moreover, $\sup_{\mathbf{z} \in \Sigma} \|u^*\|_V \leq B_u(T_1, \dots, T_N)$, with

$$B_u(T_1, \dots, T_N) = \frac{\|f\|_{V'}}{a_0 + \sum_{n=1}^N \gamma_n \left(\frac{1 - |T_n|}{2} (y_{\min} - y_{\max}) + y_{\min} \right)}.$$

Proof See [8, Lemma 23]. \square

Corollary 1 *The solution u corresponding to a diffusion coefficient (34) satisfies Assumption A1 with $\delta_n^* = |T_n^*| + \sqrt{|T_n^*|^2 - 1}$, T_n^* as in Lemma 9.*

Proof We only need to compute the parameter δ_n^* corresponding to the largest Bernstein ellipse contained in the analyticity region Σ . This can be done by enforcing $(\delta_n^* + \delta_n^{*-1})/2 = |T_n^*|$. \square

6.1 On the choice of collocation points

Before proceeding further with the description of the numerical tests performed, we specify here the families of collocation points used to build the sparse grids considered in the following, as well as the values of the corresponding quantities \mathbb{M}_n^q , cf. Definition 2. Finally, we will verify that such families of points satisfy the Assumptions A2 and A3 needed for the convergence Theorems to hold true.

As for nested points, we use the Clenshaw–Curtis rule (see e.g. [33]),

$$y_j^i = \cos\left(\frac{(j-1)\pi}{m(i)-1}\right), 1 \leq j \leq m(i),$$

that are nested when using the following level-to-nodes relation:

$$m_{db}(i) = \begin{cases} 0 & \text{if } i = 0 \\ 1 & \text{if } i = 1 \\ 2^{i-1} + 1, & \text{if } i > 1. \end{cases} \quad \Rightarrow \quad d_{db}(i) = \begin{cases} 1 & \text{if } i = 1 \\ 2 & \text{if } i = 2 \\ 2^{i-2}, & \text{if } i > 2, \end{cases} \quad (35)$$

while for non-nested points we will use the classical Gauss–Legendre (see e.g. [32]), together with the level-to-nodes relation

$$m_{lin}(i) = i, \quad \Rightarrow \quad d_{lin}(i) = 1. \quad (36)$$

Other families of points that are commonly considered for problems with uniform random variables are the Gauss–Patterson and the Leja points, see e.g. [28,10,29] and references therein. We discuss here these different choices of points.

Gauss–Legendre. As stated in Lemma 3, for Gauss–Legendre nodes there holds $\mathbb{M}_n^q = 1$. This, together with equation (36), implies that Assumption A2 holds with $C_{\mathbb{M}} = 1$. Assumption A3 holds with $C_m = 1$ due to equation (36).

Clenshaw–Curtis. In the case of nested Clenshaw–Curtis nodes, we use the standard estimate of the “ L^∞ ” Lebesgue constant (see e.g. [17,18]) as a bound for \mathbb{M}_n^q ,

$$\mathbb{M}_n^q \leq \mathbb{M}_{n,est}^q, \quad \mathbb{M}_{n,est}^q = \begin{cases} 1 & \text{for } q = 1 \\ \frac{2}{\pi} \log(q-1) + 1 & \text{for } q \geq 2. \end{cases}$$

Combining this estimate and equation (35) we obtain that Assumption A2 holds with $C_{\mathbb{M}} = 1$. Assumption A3 holds with $C_m = 2$, due again to equation (35).

Leja. Given a compact set X and an initial value $x_0 \in X$, Leja sequences are recursively defined as $x_k = \operatorname{argmax}_{y \in X} \left| \prod_{j=0}^{k-1} (y - x_j) \right|$. Choosing $x_0 = 1$ and $X = [-1, 1]$ results in the so-called *standard Leja sequence*, while choosing as X the unit disk in the complex domain together with $x_0 = 1$ and projecting the resulting sequence on $X = [-1, 1]$ one obtains the so-called *R-Leja sequence*, see [10,29] and references therein for details. Here we focus on the so-called

symmetrized Leja sequence (see again [29]), which at level i includes $2i + 1$ points defined as

$$x_0^i = 0, \quad x_1^i = 1, \quad x_2^i = -1,$$

$$x_k^i = \begin{cases} \operatorname{argmax}_{y \in X} \left| \prod_{j=0}^{k-1} (y - x_j) \right| & \text{if } k \text{ is odd, } k \leq 2i \\ -x_{k-1} & \text{if } k \text{ is even, } k \leq 2i + 1 \end{cases}$$

so that the level-to-nodes function is defined as

$$m_{SL}(i) = \begin{cases} 0 & \text{if } i = 0 \\ 1 & \text{if } i = 1 \\ 2i + 1, & \text{if } i > 1. \end{cases} \quad \Rightarrow \quad d_{SL}(i) = \begin{cases} 1 & \text{if } i = 1 \\ 2 & \text{if otherwise.} \end{cases}$$

Thus, Assumption A3 trivially holds with $C_m = 2$. The validity of Assumption A2 can be verified numerically, by computing lower and upper bounds for the Lebesgue constant. The upper bound is obtained as in the case of Clenshaw–Curtis points, by bounding $\mathbb{M}_n^{m_{SL}(i_n)}$ with the standard Lebesgue constant, that can be computed numerically. The lower bound is also established numerically, by solving approximately the maximization problem appearing in the definition of $\mathbb{M}_n^{m(i_n)}$ (see Definition 2), that can be recast into a constrained quadratic optimization problem.

