



HAL
open science

Generalized spectral decomposition method for solving stochastic finite element equations: invariant subspace problem and dedicated algorithms

Anthony Nouy

► **To cite this version:**

Anthony Nouy. Generalized spectral decomposition method for solving stochastic finite element equations: invariant subspace problem and dedicated algorithms. *Computer Methods in Applied Mechanics and Engineering*, 2008, 197 (51-52), pp.4718-4736. 10.1016/j.cma.2008.06.012 . hal-00366613

HAL Id: hal-00366613

<https://hal.science/hal-00366613v1>

Submitted on 9 Mar 2009

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Generalized spectral decomposition method for solving stochastic finite element equations: invariant subspace problem and dedicated algorithms

Anthony Nouy^{*}

*Research Institute in Civil Engineering and Mechanics (GeM), University of
Nantes, Ecole Centrale Nantes, CNRS, 2 rue de la Houssinière, B.P. 92208,
44322 Nantes Cedex 3, FRANCE*

Abstract

Stochastic Galerkin methods have become a significant tool for the resolution of stochastic partial differential equations (SPDE). However, they suffer from prohibitive computational times and memory requirements when dealing with large scale applications and high stochastic dimensionality. Some alternative techniques, based on the construction of suitable reduced deterministic or stochastic bases, have been proposed in order to reduce these computational costs. Recently, a new approach, based on the concept of generalized spectral decomposition (GSD), has been introduced for the definition and the automatic construction of reduced bases. In this paper, the concept of GSD, initially introduced for a class of linear elliptic SPDE, is extended to a wider class of stochastic problems. The proposed definition of the GSD leads to the resolution of an invariant subspace problem, which is interpreted as an eigen-like problem. This interpretation allows the construction of efficient numerical algorithms for building optimal reduced bases, which are associated with dominant generalized eigenspaces. The proposed algorithms, by separating the resolution of reduced stochastic and deterministic problems, lead to drastic computational savings. Their efficiency is illustrated on several examples, where they are compared to classical resolution techniques.

Key words: Computational Stochastic Mechanics, Stochastic partial differential equations, Stochastic Finite Element, Generalized Spectral Decomposition, Invariant Subspace problem, Stochastic Model Reduction

^{*} Corresponding author. Tel.: +33(0)2-51-12-55-20; Fax: +33(0)2-51-12-52-52
Email address: anthony.nouy@univ-nantes.fr (Anthony Nouy).

1 Introduction

Computer simulations have become an essential tool for the quantitative prediction of the response of physical models. The need to improve the reliability of numerical predictions often requires taking into account uncertainties inherent to these models.

Uncertainties, either epistemic or aleatory, are commonly modeled within a probabilistic framework. For many physical models, it leads to the resolution of a stochastic partial differential equation (SPDE) where the operator, the right-hand side, the boundary conditions or even the domain, depend on a set of random variables. Many numerical methods have been proposed for the approximation of such SPDEs. In particular, stochastic Galerkin methods [1–4] have received a growing interest in the last decade. They allow the obtention of a decomposition of the solution on a suitable approximation basis, the coefficients of the decomposition being obtained by solving a large system of equations. These methods, which lead to high quality predictions, rely on a strong mathematical basis. That allows deriving *a priori* error estimators [5–7] but also *a posteriori* error estimators [8,9] and therefore to develop adaptive approximation techniques. However, many complex applications require a fine discretization at both deterministic and stochastic levels. This dramatically increases the dimension of approximation spaces and therefore of the resulting system of equations. The use of classical solvers in a black box fashion generally leads to prohibitive computational times and memory requirements. The reduction of these computational costs has now become a key question for the development of stochastic Galerkin methods and their transfer towards large scale and industrial applications.

Some alternative resolution techniques have been investigated over the last years in order to drastically reduce computational costs induced by the use of Galerkin approximation schemes. Some of these works rely on the construction of reduced deterministic bases or stochastic bases (sets of random variables) in order to decrease the size of the problem [3,10,11]. These techniques usually start from the assertion that optimal deterministic and stochastic bases can be obtained by using a classical spectral decomposition of the solution (namely a Karhunen-Loève or Hilbert Karhunen-Loève expansion). The solution being not known *a priori*, the basic idea of these techniques is to compute an approximation of the “ideal” spectral decomposition by ad hoc numerical strategies. The obtained set of deterministic vectors (resp. random variables) is then considered as a good candidate for a reduced deterministic (resp. stochastic) basis on which the initial stochastic problem can be solved at a lower cost. Let us here mention that this kind of decomposition has already been introduced in various domains of application such as functional data analysis [12], image analysis [13], dynamical model reduction [14,15], etc. In other contexts, it is

also known as Principal Component Analysis, Proper Orthogonal Decomposition or Singular Value Decomposition.

In [16], a new approach has been proposed to define and compute suitable reduced bases, without *a priori* knowing the solution nor an approximation of it. This method, which is inspired by a technique for solving deterministic evolution equations [17–19], is based on the concept of generalized spectral decomposition (GSD). It consists in defining an optimality criterion of the decomposition based on the operator and right-hand side of the stochastic problem. In the case of a linear elliptic symmetric SPDE, the obtained decomposition can be interpreted as a generalized Karhunen-Loève expansion of the right-hand side in the metric induced by the operator. In [16], it has been shown that corresponding optimal reduced bases were solution of an optimization problem on a functional which can be interpreted as an extended Rayleigh quotient associated with an eigen-like problem. In order to solve this problem, a power-type algorithm has been proposed. This algorithm, by separating the resolution of reduced deterministic problems and reduced stochastic problems, has led to significant computational savings.

The aim of this paper is to extend the concept of generalized spectral decomposition to a wider class of stochastic problems and to provide ad hoc efficient numerical strategies for its construction. The proposed definition of the GSD leads to the resolution of an invariant subspace problem, which in fact can be interpreted as an eigen-like problem. This interpretation allows the development of suitable algorithms for the construction of the decomposition. Algorithms are inspired by resolution techniques for classical eigenproblems, such as subspace iterations or Arnoldi techniques [20]. Significant computational savings are obtained with these new algorithms, in comparison with classical resolution techniques but also with previous GSD algorithms proposed in [16].

The proposed method will be presented on a generic discretized linear problem, encountered in many physical situations, without taking care of the initial “continuous problem” and of the discretization techniques at the deterministic and stochastic levels. In this paper, we consider that the solution of the fully discretized problem is our reference solution. The proposed method then leads to an approximation of this reference approximate solution. The study of approximation error, *i.e.* the distance between the reference solution and the solution of the continuous problem, is beyond the scope of this paper. For details, the reader can refer to [4–9].

The outline of the paper is as follows. In section 2, we briefly recall the principles of stochastic Galerkin methods leading to the definition of a fully discretized version of the stochastic problem. Section 3 introduces some possible strategies for building deterministic or stochastic reduced bases. In section 4, the principles of the generalized spectral decomposition method (GSD) are in-

troduced. In particular, some mathematical considerations allow us to exhibit the underlying eigen-like problem that defines the GSD. Section 5 is devoted to the presentation of different algorithms for building the GSD. In sections 6 and 7, the method is applied to two model problems: the first one is a linear elasticity problem and the second one is based on transient heat equation. Those model problems illustrate the capabilities of the method respectively for elliptic and parabolic stochastic partial differential equations.

2 Stochastic Galerkin methods

2.1 Stochastic modeling and discretization

We adopt a probabilistic modeling of the uncertainties. We consider that the probabilistic content of the stochastic problem can be represented by a finite dimensional probability space (Θ, \mathcal{B}, P) . $\Theta \subset \mathbb{R}^m$ is the space of elementary events, \mathcal{B} an associated σ -algebra and P the probability measure. We consider that a preliminary approximation step has been performed at the deterministic level and that the stochastic problem reduces to the resolution of the following system of stochastic equations: find a random vector $\mathbf{u} : \theta \in \Theta \mapsto \mathbf{u}(\theta) \in \mathbb{R}^n$ such that we have P-almost surely

$$\mathbf{A}(\theta)\mathbf{u}(\theta) = \mathbf{b}(\theta), \quad (1)$$

where $\mathbf{A} : \Theta \rightarrow \mathbb{R}^{n \times n}$ is a random matrix and $\mathbf{b} : \Theta \rightarrow \mathbb{R}^n$ is a random vector. For the sake of clarity and generality, we do not focus on the way to obtain this semi-discretized problem. In the following, we will admit that the continuous and discretized problems are well-posed, which means that the continuous problem and the approximation technique have “good mathematical properties”. Sections 6 and 7 will illustrate two continuous model problems and associated approximation techniques that lead to a system of type (1) (by introducing usual spatial and temporal discretizations). Now, we introduce an ad-hoc real-valued random function space \mathcal{S} , classically the space of second order random variables $L^2(\Theta, dP)$, such that a weak formulation of the stochastic problem (1) can be introduced. This weak formulation, whose solution is not necessarily solution of (1), reads: find $\mathbf{u} \in \mathbb{R}^n \otimes \mathcal{S} \cong (\mathcal{S})^n$ such that

$$E(\mathbf{v}^T \mathbf{A} \mathbf{u}) = E(\mathbf{v}^T \mathbf{b}) \quad \forall \mathbf{v} \in \mathbb{R}^n \otimes \mathcal{S}. \quad (2)$$

Approximation technique at the stochastic level consists in introducing a suitable finite dimensional approximation space

$$\mathcal{S}_P = \{v(\theta) = \sum_{\alpha \in \mathcal{J}_P} v_\alpha H_\alpha(\theta), v_\alpha \in \mathbb{R}, H_\alpha \in \mathcal{S}\}, \quad (3)$$

where $\{H_\alpha\}_{\alpha \in \mathcal{J}_P}$ is a basis of \mathcal{S}_P , and $\mathcal{J}_P = \{\alpha_i, i = 1 \dots P\}$ is a set of P indices. The approximate solution $\mathbf{u} \in \mathbb{R}^n \otimes \mathcal{S}_P$ then reads

$$\mathbf{u}(\theta) = \sum_{\alpha \in \mathcal{J}_P} \mathbf{u}_\alpha H_\alpha(\theta). \quad (4)$$

A classical way to define the approximation is to use a Galerkin orthogonality criterion reading

$$E(\mathbf{v}^T \mathbf{A} \mathbf{u}) = E(\mathbf{v}^T \mathbf{b}) \quad \forall \mathbf{v} \in \mathbb{R}^n \otimes \mathcal{S}_P, \quad (5)$$

where E denotes the mathematical expectation. System (5) is equivalent to the following system of $n \times P$ equations:

$$\sum_{\beta \in \mathcal{J}_P} E(\mathbf{A} H_\alpha H_\beta) \mathbf{u}_\beta = E(H_\alpha \mathbf{b}) \quad \forall \alpha \in \mathcal{J}_P. \quad (6)$$

Several choices have been proposed for the construction of a stochastic approximation basis in $L^2(\Theta, dP)$: polynomial chaos [1], generalized polynomial chaos [21,22], finite elements [6,4], or multi-wavelets [23,24]. Such a choice depends on the regularity of the solution at the stochastic level. Several techniques have been investigated for the adaptive choice of this basis, based on *a posteriori* error estimation with respect to the continuous model [25,7–9]. For well-posed approximate problems, the solution of (5) weakly converges with P (in a mean-square sense) towards the solution of problem (2). In this paper, we will consider that this approximation basis is given (fixed P). The approximate solution of the fully discretized problem (5) will then be considered as our reference solution. The study of the stochastic approximation error, *i.e.* the distance between solutions of equations (5) and (2), is beyond the scope of this article.

2.2 Classical techniques to solve the discretized problem

System (6) can be written in the following block-matrix form:

$$\begin{pmatrix} E(\mathbf{A} H_{\alpha_1} H_{\alpha_1}) & \dots & E(\mathbf{A} H_{\alpha_1} H_{\alpha_P}) \\ \vdots & \ddots & \vdots \\ E(\mathbf{A} H_{\alpha_P} H_{\alpha_1}) & \dots & E(\mathbf{A} H_{\alpha_P} H_{\alpha_P}) \end{pmatrix} \begin{pmatrix} \mathbf{u}_{\alpha_1} \\ \vdots \\ \mathbf{u}_{\alpha_P} \end{pmatrix} = \begin{pmatrix} E(\mathbf{b} H_{\alpha_1}) \\ \vdots \\ E(\mathbf{b} H_{\alpha_P}) \end{pmatrix} \quad (7)$$

System (7) is a huge system of $n \times P$ equations. Krylov-type iterative techniques are classically used to solve this system [2,26–28], such as Preconditioned Conjugate Gradient for symmetric problems (PCG), Conjugate Gradient Square (CGS), etc. These algorithms take advantage of the sparsity of the system, coming both from the sparsity of random matrix \mathbf{A} and from classical orthogonality properties of the stochastic approximation basis [2]. These algorithms are quite efficient. However, when dealing with large scale applications (large n), and when working with high stochastic dimension, with a fine discretization at the stochastic level (requiring large P), computational costs and memory requirements induced by these techniques increase dramatically.

3 Construction of reduced approximation basis

A rising tendency in the context of computational stochastic methods consists in trying to obtain pertinent reduced models in order to drastically decrease the size of the problem when dealing with large scale applications and high stochastic dimension. The idea is to build a small set of M deterministic vectors $\mathbf{U}_i \in \mathbb{R}^n$ (or M random variables $\lambda_i \in \mathcal{S}_P$), with $M \ll n$ (or $M \ll P$), and then to compute the associated random variables λ_i (or deterministic vectors \mathbf{U}_i). The approximate solution of problem (5) can then be written:

$$\mathbf{u}(\theta) \approx \sum_{i=1}^M \lambda_i(\theta) \mathbf{U}_i. \quad (8)$$

In the following, we will denote by $\mathbf{W} = (\mathbf{U}_1 \dots \mathbf{U}_M) \in \mathbb{R}^{n \times M}$ the matrix whose columns are the deterministic vectors and $\mathbf{\Lambda} = (\lambda_1 \dots \lambda_M)^T \in \mathbb{R}^M \otimes \mathcal{S}_P$ the random vector whose components are the random variables. Decomposition (8) can then be written in a matrix form

$$\mathbf{u}(\theta) \approx \mathbf{W} \mathbf{\Lambda}(\theta). \quad (9)$$

3.1 Working on a reduced deterministic basis

Let us first suppose that a reduced deterministic basis has been computed. Then, \mathbf{W} being fixed, a natural definition of $\mathbf{\Lambda} \in \mathbb{R}^M \otimes \mathcal{S}_P$ arises from the following Galerkin orthogonality criterion:

$$E(\tilde{\mathbf{\Lambda}}^T (\mathbf{W}^T \mathbf{A} \mathbf{W}) \mathbf{\Lambda}) = E(\tilde{\mathbf{\Lambda}}^T \mathbf{W}^T \mathbf{b}) \quad \forall \tilde{\mathbf{\Lambda}} \in \mathbb{R}^M \otimes \mathcal{S}_P. \quad (10)$$

Problem (10) defines the approximation of problem (5) in the approximation subspace $\text{span}(\{\mathbf{U}_i\}_{i=1}^M) \otimes \mathcal{S}_P$. It can be interpreted as a classical stochastic

Galerkin problem on the reduced deterministic basis which is spanned¹ by the \mathbf{U}_i . As problem (5), problem (10) can be written in the following block-matrix form:

$$\begin{pmatrix} E(\mathbf{W}^T \mathbf{A} \mathbf{W} H_{\alpha_1} H_{\alpha_1}) & \dots & E(\mathbf{W}^T \mathbf{A} \mathbf{W} H_{\alpha_1} H_{\alpha_P}) \\ \vdots & \ddots & \vdots \\ E(\mathbf{W}^T \mathbf{A} \mathbf{W} H_{\alpha_P} H_{\alpha_1}) & \dots & E(\mathbf{W}^T \mathbf{A} \mathbf{W} H_{\alpha_P} H_{\alpha_P}) \end{pmatrix} \begin{pmatrix} \boldsymbol{\Lambda}_{\alpha_1} \\ \vdots \\ \boldsymbol{\Lambda}_{\alpha_P} \end{pmatrix} = \begin{pmatrix} E(\mathbf{W}^T \mathbf{b} H_{\alpha_1}) \\ \vdots \\ E(\mathbf{W}^T \mathbf{b} H_{\alpha_P}) \end{pmatrix}, \quad (11)$$

which is a system of $M \times P$ equations. Let us note that the reduced random matrix $\mathbf{W}^T \mathbf{A} \mathbf{W}$ is generally full. However, system (11) keeps its block sparsity pattern coming from orthogonality properties of the stochastic basis. This system can then be solved by classical Krylov-type iterative techniques mentioned in section 2.2.

3.2 Working on a reduced stochastic basis

Let us now suppose that a reduced stochastic basis has been computed. Then, $\boldsymbol{\Lambda}$ being fixed, a natural definition of $\mathbf{W} \in \mathbb{R}^{n \times M}$ arises from the following Galerkin orthogonality criterion:

$$E(\boldsymbol{\Lambda}^T \widetilde{\mathbf{W}}^T \mathbf{A} \mathbf{W} \boldsymbol{\Lambda}) = E(\boldsymbol{\Lambda}^T \widetilde{\mathbf{W}}^T \mathbf{b}) \quad \forall \widetilde{\mathbf{W}} \in \mathbb{R}^{n \times M}. \quad (12)$$

Problem (12) defines the approximation of problem (5) in the approximation subspace $\mathbb{R}^n \otimes \text{span}(\{\lambda_i\}_{i=1}^M)$, *i.e.* on the reduced basis of \mathcal{S}_P which is spanned² by the λ_i . It can be interpreted as a deterministic problem that can be written in the following block-matrix form:

$$\begin{pmatrix} E(\mathbf{A} \lambda_1 \lambda_1) & \dots & E(\mathbf{A} \lambda_1 \lambda_M) \\ \vdots & \ddots & \vdots \\ E(\mathbf{A} \lambda_M \lambda_1) & \dots & E(\mathbf{A} \lambda_M \lambda_M) \end{pmatrix} \begin{pmatrix} \mathbf{U}_1 \\ \vdots \\ \mathbf{U}_M \end{pmatrix} = \begin{pmatrix} E(\mathbf{b} \lambda_1) \\ \vdots \\ E(\mathbf{b} \lambda_M) \end{pmatrix}, \quad (13)$$

which is a system of $M \times n$ equations. Let us note that this block-system is generally full in the block sense but that each block inherits from the sparsity pattern of random matrix \mathbf{A} . This system can then be solved by classical direct or iterative solvers (or block solvers), the choice depending on its size.