To do this, denote by $t_1, t_2, \dots, t_{m(i_n)}$ the collocation points of $\mathcal{U}_n^{m(i_n)}$ and by $\ell_1, \ell_2, \dots, \ell_{m(i_n)}$ the associated Lagrangian polynomials, and expand

$$\begin{aligned} \int_{\Gamma_n} (\mathcal{U}_n^{m(i_n)}[f])^2 \varrho_n(y_n) dy_n &= \int_{\Gamma_n} \left(\sum_{j=1}^{m(i_n)} f(t_j) \ell_j(y_n) \right)^2 \varrho_n(y_n) dy_n \\ &= \sum_{\kappa, j}^{m(i_n)} f(t_\kappa) f(t_j) \int_{\Gamma_n} \ell_\kappa(y_n) \ell_j(y_n) \varrho_n(y_n) dy_n. \end{aligned}$$

Observe now that, being $\ell_\kappa(y_n) \ell_j(y_n)$ a polynomial of degree $2(m(i_n) - 1)$, we can integrate it exactly with a Gaussian quadrature formula with $m(i_n)$ quadrature points: we further denote by $\xi_1, \xi_2, \dots, \xi_{m(i_n)}$ the quadrature points and by $\alpha_1, \alpha_2, \dots, \alpha_{m(i_n)}$ the associated quadrature weights, and we let \mathbf{f} be the vector containing the nodal values $f_j = f(t_j)$. Computing $\mathbb{M}_n^{m(i_n)}$ amounts then to solving the quadratic optimization problem

$$\mathcal{L} = \max_{\mathbf{f} \in \mathcal{R}} \mathbf{f}^T A \mathbf{f}, \quad \mathcal{R} = \{ \mathbf{f} \in \mathbb{R}^{m(i_n)} \text{ s.t. } -1 \leq f_n \leq 1, \quad \forall n = 1, \dots, N \}$$

with $A_{ij} = \sum_{q=1}^{m(i_n)} \ell_i(\xi_q) \ell_j(\xi_q) \alpha_q$, and setting $\mathbb{M}_n^{m(i_n)} = \sqrt{\mathcal{L}}$. Being A is positive definite, the solutions of the optimization problem are located in the corners of the feasible region \mathcal{R} , and multiple maxima are possible. By repeatedly running an optimization algorithm for quadratic optimization problems

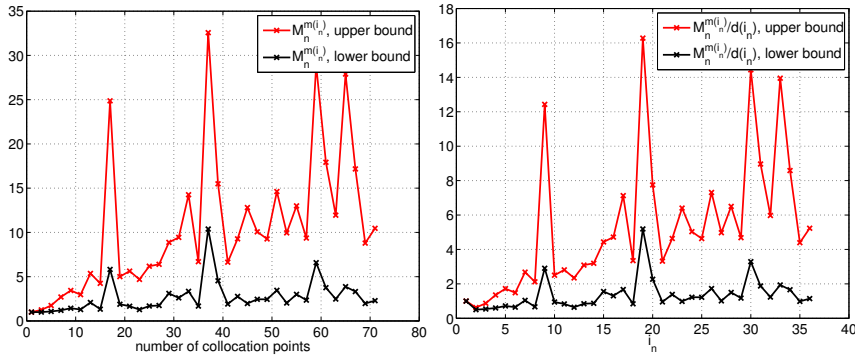


Fig. 4: Left: the lower and upper bounds for $\mathbb{M}_n^{m_{SL}(i_n)}$ for symmetric Leja points show that $\mathbb{M}_n^{m_{SL}(i_n)}$ grows polynomially. Right: the ratio $\mathbb{M}_n^{m_{SL}(i_n)}/d(i_n)$ grows polynomially as well.

(here we use the active set method, see e.g. [27]) with different initial guesses, we obtain a suboptimal solution, which is however sufficient for our purposes.

Results are reported in Figure 4, and show that we can assume a polynomial growth for $\mathbb{M}_n^{m_{SL}(i_n)}$ and for $\mathbb{M}_n^{m_{SL}(i_n)}/d(i_n)$; yet, the theory previously developed could be still be applied, at the price of modifying the auxiliary profits P^b introduced in Lemma 7 by changing the rates g_n with $g_n^b = g_n(1 - \epsilon)$ for every $\epsilon > 0$ and the constant C_M with another constant $C_M(\epsilon)$ increasing as $\epsilon \rightarrow 0$.

The same conclusion can be deduced for the R-Leja sequence, thanks to the results stated in [10], where a polynomial growth is proved for the standard “ L^∞ ” constant for this sequence of points.

Gauss–Patterson. These nodes are tabulated (see [28]). In particular, there holds

$$m_{GP}(0) = 0, \quad m_{GP}(i_n) = \sum_{k=0}^{i_n-1} 2^k$$

therefore $d_{GP}(i_n) = 2^{i_n-1}$, hence Assumption A3 holds with $C_m = 2$, while we can verify again the validity of Assumption A2 numerically.

The numerical results are shown in Figure 5, and suggest that both $\mathbb{M}_{n,est}^{m_{GP}(i_n)}$ and the ratio $\mathbb{M}_n^{m_{GP}(i_n)}/d(i_n)$ may asymptotically grow more than polynomially, hence we cannot use the results obtained in the previous Sections to derive the convergence of the quasi-optimal sparse built with Gauss–Patterson knots.

6.2 Numerical results

In this Section we consider three different “inclusions” geometries, with $N = 2, 4$ and 8 inclusions respectively (see Figure 3), with $y_{min} = -0.99$ and $y_{max} =$

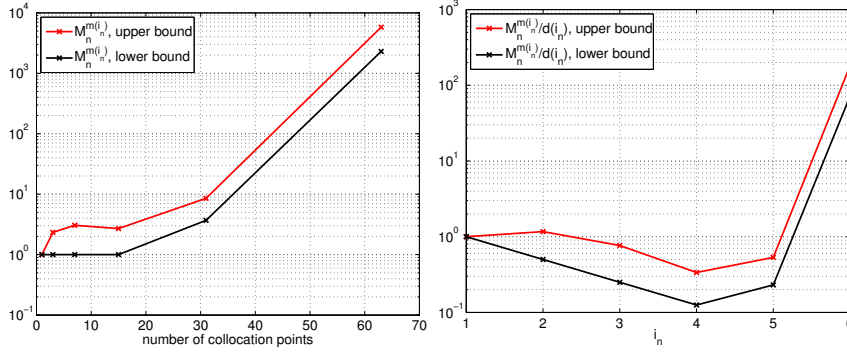


Fig. 5: The lower and upper bounds for $M_n^{m_{SL}(i_n)}$ for Gauss–Patterson points (left plot) and the corresponding ratios $M_n^{m_{SL}(i_n)}/d(i_n)$ are not bounded.

	γ_1	γ_2	γ_3	γ_4	γ_5	γ_6	γ_7	γ_8
test $N = 2$	1	0.0035						
test $N = 4$	1	0.06	0.0035	0.0002				
test $N = 8$	1	0.25	0.06	0.015	0.0035	0.0009	0.0002	0.00005

Table 1: Values of the coefficients γ_n for the anisotropic settings. Inclusions are numbered anticlockwise, starting from the bottom-left (south-west) corner.