¹ $\text{span}(\{\mathbf{U}_i\}_{i=1}^M) = \{\sum_{i=1}^M a_i \mathbf{U}_i \in \mathbb{R}^n; a_i \in \mathbb{R}\}$ is the linear subspace of \mathbb{R}^n spanned by vectors $\{\mathbf{U}_i\}_{i=1}^M$.

² $\text{span}(\{\lambda_i\}_{i=1}^M) = \{\sum_{i=1}^M a_i \lambda_i \in \mathcal{S}_P; a_i \in \mathbb{R}\}$ is the linear subspace of \mathcal{S}_P spanned by random variables $\{\lambda_i\}_{i=1}^M$.

3.3 How to define pertinent reduced bases ?

Now, the key question is: how can we define an optimal reduced basis (of deterministic vectors or random variables) which leads to the best accuracy for a given order M of decomposition ?

A possible answer is based on the following property: “the classical spectral decomposition, *i.e.* the Karhunen-Loève expansion, is the optimal reduced decomposition of the solution \mathbf{u} with respect to the natural norm in $L^2(\Theta, dP; \mathbb{R}^n)$ ”. This norm is defined as follows:

$$\|\mathbf{u}\|^2 = E(\mathbf{u}^T \mathbf{u}). \quad (14)$$

Of course, changing the norm on the solution leads to another optimal spectral decomposition, called Hilbert-Karhunen Loève decomposition in the continuous framework [13,11]. Then, if only we could compute the classical spectral decomposition of the solution, we could consider the obtained random variables (resp. deterministic vectors) as good candidates for building a reduced stochastic basis (resp. deterministic basis). The problem is that the solution, and *a fortiori* its correlation structure, is not known *a priori*. Several techniques have already been introduced to get an approximation of this spectral decomposition. In [3], the authors proposed to compute an approximation of the correlation matrix, by using Neumann expansion of \mathbf{A} , and to compute its M dominant eigenvectors. The obtained vectors, considered as a reduced deterministic basis, can then be interpreted as an approximation of vectors of the exact spectral decomposition of \mathbf{u} . In [10,11], the authors propose to first introduce a coarse approximation at the deterministic level (*e.g.* by using a coarse finite element mesh), leading to the resolution of a coarse stochastic problem $\mathbf{A}_c(\theta)\mathbf{u}_c(\theta) = \mathbf{b}_c(\theta)$, with $\mathbf{u}_c \in \mathbb{R}^{n_c} \otimes \mathcal{S}_P$, $n_c \ll n$. Then, a Karhunen-Loève (or Hilbert-Karhunen Loève) expansion of \mathbf{u}_c can be performed. After truncation at order M , it leads to the M desired random variables $\lambda_i \in \mathcal{S}_P$, considered as a reduced stochastic basis which can be used to solve the initial fine stochastic problem.

Another possible answer consists in defining another optimality criterion of the decomposition (*i.e.* of the reduced basis) which could allow its computation without *a priori* knowing the solution nor even an approximation of it. This answer has been formulated in [16] by introducing the concept of “generalized spectral decomposition” (GSD). The arising technique can be seen as a general technique for the automatic construction of both deterministic and stochastic bases simultaneously, the obtained bases being optimal with respect to operator and right-hand side of the problem. In this paper, this concept of generalized spectral decomposition will be presented in a more general context. The definition of the decomposition as a solution of an eigen-like problem will be clarified in section 4. This interpretation will help us to propose improved

algorithms for computing the generalized spectral decomposition in section 5.

Finally, let us mention another technique for the *a priori* construction of reduced bases, called the Stochastic Reduced Basis Method [29,30]. In this method, the reduced basis composed by the \mathbf{U}_i is chosen in the M-dimensional Krylov subspace of random matrix \mathbf{A} , associated with the right-hand side \mathbf{b} : $\mathbf{U}_i = \mathbf{A}\mathbf{U}_{i-1}$, for $i = 2 \dots M$, with $\mathbf{U}_1 = \mathbf{b}$. This method differs from the above techniques because the \mathbf{U}_i are also random. Its main drawback is that the computation of \mathbf{U}_i from \mathbf{U}_{i-1} requires (in practice) a stochastic projection on \mathcal{S}_P , which induces a loss of accuracy and then restricts the use of this technique to a low dimensional Krylov subspace. Another drawback, compared to the above techniques, is that it does not circumvent the problem of memory requirements since random vectors $\mathbf{U}_i \in \mathbb{R}^n \otimes \mathcal{S}_P$ have to be stored.

4 Generalized spectral decomposition

The idea of the generalized spectral decomposition method (GSD), introduced in [16], is to try to find an optimal approximation of problem (5) in the following form :

$$\mathbf{u}(\theta) \approx \sum_{i=1}^M \lambda_i(\theta) \mathbf{U}_i, \quad (15)$$

where the $\lambda_i \in \mathcal{S}_P$ are random variables and the $\mathbf{U}_i \in \mathbb{R}^n$ are deterministic vectors, none of these quantities being known *a priori*. A decomposition of this type is said optimal if the number of terms M is minimum for a given quality of approximation. The set of deterministic vectors (resp. random variables) is then considered as an optimal deterministic (resp. stochastic) reduced basis. In this section, we introduce a natural definition of this decomposition and show that the deterministic vectors (resp. random variables) are solution of an invariant subspace problem, which is interpreted as an eigen-like problem. We also recall and revisit the results obtained in [16], which corresponds to a particular case of the present article. In this section, we do not focus on algorithms allowing the computation of this decomposition, which is the aim of section 5.

4.1 Preliminary remarks

Let $\mathbf{W} = (\mathbf{U}_1 \dots \mathbf{U}_M) \in \mathbb{R}^{n \times M}$ be the matrix whose columns are the deterministic vectors and $\mathbf{\Lambda} = (\lambda_1 \dots \lambda_M)^T \in \mathbb{R}^M \otimes \mathcal{S}_P$ the random vector whose components are the random variables. Decomposition (15) is then written in

a matrix form

$$\mathbf{u}(\theta) \approx \mathbf{W}\mathbf{\Lambda}(\theta). \quad (16)$$

Due to our definition of optimality, it is natural to impose on the \mathbf{U}_i to be linearly independent, *i.e.* to span a M -dimensional linear subspace of \mathbb{R}^n : $\dim(\text{span}(\{\mathbf{U}_i\}_{i=1}^M)) = M$. Equivalently, we impose on \mathbf{W} to belong to the set of n -by- M full rank matrices, called the noncompact Stiefel manifold [31]:

$$\mathbb{S}_{n,M} = \{\mathbf{W} \in \mathbb{R}^{n \times M}; \text{rank}(\mathbf{W}) = M\}. \quad (17)$$

Indeed, if \mathbf{W} were not full rank, decomposition (15) could be equivalently rewritten as a decomposition of order $M' < M$. Then, the order M approximation would not be optimal since it would exist an order M' approximation, with $M' < M$, leading to the same approximation. For the same reason, it is natural to look for “linearly independent” λ_i , *i.e.* such that they span a M -dimensional linear subspace of \mathcal{S}_P . We then introduce the following space:

$$\mathbb{S}_{P,M}^* = \{\mathbf{\Lambda} = (\lambda_1 \dots \lambda_M)^T \in \mathbb{R}^M \otimes \mathcal{S}_P; \dim(\text{span}(\{\lambda_i\}_{i=1}^M)) = M\}. \quad (18)$$

In the following, we will use the abuse of notation: $\text{span}(\mathbf{W}) \equiv \text{span}(\{\mathbf{U}_i\}_{i=1}^M)$ and $\text{span}(\mathbf{\Lambda}) \equiv \text{span}(\{\lambda_i\}_{i=1}^M)$.

Remark 1 A function $\lambda_i \in \mathcal{S}_P$ can be identified with a vector $\boldsymbol{\lambda}_i \in \mathbb{R}^P$, whose components are the coefficients of λ_i on the basis $\{H_\alpha\}$ of \mathcal{S}_P , *i.e.* $\boldsymbol{\lambda}_i = (\dots \lambda_{i,\alpha} \dots)^T$. In the same way, a random vector $\mathbf{\Lambda} = (\lambda_1 \dots \lambda_M)^T \in \mathbb{R}^M \otimes \mathcal{S}_P$ can be identified with a matrix $\mathbf{L} \in \mathbb{R}^{P \times M}$ whose column vectors are the vectors $\boldsymbol{\lambda}_i$, *i.e.* $\mathbf{L} = (\boldsymbol{\lambda}_1 \dots \boldsymbol{\lambda}_M) = (\dots \boldsymbol{\Lambda}_\alpha \dots)^T$. The property “the λ_i are linearly independent” is simply equivalent to “the vectors $\boldsymbol{\lambda}_i$ are linearly independent” or “ $\text{rank}(\mathbf{L})=M$ ”. We then clearly have the following isomorphism: $\mathbb{S}_{P,M}^* \cong \mathbb{S}_{P,M}$.

4.2 Definition of the generalized spectral decomposition

On one hand, if \mathbf{W} were fixed, a natural definition of $\mathbf{\Lambda}$ would arise from the resolution of problem (5) in the approximation subspace $\text{span}(\mathbf{W}) \otimes \mathcal{S}_P$: find $\mathbf{\Lambda} \in \mathbb{R}^M \otimes \mathcal{S}_P$ such that

$$E(\tilde{\mathbf{\Lambda}}^T (\mathbf{W}^T \mathbf{A} \mathbf{W}) \mathbf{\Lambda}) = E(\tilde{\mathbf{\Lambda}}^T \mathbf{W}^T \mathbf{b}) \quad \forall \tilde{\mathbf{\Lambda}} \in \mathbb{R}^M \otimes \mathcal{S}_P. \quad (19)$$

The associated system of equations, written in a block-matrix form, is given in (11). We denote by $\mathbf{\Lambda} = \mathbf{f}(\mathbf{W})$ its solution, where \mathbf{f} is a mapping defined as follows:

$$\mathbf{f} : \mathbf{W} \in \mathbb{S}_{n,M} \mapsto \mathbf{\Lambda} = \mathbf{f}(\mathbf{W}) \in \mathbb{R}^M \otimes \mathcal{S}_P. \quad (20)$$

On the other hand, if $\mathbf{\Lambda}$ were fixed, a natural definition of \mathbf{W} would arise from the resolution of problem (5) in the approximation subspace $\mathbb{R}^n \otimes \text{span}(\mathbf{\Lambda})$: find $\mathbf{W} \in \mathbb{R}^{n \times M}$ such that

$$E(\mathbf{\Lambda}^T \widetilde{\mathbf{W}}^T \mathbf{A} \mathbf{W} \mathbf{\Lambda}) = E(\mathbf{\Lambda}^T \widetilde{\mathbf{W}}^T \mathbf{b}) \quad \forall \widetilde{\mathbf{W}} \in \mathbb{R}^{n \times M}. \quad (21)$$

The associated system of equations, written in a block-matrix form, is given in (13). Let us denote by $\mathbf{W} = \mathbf{F}(\mathbf{\Lambda})$ the solution of equation (21), where \mathbf{F} is the following mapping:

$$\mathbf{F} : \mathbf{\Lambda} \in \mathbb{S}_{P,M}^* \mapsto \mathbf{W} = \mathbf{F}(\mathbf{\Lambda}) \in \mathbb{R}^{n \times M}. \quad (22)$$

When neither \mathbf{W} nor $\mathbf{\Lambda}$ are fixed, it is then natural to look for a couple $(\mathbf{W}, \mathbf{\Lambda})$ that verifies both equations (19) and (21) simultaneously. The problem then reads: find $(\mathbf{W}, \mathbf{\Lambda}) \in \mathbb{S}_{n,M} \times \mathbb{S}_{P,M}^*$ such that

$$\mathbf{W} = \mathbf{F}(\mathbf{\Lambda}) \quad \text{and} \quad \mathbf{\Lambda} = \mathbf{f}(\mathbf{W}). \quad (23)$$

We will see in the following that problem (23) can be interpreted as an invariant subspace problem, which can be interpreted as an eigen-like problem.

4.3 Non-uniqueness of the decomposition - equivalence class of solutions

The couples $(\mathbf{W}, \mathbf{\Lambda})$ and $(\mathbf{W}\mathbf{P}, \mathbf{P}^{-1}\mathbf{\Lambda})$ clearly lead to the same decomposition for all $\mathbf{P} \in \mathbb{GL}_M$.³ We can then define an equivalence class of couples in $\mathbb{S}_{n,M} \times \mathbb{S}_{P,M}^*$ that leads to the same approximation:

$$\begin{aligned} & (\mathbf{W}_1, \mathbf{\Lambda}_1) \sim (\mathbf{W}_2, \mathbf{\Lambda}_2) \\ \Leftrightarrow & \{ \mathbf{W}_1 = \mathbf{W}_2 \mathbf{P}, \mathbf{\Lambda}_1 = \mathbf{P}^{-1} \mathbf{\Lambda}_2, \mathbf{P} \in \mathbb{GL}_M \}. \end{aligned} \quad (24)$$

If $(\mathbf{W}, \mathbf{\Lambda})$ verifies problem (23), the obtained decomposition can be equivalently written in terms of $\mathbf{\Lambda}$ or \mathbf{W} :

$$\mathbf{u} \approx \mathbf{F}(\mathbf{\Lambda})\mathbf{\Lambda} \quad \text{or} \quad \mathbf{u} \approx \mathbf{W}\mathbf{f}(\mathbf{W}) \quad (25)$$

Proposition 2 *The mappings \mathbf{f} and \mathbf{F} verify the following homogeneity properties: $\forall \mathbf{P} \in \mathbb{GL}_M$,*

$$\begin{aligned} \mathbf{f}(\mathbf{W}\mathbf{P}) &= \mathbf{P}^{-1}\mathbf{f}(\mathbf{W}) \\ \mathbf{F}(\mathbf{P}\mathbf{\Lambda}) &= \mathbf{F}(\mathbf{\Lambda})\mathbf{P}^{-1} \end{aligned}$$

Proposition 2 allows us to introduce equivalence classes for $\mathbf{\Lambda}$ and \mathbf{W} separately, defined by the following proposition.

³ \mathbb{GL}_M denotes the linear group of invertible matrices in $\mathbb{R}^{M \times M}$

Proposition 3 *The obtained decomposition $\mathbf{F}(\Lambda)\Lambda$ (resp. $\mathbf{Wf}(\mathbf{W})$) is unique on the equivalence class on $\mathbb{S}_{P,M}^*$ (resp. $\mathbb{S}_{n,M}$) defined by the equivalence relation $\overset{\Lambda}{\sim}$ (resp. $\overset{\mathbf{W}}{\sim}$) where*

$$\Lambda_1 \overset{\Lambda}{\sim} \Lambda_2 \Leftrightarrow \{\Lambda_1 = \mathbf{P}\Lambda_2, \mathbf{P} \in \mathbb{GL}_M\} \quad (26)$$

$$\mathbf{W}_1 \overset{\mathbf{W}}{\sim} \mathbf{W}_2 \Leftrightarrow \{\mathbf{W}_1 = \mathbf{W}_2\mathbf{P}, \mathbf{P} \in \mathbb{GL}_M\} \quad (27)$$

$\mathbf{W}_1 \overset{\mathbf{W}}{\sim} \mathbf{W}_2$ implies that $(\mathbf{W}_1, \mathbf{f}(\mathbf{W}_1)) \sim (\mathbf{W}_2, \mathbf{f}(\mathbf{W}_2))$ and $\Lambda_1 \overset{\Lambda}{\sim} \Lambda_2$ implies that $(\mathbf{F}(\Lambda_1), \Lambda_1) \sim (\mathbf{F}(\Lambda_2), \Lambda_2)$

Remark 4 *The non-uniqueness of the decomposition offers a flexibility in the choice of deterministic vectors or random variables. For example, it is possible to choose the particular solution corresponding to orthonormal deterministic vectors (or random variables), which can be interesting from a computational point of view.*

4.4 Interpretation as an eigen-like problem

For the interpretation of the generalized spectral decomposition, we will focus on a formulation on \mathbf{W} . Problem (23) can be rewritten as a problem on \mathbf{W} :

$$\mathbf{W} = \mathbf{F} \circ \mathbf{f}(\mathbf{W}) \quad (28)$$

Let us now introduce the following mapping:

$$\mathbf{T} : \mathbf{W} \in \mathbb{S}_{n,M} \mapsto \mathbf{F} \circ \mathbf{f}(\mathbf{W}) \in \mathbb{R}^{n \times M} \quad (29)$$

Equation (28) then reads

$$\mathbf{W} = \mathbf{T}(\mathbf{W}) \quad (30)$$

From mappings homogeneity properties (proposition 2), we deduce the following homogeneity property for \mathbf{T} :

Proposition 5 *The mapping \mathbf{T} verifies the following homogeneity property : $\forall \mathbf{P} \in \mathbb{GL}_M$,*

$$\mathbf{T}(\mathbf{W}\mathbf{P}) = \mathbf{T}(\mathbf{W})\mathbf{P}$$

Regarding proposition 5, if $\mathbf{W} \in \mathbb{S}_{n,M}$ verifies equation (30), all matrices in its equivalence class, defined by (27), also verifies this equation. The problem can then be reformulated in the quotient space $Gr_{n,M} = (\mathbb{S}_{n,M} / \overset{\mathbf{W}}{\sim})$, which can be identified with the set of M -dimensional linear subspace of \mathbb{R}^n , the so

called Grassmann manifold (see *e.g.* [31,32]). An element $\mathbf{W} \in Gr_{n,M}$ can be associated with all matrices $\mathbf{W} \in \mathbb{R}^{n \times M}$ such that $span(\mathbf{W}) = \mathbf{W}$, *i.e.* whose column vectors span \mathbf{W} . The problem can then be interpreted as follows: find a M -dimensional linear subspace of \mathbb{R}^n such that for all $\mathbf{W} \in \mathbb{S}_{n,M}$ that spans this subspace, equation (30) holds. Equation (30) can then be rewritten:

$$\mathbf{W} = \mathcal{J}(\mathbf{W}) \quad (31)$$

where \mathcal{J} is the following mapping :

$$\mathcal{J} : \mathbf{W} = span(\mathbf{W}) \in Gr_{n,M} \mapsto span(\mathbf{T}(\mathbf{W})) \quad (32)$$

Equation (31) means that we look for a linear subspace \mathbf{W} that is invariant by the mapping \mathcal{J} . This is a fixed point problem on the Grassmann manifold.