0.99. We set homogeneous Dirichlet boundary conditions and use a constant forcing term defined on the square subdomain F located in the center of the physical domain D (see again Figure 3), i.e. $f(\mathbf{x}) = 100\chi_F(\mathbf{x})$. Each geometry is considered both in an isotropic setting (i.e. γ_n in (32) are set to 1 for all $n = 1, \dots, N$) and an anisotropic setting (see Table 1 for the values of γ_n).

As already mentioned, we will test the performances of two different versions of the quasi-optimal sparse grid proposed in the previous section: one using nested points, which implies using Clenshaw–Curtis points and the profits in (14), and one with non-nested points, which implies using Gauss–Legendre points and the profits in (17); from here on, we will denote these two grids as “OPT-N” and “OPT-NN”. In particular, we will verify the accuracy of the convergence ansatz

$$\left\| u - \mathcal{S}_{\mathcal{I}(w)}^m[u] \right\|_{V \otimes L^2_\varepsilon(\Gamma)} \leq \alpha_n \exp\left(-\beta_n N \sqrt[N]{W_{\mathcal{I}(w),m}}\right), \quad (37)$$

for nested points and of

$$\left\| u - \mathcal{S}_{\mathcal{I}(w)}^m[u] \right\|_{V \otimes L^2_\varepsilon(\Gamma)} \leq \alpha_{nn} \exp\left(-\beta_{nn} N \sqrt[2N]{W_{\mathcal{I}(w),m}}\right) \quad (38)$$

for non-nested points, where $\alpha_n, \beta_n, \alpha_{nn}, \beta_{nn}$ are numerical values. In practice, we introduce the bounded linear functional $\Theta : V \rightarrow \mathbb{R}$,

$$\Theta(u) = \int_F u(\mathbf{x}) d\mathbf{x}$$

	test $N = 2$	test $N = 4$	test $N = 8$
y_1	1.41 (1)	1.41 (1)	1.41 (1)
y_2	7.16 (5)	4.31 (3)	2.88 (2)
y_3		7.16 (5)	4.31 (3)
y_4		10.01 (7)	5.70 (4)
y_5			7.16 (5)
y_6			8.51 (6)
y_7			10.01 (7)
y_8			11.40 (8)

Table 2: Absolute and normalized values of the anisotropy rates g_n for the anisotropic settings. The normalized values are shown in parenthesis, and are defined as g_n/g_1 . Inclusions are numbered anticlockwise, starting from the bottom-left (south-west) corner. Finally, in the isotropic setting the value of g is $g = 1.41$ for all variables.

and we monitor the convergence of the quantity

$$\varepsilon = \sqrt{\mathbb{E} \left[\left(\Theta(\mathcal{S}_{\mathcal{I}(w)}^m[u]) - \Theta(u) \right)^2 \right]}, \quad (39)$$

with respect to number of sparse grid points, that will converge with the same rate as the full error $\|u - \mathcal{S}_{\mathcal{I}(w)}^m[u]\|_{V \otimes L^2_\rho(\Gamma)}$, given the linearity of Θ .

Moreover, from a practical point of view, it is important to observe that using the logarithm of the Bernstein radii δ_n to estimate the quantities g_n needed in the sparse grid construction (cf. equations (14) and (17)) turns out to be very pessimistic, as observed in [4, 6, 8]. Such quantities are better assessed with the numerical procedure described in [4, 6]: the results are shown in Table 2.

Remark 7 In our numerical experiments, the set of the w largest profits is always downward closed, i.e. we never have to explicitly enforce the admissibility condition (2), both in the nested and non-nested case. By following closely the argument in [31, Chapter 6, Lemma 19] it is actually easy to show that, regardless of the values of g_1, \dots, g_N , the set of the w largest profits is necessarily downward closed when considering Gauss–Legendre points. Conversely, the set of the w largest profits is downward closed when considering Clenshaw–Curtis points only under the assumption that the rates g_n are sufficiently large. However, such condition is very mild ($g_n \geq \bar{g} \approx 0.13$) and satisfied by the values in Table 2.

We will furthermore compare the performances of the quasi-optimal sparse grids “OPT-N” and “OPT-NN”, with that of a number of different sparse grid schemes. In particular, we will consider:

1. A standard sparse grid (labeled “SM”) built with the classical Clenshaw–Curtis abscissas, together with the level-nodes relation $m(i_n) = m_{db}(i)$ and

using the multi-index set

$$\mathcal{I}_{SM}(w) = \left\{ \mathbf{i} \in \mathbb{N}_+^N : \sum_{n=1}^N (i_n - 1) \leq w \right\},$$

and its anisotropic counterpart (“aSM”)

$$\mathcal{I}_{aSM}(w) = \left\{ \mathbf{i} \in \mathbb{N}_+^N : \sum_{n=1}^N g_n (i_n - 1) \leq w \right\},$$

proposed in [26,4]; the rates g_n used here are those listed in Table 2.

2. The (anisotropic) Total Degree sparse grid (labeled respectively “TD” and “aTD”) proposed in [4] with Gauss–Legendre points, $m(i_n) = m_{i_n}(i)$ and

$$\mathcal{I}_{TD}(w) = \left\{ \mathbf{i} \in \mathbb{N}_+^N : \sum_{n=1}^N (i_n - 1) \leq w \right\},$$

$$\mathcal{I}_{aTD}(w) = \left\{ \mathbf{i} \in \mathbb{N}_+^N : \sum_{n=1}^N g_n (i_n - 1) \leq w \right\},$$

again, the rates g_n used here are those listed in Table 2.