In fact, problem (31) can be interpreted as an eigen-like problem. If $\mathbf{W} = span(\mathbf{W})$ is one of its solutions, then \mathbf{W} is interpreted as a generalized eigenspace of operator \mathbf{T} . The interpretation of (31) as an eigen-like problem is crucial since it allows characterizing the best invariant subspace, regarding the decomposition of the solution. This best invariant subspace appears to be the “dominant eigenspace” of operator \mathbf{T} . In some particular cases, problem (31) exactly coincides with a classical eigenproblem, associated with a classical spectral decomposition of the solution (see section 4.5 and [16]). In the general case, this interpretation is motivated by the observed properties of the problem. It naturally leads to the introduction of dedicated algorithms, inspired by classical algorithms for solving eigenproblems. These algorithms, introduced in section (5), present similar behaviors when applied to the present eigen-like problem or to a classical eigenproblem. Numerical examples will illustrate this classical behavior of algorithms and the soundness of the interpretation as an eigen-like problem.

Remark 6 *In fact, due to the definition of mappings \mathbf{f} and \mathbf{F} , the mapping \mathbf{T} is only defined on a subset $\overline{\mathbb{S}}_{n,M} \subset \mathbb{S}_{n,M}$ such that $\mathbf{f}(\overline{\mathbb{S}}_{n,M}) \subset \mathbb{S}_{P,M}^*$. For the same reason, mapping \mathcal{J} is only defined on the quotient space $\overline{Gr}_{n,M} = (\overline{\mathbb{S}}_{n,M} / \mathcal{W})$, which is a subset of the Grassmann manifold $Gr_{n,M}$.*

Remark 7 *We could have equivalently written the problem in terms of $\mathbf{\Lambda}$:*

$$\mathbf{\Lambda} = \mathbf{f} \circ \mathbf{F}(\mathbf{\Lambda}) \equiv \mathbf{T}^*(\mathbf{\Lambda})$$

This equation is an invariant subspace problem on $Gr_{P,M}$, which can be interpreted as an eigen-like problem on mapping \mathbf{T}^ .*

4.5 The case of a coercive symmetric bounded random matrix

We here briefly recall and comment some mathematical results obtained in [16] for the case where \mathbf{A} is a bounded linear coercive symmetric operator from $\mathbb{R}^n \otimes \mathcal{S}$ to $\mathbb{R}^n \otimes \mathcal{S}$. These results clarify in which sense the decomposition is optimal and generalizes the concept of Rayleigh quotient in the case of our eigen-like problem. We first recall that in this case, \mathbf{A} defines a norm on $\mathbb{R}^n \otimes \mathcal{S}_P$, equivalent to the L^2 norm (14), and defined by:

$$\|\mathbf{u}\|_{\mathbf{A}}^2 = E(\mathbf{u}^T \mathbf{A} \mathbf{u}) \quad (33)$$

Proposition 8 *In the case of a bounded coercive symmetric random matrix \mathbf{A} , the optimal decomposition $\mathbf{W}\mathbf{f}(\mathbf{W})$ with respect to the \mathbf{A} -norm is such that \mathbf{W} verifies the following optimization problem:*

$$\mathbf{W} = \underset{\mathbf{W} \in \mathcal{S}_{n,M}}{\operatorname{argmax}} R(\mathbf{W}) \quad (34)$$

where $R(\mathbf{W})$ is a functional defined by:

$$R(\mathbf{W}) = \operatorname{Trace}(\mathbf{R}(\mathbf{W})) \quad (35)$$

$$\text{with } \mathbf{R}(\mathbf{W}) = E(\mathbf{f}(\mathbf{W})\mathbf{b}^T \mathbf{W}) \quad (36)$$

The error in \mathbf{A} -norm then verifies:

$$\|\mathbf{u} - \mathbf{W}\mathbf{f}(\mathbf{W})\|_{\mathbf{A}}^2 = \|\mathbf{u}\|_{\mathbf{A}}^2 - R(\mathbf{W}) \quad (37)$$

In the case where \mathbf{A} is deterministic, functional $\mathbf{R}(\mathbf{W})$ reduces to

$$\mathbf{R}(\mathbf{W}) = (\mathbf{W}^T \mathbf{A} \mathbf{W})^{-1} \mathbf{W}^T E(\mathbf{b}\mathbf{b}^T) \mathbf{W}, \quad (38)$$

and the mapping \mathbf{T} reads

$$\mathbf{T}(\mathbf{W}) = \mathbf{A}^{-1} E(\mathbf{b}\mathbf{b}^T) \mathbf{W} \mathbf{R}(\mathbf{W})^{-1}. \quad (39)$$

Then, problem (30) reads:

$$\mathbf{A} \mathbf{W} \mathbf{R}(\mathbf{W}) = E(\mathbf{b}\mathbf{b}^T) \mathbf{W} \quad (40)$$

which is a classical generalized eigenproblem. Moreover, if \mathbf{A} is symmetric, \mathbf{R} appears to be the classical associated matrix Rayleigh quotient (see *e.g.* [33]). The obtained decomposition is then a classical spectral decomposition of $\mathbf{A}^{-1}\mathbf{b}$ in the metric induced by \mathbf{A} . In the case of a random matrix, we also have the following properties of classical Rayleigh quotients [34].

Proposition 9 *Functionals $\mathbf{R}(\mathbf{W})$ and $R(\mathbf{W})$, defined by equation (36) and (35), verify the following properties:*

(i) *Homogeneity*: $\forall \mathbf{P} \in \mathbb{GL}_M$,

$$\mathbf{R}(\mathbf{W}\mathbf{P}) = \mathbf{P}^{-1}\mathbf{R}(\mathbf{W})\mathbf{P} \quad \text{and} \quad R(\mathbf{W}\mathbf{P}) = R(\mathbf{W}).$$

(ii) *Stationarity* : \mathbf{W} verifies eigen-like problem (30) if and only if it is a stationarity point of $R(\mathbf{W})$.

Regarding proposition 9, functional \mathbf{R} (resp. R) can still be interpreted as a generalized matrix (resp. scalar) Rayleigh quotient associated with eigen-like problem (30). Vectors that makes R stationary are then naturally called generalized eigenvectors, the value of R being interpreted as a generalized eigenvalue. This functional allows us to quantify the quality of generalized eigenspaces and eigenvectors, regarding (37). The best eigenspace (resp. eigenvector), which maximizes R , will then be called the dominant eigenspace (resp. eigenvector).

Remark 10 *Due to homogeneity property of the generalized Rayleigh quotient, the optimization problem (34) can be interpreted as an optimization problem on Grasmann Manifold $Gr_{n,M}$ [33].*

Remark 11 *In this particular case of a coercive symmetric bounded random matrix, the generalized spectral decomposition can be thought as a spectral decomposition of $\mathbf{A}^{-1}\mathbf{b}$ in the metric induced by random matrix \mathbf{A} , i.e. associated with the inner product*

$$((\mathbf{u}, \mathbf{v}))_{\mathbf{A}} = E(\mathbf{u}^T \mathbf{A} \mathbf{v}) \tag{41}$$

However, this decomposition is not classical and does not lead to a classical eigenproblem since inner product (41) is not a natural inner product on tensor product space $\mathbb{R}^n \otimes \mathcal{S}_P$, usually built by tensorisation of inner products on \mathbb{R}^n and \mathcal{S}_P . It only coincides with a classical spectral decomposition for the case of a deterministic symmetric positive definite matrix \mathbf{A} . In the continuous framework, this classical decomposition is called a Hilbert Karhunen Loève decomposition [11].

5 Algorithms for the construction of the generalized spectral decomposition

We have seen in the previous section that the construction of the generalized spectral decomposition (GSD) consists in solving a fixed point problem in the non-compact Stiefel manifold $\mathbb{S}_{n,M}$, i.e. to find $\mathbf{W} \in \mathbb{S}_{n,M}$ such that

$$\mathbf{T}(\mathbf{W}) = \mathbf{W}, \tag{42}$$

where \mathbf{T} is defined in (29). From properties of operator \mathbf{T} , this problem has been interpreted as an eigen-like problem, the optimal GSD of order M being associated with the M -dimensional dominant eigenspace of operator \mathbf{T} . This interpretation naturally leads to the introduction of the following algorithms, inspired by classical algorithms that allow the capture of dominant eigenspaces of linear operators. Numerical examples in Sections 6 and 7 will illustrate the ability of the proposed algorithms to construct the GSD.

5.1 Subspace Iteration type algorithm (SI-GSD)

For classical eigenproblems, the basic algorithm for finding the dominant eigenspace of a linear operator \mathbf{T} is the subspace iteration method. Subspace iterations consist in building the series $\mathbf{W}^{(k+1)} = \mathbf{T}(\mathbf{W}^{(k)})$. In the case of classical eigenproblems, this series converges towards the dominant eigenspace. The extension in the case of our eigen-like problem (42) is straightforward and leads to algorithm 1. In practice, as shown in examples, this algorithm converges very quickly towards the dominant generalized eigenspace.

Algorithm 1 SUBSPACE ITERATION (SI-GSD)

- 1: Initialize $\mathbf{W}^{(0)} \in \mathbb{S}_{n,M}$
- 2: **for** $k = 1$ to k_{max} **do**
- 3: Compute $\mathbf{W}^{(k)} = \mathbf{T}(\mathbf{W}^{(k-1)})$
- 4: Orthonormalize $\mathbf{W}^{(k)}$ (e.g. by QR factorization)
- 5: **end for**
- 6: Set $\mathbf{W} = \mathbf{W}^{(k)}$ and compute $\mathbf{\Lambda} = \mathbf{f}(\mathbf{W})$

Remark 12 Algorithm 1 with $M = 1$ corresponds to a power-type algorithm, which leads to the construction of the dominant generalized eigenvector of \mathbf{T} .

At each iteration k , the computation of $\mathbf{T}(\mathbf{W}^{(k-1)})$ (step 3 of (SI-GSD)) can be decomposed into two steps:

$$\mathbf{\Lambda}^{(k-1)} = \mathbf{f}(\mathbf{W}^{(k-1)}) \quad \text{and} \quad \mathbf{W}^{(k)} = \mathbf{F}(\mathbf{\Lambda}^{(k-1)}) \quad (43)$$

The first step consists in solving a stochastic problem on a fixed deterministic basis (problem of type (10)), which is a problem of size $M \times P$. The second one consists in solving a deterministic problem on a fixed stochastic basis (problem of type (12)), which is a problem of size $M \times n$. Iterations of (SI-GSD) then asks at most for the resolution of k_{max} problems of size $M \times P$ and k_{max} problems of size $M \times n$.

We observe in practice that the initialization step has a low influence on the convergence of the algorithm. A simple and efficient choice consists in taking an initial vector $\mathbf{\Lambda}^{(0)}(\theta) = \sum_{\alpha} \mathbf{\Lambda}_{\alpha}^{(0)} H_{\alpha}(\theta) \in \mathbb{R}^n \otimes \mathbb{S}_P$ with random coefficients

$\Lambda_\alpha^{(0)}$. Then, we simply compute $\mathbf{W}^{(0)} = \mathbf{f}(\Lambda^{(0)})$ by solving a deterministic problem of size $n \times M$ (problem of type (12)).

5.2 Arnoldi type algorithm (A-GSD)

Subspace iteration algorithm 1 can be considered as the reference algorithm, which leads to the “ideal” generalized spectral decomposition associated with the dominant eigenspace. Here, we present an algorithm which gives an approximation of this decomposition and leads to significant computational savings. It is inspired by the Arnoldi technique for solving classical eigenproblems (see *e.g.* [20]). Here, the idea is to find a matrix \mathbf{W} whose columns span a M -dimensional “generalized Krylov subspace” of operator \mathbf{T} and then to build the associated random vector $\Lambda = \mathbf{f}(\mathbf{W})$. The obtained linear subspace can then be considered as a “Ritz approximate” of the dominant generalized eigenspace of \mathbf{T} . The “generalized Krylov subspace” $\mathcal{K}_M(\mathbf{T}, \mathbf{U}_1)$ associated with operator \mathbf{T} and an initial deterministic vector \mathbf{U}_1 can be defined as follows:

$$\begin{aligned} \mathcal{K}_M(\mathbf{T}, \mathbf{U}_1) &= \text{span}(\{\mathbf{U}_i\}_{i=1}^M), \\ \text{with } \mathbf{U}_i &= \mathbf{T}(\mathbf{U}_{i-1}), \quad i = 2 \dots M. \end{aligned} \quad (44)$$

That leads to the following algorithm.

Algorithm 2 ARNOLDI TYPE ALGORITHM (A-GSD)

- 1: Initialize $\mathbf{U}_1 \in \mathbb{R}^n$ and set $\mathbf{U}_1 = \frac{\mathbf{U}_1}{\|\mathbf{U}_1\|}$
- 2: **for** $i = 1$ to M **do**
- 3: Compute $\mathbf{U} = \mathbf{T}(\mathbf{U}_i)$
- 4: Compute $\mathbf{U}_{i+1} = \mathbf{U} - \sum_{j=1}^i (\mathbf{U}_j^T \mathbf{U}) \mathbf{U}_j$ (Orthogonalization step)
- 5: **if** $\|\mathbf{U}_{i+1}\| < \epsilon \|\mathbf{U}\|$ **then**
- 6: *break*
- 7: **end if**
- 8: $\mathbf{U}_{i+1} = \frac{\mathbf{U}_{i+1}}{\|\mathbf{U}_{i+1}\|}$
- 9: **end for**
- 10: Set $\mathbf{W} = (\mathbf{U}_1 \dots \mathbf{U}_i)$ and compute $\Lambda = \mathbf{f}(\mathbf{W})$

The construction of a M -dimensional Krylov subspace (steps 2 to 9) requires to apply $(M - 1)$ times the mapping \mathbf{T} to a vector. Then, it requires only the resolution of $(M - 1)$ problems of size P and $(M - 1)$ problems of size n . Finally, the calculation of Λ requires to solve a problem of size $M \times P$. That leads to significant computational savings compared to algorithm 1. In practice, for the initialization, we introduce an initial $\lambda_0(\theta) = \sum_\alpha \lambda_{0,\alpha} H_\alpha(\theta) \in \mathcal{S}_P$ with random coefficients $\lambda_{0,\alpha}$. Then, we compute the initial vector $\mathbf{U}_1 = \mathbf{f}(\lambda_0)$ by solving a simple deterministic problem of size n (problem of type (12) with $M = 1$).

By using a M -dimensional Krylov subspace, we only get an approximation of the “ideal” generalized spectral decomposition of order M . If we want to improve the quality of this decomposition, we can of course generate a $(M+k)$ -dimensional Krylov subspace $\mathcal{K}^{M+k}(\mathbf{T}, \mathbf{U}_1)$. That leads to a decomposition of order $(M+k)$, from which we can select the M most significant modes with respect to a given metric. This selection step consists in a classical spectral decomposition with respect to this metric (see appendix A). This methodology will be denoted by $(A^{M+k}\text{-GSD})$. The computational time required for the selection is very low.

5.3 Restarting algorithms by “operator deflation”

Of course, when using algorithms 1 or 2 to solve a stochastic problem, we do not know *a priori* the required order M to reach a given accuracy. Suppose that we have built a first decomposition of a given order M and that we have estimated the residual error. If the current decomposition does not reach the required accuracy, we can of course reuse algorithm 1 or 2 with a higher order M .

Remark 13 *A basic way to estimate the error is to compute the norm of the residual of the stochastic problem. For the case of a symmetric bounded coercive random matrix \mathbf{A} , another estimator (computationally cheaper) has been proposed in [16]. Other strategies to estimate the error of the decomposition in a more general context will be introduced in a subsequent paper.*

Another possibility consists in building the subsequent deterministic vectors and random variables by using a “deflation” of operator \mathbf{T} . Let us suppose that we have built a first decomposition $\mathbf{W}_r \mathbf{\Lambda}_r(\theta)$. The deflated operator can be defined by $\mathbf{T}^{(r)} = \mathbf{F}^{(r)} \circ \mathbf{f}^{(r)}$, where mappings $\mathbf{f}^{(r)}$ and $\mathbf{F}^{(r)}$ are defined as mappings \mathbf{f} and \mathbf{F} by replacing the right-hand side \mathbf{b} by the residual of the stochastic problem

$$\mathbf{b}^{(r)} = \mathbf{b} - \mathbf{A} \mathbf{W}_r \mathbf{\Lambda}_r. \quad (45)$$

The subsequent deterministic vectors can then be found by computing the dominant eigenspace of $\mathbf{T}^{(r)}$ with algorithm 1 or 2. This strategy leads to the global algorithm 3.

Algorithm 3 RESTARTED ALGORITHM BY OPERATOR DEFLATION

- 1: Set $\mathbf{b}^{(0)} = \mathbf{b}$, $\mathbf{W}_0 = \emptyset$, $\mathbf{\Lambda}_0 = \emptyset$
- 2: **for** $r = 1$ to r_{max} **do**
- 3: Compute the M_r -dimensional dominant eigenspace $\widehat{\mathbf{W}}_r$ of the deflated operator $\mathbf{T}^{(r-1)} = \mathbf{F}^{(r-1)} \circ \mathbf{f}^{(r-1)}$
- 4: Set $\mathbf{W}_r = (\mathbf{W}_{r-1} \widehat{\mathbf{W}}_r)$

- 5: (without global updating) Compute $\widehat{\mathbf{\Lambda}}_r = \mathbf{f}^{(r-1)}(\widehat{\mathbf{W}}_r)$ and set $\mathbf{\Lambda}_r = (\mathbf{\Lambda}_{r-1}^T \widehat{\mathbf{\Lambda}}_r^T)^T$
 (with global updating) Compute $\mathbf{\Lambda}_r = \mathbf{f}(\mathbf{W}_r)$
- 6: Check convergence
- 7: **end for**

Algorithm 3 introduces two variants which only differ at step 5, where random variables are or are not updated with respect to all previously computed deterministic vectors. In the case of a deterministic definite matrix \mathbf{A} (symmetric or not), we can prove (see proposition 19 in appendix B) that these two variants are equivalent and that they correspond to a usual deflation procedure to solve the associated classical generalized eigenproblem written in (40). In the general case of a random operator, proposition 19 is not true, even with global updating. Indeed, the obtained linear subspace $\text{span}(\mathbf{W})$, being the sum of invariant subspaces of subsequent deflated operators $\mathbf{T}^{(r)}$, is not necessarily an invariant subspace of the initial operator \mathbf{T} . In practice, the global updating step leads to a better accuracy for a given order of decomposition.

Remark 14 *In fact, the global updating can be interpreted as a subspace iteration on the initial eigen-like problem. Of course, the accuracy of the GSD could be improved by performing additional subspace iterations. However, we observe in practice that these additional iterations do not significantly improve the accuracy. Then, since they are computationally expensive, they should be avoided. In the general case, an important question concerns the study of the relationship between generalized eigenspaces of deflated operators $\mathbf{T}^{(r)}$ and generalized eigenspaces of the initial operator \mathbf{T} . In particular, that could allow to modify the proposed algorithms in order to directly build the optimal GSD, without performing any global updating. This question is currently under investigation.*

We have seen that algorithm 1 for $M = 1$ corresponds to a power-type algorithm. If we use $M_r = 1$ in algorithm 3, the use of the power-type algorithm allows the construction of the dominant eigenvector of the deflated operator $\mathbf{T}^{(r)}$. It leads to a power-type algorithm with deflation, which was first introduced in [16]. This algorithm was called (P-GSD) when no global updating is performed and (PU-GSD) when the global updating is performed. The effect of global updating was first illustrated in [16].

Remark 15 *The proposed deflation procedure can be also interesting when dealing with large scale applications (needing large n and large P) and when a high order M of decomposition is required to reach a good accuracy. Indeed, in this case, algorithms 1 or 2 require the resolution of “not so reduced” problems of size $M \times P$ or $M \times n$. When using restarted algorithm 3, algorithm 1 (or 2), which is used at step 3, requires to solve problems of size $M_r \times n$ or $M_r \times P$, with $M_r < M$.*

5.4 Ability to capture a solution with “low dimensionality”

An “ideal” GSD algorithm should be able to automatically capture the optimal reduced basis. In other words, it should be able to capture an exact solution $\mathbf{u} \in \mathbb{R}^n \otimes \mathcal{S}_P$ with a decomposition order M equal to the ideal decomposition order. This ideal order can be defined as the minimum order M that leads to an exact decomposition of \mathbf{u} under the form (15), *i.e.*

$$M_{\mathbf{u}} = \min\{M \in \mathbb{N}; \mathbf{u} = \sum_{i=1}^M \lambda_i \mathbf{U}_i, \lambda_i \in \mathcal{S}_P, \mathbf{U}_i \in \mathbb{R}^n\} \quad (46)$$

In fact, $M_{\mathbf{u}}$ is the number of terms in the exact spectral decomposition of \mathbf{u} . It can be called the dimensionality of random vector \mathbf{u} . As $\mathbb{R}^n \otimes \mathcal{S}_P$ is isomorphic to $\mathbb{R}^n \otimes \mathbb{R}^P$, the dimensionality is clearly finite and verifies $M_{\mathbf{u}} \leq \min(n, P)$. Of course, in general, the dimensionality of a solution verifies $M_{\mathbf{u}} = \min(n, P)$. However, for some problems, the solution may have a low dimensionality. For example, for the case where matrix \mathbf{A} is deterministic, the dimensionality of the solution is the dimensionality of the right-hand side \mathbf{b} , *i.e.* $M_{\mathbf{u}} = M_{\mathbf{b}}$. In this particular case, we can prove that all the proposed algorithms (P-GSD, PU-GSD, A-GSD, SI-GSD) allow the capture of the exact solution \mathbf{u} with exactly $M_{\mathbf{u}}$ modes. This good property has been also verified in [16] for the power-type algorithm when matrix \mathbf{A} has a very specific structure (product of a random variable by a deterministic matrix). In more general cases, the proposed subspace iteration algorithm (SI-GSD) and Arnoldi-type algorithm (\mathbf{A}^{M+k} -GSD) have also this property. This will be illustrated in example 2 (section 7.7).

Remark 16 $M_{\mathbf{u}}$ can also be interpreted as the unique integer such that there exists $\mathbf{W}_{\mathbf{u}} \in \mathcal{S}_{n, M_{\mathbf{u}}}$ and $\mathbf{\Lambda}_{\mathbf{u}} \in \mathcal{S}_{P, M_{\mathbf{u}}}^*$ such that $\mathbf{u} = \mathbf{W}_{\mathbf{u}} \mathbf{\Lambda}_{\mathbf{u}}$. Mapping \mathcal{J} , defined in (32), seems to have a unique fixed point on $\overline{\text{Gr}}_{n, M_{\mathbf{u}}}$ (see remark 6), which is the span of $\mathbf{W}_{\mathbf{u}}$ (observed in practice). If $M > M_{\mathbf{u}}$, problem (30) has no solution. In fact, in this case, $\mathcal{S}_{n, M}$ (resp. $\overline{\text{Gr}}_{n, M}$) is empty, which means that we can not find a full rank matrix in $\mathcal{S}_{n, M}$ such that the associated random variables are linearly independent. One could interpret this property as follows: if M were greater than $M_{\mathbf{u}}$, a part of the generalized eigenspace would be associated with zero eigenvalue. All these remarks are easily proven in the case of a deterministic symmetric matrix \mathbf{A} but need for more mathematical investigations in the general case.

6 Example 1: a linear elasticity problem

In this first model problem, the mathematical framework is the one of section 4.5, already illustrated in [16]. The aim of this section is to validate the new algorithms proposed in this article. In particular, we show that the subspace iteration algorithm 1 (SI-GSD) allows constructing the ideal generalized spectral decomposition and that Arnoldi algorithm 2 (A-GSD) leads to a rather good approximation of this ideal decomposition, which converges towards the ideal decomposition when using higher dimensional generalized Krylov subspaces. The computational costs of these algorithms are then illustrated.

6.1 Formulation of the problem, stochastic modeling and approximation

6.1.1 Formulation of the problem and spatial discretization

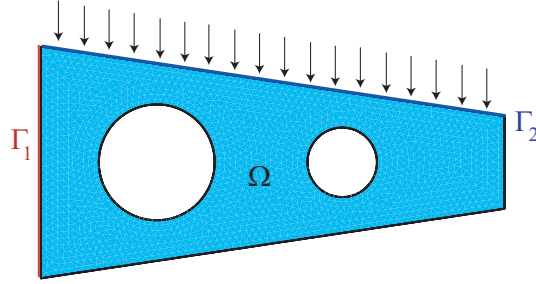


Fig. 1. Model problem 1: bending of an elastic structure

We consider a classical linear elasticity problem on a domain $\Omega \subset \mathbb{R}^2$ (Figure 1). We work under plane strain assumption. We denote by $\mathbf{u}(\mathbf{x}, \theta)$ the displacement field. We denote by $\mathbf{g}(\mathbf{x}, \theta)$ the surface load applied on a part Γ_2 of the boundary. Homogeneous Dirichlet boundary conditions are applied on another part of the boundary, denoted by Γ_1 . We consider the complementary part of $\Gamma_1 \cup \Gamma_2$ in $\partial\Omega$ as a free boundary. A classical weak formulation of this problem reads: find $\mathbf{u} \in \mathcal{V} \otimes \mathcal{S}$ such that

$$A(\mathbf{u}, \mathbf{v}) = B(\mathbf{v}) \quad \forall \mathbf{v} \in \mathcal{V} \otimes \mathcal{S} \quad (47)$$

where

$$\begin{aligned} A(\mathbf{u}, \mathbf{v}) &= E \left(\int_{\Omega} \boldsymbol{\varepsilon}(\mathbf{v}) : \mathbf{C} : \boldsymbol{\varepsilon}(\mathbf{u}) \, dx \right) \\ B(\mathbf{v}) &= E \left(\int_{\Gamma_2} \mathbf{g} \cdot \mathbf{v} \, ds \right) \end{aligned}$$

$\boldsymbol{\varepsilon}(\mathbf{u})$ is the symmetric part of the displacement gradient (or stain tensor) and \mathbf{C} the Hooke fourth-order tensor. Under classical regularity assumptions on the data (see [4,35]), *i.e.* material properties and loadings, an ad-hoc choice

for function spaces consists in taking $\mathcal{V} = \{\mathbf{v}(\mathbf{x}) \in (H^1(\Omega))^2; \mathbf{v}|_{\Gamma_1} = 0\}$ and $\mathcal{S} = L^2(\Theta, dP)$. A classical finite element approximation at the space level can be introduced. Let us denote by $\mathcal{V}_n = \{\mathbf{v}(\mathbf{x}) = \sum_{i=1}^n v_i \boldsymbol{\varphi}_i(\mathbf{x}), \boldsymbol{\varphi}_i \in \mathcal{V}\} \subset \mathcal{V}$ the finite element approximation space. A function $\mathbf{v} \in \mathcal{V}_n \otimes \mathcal{S}$ will then be associated with the random vector $\mathbf{v}(\theta) = (v_1(\theta), \dots, v_n(\theta))^T \in \mathbb{R}^n \otimes \mathcal{S}$. Random matrix \mathbf{A} and random vector \mathbf{b} in equation (1) are then defined as follows: $\forall \mathbf{u}, \mathbf{v} \in \mathcal{V}_n$,

$$\begin{aligned} A(\mathbf{u}, \mathbf{v}) &= E(\mathbf{v}^T \mathbf{A} \mathbf{u}), \\ B(\mathbf{v}) &= E(\mathbf{v}^T \mathbf{b}). \end{aligned}$$

Here, we use a mesh composed by 3-nodes triangles, illustrated on Figure 2. The dimension of the approximation space \mathcal{V}_n is $n = 1624$.

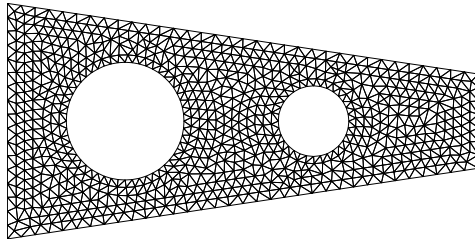


Fig. 2. Finite element mesh: 3-nodes triangles

6.1.2 Stochastic modeling and stochastic discretization

We consider an isotropic material, with a Poisson coefficient $\nu = 0.3$ and a Young modulus $\kappa(\mathbf{x}, \theta)$ which is a lognormal random field reading

$$\kappa(\mathbf{x}, \theta) = \exp(\mu + \sigma \gamma(\mathbf{x}, \theta))$$

where $\gamma(\mathbf{x}, \theta)$ is a homogeneous Gaussian random field with a zero mean, a unitary standard deviation and an exponential square correlation function with a correlation length equal to 1 (the horizontal length of the structure is 2): $E(\gamma(\mathbf{x}, \theta) \gamma(\mathbf{x}', \theta)) = \exp(-\|\mathbf{x} - \mathbf{x}'\|^2)$. μ and σ are chosen such that the marginal distribution of κ has a mean equal to 1 and a standard deviation equal to 0.25. The stochastic field κ is discretized as follows: we first perform a truncated Karhunen-Loève decomposition of γ on 9 modes: $\gamma(\mathbf{x}, \theta) \approx \sum_{i=1}^9 \gamma_i(\mathbf{x}) \xi_i(\theta)$, where the $\xi_i \in N(0, 1)$ ⁴ are statistically independent standard Gaussian random variables. We then decompose κ on a Hermite polynomial chaos of degree 3 in dimension 9: $\kappa \approx \sum_{\alpha} \kappa_{\alpha}(\mathbf{x}) H_{\alpha}(\boldsymbol{\xi}(\theta))$. Coefficients κ_{α} of the decomposition can be obtained analytically (see *e.g.* [10]).

⁴ $\xi \in N(\mu, \sigma)$ is a Gaussian random variable with mean μ and standard deviation σ

Finally, a truncated Karhunen Loève decomposition of order 9 is performed in order to reduce the number of space functions for the representation of κ : $\kappa \approx \sum_{i=1}^9 \kappa_i(\mathbf{x})e_i(\theta)$. The first 9 modes of this decomposition are shown in Figure 3. The overall discretization procedure of κ leads to a relative L^2 error of 10^{-2} between the discretized stochastic field and the initial stochastic field.

We consider that the surface load \mathbf{g} is vertical, uniformly distributed on Γ_2 : $\mathbf{g}(\mathbf{x}, \theta) = -\xi_{10}(\theta)\mathbf{e}_y$, where $\xi_{10} \in N(1, 0.2)$ is a Gaussian random variable, statistically independent of the previous random variables $\{\xi_i\}_{i=1}^9$.

The probabilistic content is then represented by $m = 10$ random variables $\{\xi_i\}_{i=1}^m$. For the approximation space \mathcal{S}_P , we first choose a polynomial chaos of degree $p = 3$ in dimension $m = 10$, thus leading to a dimension $P = 286$ of the stochastic approximation space \mathcal{S}_P .

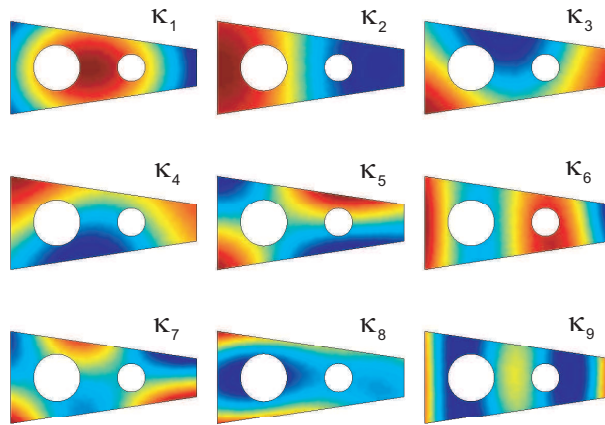


Fig. 3. First modes $\{\kappa_i(\mathbf{x})\}_{i=1}^9$ of the Karhunen-Loève decomposition of the lognormal field $\kappa(\mathbf{x}, \theta)$.

6.2 Reference solution and definition of errors

The reference solution, denoted by \mathbf{u} , is the solution of the initial discretized problem (5). It is computed by a classical Preconditioned Conjugate Gradient algorithm (PCG). The preconditioner is a block preconditioner based on the expectation $E(\mathbf{A})$ of the random matrix (see [26] for its definition). We denote by (GSD) the generalized spectral decomposition and (SD) the classical spectral decomposition of the reference solution. We denote by $\mathbf{u}^{(M)}$ a spectral decomposition of order M . In order to estimate the quality of approximate solutions, we will use the following relative errors, respectively in L^2 -norm and \mathbf{A} -norm:

$$\varepsilon_{L^2}^{(M)} = \frac{\|\mathbf{u} - \mathbf{u}^{(M)}\|}{\|\mathbf{u}\|}, \quad \varepsilon_{\mathbf{A}}^{(M)} = \frac{\|\mathbf{u} - \mathbf{u}^{(M)}\|_{\mathbf{A}}}{\|\mathbf{u}\|_{\mathbf{A}}}$$

The norms $\|\cdot\|$ and $\|\cdot\|_{\mathbf{A}}$ are defined in equations (14) and (33) respectively. In this example, we can notice that the \mathbf{A} -norm is in fact equivalent to an energy norm:

$$\begin{aligned}\|\mathbf{v}\|_{\mathbf{A}}^2 &= E(\mathbf{v}^T \mathbf{A} \mathbf{v}) \\ &= \int_{\Theta} \int_{\Omega} \boldsymbol{\varepsilon}(\mathbf{v}(\mathbf{x}, \theta)) : \mathbf{C} : \boldsymbol{\varepsilon}(\mathbf{v}(\mathbf{x}, \theta)) dx dP(\theta)\end{aligned}$$

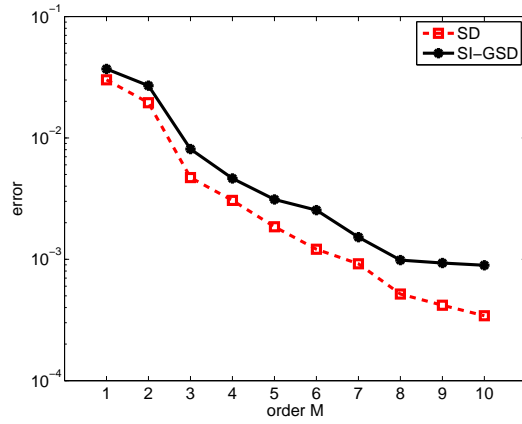
6.3 Comparison between Generalized Spectral Decomposition (GSD) and classical Spectral Decomposition (SD)

Here we compare the generalized spectral decomposition (GSD) with the classical spectral decomposition (SD) of the reference solution. The reference GSD is obtained by algorithm 1 (SI-GSD). Figure 4 shows the convergence of SD and GSD with respect to the order M of decomposition. We clearly observe that the GSD is better than the SD with respect to the \mathbf{A} -norm and that the SD is better than the GSD with respect to the L^2 -norm. This result was expected regarding the definition of spectral decompositions. For this model problem, we recall that the \mathbf{A} -norm is equivalent to an energy norm and then, GSD leads to a better decomposition with respect to the energy norm.

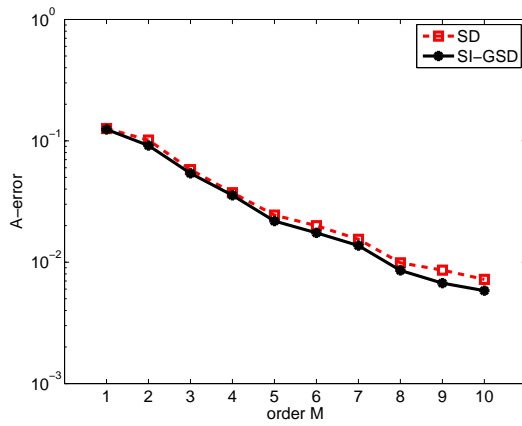
In Figure 5 (resp. 6), we plot the first 9 deterministic functions of the GSD (resp. SD).