3. The adaptive strategy proposed by [19], in the implementation provided by [23] and available at <http://www.ians.uni-stuttgart.de/spinterp> (labeled “KL”). As already mentioned in the introduction and in Section 3, this is an adaptive algorithm that explores the set of admissible hierarchical surpluses and adds to the sparse grid approximation the most “profitable” ones, according to suitable a-posteriori estimates. The implementation considered here has a tunable parameter $\tilde{\omega}$ that allows one to move continuously from the standard sparse grid just described ($\tilde{\omega} = 0$) to the fully adaptive algorithm ($\tilde{\omega} = 1$). Following [23], in the present work we have set $\tilde{\omega} = 0.9$, that numerically has been proved to be a well performing choice. Clearly, this strategy is bounded to use nested (i.e. Clenshaw–Curtis) points, and (at least in the implementation considered here) works only on problems with a finite number of dimensions. We will measure the convergence of this algorithm in terms of the *total* number of points, i.e. including also those necessary to explore the set of hierarchical surpluses.
4. Two “brute force” approximations of the quasi-optimal sparse grid (one for the nested points case and one for the non-nested points case, labeled “BF-N” and “BF-NN” respectively), that we obtain by first computing numerically the profits of all the hierarchical surpluses in a sufficiently large “universe” $\mathbb{U} \subset \mathbb{N}_+^N$ (see Table 3) and then sorting them in decreasing order. Whenever such ordering does not satisfy the admissibility condition (2), all the hierarchical surpluses needed are added (this approach is equivalent to modifying the profits according to (9)). These grids were computed for $N = 2$ and $N = 4$ only.

test case	U-nested	U-non nested	MC samples
iso2D	TD(8)	TD(8)	6000
iso4D	TD(8)	TD(8)	25000
iso8D	TD(10)	TD(10)	50000
aniso2D	TD(6)	TD(10)	5000
aniso4D	TD(6)	TD(13)	15000
aniso8D	TD(10)	TD(13)	25000

Table 3: Universe \mathbb{U} and size of Monte Carlo sample considered in each computational test. Due to the different level-to-nodes relations used, we use two different sets \mathbb{U} for the nested and non-nested case.

The error (39) has been computed with a Monte Carlo sampling, see Table 3; we emphasize that the number of Monte Carlo samples used has been verified to be sufficient for our purposes. The same sampling strategy has also been employed for the computation of the profits needed to build the brute force sparse grids “BF-N” and “BF-NN”.

We show the results for the isotropic and anisotropic setting in Figures 6 and 7 respectively; from the analysis of the numerical results, several conclusions can be drawn. First, the proposed profit estimates are quite sharp, both in the nested and non-nested case, since the convergence curves for the “brute force” sparse grids “BF-N”/“BF-NN” and their estimated counterparts “OPT-N”/“OPT-NN” are very close in every test. Observe that while this was expected for “BF-N” and “OPT-N”, given the corresponding results presented in previous works [6, 7], this was not obvious for “BF-NN” and “OPT-NN”, given the pessimistic approach adopted to estimate the work contribution of each hierarchical surplus. The non-monotone convergence curve for the “OPT-NN” scheme can be explained with the fact, already pointed out, that increasing the number of multi-indices does not necessarily lead to an increase of the number of total points in a sparse grid when using non-nested points (cf. Figure 2 and Example 1).

Second, “OPT-N” is found to be more efficient than “OPT-NN”, as expected given the convergence Theorems 2 and 3; “OPT-N” is furthermore found to be competitive with the a-posteriori sparse grid construction (“KL”), again in agreement with the previous work [6].

Third, comparing the performance of “OPT-N” and “OPT-NN” with that of non-optimized sparse grids, like the Smolyak and Total Degree ones (“SM”/“aSM” and “TD”/“aTD”), we see that the “TD”/“aTD” convergence behavior closely resembles that of “OPT-NN”, while the same does not hold true for the nested-points corresponding grids, i.e. “SM”/“aSM” versus “OPT-N”. Indeed, if on the one hand in the isotropic setting “SM” is competitive with the nested quasi-optimal grid (although less efficient for low tolerances), on the other hand “aSM” is instead quite less efficient than “OPT-N”, and even less efficient than the isotropic “SM” for low values of N , i.e. $N = 2, 4$. Observe also the significant loss of efficiency caused by using isotropic approximations in the anisotropic setting as the number of variables increases, highlighting the need for anisotropic approximation schemes.

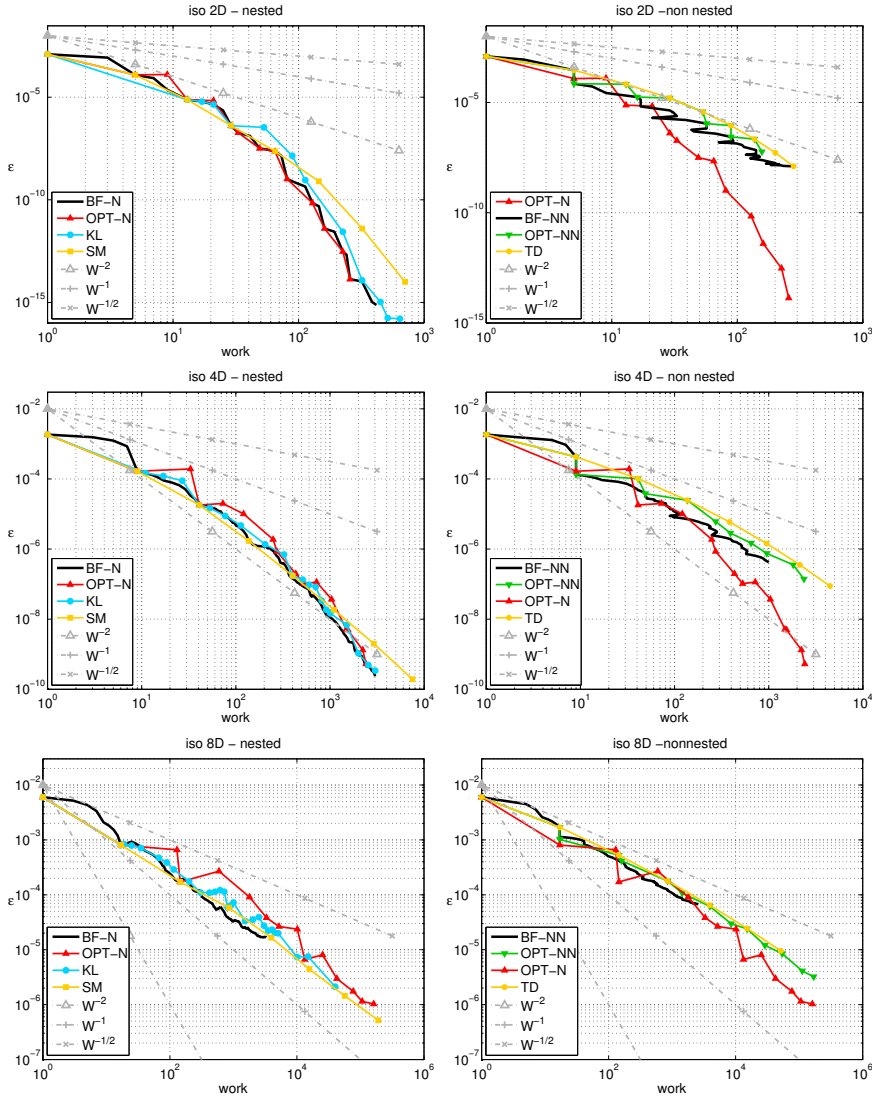


Fig. 6: Results for the **isotropic** setting. Top row: case $N = 2$; middle row: case $N = 4$; bottom row: case $N = 8$. Left column: sparse grids with nested points; right column: sparse grid with non-nested points.