For SD, the deterministic vectors are uniquely defined as the M dominant eigenvectors of correlation matrix $E(\mathbf{u}\mathbf{u}^T)$, sorted by decreasing eigenvalues. This corresponds to a sorting regarding their contributions to the spectral decomposition with respect to the L^2 -norm. The GSD vectors are not uniquely defined. Every set of vectors which spans the same linear subspace of \mathbb{R}^n leads the same GSD. Then, in order to compare the deterministic vectors with those of SD, we perform a classical spectral decomposition of the GSD with respect to the L^2 -norm. This is just a rewriting of the GSD, the initial and final deterministic vectors belonging to the same equivalence class, *i.e.* spanning the same linear subspace (see appendix A for the definition of the sorting procedure). The obtained vectors, shown in Figure 7, are very similar to the one obtained by SD.

Finally, in Figure 8 we compare the GSD decomposition with a classical spectral decomposition SD in the metric associated with $E(\mathbf{A})$, which is in fact a Hilbert-Karhunen Loève decomposition of the solution with respect to the inner product $\langle \mathbf{u}, \mathbf{v} \rangle = E(\mathbf{u}^T E(\mathbf{A}) \mathbf{v})$ in $\mathbb{R}^n \otimes \mathcal{S}_P$. We observe that such a SD decomposition is more similar to the GSD decomposition.



(a) L^2 -norm



(b) \mathbf{A} -norm

Fig. 4. Classical SD in the natural metric of $L^2(\Theta, dP; \mathbb{R}^n)$ versus SI-GSD: convergence in L^2 -norm (a) and \mathbf{A} -norm (b)

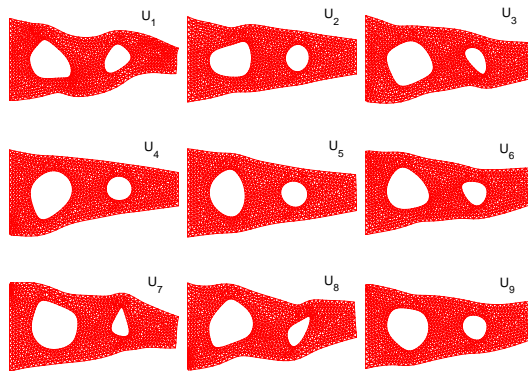


Fig. 5. First modes $\{\mathbf{U}_i\}_{i=1}^9$ of the GSD obtained by (SI-GSD)

6.4 Convergence of algorithm (SI-GSD)

In order to evaluate the convergence of the subspace iteration algorithm (SI-GSD), we compare successive linear subspaces $\text{span}(\mathbf{W}^{(k)})$. The comparison

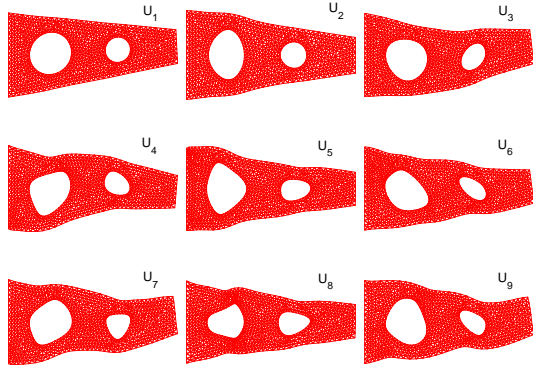


Fig. 6. First modes $\{\mathbf{U}_i\}_{i=1}^9$ of (SD)

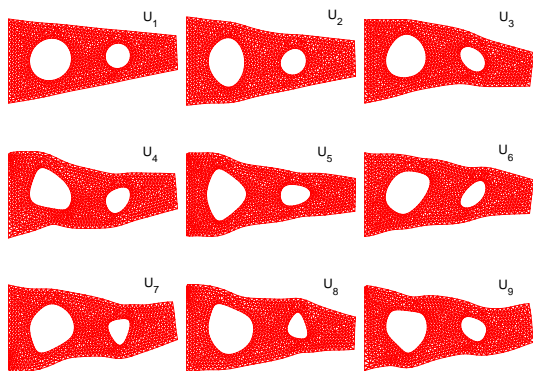
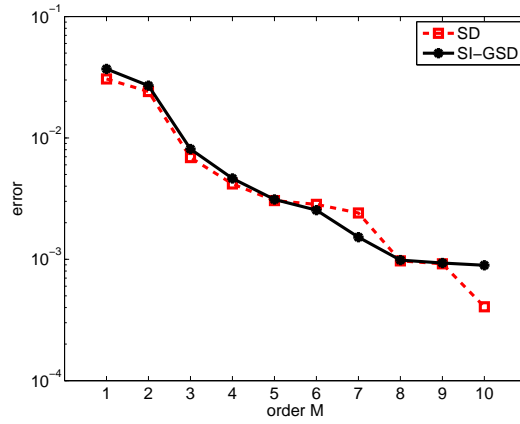


Fig. 7. First modes $\{\mathbf{U}_i\}_{i=1}^9$ of the GSD, sorted with respect to the L^2 norm

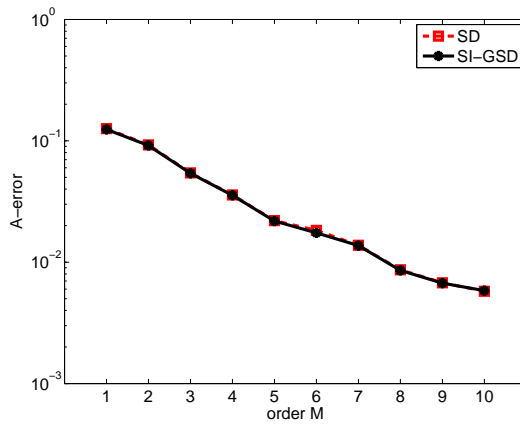
between two linear subspaces can be made by computing the largest principal angle between these linear subspaces. Two matrices \mathbf{W} and $\widetilde{\mathbf{W}}$ being given, the largest principal angle between the associated linear subspaces $span(\mathbf{W})$ and $span(\widetilde{\mathbf{W}})$ is defined by

$$\angle(\mathbf{W}, \widetilde{\mathbf{W}}) = \max_{\mathbf{U} \in span(\mathbf{W})} \min_{\widetilde{\mathbf{U}} \in span(\widetilde{\mathbf{W}})} \angle(\mathbf{U}, \widetilde{\mathbf{U}})$$

where $\angle(\mathbf{U}, \widetilde{\mathbf{U}})$ is the classical angle between the two vectors (see *e.g.* [36]). Figure 9 shows the cosine of angle $\angle(\mathbf{W}^{(k)}, \mathbf{W}^{(k+1)})$ between two successive iterates of algorithm SI-GSD for different orders M of decomposition. Figure 10 shows, for different orders M , the cosine of angle $\angle(\mathbf{W}^{(k)}, \mathbf{W}^{(ref)})$ between an iterate $\mathbf{W}^{(k)}$ of algorithm SI-GSD and the reference subspace $\mathbf{W}^{(ref)}$. Finally, we observe the convergence of SI-GSD in Figure 11 by estimating the distance (in \mathbf{A} -norm) between the reference GSD $\mathbf{u}^{(M)}$ and $\mathbf{W}^{(k)}\mathbf{f}(\mathbf{W}^{(k)})$, which is the approximate GSD obtained at iteration k . On all these figures, we observe a very fast convergence of SI-GSD towards the dominant generalized M -dimensional eigenspace, whatever the dimension M is.



(a) L^2 -norm



(b) \mathbf{A} -norm

Fig. 8. classical SD in the metric induced by $E(\mathbf{A})$ versus SI-GSD: convergence in L^2 -norm (a) and \mathbf{A} -norm (b)

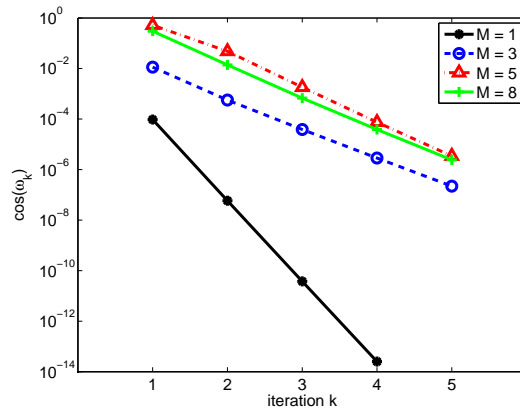


Fig. 9. Convergence of (SI-GSD) for different orders M of decomposition: cosine of the largest principal angle between two iterates $\omega_k = \angle(\mathbf{W}^{(k)}, \mathbf{W}^{(k+1)})$

6.5 Arnoldi-type algorithm (A-GSD) versus Subspace Iterations (SI-GSD)

Here, we will evaluate the quality of the decomposition obtained by the Arnoldi-type algorithm 2 (A-GSD). The obtained decomposition is compared with

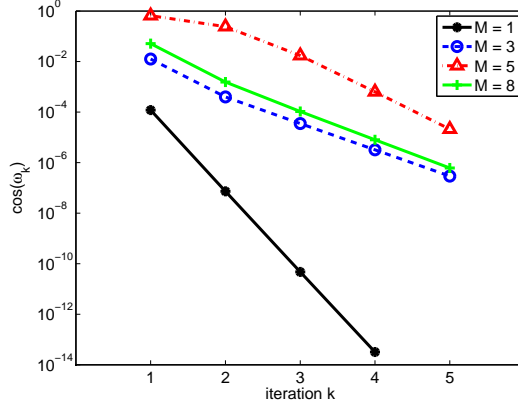


Fig. 10. Convergence of (SI-GSD) for different orders M of decomposition: cosine of the largest principal angle between iterates and reference subspace $\mathbf{W}^{(ref)}$: $\omega_k = \angle(\mathbf{W}^{(k)}, \mathbf{W}^{(ref)})$

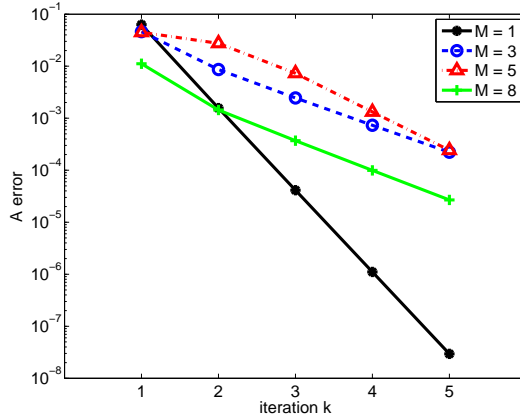


Fig. 11. Convergence of (SI-GSD) for different orders M of decomposition: error $\|\mathbf{W}^{(k)}\mathbf{f}(\mathbf{W}^{(k)}) - \mathbf{u}^{(M)}\|_{\mathbf{A}}$ with respect to k

the decomposition obtained by algorithm (SI-GSD), considered as the exact GSD. Figure 12 shows the error in \mathbf{A} -norm with respect to the order M of the decomposition for both algorithms. Here, A-GSD consists in generating a M -dimensional Krylov subspace. We observe that A-GSD leads to a relatively good approximate decomposition.

In Figure 13, we compare the convergence of the GSD decompositions obtained by SI-GSD and \mathbf{A}^{M+k} -GSD. We recall that \mathbf{A}^{M+k} -GSD consists in building a generalized spectral decomposition of order $M+k$, by generating a $(M+k)$ -dimensional generalized Krylov subspace and then to select the M most significant modes (see section 5.2). In Figure 13(a) (resp. 13(b)), the M modes for \mathbf{A}^{M+k} -GSD are selected with respect to the natural inner product in $L^2(\Theta, dP; \mathbb{R}^n)$ (resp. to the inner product induced by $E(\mathbf{A})$). We observe that by increasing the dimension of the generalized Krylov subspace, we rapidly converge towards an optimal spectral decomposition, whatever the metric used for the selection is.

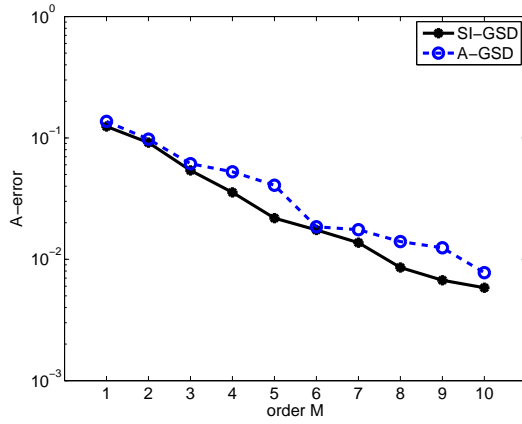
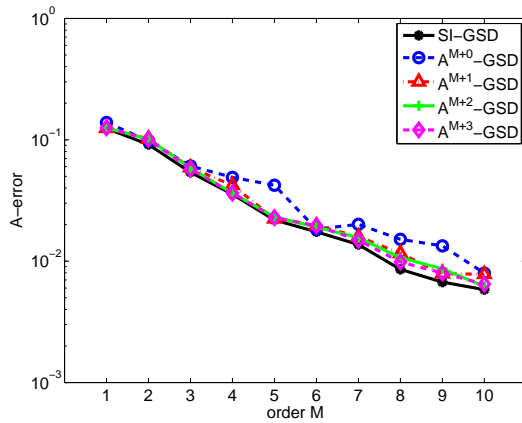
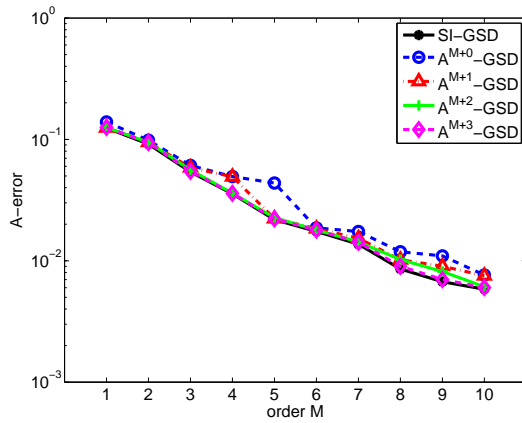


Fig. 12. SI-GSD versus A-GSD: convergence of the decomposition \mathbf{A} -norm



(a) Selection in L^2 metric



(b) Selection in $E(\mathbf{A})$ metric

Fig. 13. SI-GSD versus A^{M+k} -GSD: selection of the M most significant modes with respect to the natural metric in $L^2(\Theta, dP; \mathbb{R}^n)$ (a) or with respect to the metric induced by $E(\mathbf{A})$ (b)

6.6 Computational costs: comparison with a classical resolution technique

6.6.1 Comparison between GSD algorithms and a standard PCG

Here, we compare the computational times required by the classical Preconditioned Conjugate Gradient (PCG) and by the (GSD) algorithms. Figure 14

shows the error with respect to computational time for different algorithms.

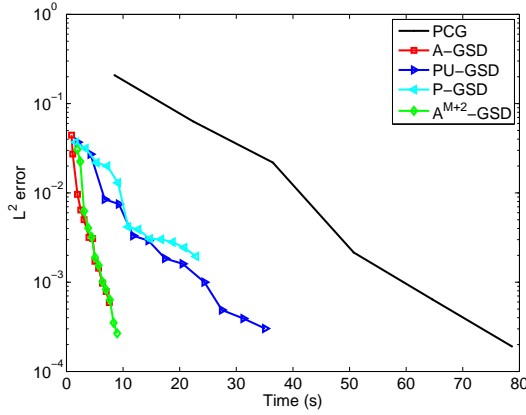


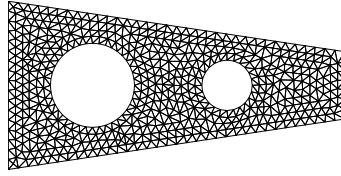
Fig. 14. Comparison between resolution techniques: L^2 -error versus computational time (reference discretization)

Arnoldi-type algorithm (A-GSD) leads to significant computational savings compared to PCG but also compared to algorithms P-GSD and PU-GSD, which were proposed in a previous paper [16]. PU-GSD (resp. P-GSD) is a power-type algorithm with deflation and with updating (resp. without updating). We recall that these power-type algorithms are equivalent to a subspace iteration algorithm to build the generalized spectral decomposition. We observe that PU-GSD and P-GSD lead to similar computational times. In fact, the computational time required by the updating step in PU-GSD is balanced by the fact that PU-GSD needs for a lower order of decomposition than P-GSD (the computed modes are more pertinent). We also observe that algorithms A^{M+k} -GSD and A-GSD lead to similar computational times. Of course, for a given order of decomposition, A^{M+k} -GSD requires more computational times. However, this decomposition is more accurate than with A-GSD since the computed modes are more pertinent.

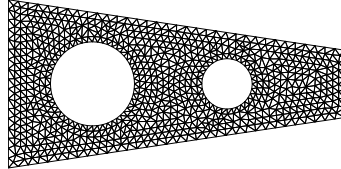
6.6.2 Influence of the dimension of approximation spaces

To go further in the comparison of computational costs, we will analyze the influence of the dimensions P and n of stochastic and deterministic approximation spaces. Here, we only compare PCG with the most efficient GSD algorithm, namely the Arnoldi-type algorithm (A-GSD). Four meshes, shown in Figure 15, are used to analyze the influence of n . Meshes 1 to 4 correspond respectively to $n = 1150, 1624, 3590$ and 6166 . To analyze the influence of P , we simply increase the order p of the polynomial chaos expansion of the solution. We will use $p = 2, 3$ or 4 , corresponding to $P = 66, 286$ or $P = 1001$ respectively.