Next, we verify the sharpness of the theoretical bounds provided in Theorems 2 and 3, and in particular of the convergence ansatz (37)-(38). To this end, we plot in semi-logarithmic scale the quantities $N \sqrt[N]{W_{\mathcal{I}(w),m}}$, $N^{2N} \sqrt[2N]{W_{\mathcal{I}(w),m}}$ versus the sparse grid error for “OPT-N” and “OPT-NN” sparse grids. Results are shown in Figure 8: we correctly get straight lines, thus suggesting that the ansatz can be considered quite effective.

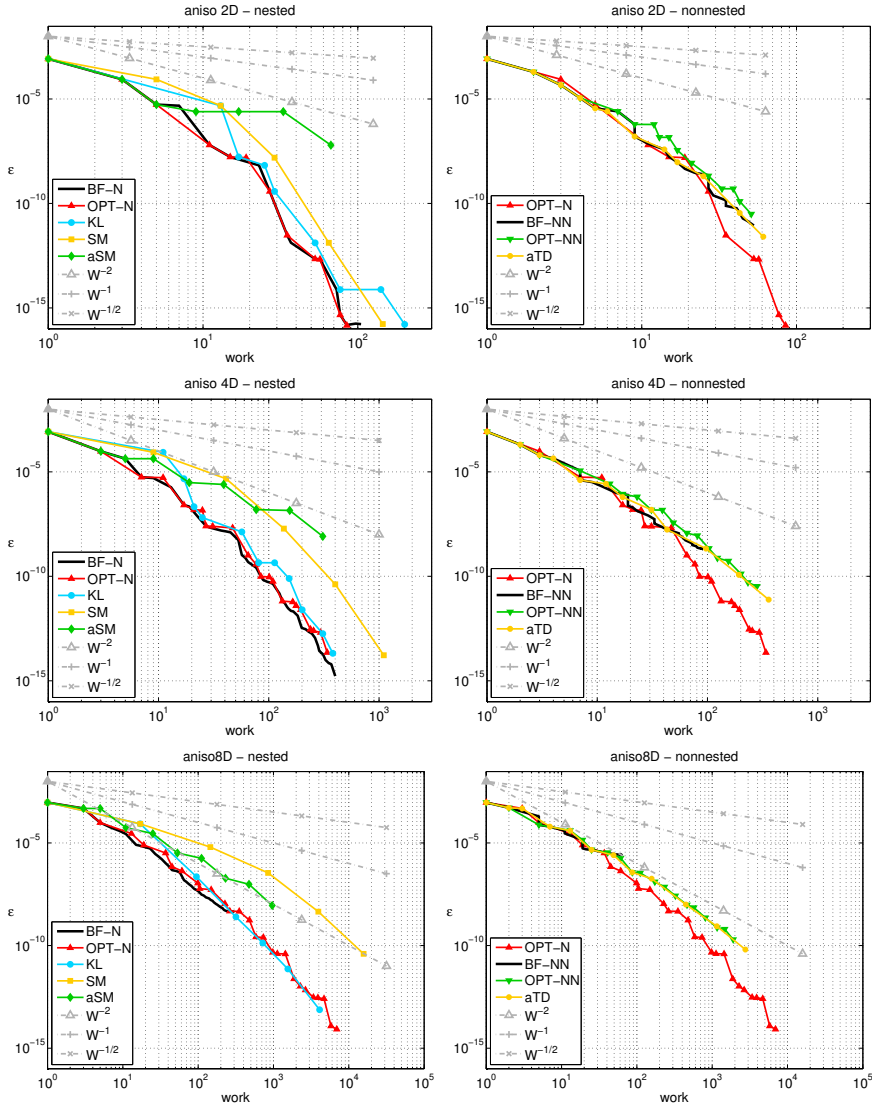


Fig. 7: Results for the **anisotropic** setting. Top row: case $N = 2$; middle row: case $N = 4$; bottom row: case $N = 8$. Left column: sparse grids with nested points; right column: sparse grid with non-nested points.

Finally, we investigate the convergence of the expected value of $\Theta(u)$, $|\mathbb{E}[\mathcal{S}_{\mathcal{I}(w)}^n[\Theta(u)] - \mathbb{E}[\Theta(u)]]|$, see Figures 9 and 10. As expected, the convergence in this case is faster than the convergence of the full $V \otimes L^2_\theta(\Gamma)$ norm inspected previously. Moreover, the convergence is significantly less smooth than the previous case, due to cancellations among hierarchical surpluses. Finally, in the