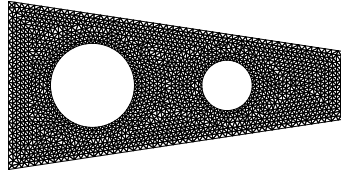
Figure 16 (resp. 17) shows the computational times for different p and for the fixed finite element mesh 2 (resp. mesh 4). We can observe that the PCG



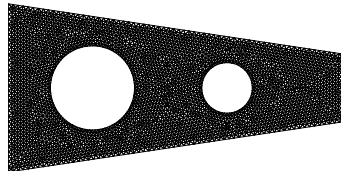
(a) mesh 1 : $n = 1150$



(b) mesh 2 : $n = 1624$



(c) mesh 3 : $n = 3590$



(d) mesh 4 : $n = 6166$

Fig. 15. Different finite element meshes

computational times drastically increase with p while the A-GSD computational times are almost independent of p . More precisely, we notice that when $P \ll n$, the computational times required by A-GSD are almost independent of P . This result is clearly observed when we use the fine mesh 4. This is due to the fact that deterministic and stochastic problems are uncoupled. Then, for large n , P has a low influence on computational times, which comes essentially from the resolution of deterministic equations.

Figure 18 (resp. 19) shows the computational times for different n and for a fixed $p = 2$ (resp. $p = 4$). We can observe that the PCG computational times drastically increase with n while the A-GSD computational times are almost

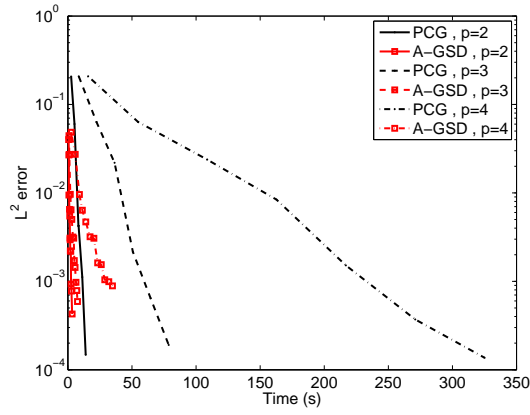


Fig. 16. Influence of stochastic dimension (variable p) with fixed mesh 2

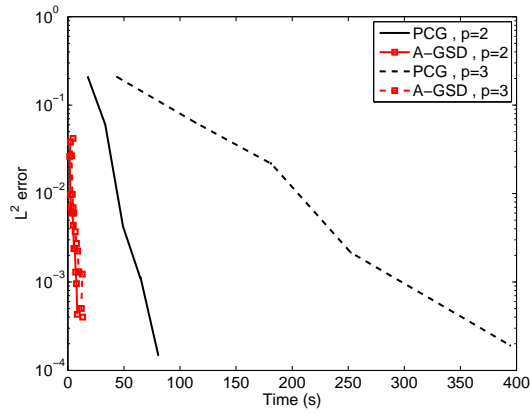


Fig. 17. Influence of stochastic dimension (variable p) with fixed mesh 4

independent of n (especially for large P). This result is clearly observed when we use $p = 4$ ($P = 1001$). This is still due to the fact that deterministic and stochastic problems are uncoupled. Then, for large P , n has a low influence on computational times, which comes essentially from the resolution of stochastic equations.

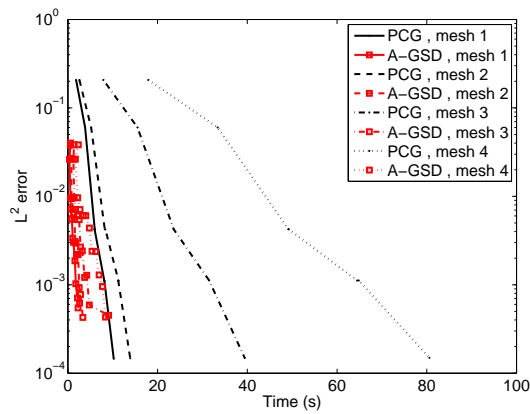


Fig. 18. Influence of deterministic dimension (different meshes) with fixed stochastic discretization ($p = 2$)

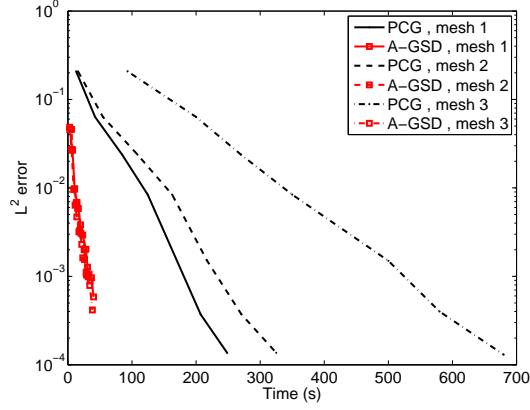


Fig. 19. Influence of deterministic dimension (different meshes) with fixed stochastic discretization ($p = 4$)

Finally, Table 1 shows the gains in terms of computational times to reach a given relative error of 10^{-2} . It also shows the gains in terms of memory requirements to store the corresponding approximate solution. For the largest P and n , we observe that the computational time with A-GSD is 50 times lower than with PCG and that memory requirements are about 200 times lower.

	P=66 (p=2)		P=286 (p=3)		P=1001 (p=4)	
	T_g	M_g	T_g	M_g	T_g	M_g
n=1150	9.3	15.3	15.4	57.3	11.2	133.8
n=1624	8.8	15.9	21.6	60.8	17.2	154.8
n=3590	14.8	16.2	42.6	66.2	33.1	195.7
n=6166	20.2	16.3	51.9	68.3	47.2	215.3

Table 1

Comparison between PCG and A-GSD: computational time gain factor $T_g = \frac{\text{time}(PCG)}{\text{time}(A-GSD)}$ to reach a relative error of 10^{-2} and memory gain factor $M_g = \frac{\text{memory}(PCG)}{\text{memory}(A-GSD)}$ to store the approximate solution

7 Model problem 2: a transient heat diffusion problem

7.1 Formulation of the problem and semi-discretization

We consider a transient heat diffusion equation as a model problem for parabolic stochastic partial differential equations (see Figure 20). The problem reads:

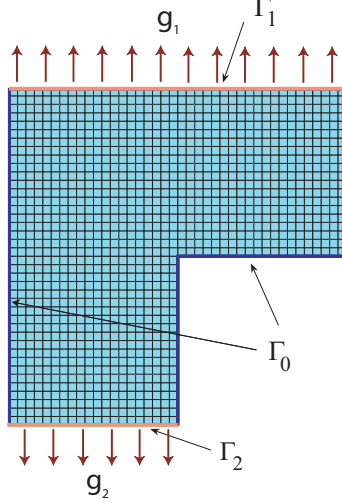


Fig. 20. Model problem 2 and associated finite element mesh

find a temperature field $u : \Omega \times (0, T) \times \Theta \rightarrow \mathbb{R}$ such that

$$\begin{aligned}
 c\partial_t u - \nabla \cdot (\kappa \nabla u) &= f & \text{on } \Omega \times (0, T) \\
 -\kappa \nabla u \cdot n &= g_i & \text{on } \Gamma_i \times (0, T), \quad i = 1, 2 \\
 u &= 0 & \text{on } \Gamma_0 \times (0, T) \\
 u|_{t=0} &= u_0 & \text{on } \Omega
 \end{aligned} \tag{48}$$

where Ω denotes the spatial domain, $(0, T)$ the time domain, c and κ the material parameters, f the volumic heat source and g_1 (resp. g_2) a normal flux on a part Γ_1 (resp. Γ_2) of the boundary. Homogeneous Dirichlet boundary conditions are applied on a part Γ_0 of the boundary, which is the complementary part of $\Gamma_1 \cup \Gamma_2$. For numerical examples, we will take $u_0 = 0$. However, for the sake of generality, this term is kept in the presentation of the method.

For space discretization, we use four-nodes linear finite elements (see mesh in Figure 20). It classically leads to the following system of stochastic differential equations in time: find $\mathbf{u} : (0, T) \times \Theta \rightarrow \mathbb{R}^{n_x}$ such that

$$\mathbf{M}(\theta) \dot{\mathbf{u}}(t, \theta) + \mathbf{B}(\theta) \mathbf{u}(t, \theta) = \mathbf{c}(t, \theta) \tag{49}$$

$$\mathbf{u}(0, \theta) = \mathbf{u}_0(\theta) \tag{50}$$

When using a classical time integration scheme, the resulting semi-discretized stochastic problem can formally be written as in equation (1). Let us illustrate this for the case of a standard backward Euler scheme. We denote by $(0 = t_0, t_1, \dots, t_{n_t} = T)$ the time grid. For the sake of simplicity, we suppose that we have a uniform time step δt . Denoting the solution by $\mathbf{u}(\theta) \equiv (\mathbf{u}(t_1, \theta)^T, \mathbf{u}(t_2, \theta)^T, \dots, \mathbf{u}(t_{n_t}, \theta)^T)^T \in \mathbb{R}^n \otimes \mathcal{S}$, with $n = n_x \times n_t$, problem ((49),(50)) can then be written under the form (1) where random matrix \mathbf{A}

and random vector \mathbf{b} are defined by:

$$\mathbf{A} = \begin{pmatrix} \mathbf{M} + \mathbf{B}\delta t & 0 & \dots & 0 \\ -\mathbf{M} & \mathbf{M} + \mathbf{B}\delta t & \ddots & \vdots \\ \vdots & \ddots & \ddots & 0 \\ 0 & \dots & -\mathbf{M} & \mathbf{M} + \mathbf{B}\delta t \end{pmatrix} \quad (51)$$

$$\mathbf{b} = \begin{pmatrix} \mathbf{c}(t_1)\delta t + \mathbf{M}\mathbf{u}_0 \\ \mathbf{c}(t_2)\delta t \\ \vdots \\ \mathbf{c}(t_{n_t})\delta t \end{pmatrix} \quad (52)$$

The present model problem is then equivalent to our generic problem (1), whose weak formulation writes (5). For the reference problem, we consider $n_x = 1179$ and $n_t = 30$.

The generalized spectral decomposition technique can then be directly applied to problem (5). Here, the GSD decomposition (15) reads

$$\mathbf{u}(t, \theta) \approx \mathbf{W}(t)\mathbf{\Lambda}(\theta) \Leftrightarrow \mathbf{u}(\theta) \approx \mathbf{W}\mathbf{\Lambda}(\theta)$$

where $\mathbf{W} \equiv (\mathbf{W}(t_1)^T, \dots, \mathbf{W}(t_{n_t})^T)^T \in \mathbb{R}^{n \times M}$. With the proposed algorithms, the construction of this decomposition requires the resolution of reduced deterministic problems (13) and reduced stochastic problems (11). Of course, this writing is just formal. In practice, those problems are re-interpreted with respect to the initial evolution problem. This interpretation and computational aspects are detailed in appendix C. In particular, it is shown that GSD takes into account initial conditions in a weak sense.

7.2 Stochastic modeling and approximation

Material parameters are considered as simple random variables, independent of space and time. We take $c(\theta) = \xi_1(\theta) + \xi_2(\theta)$ and $\kappa = \xi_3(\theta) + \xi_4(\theta)$, where $\xi_1, \dots, \xi_4 \in U(0.7, 1.3)$ ⁵ are four independent identically distributed uniform random variables. The volumic heat source is considered as a simple Gaussian random variable, independent of time and space: $f = \xi_5(\theta) \in N(1, 0.2)$. Normal fluxes are taken as: $g_1(t, \theta) = \xi_6(\theta)\frac{t}{T}$ and $g_2(t, \theta) = \xi_7(\theta)\frac{t}{T}$, where $\xi_6, \xi_7 \in N(1, 0.2)$. The probabilistic content is then represented by $m = 7$

⁵ $\xi \in U(a, b)$ is a uniform random variable on (a, b)

random variables $\{\xi_i\}_{i=1}^7$, which are considered statistically independent. For the approximation at the stochastic level, we use a generalized polynomial chaos approximation of order p [21,22]. Basis functions of \mathcal{S}_P are then multi-dimensional polynomials, which are the product of Legendre polynomials in ξ_1, ξ_2, ξ_3 and ξ_4 and Hermite polynomials in ξ_5, ξ_6 and ξ_7 . For the reference problem, we take $p = 4$, corresponding to $P = \frac{(p+m)!}{m!p!} = 330$ basis functions in \mathcal{S}_P .

7.3 Reference solution and error indicator

The reference solution, denoted by \mathbf{u} , is the solution of the initial discretized problem (5). It is computed by a time incremental resolution. At each t_i , $i = 1, \dots, n_t$, $\mathbf{u}(t_i, \theta) \in \mathbb{R}^{n_x} \otimes \mathcal{S}_P$ verifies the problem: $\forall \mathbf{v} \in \mathbb{R}^{n_x} \otimes \mathcal{S}_P$,

$$E(\mathbf{v}^T (\mathbf{M}(\theta) + \mathbf{B}(\theta)\delta t)\mathbf{u}(t_i, \theta)) = E(\mathbf{v}^T (\mathbf{M}(\theta)\mathbf{u}(t_{i-1}, \theta) + \mathbf{c}(t_i, \theta)\delta t)) \quad (53)$$

Problem (53), which is a system of $n_x \times P$ equations, is solved as in model problem 1 using a classical Preconditioned Conjugate Gradient algorithm (PCG) with a tolerance 10^{-5} . The preconditioner is a block diagonal preconditioner based on the expectation $E(\mathbf{M} + \mathbf{B}\delta t)$ (see [26] for its definition). In the following, this reference resolution technique is simply denoted by PCG. We denote by $\mathbf{u}^{(M)}$ a spectral decomposition of order M . In order to compare approximate solutions, we use the following relative error:

$$\varepsilon^{(M)} = \frac{\|\mathbf{u} - \mathbf{u}^{(M)}\|}{\|\mathbf{u}\|} \quad (54)$$

where the norm $\|\cdot\|$ is the following L^2 -norm on the approximation space :

$$\|\mathbf{u}\|^2 = \sum_{i=1}^{n_t} \delta t E(\mathbf{u}(t_i, \theta)^T \mathbf{u}(t_i, \theta)) \quad (55)$$

Remark 17 *The L^2 -norm defined in (55) coincides with the natural norm in $L^2(\Theta, dP; L^2((0, T); \mathbb{R}^n))$ if we consider that the approximation is piecewise constant on each time interval and lower semi-continuous:*

$$\|\mathbf{u}\|^2 = \int_{\Theta} \int_0^T \mathbf{u}(t, \theta)^T \mathbf{u}(t, \theta) dt dP(\theta)$$

7.4 Comparison between Generalized Spectral Decomposition (GSD) and classical Spectral Decomposition (SD)

Here, we compare the generalized spectral decomposition (GSD) with the classical spectral decomposition (SD) of the reference solution. The SD of the

reference solution $\mathbf{u} \in \mathbb{R}^n \otimes \mathcal{S}_P$ is defined as the best decomposition of order M with respect to the norm defined in (55). The time grid being uniform, this decomposition can be obtained by applying a classical Karhunen-Loève decomposition to random vector \mathbf{u} . The GSD is obtained by algorithm 2 (A-GSD). Figure 21 shows the convergence of SD and GSD with respect to the order M of decomposition. We observe that A-GSD leads to a rather good decomposition with respect to the L^2 -norm.

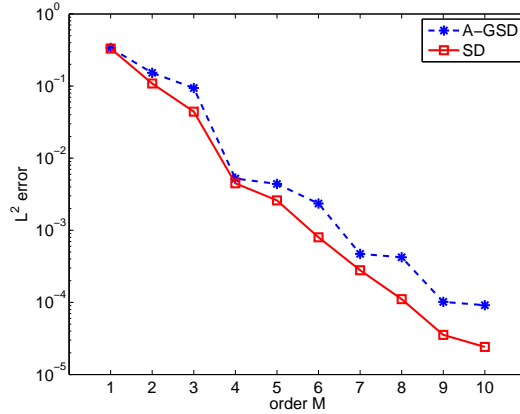


Fig. 21. Convergence of GSD, obtained by algorithm (A-GSD), versus convergence of classical SD (in the L^2 metric)

We recall that A-GSD consists in generating a M -dimensional generalized Krylov subspace. We now use algorithm A^{M+k} -GSD in order to improve the quality of the decomposition of order M (see section 5.2). We recall that this algorithm consists in building a GSD of order $(M + k)$ (by building a $(M + k)$ -dimensional generalized Krylov subspace) and then in selecting the M most significant modes with respect to a given metric. Here, we use the L^2 metric for the selection. In Figure 22, we compare the convergence of the GSD decomposition obtained by A^{M+k} -GSD with the convergence of the classical SD of the reference solution. We still observe that by increasing the dimension of the Krylov subspace, the GSD quickly converges towards an optimal spectral decomposition, which is very similar to L^2 -optimal classical SD.

7.5 Quality of the generalized spectral decomposition

In Figure 23 (resp. 24), we can see the quantiles (5% and 95%) of the GSD approximation at a given point (resp. at a given time on a vertical line). We observe that the approximation of these quantiles converges very fast with the order M of the GSD. For $M = 4$, corresponding to an error in L^2 norm inferior to 10^{-2} , the approximation of the quantiles is very good.

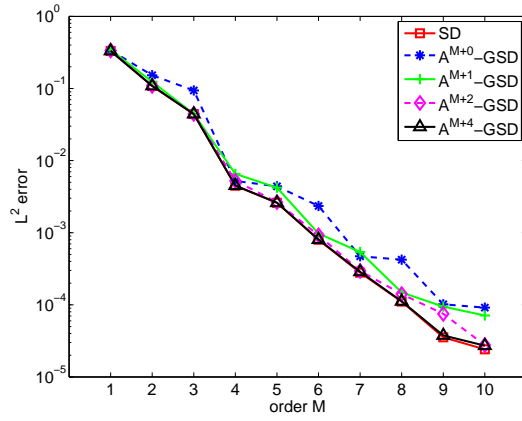


Fig. 22. (SD) versus (A^{M+k} -GSD): selection of the M most significant modes with respect to the L^2 metric

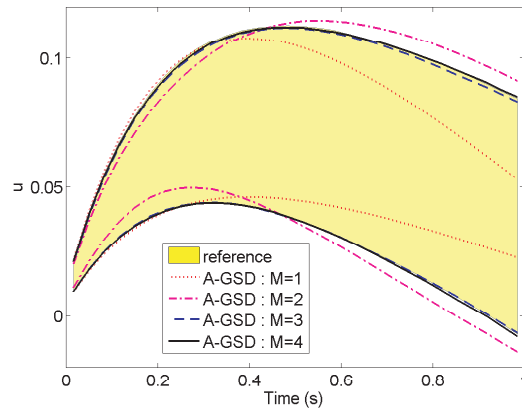


Fig. 23. Quantiles 0.05 and 0.95 of approximate solutions at point $(0.5, 1.5)$

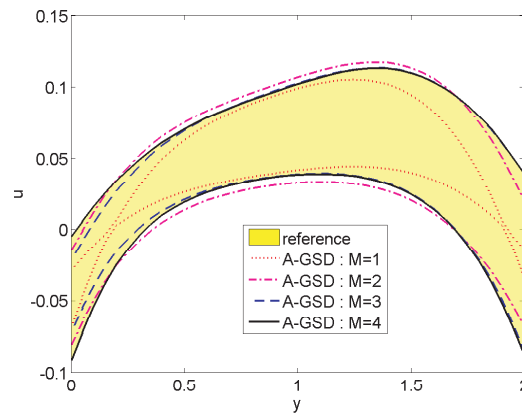


Fig. 24. Quantiles 0.05 and 0.95 of the approximate solutions on the line $x = 0.5$ at time $t = 0.65s$

7.6 Computational costs

Now, we illustrate the efficiency of the proposed algorithm A-GSD by comparing it to the classical resolution technique PCG. Figure 25 shows the evolution

of the error with respect to computational time. We can observe that the convergence rate of A-GSD is far better than the one of PCG.

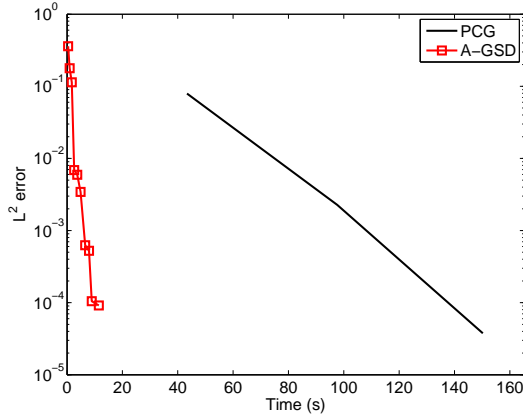


Fig. 25. Error versus computation time for PCG and A-GSD (reference discretization)

To go further in the comparison of computational costs, we will analyze the influence of the dimensions P and n of stochastic and deterministic approximation spaces. Figure 26 shows the convergence curves for different orders $p = 3, 4$ and 5 of polynomial chaos expansions, respectively corresponding to $P = 120, 330$ and 792 . We clearly observe that an increasing P has a very low influence on the convergence rate of A-GSD while it drastically deteriorates the convergence rate of PCG.

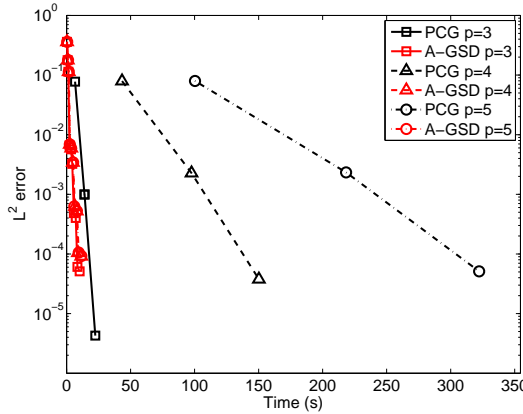


Fig. 26. Influence of the dimension of the stochastic approximation space for PCG and A-GSD (for $p=3, 4$ or 5)

Table 2 shows the gains in terms of computational times for different discretizations at stochastic level (different orders p of polynomial chaos) and deterministic level (different numbers of spatial degrees of freedom n_x and time steps n_t). The gain is computed by comparing computational times of PCG and A-GSD to reach a given relative error of 10^{-2} . We observe that for the finest discretizations, computational times can be divided by 100.

T_g			
p=3	$n_t=10$	$n_t=30$	$n_t=50$
$n_x=289$	7	15	19
$n_x=659$	8	15	19
$n_x=1179$	6	14	19
p=4	$n_t=10$	$n_t=30$	$n_t=50$
$n_x=289$	17	39	49
$n_x=659$	18	39	50
$n_x=1179$	13	33	43
p=5	$n_t=10$	$n_t=30$	$n_t=50$
$n_x=289$	35	76	103
$n_x=659$	36	85	103
$n_x=1179$	27	64	87

Table 2

Time gain factor $T_g = \frac{\text{time}(PCG)}{\text{time}(A-GSD)}$

In Figure 27, we illustrate the computational time required by (PCG) to reach a relative error of 10^{-2} with respect to the total dimension $n \times P$ of the approximation space. We observe that computational time grows approximately linearly with respect to $n \times P$ (unitary slope in log-log plot). Results of Table 2 are plotted in Figure 28, where we can observe a quasi linearity between the gain and the total dimension of the approximation space (unitary slope in log-log plot).

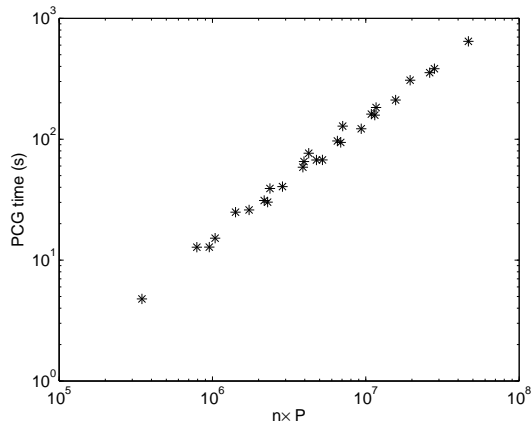


Fig. 27. $\text{time}(PCG)$ with respect to $n \times P$

Finally, Table 3 shows the gains in terms of memory requirements for the storage of the solution, for different discretizations. For the finest discretizations,

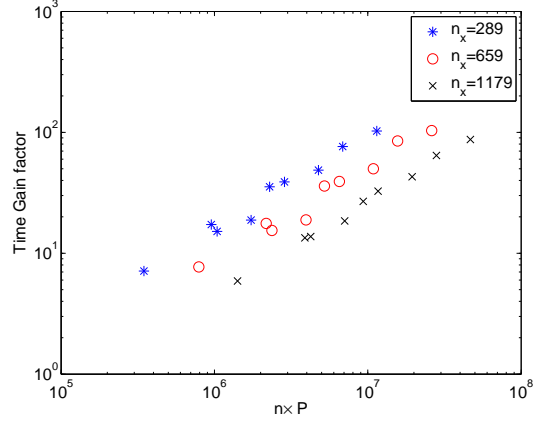


Fig. 28. Time gain factor $T_g = \frac{\text{time}(PCG)}{\text{time}(A-GSD)}$ with respect to $n \times P$

the memory requirements are divided by 200.

M_g			
p=3	$n_t=10$	$n_t=30$	$n_t=50$
$n_x=289$	29	30	30
$n_x=659$	30	30	30
$n_x=1179$	30	30	30
p=4	$n_t=10$	$n_t=30$	$n_t=50$
$n_x=289$	74	80	81
$n_x=659$	79	81	82
$n_x=1179$	80	82	82
p=5	$n_t=10$	$n_t=30$	$n_t=50$
$n_x=289$	155	181	188
$n_x=659$	177	190	193
$n_x=1179$	186	194	195

Table 3

Memory gain factor $M_g = \frac{\text{memory}(PCG)}{\text{memory}(A-GSD)}$ for the storage of the solution

7.7 Manufactured problem with low dimensionality solution

The aim of this last section is to illustrate the ability of the proposed algorithms to capture an exact solution $\mathbf{u} \in \mathbb{R}^n \otimes \mathcal{S}_P$ with “low dimensionality” (see section 5.4). For the case where matrix \mathbf{A} is deterministic, the dimensionality of the solution is the dimensionality of the right-hand side \mathbf{b} , *i.e.* $M_{\mathbf{u}} = M_{\mathbf{b}}$.

In this particular case, all the proposed algorithms allow the capture of the exact solution \mathbf{u} with exactly $M_{\mathbf{u}}$ modes. For the present example, it has been verified by taking \mathbf{M} and \mathbf{B} deterministic and by keeping unchanged the right-hand side, whose dimensionality is 3.

Here, in a more general case, we want to illustrate that the Arnoldi-type algorithm still has this ability to construct the ideal decomposition, *i.e.* to construct the optimal $M_{\mathbf{u}}$ -dimensional reduced basis. For that purpose, we give us a solution \mathbf{u} with a desired dimensionality $M_{\mathbf{u}}$ and define a manufactured right-hand side $\mathbf{c} = \mathbf{M}\dot{\mathbf{u}} + \mathbf{B}\mathbf{u}$ such that the solution of the new evolution problem is \mathbf{u} . For \mathbf{u} , we consider the truncation at order 3 of the SD of the solution of the previous reference problem. We use A^{M+k} -GSD with $M = 3$. That means that we only look for an order $M = 3$ generalized spectral decomposition $\mathbf{u}^{(3)}(t, \theta) = \mathbf{W}(t)\mathbf{\Lambda}(\theta)$. The 3 “spectral modes” can be interpreted as the Ritz approximate of the 3 first “exact” dominant generalized eigenmodes using a generalized Krylov subspace of dimension $(3 + k)$. In Figure 29, we observe that the relative error in L^2 norm between $\mathbf{u}^{(3)}$ and the manufactured solution \mathbf{u} decreases very rapidly towards the machine precision, which means that the linear subspace which is spanned by the 3 generalized eigenmodes quickly converges towards the expected dominant 3-dimensional generalized eigenspace of operator \mathbf{T} .

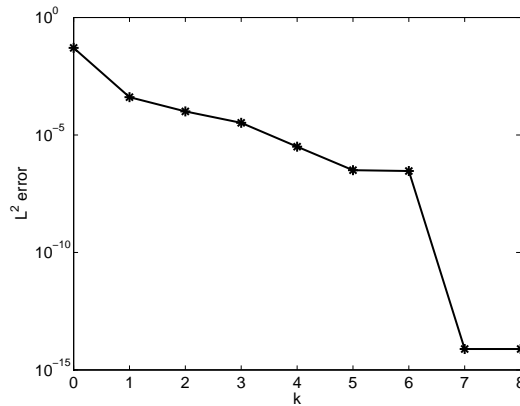


Fig. 29. Manufactured problem with a solution of dimensionality 3: L^2 error obtained by (A^{M+k} -GSD) with $M = 3$

Remark 18 *The optimal 3-dimensional eigenspace can also be captured with a very high precision by the SI-GSD in a few iterations.*

8 Conclusion

An efficient alternative approach has been proposed for the resolution of equations arising from stochastic Galerkin schemes. This approach generalizes the

concept of spectral decomposition for the resolution of stochastic problems. Deterministic vectors and random variables that appear in this decomposition are obtained by solving an invariant subspace problem, which can be interpreted as an eigen-like problem. It has been shown that the associated dominant generalized eigenspaces define pertinent reduced bases, allowing the obtention of a very accurate low-order spectral decomposition of the solution. This interpretation has allowed us to derive efficient algorithms, inspired by Subspace Iteration and Arnoldi algorithms for classical eigenproblems. These algorithms have the great advantage to only require the resolution of separated reduced stochastic problems and reduced deterministic problems. The method has been illustrated on two model problems: a linear elasticity problem and a transient heat diffusion equation. The proposed algorithms have been systematically compared to classical resolution techniques. The examples clearly show that the proposed GSD algorithms lead to significant savings in terms of computational times and memory requirements. The method has been developed on a generic discretized problem, arising from spatial, temporal and stochastic discretizations of SPDEs. Further mathematical and numerical investigations will be necessary to rigorously define the class of problems for which the GSD can be applied. The introduction of pertinent error estimators, based on *a posteriori* error estimation techniques, is also an important issue for the proposed methodology. Indeed, it could lead to a more pertinent and adaptive construction of the GSD. Lastly, the methodology and algorithms can be naturally extended to the case of stochastic non-linear problems. This will be developed and illustrated in a forthcoming paper.

Acknowledgement

This work is supported by the French National Research Agency (grant ANR-06-JCJC-0064).

Appendix

A Basic operations on spectral decompositions

A.1 Orthogonalization of deterministic vectors and random variables

Let us suppose that a spectral decomposition $\mathbf{u}(\theta) = \sum_{i=1}^M \lambda_i(\theta) \mathbf{U}_i = \mathbf{W} \mathbf{\Lambda}(\theta)$ has been found, with $\mathbf{W} \in \mathbb{S}_{n,M}$ and $\mathbf{\Lambda} \in \mathbb{S}_{P,M}^*$. Let us denote by $\mathbf{D}_{\mathbf{W}} \in \mathbb{R}^{M \times M}$

and $\mathbf{D}_\Lambda \in \mathbb{R}^{M \times M}$ the matrices such that $\widetilde{\mathbf{W}} = \mathbf{W}\mathbf{D}_\mathbf{W}^{-1}$ and $\widetilde{\Lambda} = \mathbf{D}_\Lambda^{-T}\Lambda$ are orthonormal in the following sense:

$$\widetilde{\mathbf{W}}^T \mathbf{M} \widetilde{\mathbf{W}} = \mathbf{I}_M, \quad E(\widetilde{\Lambda} \widetilde{\Lambda}^T) = \mathbf{I}_M, \quad (\text{A.1})$$

where \mathbf{M} is a symmetric definite positive matrix. Matrix $\mathbf{D}_\mathbf{W}$ can for example be obtained using a cholesky factorization of \mathbf{M} and a QR decomposition of $\mathbf{W} \in \mathbb{R}^{n \times M}$. Matrix \mathbf{D}_Λ can be obtained by a QR factorization of the matrix in $\mathbb{R}^{P \times M}$ whose columns represent the components of the λ_i on the orthonormal basis of \mathcal{S}_P (see remark 1). Then, denoting by $\mathbf{D} = \mathbf{D}_\mathbf{W} \mathbf{D}_\Lambda^T$, the initial decomposition can be rewritten

$$\mathbf{u} = \widetilde{\mathbf{W}} \mathbf{D} \widetilde{\Lambda}. \quad (\text{A.2})$$

A.2 Sorting and truncating a spectral decomposition

Suppose that we have obtained a spectral decomposition $\mathbf{u} = \mathbf{W}\Lambda$. We have seen that there exists an infinite number of matrices and random vectors leading to the same decomposition (see section 4.3). It could be interesting to extract some deterministic vectors and random variables, respectively in $\text{span}(\mathbf{W})$ and $\text{span}(\Lambda)$, which are sorted by decreasing contribution in the decomposition. The simplest way to proceed consists in performing a Karhunen-Loève decomposition of \mathbf{u} with respect to a given metric. Let us explain how to build this decomposition in practice. Let us denote by $((\mathbf{u}, \mathbf{v}))_{\mathbf{M}} = E(\mathbf{u}^T \mathbf{M} \mathbf{v})$ the chosen inner product, where \mathbf{M} is a deterministic symmetric positive definite matrix. We start to compute the decomposition (A.2), leading to orthonormal deterministic vectors and random variables. We then apply a singular value decomposition of matrix \mathbf{D} :

$$\mathbf{D} = \mathbf{U}_D \Sigma_D \mathbf{V}_D^T,$$

where $\mathbf{U}_D^T \mathbf{U}_D = \mathbf{I}_M$, $\mathbf{V}_D^T \mathbf{V}_D = \mathbf{I}_M$ and $\Sigma_D = \text{diag}(\sigma_1, \dots, \sigma_M) \in \mathbb{R}^{M \times M}$ is diagonal, with $\sigma_1 \geq \dots \geq \sigma_M$. The decomposition can then be rewritten

$$\mathbf{u} = \mathbf{W} \Sigma_D \Lambda = \sum_{i=1}^M \sigma_i \mathbf{U}_i \lambda_i,$$

where $\mathbf{W} = \widetilde{\mathbf{W}} \mathbf{U}_D$ and $\Lambda = \mathbf{V}_D^T \widetilde{\Lambda}$ are orthonormal in the sense of (A.1). We can easily verify that

$$E(\mathbf{u} \mathbf{u}^T) = \sum_{i=1}^M \sigma_i^2 \mathbf{U}_i \mathbf{U}_i^T,$$

$$\|\mathbf{u}\|_{\mathbf{M}}^2 = E(\mathbf{u}^T \mathbf{M} \mathbf{u}) = \sum_{i=1}^M \sigma_i^2.$$

A truncation of the spectral decomposition can then be simply performed. If we want a precision ϵ on the relative error, a simple criterion is to choose the minimum order $M' \leq M$ such that

$$\sum_{i=M'+1}^M \sigma_i^2 \leq \epsilon^2 \sum_{i=1}^M \sigma_i^2$$

B Restarted algorithm in the case of a deterministic definite matrix

Proposition 19 *In the case where \mathbf{A} is a deterministic definite matrix, restarted algorithm 3 corresponds to a classical deflation technique. At stage r , $\text{span}(\mathbf{W}_r)$ is the $(\sum_{i=1}^r M_i)$ -dimensional dominant invariant subspace of \mathbf{T} . The approximation $\mathbf{W}_r \mathbf{\Lambda}_r$ is the same as the one obtained by computing the $(\sum_{i=1}^r M_i)$ -dimensional dominant invariant subspace of \mathbf{T} .*

PROOF. We will first prove this result for the algorithm with global updating (i) and then show that algorithms with or without global updating are equivalent in the case where \mathbf{A} is a deterministic definite matrix (ii).