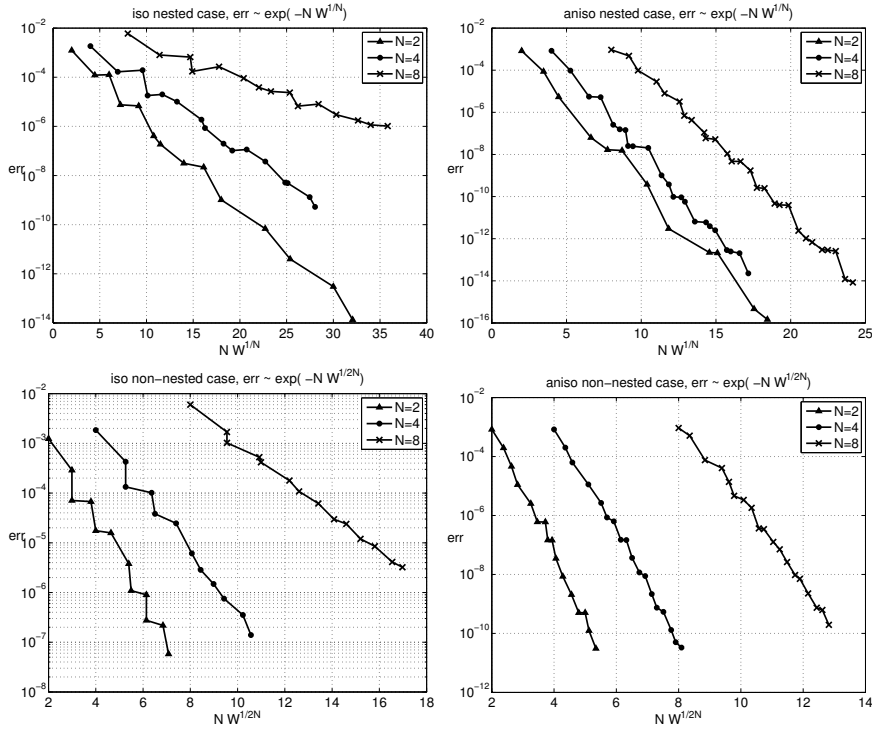


Fig. 8: Verification of the quasi-optimal sparse grid convergence estimates. Top row: nested case; bottom row: non-nested case. Left column: isotropic setting; right column: anisotropic setting.

anisotropic setting the non-nested grids are surprisingly competitive with the nested schemes. A possible explanation for this is that in this test we are actually considering the quadrature capabilities of the sparse grids rather than the interpolation ones, for which Gaussian collocation points are particularly suitable. Beside this aspect, the other observations on the performances of the sparse grids schemes are the same as the previous case.

7 Conclusions

In this work we have proved an error estimate for the stochastic collocation based on quasi-optimal sparse grid constructed by choosing the w most profitable hierarchical surpluses, that is the w surpluses with the highest error reduction / cost ratio. The convergence of such grid is proved in terms of weighted τ -summability of such profits: as the true profits are unknown, we propose to build the quasi-optimal sparse grid introducing a-priori estimates on the decay of the profits. We have then considered the application of such quasi-optimal sparse grid to Hilbert-valued functions which are analytic in

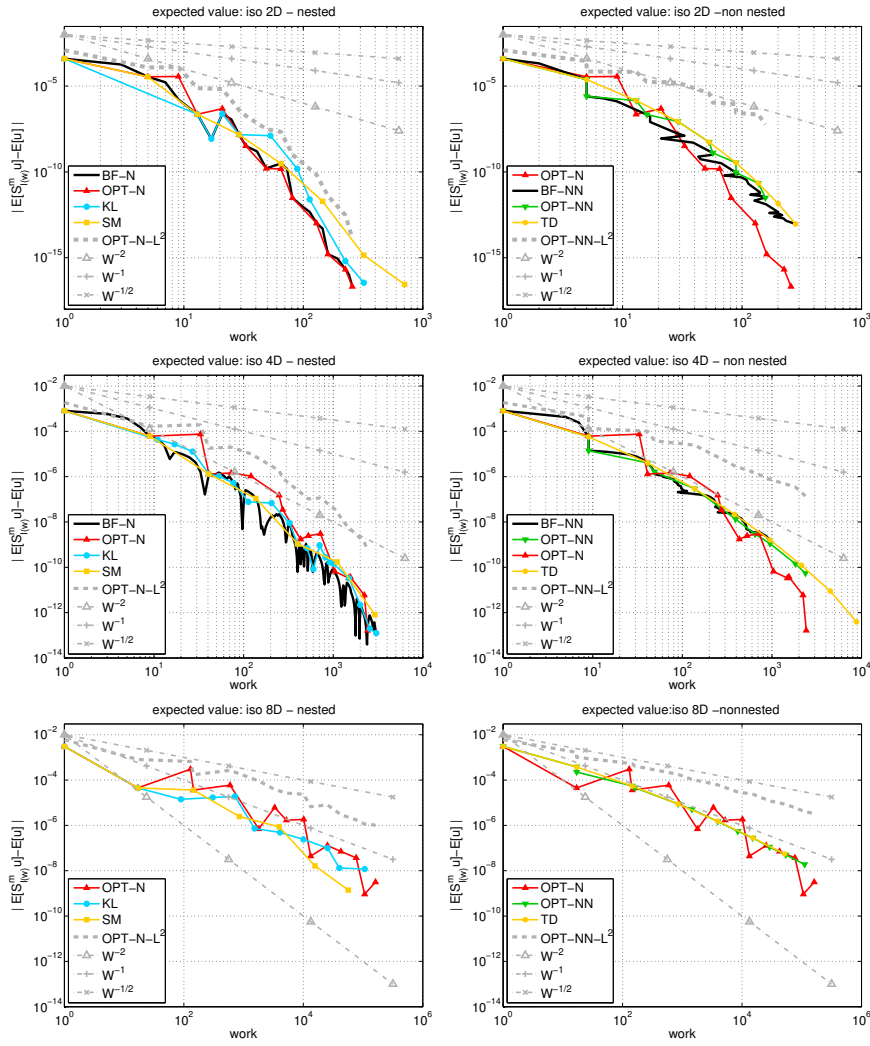


Fig. 9: Results for the **isotropic** setting, convergence of the expected value of $\Theta(u)$. Top row: case $N = 2$; middle row: case $N = 4$; bottom row: case $N = 8$. Left column: sparse grids with nested points; right column: sparse grid with non-nested points.

certain polyellipses. We have considered two variations of the scheme, one using nested collocation points and the other one using non-nested points; in both cases we were able to derive profit bounds and prove the corresponding τ -summability. After having verified that the solution of the so-called “inclusion problem” satisfied the above-mentioned analyticity Assumption, we have shown with some numerical tests that the profit estimates are quite sharp and

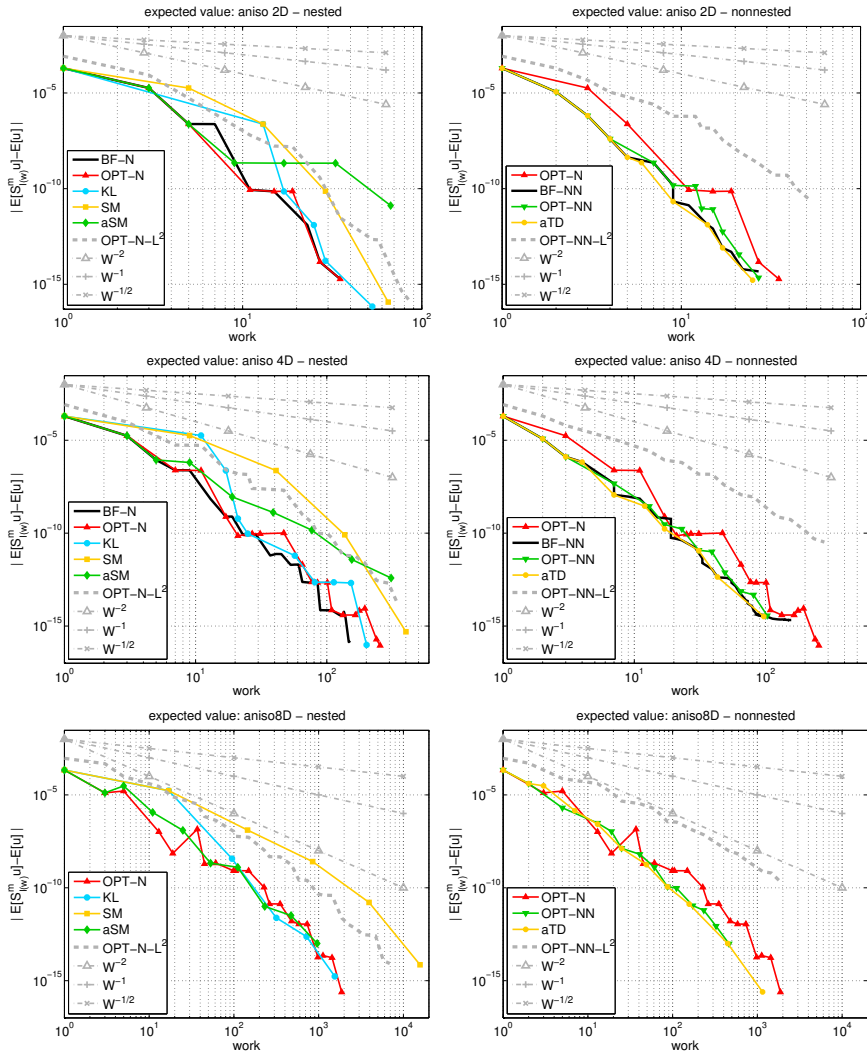


Fig. 10: Results for the **anisotropic** setting, convergence of the expected value of $\Theta(u)$. Top row: case $N = 2$; middle row: case $N = 4$; bottom row: case $N = 8$. Left column: sparse grids with nested points; right column: sparse grid with non-nested points.

that the convergence results provide the correct ansatz for the error decay, though with constants fitted numerically. The proposed method is therefore competitive with the a-posteriori adaptive scheme [23], and possibly outperforms the previously proposed anisotropic sparse grids.

Acknowledgements The authors would like to recognize the support of King Abdullah University of Science and Technology (KAUST) AEA project “Predictability and Uncertainty Quantification for Models of Porous Media” and University of Texas at Austin AEA Rnd 3 “Uncertainty quantification for predictive modeling of the dissolution of porous and fractured media”. F. Nobile and L. Tamellini have been supported by the Italian grant FIRB-IDEAS (Project n. RBID08223Z) “Advanced numerical techniques for uncertainty quantification in engineering and life science problems”. They also received support from the Center for Advanced Modeling Science (CADMOS). R. Tempone is a member of the KAUST SRI Center for Uncertainty Quantification in Computational Science and Engineering. We acknowledge the usage of the Matlab[®] functions `patterson_rule.m` by J. Burkardt⁴ for the computation of Gauss–Patterson points and `lejapoints.m` by M. Caliari⁵ for the computation of Symmetrized Leja Points.

References

1. K. I. Babenko. Approximation by trigonometric polynomials in a certain class of periodic functions of several variables. *Soviet Math. Dokl.*, 1:672–675, 1960.
2. I. 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.
3. I. Babuška, F. Nobile, and R. Tempone. A stochastic collocation method for elliptic partial differential equations with random input data. *SIAM Review*, 52(2):317–355, June 2010.
4. 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.
5. V. Barthelmann, E. Novak, and K. Ritter. High dimensional polynomial interpolation on sparse grids. *Adv. Comput. Math.*, 12(4):273–288, 2000.
6. J. Beck, F. Nobile, L. Tamellini, and R. Tempone. On the optimal polynomial approximation of stochastic PDEs by Galerkin and collocation methods. *Mathematical Models and Methods in Applied Sciences*, 22(09), 2012.
7. J. Beck, F. Nobile, L. Tamellini, and R. Tempone. A quasi-optimal sparse grids procedure for groundwater flows. MATHICSE Report 46/2012, Ecole Polytechnique Fédérale de Lausanne, 2012. To appear in M. Azaiez, H. El Fekih, and J. S. Hesthaven editors, *Spectral and High Order Methods for Partial Differential Equations*. Selected papers from the ICOSAHOM '12 conference.
8. J. Beck, F. Nobile, L. Tamellini, and R. Tempone. Convergence of quasi-optimal Stochastic Galerkin methods for a class of PDEs with random coefficients. *Computers & Mathematics with Applications*, 2013. In press. Available online, DOI information: 10.1016/j.camwa.2013.03.004.
9. H.J Bungartz and M. Griebel. Sparse grids. *Acta Numer.*, 13:147–269, 2004.
10. A. Chkifa. On the lebesgue constant of leja sequences for the complex unit disk and of their real projection. *Journal of Approximation Theory*, 166(0):176 – 200, 2013.
11. A. Chkifa, A. Cohen, R. Devore, and C. Schwab. Sparse adaptive Taylor approximation algorithms for parametric and stochastic elliptic PDEs. *ESAIM: Mathematical Modelling and Numerical Analysis*, 47(1):253–280, 2013.
12. A. Chkifa, A. Cohen, and C. Schwab. High-dimensional adaptive sparse polynomial interpolation and applications to parametric PDEs. *Foundations of Computational Mathematics*, pages 1–33, 2013.