(i) *Algorithm with global updating*

Let us use a proof by induction. Proposition 19 is trivially verified for $r = 1$. Let us now suppose that it is verified at stage r . $\text{span}(\mathbf{W}_r)$ is the dominant invariant subspace of \mathbf{T} , with dimension $M = \sum_{i=1}^r M_i$. Then, \mathbf{W}_r verifies the following generalized eigenproblem:

$$\mathbf{A} \mathbf{W}_r \mathbf{R}(\mathbf{W}_r) = E(\mathbf{b} \mathbf{b}^T) \mathbf{W}_r,$$

where $\mathbf{R}(\mathbf{W}_r) = (\mathbf{W}_r^T \mathbf{A} \mathbf{W}_r)^{-1} \mathbf{W}_r^T E(\mathbf{b} \mathbf{b}^T) \mathbf{W}_r$. Let $\text{span}(\widehat{\mathbf{W}}_{r+1})$ be the dominant M_{r+1} -dimensional invariant subspace of the deflated operator $\mathbf{T}^{(r)}$, which verifies the following generalized eigenproblem:

$$\mathbf{A} \mathbf{W} \mathbf{R}^{(r)}(\mathbf{W}) = E(\mathbf{b}^{(r)} \mathbf{b}^{(r)T}) \mathbf{W},$$

where $\mathbf{R}^{(r)}(\mathbf{W}) = (\mathbf{W}^T \mathbf{A} \mathbf{W})^{-1} \mathbf{W}^T E(\mathbf{b}^{(r)} \mathbf{b}^{(r)T}) \mathbf{W}$. After some algebra, we can show that

$$E(\mathbf{b}^{(r)} \mathbf{b}^{(r)T}) = E(\mathbf{b} \mathbf{b}^T) - \mathbf{A} \mathbf{W}_r \mathbf{R}(\mathbf{W}_r) (\mathbf{W}_r^T \mathbf{A} \mathbf{W}_r)^{-T} \mathbf{W}_r^T \mathbf{A}^T,$$

which is a classical deflation of matrix $E(\mathbf{b} \mathbf{b}^T)$. Now, using lemma 20, we conclude that $\text{span}((\mathbf{W}_r \widehat{\mathbf{W}}_{r+1}))$ is the dominant $(M + M_{r+1})$ -dimensional invariant subspace of the initial operator \mathbf{T} , which ends the proof for the case with global updating.

(ii) *Without global updating*

To end up with the proof, we just have to show that with or without global updating, we obtain the same random vector, *i.e.*

$$\mathbf{f}(\mathbf{W}_r) = \mathbf{f}((\mathbf{W}_{r-1} \widehat{\mathbf{W}}_r)) = \begin{pmatrix} \mathbf{f}(\mathbf{W}_{r-1}) \\ \mathbf{f}^{(r-1)}(\widehat{\mathbf{W}}_r) \end{pmatrix}.$$

Noting that $\mathbf{W}_{r-1}^T \mathbf{A} \widehat{\mathbf{W}}_r = 0$ (lemma 20), we have, with a little algebra,

$$\begin{aligned} \mathbf{f}((\mathbf{W}_{r-1} \widehat{\mathbf{W}}_r)) &= \begin{pmatrix} \mathbf{W}_{r-1}^T \mathbf{A} \mathbf{W}_{r-1} & 0 \\ \widehat{\mathbf{W}}_r^T \mathbf{A} \mathbf{W}_{r-1} & \widehat{\mathbf{W}}_r^T \mathbf{A} \widehat{\mathbf{W}}_r \end{pmatrix}^{-1} \begin{pmatrix} \mathbf{W}_{r-1}^T \mathbf{b} \\ \widehat{\mathbf{W}}_r^T \mathbf{b} \end{pmatrix} \\ &= \begin{pmatrix} (\mathbf{W}_{r-1}^T \mathbf{A} \mathbf{W}_{r-1})^{-1} \mathbf{W}_{r-1}^T \mathbf{b} \\ (\widehat{\mathbf{W}}_r^T \mathbf{A} \widehat{\mathbf{W}}_r)^{-1} \widehat{\mathbf{W}}_r^T (\mathbf{b} - \mathbf{A} \mathbf{W}_{r-1} \mathbf{f}(\mathbf{W}_{r-1})) \end{pmatrix} \end{aligned}$$

which ends the proof.

The following lemma is quite classical and can be proven with little algebra.

Lemma 20 *We consider a M_1 -dimensional subspace $\text{span}(\mathbf{W}_1)$ which verifies the generalized eigenproblem $\mathbf{A} \mathbf{W}_1 \mathbf{R}(\mathbf{W}_1) = \mathbf{B} \mathbf{W}_1$, where \mathbf{A} is a definite matrix, \mathbf{B} a (semi)definite symmetric positive matrix and $\mathbf{R}(\mathbf{W}) = (\mathbf{W}^T \mathbf{A} \mathbf{W})^{-1} (\mathbf{W}^T \mathbf{B} \mathbf{W})$. We then consider a M_2 -dimensional subspace $\text{span}(\widehat{\mathbf{W}}_2)$ which verifies the deflated generalized eigenproblem: $\mathbf{A} \widehat{\mathbf{W}}_2 \mathbf{R}^{(1)}(\widehat{\mathbf{W}}_2) = \mathbf{B}^{(1)} \widehat{\mathbf{W}}_2$, where $\mathbf{B}^{(1)}$ is the deflated matrix*

$$\mathbf{B}^{(1)} = \mathbf{B} - \mathbf{A} \mathbf{W}_1 \mathbf{R}(\mathbf{W}_1) (\mathbf{W}_1^T \mathbf{A} \mathbf{W}_1)^{-T} \mathbf{W}_1^T \mathbf{A}^T,$$

and $\mathbf{R}^{(1)}(\mathbf{W}) = (\mathbf{W}^T \mathbf{A} \mathbf{W})^{-1} (\mathbf{W}^T \mathbf{B}^{(1)} \mathbf{W})$. We have the following properties:

- (i) $\text{span}(\mathbf{W}_1)$ is an eigenspace of the deflated eigenproblem associated with zero eigenvalues.
- (ii) $\mathbf{W}_1^T \mathbf{A} \widehat{\mathbf{W}}_2 = 0$
- (iii) there exists \mathbf{W}_2 such that $\text{span}(\mathbf{W}_2) \subset \text{span}((\mathbf{W}_1 \widehat{\mathbf{W}}_2))$ is an eigenspace of the initial eigenproblem and such that

$$\mathbf{R}(\mathbf{W}_2) = \mathbf{R}^{(1)}(\widehat{\mathbf{W}}_2)$$

- (iv) $\text{span}((\mathbf{W}_1 \widehat{\mathbf{W}}_2))$ is an eigenspace of the initial generalized eigenproblem associated with the same eigenvalues as $\text{span}((\mathbf{W}_1 \mathbf{W}_2))$.

C Interpretation of GSD for stochastic time evolution equation

Here, we give an interpretation of reduced problems arising from GSD algorithms in the case of model problem 2. After having introduced space and time discretizations, we have seen that this problem is equivalent to problem (5), where random matrix \mathbf{A} and random vector \mathbf{b} are respectively given in (51) and (52).

C.1 Reduced stochastic problems

Computing $\mathbf{\Lambda} = \mathbf{f}(\mathbf{W})$ for a fixed $\mathbf{W} \in \mathbb{R}^{n \times M}$ requires to solve the following equation:

$$E(\tilde{\mathbf{\Lambda}}^T \mathbf{W}^T \mathbf{A} \mathbf{W} \mathbf{\Lambda}) = E(\tilde{\mathbf{\Lambda}}^T \mathbf{W}^T \mathbf{b}) \quad \forall \tilde{\mathbf{\Lambda}} \in \mathbb{R}^M \otimes \mathcal{S}_P$$

where

$$\begin{aligned} \mathbf{W}^T \mathbf{A} \mathbf{W} &= \delta t \sum_{i=1}^{n_t} \mathbf{W}(t_i)^T \mathbf{B} \mathbf{W}(t_i) + \mathbf{W}(t_1)^T \mathbf{M} \mathbf{W}(t_1) + \\ &\quad \sum_{i=2}^{n_t} \mathbf{W}(t_i)^T \mathbf{M} (\mathbf{W}(t_i) - \mathbf{W}(t_{i-1})) \\ \mathbf{W}^T \mathbf{b} &= \mathbf{W}(t_1)^T \mathbf{M} \mathbf{u}_0 + \delta t \sum_{i=1}^{n_t} \mathbf{W}(t_i)^T \mathbf{c}(t_i) \end{aligned}$$

In fact, the problem to be solved is a time-discretized version of the following equation, arising from a natural weak formulation in time:

$$\begin{aligned} &\left(\int_{(0,T)} \mathbf{W}(t)^T (\mathbf{M} \dot{\mathbf{W}}(t) + \mathbf{B} \mathbf{W}(t)) dt + \mathbf{W}(0^+)^T \mathbf{M} \mathbf{W}(0^+) \right) \mathbf{\Lambda} = \\ &\int_{(0,T)} \mathbf{W}(t)^T \mathbf{c}(t) dt + \mathbf{W}(0^+)^T \mathbf{M} \mathbf{u}_0 \end{aligned}$$

where, in the time-discretized version, $a(0^+)$ is identified with $a(t_1)$ and the time integrals must be interpreted as follows. For any $a : (0, T) \rightarrow \mathbb{R}$ and $b : (0, T) \rightarrow \mathbb{R}$,

$$\int_{(0,T)} a(t)b(t) dt \approx \sum_{i=1}^{n_t} a(t_i)b(t_i)\delta t$$

$$\int_{(0,T)} a(t)\dot{b}(t) dt \approx \sum_{i=2}^{n_t} a(t_i)(b(t_i) - b(t_{i-1}))$$

C.2 Reduced deterministic problems

Computing $\mathbf{W} = \mathbf{F}(\boldsymbol{\Lambda})$ for a fixed $\boldsymbol{\Lambda} \in \mathbb{R}^M \otimes \mathcal{S}_P$ requires to solve the following equation:

$$E(\mathbf{A}\mathbf{W}\boldsymbol{\Lambda}\boldsymbol{\Lambda}^T) = E(\mathbf{b}\boldsymbol{\Lambda}^T)$$

In fact, it is a time-discretized version of the following system of evolution equations:

$$E(\mathbf{M}\dot{\mathbf{W}}(t)\boldsymbol{\Lambda}\boldsymbol{\Lambda}^T) + E(\mathbf{B}\mathbf{W}(t)\boldsymbol{\Lambda}\boldsymbol{\Lambda}^T) = E(\mathbf{c}(t)\boldsymbol{\Lambda}^T) \quad (\text{C.1})$$

$$E(\mathbf{M}\mathbf{W}(0)\boldsymbol{\Lambda}\boldsymbol{\Lambda}^T) = E(\mathbf{M}\mathbf{u}_0\boldsymbol{\Lambda}^T) \quad (\text{C.2})$$

which is equivalent to the set of M coupled deterministic evolution equations: for $i = 1, \dots, M$,

$$\sum_{j=1}^M E(\mathbf{M}\lambda_i\lambda_j)\dot{\mathbf{U}}_j(t) + \sum_{j=1}^M E(\mathbf{B}\lambda_i\lambda_j)\mathbf{U}_j(t) = E(\mathbf{c}(t)\lambda_i) \quad (\text{C.3})$$

$$\sum_{j=1}^M E(\mathbf{M}\lambda_i\lambda_j)\mathbf{U}_j(0) = E(\mathbf{M}\mathbf{u}_0\lambda_i) \quad (\text{C.4})$$

Equation (C.2) (or (C.4)) is equivalent to a weak imposition of the initial condition.

Remark 21 *In the case where \mathbf{M} and \mathbf{B} are deterministic, it can be convenient to orthonormalize the λ_i . Indeed, computing \mathbf{W} then requires to solve a set of M uncoupled deterministic evolution equations: for $i = 1, \dots, M$,*

$$\begin{aligned} \mathbf{M}\dot{\mathbf{U}}_i(t) + \mathbf{B}\mathbf{U}_i(t) &= E(\mathbf{c}(t)\lambda_i) \\ \mathbf{U}_i(0) &= E(\mathbf{u}_0\lambda_i) \end{aligned}$$

References

- [1] R. Ghanem and P. Spanos. *Stochastic finite elements: a spectral approach*. Springer, Berlin, 1991.
- [2] R. Ghanem. Ingredients for a general purpose stochastic finite elements implementation. *Computer Methods in Applied Mechanics and Engineering*, 168:19–34, 1999.
- [3] H. G. Matthies and A. Keese. Galerkin methods for linear and nonlinear elliptic stochastic partial differential equations. *Computer Methods in Applied Mechanics and Engineering*, 194(12-16):1295–1331, 2005.

- [4] I. Babuška, R. Tempone, and G. E. Zouraris. Solving elliptic boundary value problems with uncertain coefficients by the finite element method: the stochastic formulation. *Computer Methods in Applied Mechanics and Engineering*, 194:1251–1294, 2005.
- [5] F. E. Benth and J. Gjerde. Convergence rates for finite element approximations of stochastic partial differential equations. *Stochastics and Stochastics Rep.*, 63(3-4):313–326, 1998.
- [6] M. Deb, I. Babuška, and J. T. Oden. Solution of stochastic partial differential equations using Galerkin finite element techniques. *Computer Methods in Applied Mechanics and Engineering*, 190:6359–6372, 2001.
- [7] P. Frauenfelder, C. Schwab, and R. A. Todor. Finite elements for elliptic problems with stochastic coefficients. *Computer Methods in Applied Mechanics and Engineering*, 194(2-5):205–228, 2005.
- [8] P. Ladevèze and E. Florentin. Verification of stochastic models in uncertain environments using the constitutive relation error method. *Computer Methods in Applied Mechanics and Engineering*, 196(1-3):225–234, 2006.
- [9] L. Mathelin and O. Le Maître. Dual-based a posteriori error estimate for stochastic finite element methods. *Communications in Applied Mathematics and Computational Science*, 2:83–116, 2007.
- [10] R. Ghanem, G. Saad, and A. Doostan. Efficient solution of stochastic systems: application to the embankment dam problem. *Structural Safety*, 29(3):238–251, 2007.
- [11] A. Doostan, R. Ghanem, and J. Red-Horse. Stochastic model reductions for chaos representations. *Computer Methods in Applied Mechanics and Engineering*, 196(37-40):3951–3966, 2007.
- [12] B. W. Silverman. Smoothed functional principal components analysis by choice of norm. *Ann. Statist.*, 24(1):1–24, 1996.
- [13] A. Levy and J. Rubinstein. Some properties of smoothed principal component analysis for functional data. *Journal of The Optical Society of America*, 16(1):28–35, 1999.
- [14] M. Kirby. Minimal dynamical systems from pdes using sobolev eigenfunctions. *Physica D: Nonlinear Phenomena*, 57(3-4):466–475, 1992.
- [15] G. Berkooz, P. Holmes, and J. L. Lumley. The proper orthogonal decomposition in the analysis of turbulent flows. *Annual review of fluid mechanics*, 25:539–575, 1993.
- [16] A. Nouy. A generalized spectral decomposition technique to solve a class of linear stochastic partial differential equations. *Computer Methods in Applied Methods in Engineering*, 196(45-48):4521–4537, 2007.
- [17] P. Ladevèze. *Nonlinear Computational Structural Mechanics - New Approaches and Non-Incremental Methods of Calculation*. Springer Verlag, 1999.

- [18] P. Ladevèze. and A. Nouy On a Multiscale Computational Strategy with Time and Space Homogenization for Structural Mechanics. *Computer Methods in Applied Methods in Engineering*, 192:3061–3087, 2003.
- [19] A. Nouy and P. Ladevèze. Multiscale computational strategy with time and space homogenization: a radial-type approximation technique for solving micro problems. *International Journal for Multiscale Coputational Engineering*, 170(2):557–574, 2004.
- [20] Y. Saad. *Numerical methods for large eigenvalue problems*. Halstead Press, New York, 1992.
- [21] D. B. Xiu and G. E. Karniadakis. The Wiener-Askey polynomial chaos for stochastic differential equations. *SIAM J. Sci. Comput.*, 24(2):619–644, 2002.
- [22] C. Soize and R. Ghanem. Physical systems with random uncertainties: chaos representations with arbitrary probability measure. *SIAM J. Sci. Comput.*, 26(2):395–410, 2004.
- [23] O. P. Le Maître, O. M. Knio, H. N. Najm, and R. G. Ghanem. Uncertainty propagation using Wiener-Haar expansions. *Journal of Computational Physics*, 197(1):28–57, 2004.
- [24] O. P. Le Maître, H. N. Najm, R. G. Ghanem, and O. M. Knio. Multi-resolution analysis of wiener-type uncertainty propagation schemes. *Journal of Computational Physics*, 197(2):502–531, 2004.
- [25] X. Wan and G.E. Karniadakis. An adaptive multi-element generalized polynomial chaos method for stochastic diffential equations. *J. Comp. Phys.*, 209:617–642, 2005.
- [26] R. G. Ghanem and R. M. Kruger. Numerical solution of spectral stochastic finite element systems. *Computer Methods in Applied Mechanics and Engineering*, 129:289–303, 1996.
- [27] M. F. Pellissetti and R. G. Ghanem. Iterative solution of systems of linear equations arising in the context of stochastic finite elements. *Advances in Engineering Software*, 31:607–616, 2000.
- [28] A. Keese and H. G. Mathhies. Hierarchical parallelisation for the solution of stochastic finite element equations. *Computer Methods in Applied Mechanics and Engineering*, 83:1033–1047, 2005.
- [29] P. B. Nair and A. J. Keane. Stochastic reduced basis methods. *AIAA Journal*, 40(8):1653–1664, 2002.
- [30] S. K. Sachdeva, P. B. Nair, and A. J. Keane. Hybridization of stochastic reduced basis methods with polynomial chaos expansions. *Probabilistic Engineering Mechanics*, 21(2):182–192, 2006.
- [31] A. Edelman, T. A. Arias, and S. T. Smith. The geometry of algorithms with orthogonality constraints. *SIAM J. Matrix Anal. Appl.*, 20(2):303353, 1998.

- [32] P.-A. Absil, R. Mahony, and R. Sepulchre. Riemannian geometry of Grassmann manifolds with a view on algorithmic computation. *Acta Appl. Math.*, 80:199220, 2004.
- [33] A. Sameh and Z. Tong. The trace minimization method for the symmetric generalized eigenvalue problem. *J. Comput. Appl. Math.*, 123:155–175, 2000.
- [34] P.-A. Absil, R. Mahony, R. Sepulchre, and P. Vandooren. A Grassmann-Rayleigh Quotient Iteration for computing invariant subspaces. *SIAM Review*, 44:57–73, 2002.
- [35] C. Soize. Non-gaussian positive-definite matrix-valued random fields for elliptic stochastic partial differential operators. *Computer Methods in Applied Mechanics and Engineering*, 195(1-3):26–64, 2006.
- [36] G. H. Golub and C. F. Van Loan. *Matrix Computations, 3rd ed.* Johns Hopkins University Press, Baltimore, MD, 1996.