⁴ http://people.sc.fsu.edu/~jburkardt/m_src/patterson_rule/patterson_rule.html

⁵ <http://profs.sci.univr.it/~caliari/software/lejapoints.m>

13. A. Cohen, R. DeVore, and C. Schwab. Analytic regularity and polynomial approximation of parametric and stochastic elliptic PDE'S. *Anal. Appl. (Singap.)*, 9(1):11–47, 2011.
14. B. A. Davey and H. A. Priestley. *Introduction to lattices and order*. Cambridge University Press, New York, second edition, 2002.
15. R. A. DeVore. Nonlinear approximation. In *Acta numerica, 1998*, volume 7 of *Acta Numer.*, pages 51–150. Cambridge Univ. Press, Cambridge, 1998.
16. R. A. DeVore and G. G. Lorentz. *Constructive Approximation*. Die Grundlehren der mathematischen Wissenschaften in Einzeldarstellungen. Springer, 1993.
17. V. K. Dzjadik and V. V. Ivanov. On asymptotics and estimates for the uniform norms of the Lagrange interpolation polynomials corresponding to the Chebyshev nodal points. *Anal. Math.*, 9(2):85–97, 1983.
18. H. Ehlich and K. Zeller. Auswertung der Normen von Interpolationsoperatoren. *Math. Ann.*, 164:105–112, 1966.
19. T. Gerstner and M. Griebel. Dimension-adaptive tensor-product quadrature. *Computing*, 71(1):65–87, 2003.
20. R. G. Ghanem and P. D. Spanos. *Stochastic finite elements: a spectral approach*. Springer-Verlag, New York, 1991.
21. M. Griebel and S. Knapek. Optimized general sparse grid approximation spaces for operator equations. *Math. Comp.*, 78(268):2223–2257, 2009.
22. W. Gui and I. Babuka. The h,p and h-p versions of the finite element method in 1 dimension - part i. the error analysis of the p-version. *Numerische Mathematik*, 49(6):577–612, 1986.
23. A. Klimke. *Uncertainty modeling using fuzzy arithmetic and sparse grids*. PhD thesis, Universität Stuttgart, Shaker Verlag, Aachen, 2006.
24. O. P. Le Maître and O. M. Knio. *Spectral methods for uncertainty quantification*. Scientific Computation. Springer, New York, 2010. With applications to computational fluid dynamics.
25. C. Lubich. *From quantum to classical molecular dynamics: reduced models and numerical analysis*. Zurich lectures in advanced mathematics. European Mathematical Society, 2008.
26. 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.
27. J. Nocedal and S. J. Wright. *Numerical Optimization*. Springer, 1999.
28. T. N. L. Patterson. The optimum addition of points to quadrature formulae. *Math. Comp.* 22 (1968), 847–856; addendum, *ibid.*, 22(104):C1–C11, 1968.
29. C. Schillings and C. Schwab. Sparse, adaptive Smolyak quadratures for Bayesian inverse problems. *Inverse Problems*, 29(6), 2013.
30. Jie Shen and Li-Lian Wang. Sparse spectral approximations of high-dimensional problems based on hyperbolic cross. *SIAM J. Numer. Anal.*, 48(3):1087–1109, 2010.
31. L. Tamellini. *Polynomial approximation of PDEs with stochastic coefficients*. PhD thesis, Politecnico di Milano, 2012.
32. L. N. Trefethen. Is Gauss quadrature better than Clenshaw-Curtis? *SIAM Rev.*, 50(1):67–87, 2008.
33. L.N. Trefethen. *Approximation Theory and Approximation Practice*. Society for Industrial and Applied Mathematics, 2013.
34. G.W. Wasilkowski and H. Wozniakowski. Explicit cost bounds of algorithms for multivariate tensor product problems. *Journal of Complexity*, 11(1):1 – 56, 1995.

Recent publications:

MATHEMATICS INSTITUTE OF COMPUTATIONAL SCIENCE AND ENGINEERING
Section of Mathematics
Ecole Polytechnique Fédérale
CH-1015 Lausanne

- 44.2013** A. KOSHAKJI, A. QUARTERONI, G. ROZZA:
Free form deformation techniques applied to 3D shape optimization problems
- 45.2013** J. E. CATRILLON-CANDAS, F. NOBILE, R. F. TEMPONE:
Analytic regularity and collocation approximation for PDEs with random domain deformations
- 01.2014** GIOVANNI MIGLIORATI:
Multivariate Markov-type and Nicolskii-type inequalities for polynomials associated with downward closed multi-index sets
- 02.2014** FEDERICO NEGRI, ANDREA MANZONI, GIANLUIGI ROZZA:
Certified reduced basis method for parametrized optimal control problems governed by the Stokes equations
- 03.2014** CEDRIC EFFENBERGER, DANIEL KRESSNER:
On the residual inverse iteration for nonlinear eigenvalue problems admitting a Rayleigh functional
- 04.2014** TAKAHITO KASHIWABARA, CLAUDIA M. COLCIAGO, LUCA DEDÈ, ALFIO QUARTERONI:
Numerical Well-posedness, regularity, and convergence analysis of the finite element approximation of a generalized Robin boundary value problem
- 05.2014** BJÖRN ADLERBORN, BO KAGSTRÖM, DANIEL KRESSNER:
A parallel QZ algorithm for distributed memory HPC systems
- 06.2014** MICHELE BENZI, SIMONE DEPARIS, GWENOL GRANDPERRIN, ALFIO QUARTERONI:
Parameter estimates for the relaxed dimensional factorization preconditioner and application to hemodynamics
- 07.2014** ASSYR ABDULLE, YUN BAI:
Reduced order modelling numerical homogenization
- 08.2014** ANDREA MANZONI, FEDERICO NEGRI:
Rigorous and heuristic strategies for the approximation of stability factors in nonlinear parametrized PDEs
- 09.2014** PENG CHEN, ALFIO QUARTERONI:
A new algorithm for high-dimensional uncertainty quantification problems based on dimension-adaptive and reduced basis methods
- 10.2014** NATHAN COLLIER, ABDUL-LATEEF HAJI-ALI, FABIO NOBILE, ERIK VON SCHWERIN, RAÚL TEMPONE:
A continuation multilevel Monte Carlo algorithm
- 11.2014** LUKA GRUBISIC, DANIEL KRESSNER:
On the eigenvalue decay of solutions to operator Lyapunov equations
- 12.2014** FABIO NOBILE, LORENZO TAMELLINI, RAÚL TEMPONE:
Convergence of quasi-optimal sparse grid approximation of Hilbert-valued functions: application to random elliptic PDEs